

Mining Signed Social Networks Using Unsupervised Learning Algorithms

by

Kewei Cheng

A Thesis Presented in Partial Fulfillment
of the Requirements for the Degree
Master of Science

Approved February 2017 by the
Graduate Supervisory Committee:

Huan Liu, Chair
Hanghang Tong
Chitta Baral

ARIZONA STATE UNIVERSITY

May 2017

ABSTRACT

Due to vast resources brought by social media services, social data mining has received increasing attention in recent years. The availability of sheer amounts of user-generated data presents data scientists both opportunities and challenges. Opportunities are presented with additional data sources. The abundant link information in social networks could provide another rich source in deriving implicit information for social data mining. However, the vast majority of existing studies overwhelmingly focus on positive links between users while negative links are also prevailing in real-world social networks such as distrust relations in Epinions and foe links in Slashdot. Though recent studies show that negative links have some added value over positive links, it is difficult to directly employ them because of its distinct characteristics from positive interactions. Another challenge is that label information is rather limited in social media as the labeling process requires human attention and may be very expensive. Hence, alternative criteria are needed to guide the learning process for many tasks such as feature selection and sentiment analysis.

To address above-mentioned issues, I study two novel problems for signed social networks mining, (1) unsupervised feature selection in signed social networks; and (2) unsupervised sentiment analysis with signed social networks. To tackle the first problem, I propose a novel unsupervised feature selection framework SignedFS. In particular, I model positive and negative links simultaneously for user preference learning, and then embed the user preference learning into feature selection. To study the second problem, I incorporate explicit sentiment signals in textual terms and implicit sentiment signals from signed social networks into a coherent model Signed-Senti. Empirical experiments on real-world datasets corroborate the effectiveness of these two frameworks on the tasks of feature selection and sentiment analysis.

This thesis is dedicated to my parents for their endless love and support.

ACKNOWLEDGMENTS

I would like to thank all the people for helping me to get this thesis finished. First, I owe particular thanks to my committee chair and advisor, Professor Huan Liu, who has been very nice and patient to me during my Master study. This thesis is impossible without the help from my advisor Dr. Huan Liu. I offer my sincere appreciation and gratitude to his help. Second, I would like to thank Professor Hanghang Tong and Professor Chitta Baral for being a part of my thesis committee and giving me helpful suggestions and insightful comments to my thesis.

I would like to thank everyone in our lab, Jiliang Tang, Robert Trevino, Fred Morstatter, Isaac Jones, Suhas Ranganath, Suhang Wang, Tahora Hossein Nazer, Jundong Li, Liang Wu, Ghazaleh Beigi, Kai Shu, Justin Sampson and Harsh Dani, for your care and nice help throughout my life and education in U.S. In particular, thanks to Jundong Li who helped me write my first paper and gave me many helpful suggestions for my thesis.

Finally yet importantly, I would like to thank my parents, without their love and support, I would not have been able to complete a graduation education.

TABLE OF CONTENTS

	Page
LIST OF TABLES	vi
LIST OF FIGURES	vii
CHAPTER	
1 INTRODUCTION	1
1.1 Research Challenges	2
1.2 Contributions	3
1.3 Organization	4
2 RELATED WORK	5
2.1 Feature Selection in Social Media	5
2.2 Sentiment Analysis in Social Media	6
2.3 Signed Social Networks Analysis	6
3 ANALYSIS ON NEGATIVE LINKS	8
3.1 Negative Links and Node Similarity	8
3.1.1 First-order Proximity	10
3.1.2 Second-order Proximity	11
3.2 Negative Links and Sentiment Similarity	12
4 UNSUPERVISED FEATURE SELECTION IN SIGNED SOCIAL NET- WORKS	15
4.1 Problem Statement	16
4.2 The Proposed Framework - SignedFS	16
4.2.1 Modeling Positive and Negative Links	17
4.2.2 Modeling Feature Information	18
4.2.3 Signed Graph Regularization	19
4.3 Optimization	22

CHAPTER	Page
4.4 Experiments.....	25
4.4.1 Experimental Setting.....	25
4.4.2 Quality of Selected Features by SignedFS	26
4.4.3 Impact of Negative Links	29
4.4.4 Parameter Analysis	30
5 UNSUPERVISED SENTIMENT ANALYSIS IN SIGNED SOCIAL NET- WORKS	33
5.1 Problem Statement	34
5.2 The Proposed Framework-SignedSenti.....	35
5.2.1 Basic Model for Unsupervised Sentiment Analysis	35
5.2.2 Sentiment Signals from Textual Terms.....	36
5.2.3 Exploiting Positive and Negative Interactions	37
5.2.4 Objective Function of SignedSenti.....	38
5.3 Optimization Algorithm for SignedSenti	39
5.4 Experiments.....	41
5.4.1 Experimental Setting.....	42
5.4.2 Sentiment Polarity Prediction Performance	43
5.4.3 Parameter Analysis	45
6 CONCLUSION AND FUTURE WORK.....	46
REFERENCES	48

LIST OF TABLES

Table	Page
3.1 Statistics of Datasets for Validation of Homophily Effect in Signed Social Networks.....	9
3.2 P -values for t -test Results	11
3.3 Statistics of Datasets for Validation of Signed Link Based Partial Order Assumption	13
4.1 Clustering Performance of Different Feature Selection Algorithms in Epinions	27
4.2 Clustering Performance of Different Feature Selection Algorithm in Wiki-rfa.....	28
5.1 Sentiment Polarity Prediction Accuracy.....	44

LIST OF FIGURES

Figure	Page
4.1 Illustration of the Proposed SignedFS Framework.	17
4.2 The Impact of Negative Links for SignedFS on Wiki-rfa	30
4.3 Parameter Analysis of SignedFS on Wikipedia	31
5.1 An Illustration of Unsupervised Sentiment Analysis With Signed Social Networks.	35
5.2 Parameter Analysis of SignedSenti on Slashdot.	44

Chapter 1

INTRODUCTION

With the rise of online social platforms such as Facebook ¹ and Twitter ², social network analysis has gained increasing attentions in recent years. The popularity of social media services greatly diversifies the way people communicate and socialize, enabling users to share and exchange opinions in different aspects. Huge volumes of data are user generated at an unprecedented speed. For example, over 500 terabyte data are generated on Facebook every day ³ and around 6000 tweets are tweeted on Twitter every second ⁴. These massive amounts of high-dimensional social media data (e.g., posts, images, videos) present challenges to traditional data mining tasks due to the curse of dimensionality [10]. The sheer volume of opinion-rich data also present great opportunities by providing rich sources in understanding individual and public opinions. For example, unveiling the opinions of customers is valuable for business advertisers in devising better targeted marketing tactics [30]; politicians could also adjust their campaign strategies according to the aggregated sentiments of tweets about election [33].

Social media data is inherently linked by various types of social relations, making it distinct from traditional independent and identically distributed, i.e., i.i.d. data. Motivated by social science theories such as social influence and homophily effect [9, 23, 32, 35], rich sources of information may exist among user interactions. Since label

¹<https://www.facebook.com/>

²<https://twitter.com/>

³<http://www.cnet.com/news/facebook-processes-more-than-500-tb-of-data-daily>

⁴<http://www.internetlivestats.com/twitter-statistics>

information (e.g., user group, sentiment polarity) is costly and labor-intensive to obtain for social media data, these social science theories could be potentially helpful to direct the learning process for a variety of social media mining tasks including feature selection and sentiment analysis.

A majority of existing methods for social media data mining mainly leverage positive interactions among users to guide the learning process. However, in addition to positive links, many real-world social networks may also contain negative links, such as distrust relations in Epinions⁵ and foes links in Slashdot⁶. Even for some platforms without explicit negative links, it is still possible to infer the attitude of a link (positive or negative) from user rating scores or reviews implicitly [42]. The social networks with both positive and negative links are often referred to signed social networks. The availability of negative interactions bring about a richer source of information and recent work shows that negative links have additional values over positive interactions, which could benefit a variety of learning tasks such as community detection [22, 28], recommendation [41] and link prediction [7, 14, 18, 24, 37]. The Recent advance of signed social network analysis motivates me to investigate if negative links can help us mining social media when label information is not available.

1.1 Research Challenges

Despite the potential opportunities from negative links, the development of a principled learning model for unsupervised learning methods for signed social networks mining is still in its infancy. The reason can mainly be attributed as follows:

- Due to the lack of label information, unsupervised methods are more appealing in practice for social data mining. Without label information, the most

⁵<http://www.epinions.com>

⁶<http://slashdot.org>

challenging part is to exploit alternative criteria to guide learning tasks.

- Existing methods for social data mining mainly extract latent representations from positive links and then employ these latent representations to guide learning tasks [5]. However, different from positive links, negative links carry out different information. For example, trust information is often a good indicator of positive emotions such as joy and altruism; while distrust relations may be indicators of negative emotions like anger and pessimism. Hence, mining signed social networks can not simply be extended from existing methods of mining unsigned social networks in a straightforward way;
- Majority of existing methods of mining unsigned social networks are based on some social theories [1, 9, 15, 23], assuming that individuals tend to be similar when they are connected. Nonetheless, these theories may not be directly applicable to signed social networks where individuals with negative links may show contrastive properties. Hence, mining signed social networks is a non-trivial problem.

1.2 Contributions

In this thesis, I study two novel problems in signed social networks, unsupervised feature selection and unsupervised sentiment analysis, which have not been studied previously. In particular, I focus on answering three questions: (1) how to employ and adapt existing social science theories on unsigned networks for signed social networks? (2) how to mathematically model both positive and negative links for feature selection? (3) how to explicitly model positive and negative interactions among users for unsupervised sentiment analysis? The main contributions of this paper are summarized as follows:

- My preliminary data analysis on signed social networks pave the way for adapting existing social science theories on unsigned networks for the mining of signed social networks;
- I propose an unsupervised feature selection framework SignedFS which aims to identify relevant features by leveraging both positive and negative links in signed social networks. In detail, I provide a principled way to mathematically model positive and negative links into a coherent latent representation and embed the latent representation into feature selection phase;
- I propose a novel framework SignedSenti to leverage implicit sentiment signals in positive and negative user interactions for unsupervised sentiment analysis. Methodologically, I propose to incorporate the signed social relations and the sentimental signals from textual terms into a unified framework because of the lack of sentiment labels;
- I evaluate the efficiency of the proposed SignedFS and SignedSenti framework on real-world signed social datasets.

1.3 Organization

The rest of the thesis is organized as follows. In Chapter 2, I briefly reviews related work. In Chapter 3, I introduce some real-world signed social networks datasets and conduct preliminary data analysis on signed social networks. In Chapter 4, I formally define the problem of unsupervised feature selection in signed social networks and introduce the details about the proposed unsupervised feature selection framework SignedFS. In Chapter 5, I study a novel problem of sentiment analysis with signed social networks under an unsupervised scenario and propose a novel framework SignedSenti. The thesis is concluded with future work in Chapter 6.

Chapter 2

RELATED WORK

In this section, I briefly review related work from three aspects: (1) feature selection in social media; (2) sentiment analysis in social media; and (3) signed social network analysis.

2.1 Feature Selection in Social Media

With existence of link information, feature selection in networked data are distinct from traditional feature selections which assumes that data is independent and identically distributed. In [13], a supervised feature selection algorithm FSNet was proposed for network data. FSNet captures the correlation between content information and class labels by a linear classifier and it incorporates link information via graph regularization. Distinct from traditional networked data, social media data present its unique characteristics with the existence of complex linkage structure such as CoPost, CoFollowing, CoFollowed and Following. Motivated by these observations, Tang and Liu [44] made the first attempt to perform feature selection for social media data. Since networked data are usually costly to label, an unsupervised feature selection framework LUFs was proposed in [45]. In particular, LUFs extracts social dimensions from link information to help select relevant features. However, link information may contain a lot of noise and itself may be incomplete. In order to alleviate the negative impacts from noisy and incomplete links, Li et al. [27] proposed a robust unsupervised feature selection framework for networked data. However, all above mentioned approaches only consider the positive interactions among networked instance, to the best of our knowledge, I am the first attempt to study unsupervised

feature selection on signed networks.

2.2 Sentiment Analysis in Social Media

Sentiment analysis in social media has been a surge of research recently. However, it faces some challenges mainly because of the bewildering combination of heterogeneous data sources and structures. Also, since labels of social media data are costly to obtain, unsupervised sentiment analysis is more desired. Recent years have witnessed some efforts in exploring external information for unsupervised sentiment analysis. As the most representative unsupervised sentiment analysis algorithms, lexicon-based methods [33, 39, 51] determine sentiment polarity of texts by exploiting sentiment signals revealed by words or phrases. In addition to rich source of text information, abundant emotional signals are widely observed in social media. In [19], the authors proposed a framework to incorporate two categories of emotional signals for unsupervised sentiment analysis. [49] made one of the first attempt to leverage social media images for unsupervised sentiment analysis. Different from above mentioned approaches, I present the first study on unsupervised sentiment analysis with both positive and negative social interactions.

2.3 Signed Social Networks Analysis

Even though mining signed graph is still in its early stage, some problems in signed networks have already been well studied, such as link prediction and community detection. Existing link prediction methods on signed social network can be broadly divided into two groups: supervised methods and unsupervised methods. Supervised links prediction can be regarded as classification problem. Like normal classification problem, the most important parts of it is to construct features. Some common used features include local topology features [24] and feature derived from long cycles [7].

Unsupervised methods predict missing links without label information. These methods mainly predict signs of links according to the topological properties of signed networks [14, 37]. Community detection is another fundamental problem for mining signed social networks. In [28], Li extends modularity maximization to signed networks which takes both the tendency of users with positive links to form communities and the tendency of users with negative links to destroy them into consideration. A spectral algorithm was proposed in [22]. It is the first attempt to define a signed laplacian matrix which can separate users with negative links and force users with positive links to be closer.

Chapter 3

ANALYSIS ON NEGATIVE LINKS

In unsigned social networks, some social science theories such as social influences and homophily [9, 23, 32, 35] are widely adopted in social network analysis to bridges the gap between learning task and network structure, especially in cases when label information is costly to obtain. In this chapter, I investigate whether negative links reveal some useful information for signed social networks mining.

3.1 Negative Links and Node Similarity

The homophily effect [32] in social science theories suggests that users are similar to each other when they are interconnected. However, it is not appropriate to directly apply the homophily effect on signed social network analysis [43] as instances may also be negatively connected. In this subsection, I revisit the homophily effect in signed social networks.

I first introduce two real-world signed networks used in this study. I used two real-world signed social networks datasets from Epinions¹ and Wiki-rfa².

Epinions: Epinions is a consumer review website in which users share their reviews about products. Users can either trust or distrust other users. They can also write reviews for products from various categories. For users, I collect their positive and negative links as well as their reviews comments. Features are formed by the bag-of-words model based on the reviews comments. The major categories of reviews by users are taken as the ground truth of class labels.

¹<http://jiliang.xyz/trust.html>

²<https://snap.stanford.edu/data/wiki-RfA.html>

Table 3.1: Statistics of Datasets for Validation of Homophily Effect in Signed Social Networks

Datasets	Epinions	Wiki-rfa
# of Users	7,140	7,096
# of Features	15,069	10,608
# of Classes	24	2
# of Positive Links	13,569	104,555
Density of Positive Links	2.7e(-4)	2.1e(-3)
# of Negative Links	3,010	23,516
Density of Negative Links	5.9e(-5)	4.7e(-4)

Wiki-rfa: Wikipedia Requests for Adminship is a who-votes-for-whom network where a signed link indicates a positive or a negative vote by one user on the promotion of another. Each vote is typically accompanied by a short comment which is used to construct features by the bag-of-words model. The person voted by the user could be rejected or accepted, which is taken as ground truth.

Detailed statistics of these two datasets are presented in Table 3.1. I notice that positive links are denser than negative links in both two datasets. With these properties, I now study the first-order proximity and the second-order proximity in signed social networks.

In social sciences, some theories such as homophily effect [32] and balance theory [17] suggest the correlations between user similarity and positive/negative links. These theories bridges the gap between user features and network structure, and is widely adopted in social network analysis. Two kinds of network structures have been investigated in social theories. One is represented by the observed links in the networks, which reveals the first-order proximity between the users. For example, the homophily effect explores the first-order proximity between users in social network-

s. The other is represented by two users with shared neighborhoods. For instance, balance theory suggests the second-order proximity between the users. In this subsection, I would like to explore the first-order and the second-order proximity in signed social networks.

3.1.1 First-order Proximity

The homophily effect in social science theories suggests that users are similar to each other when they are interconnected. However, it is not appropriate to directly apply the homophily effect on signed social network analysis [43] as instances may also be negatively connected. To explore the first-order proximity in signed social networks, I revisit the homophily effect in signed social networks by attempting to answering the following questions: are users with positive relations tend to be more similar than users with negative relations?

To answer these questions, first, I define the user similarity score between two users u_i and u_j as $sim_{ij} = \|\mathbf{y}_i - \mathbf{y}_j\|_2$, where $\mathbf{y}_i \in \mathbb{R}^{1 \times c}$ and $\mathbf{y}_j \in \mathbb{R}^{1 \times c}$ are the ground truth of user labels for user u_i and u_j , respectively. k denotes the number of user labels.

With the definition of user similarity, I construct two vectors \mathbf{p}_1 and \mathbf{n}_1 with the same length to denote the user similarity between positively connected users and negatively connected users, respectively. To be specific, elements in \mathbf{p}_1 denotes the similarity score between two users (u_i, u_j) with positive relation. Elements in \mathbf{n}_1 denotes the similarity score between two users (u_i, u_j) with negative relation. To see the significance, I sample 500 pairs of users for each of vectors \mathbf{p}_1 and \mathbf{n}_1 and conduct two samples t-tests on these three vectors. The null hypothesis is rejected at significance level $\alpha = 0.01$ with p -values shown in Table 3.2. Therefore, we verify that users with positive relations tend to be more similar than users with negative

relations.

3.1.2 Second-order Proximity

Balance theory in signed social networks suggests that "the friend of my friend is my friend" and "the enemy of my enemy is my friend". Based on balance theory, I would like to investigate the second-order proximity in signed social networks. Specifically, I aim to answer the following two questions: (1) Is friend of my friend tend to be similar with me? (2) Is enemy of my enemy more likely to be similar with me?

With user similarity and vector \mathbf{r} defined in section 3.1.1, I construct another three vectors \mathbf{p}_2 , \mathbf{n}_2 and \mathbf{r} to denote the user similarity between two users with a shared friend, two users with a shared enemy and randomly chosen users, respectively. For example, each element in \mathbf{p}_2 denotes the similarity score between two users u_i and u_k . Both u_i and u_k have a friend u_j . The element in \mathbf{n}_2 denotes the similarity score between two users u_i and u_k . Both u_i and u_k have an enemy u_j . And the element of \mathbf{r} represents the similarity score between u_i and another randomly selected user u_r . I also sample 500 pairs of users for each of vectors \mathbf{p}_2 , \mathbf{n}_2 and \mathbf{r} and conduct two samples t-tests on these three vectors. The null hypothesis is rejected at significance level $\alpha = 0.01$ with p -values shown in Table 3.2. From the table, we observe that both of the friend of my friend and the enemy of my enemy are more likely to be similar to me.

Table 3.2: P -values for t -test Results

Hypothesis	Epinions	Wiki-rfa
$H_0 : \mathbf{p}_1 \geq \mathbf{n}_1$ $H_1 : \mathbf{p}_1 < \mathbf{n}_1$	$2.3974e(-7)$	$8.3255e(-4)$
$H_0 : \mathbf{p}_2 \geq \mathbf{r}$ $H_1 : \mathbf{p}_2 < \mathbf{r}$	$1.3614e(-6)$	$9.8577e(-7)$
$H_0 : \mathbf{n}_2 \geq \mathbf{r}$ $H_1 : \mathbf{n}_2 < \mathbf{r}$	$5.5854e(-5)$	$1.3126e(-12)$

3.2 Negative Links and Sentiment Similarity

Above results show that users are likely to be more similar to their friends than their foes. Hence, it motivates me to investigate if friends are more likely to exhibit similar sentiments than foes on the same item which I conclude as the signed link based partial order assumption. To introduce signed link based partial order assumption, I first define the concepts of positive linked set, negative linked set.

Definition 1. Positive Linked Set:

For a specific text post t_i on the item o_r posted by user u_a , its positive linked set $\mathcal{P}(t_i)$ is defined as the whole set of text posts t_j on the same item o_r that are posted by user u_b , where u_b is positively connected from u_a , i.e., $\mathcal{P}(t_i) = \{t_j | \forall (j, r, a, b) \text{ s.t. } \mathbf{O}_{ir} = 1, \mathbf{O}_{jr} = 1, \mathbf{T}_{ai} = 1, \mathbf{T}_{bj} = 1, \mathbf{A}_{ab} = 1\}$.

Definition 2. Negative Linked Set:

For a specific text post t_i on the item o_r posted by user u_a , its negative linked set $\mathcal{N}(t_i)$ is defined as the whole set of text posts t_k on the same item o_r that are posted by user u_b , where u_b is negatively connected from u_a , i.e., $\mathcal{N}(t_i) = \{t_k | \forall (k, r, a, b) \text{ s.t. } \mathbf{O}_{ir} = 1, \mathbf{O}_{kr} = 1, \mathbf{T}_{ai} = 1, \mathbf{T}_{bk} = 1, \mathbf{A}_{ab} = -1\}$.

With the concepts of positive linked set, negative linked set, the signed link based partial order assumption is defined as following:

Assumption 1. Signed Link Based Partial Order:

For text post t_j in the positive linked set of t_i and text post t_k in the negative linked set of t_i , sentiment polarity of t_i is usually more similar to the sentiment polarity of t_j than t_k . I denote such property as signed link partial order which can be formulated as follows:

$$\text{sim}(t_i, t_j) > \text{sim}(t_i, t_k), t_j \in \mathcal{P}(t_i), t_k \in \mathcal{N}(t_i) \quad (3.1)$$

Table 3.3: Statistics of Datasets for Validation of Signed Link Based Partial Order Assumption

Statistics	Epinions	Slashdot
# of posts	1,559,803	133,335
# of items	200,952	72,241
# of users	326,978	7,897
# of positive links	717,667	52,639
# of negative links	123,705	17,535

Similarly, I first introduce two real-world signed social networks datasets from Epinions³ and Slashdot⁴ used in validating the signed link based partial order assumption. Detailed statistics of these two datasets are shown in Table 3.3.

Epinions: As shown in Section 2.2, Epinions is a product review website where users share their reviews about products. I crawled a set of reviews, products and users as well as their interactions. The unigram model is employed on product reviews to construct the feature space, and term frequency is used as feature weight. For the evaluation purpose, I take the rating scores of reviews as the ground truth of sentiment labels. In particular, the ratings of 4, 5 and 6 are considered as positive labels while the ratings of 1,2 and 3 are taken as negative labels.

Slashdot: Slashdot is a technology news website for users to share and comment new articles on science and technology. Users can tag others as friends or foes. Likewise, I crawled and collect comments, articles, users and their relations. The feature space is also built with unigram model and the ratings of comments are employed to establish ground truth in the same way as Epinions.

With these two datasets, I start to validate whether the signed link based partial

³<http://jiliang.xyz/trust.html>

⁴<https://slashdot.org/>

order assumption holds for text posts in real-world signed networks.

First, I define the sentiment similarity between two text posts t_i and t_j as $sim(t_i, t_j) = \|\mathbf{y}_i - \mathbf{y}_j\|_2$, where $\mathbf{y}_i \in \mathbb{R}^{1 \times k}$ and $\mathbf{y}_j \in \mathbb{R}^{1 \times k}$ are the ground truth of sentiment labels for text posts x_i and x_j , respectively. k denotes the number of sentiment labels. With the definition of text post sentiment similarity, to verify if the signed link based partial order assumption holds, I construct two vectors \mathbf{s}_p and \mathbf{s}_n of the same length. Elements in \mathbf{s}_p denote the sentiment similarity of two text posts t_i and t_j , where t_j is from the positive linked set of t_i . Elements in \mathbf{s}_n indicate the sentiment similarity between two text posts t_i and t_k where t_k is from the negative linked set of t_i . To validate the assumption, I first sample 500 pairs in each group to construct \mathbf{s}_p and \mathbf{s}_n , and then conduct two sample t-test on these two vectors. The null hypothesis is $H_0 : \mathbf{c}_p \geq \mathbf{c}_n$ while the alternative hypothesis is $H_1 : \mathbf{c}_p < \mathbf{c}_n$. In the formulations, \mathbf{c}_p and \mathbf{c}_n represent the sample means in these two groups \mathbf{s}_p and \mathbf{s}_n , respectively. The null hypothesis is rejected at the significant level $\alpha = 0.01$ with p -values of $4.3e(-7)$ and $7.2e(-4)$ in Epinions and Slashdot, respectively. It indicates that the signed link based partial order assumption indeed holds in real-world signed social networks. In other words, it suggests the existence of implicit sentiment signals among positive and negative user interactions, which paves way for unsupervised sentiment analysis.

UNSUPERVISED FEATURE SELECTION IN SIGNED SOCIAL NETWORKS

The rapid growth of social media services brings large amounts of high-dimensional social media data at an unprecedented rate. Feature selection has shown to be powerful to prepare high-dimensional data for effective machine learning tasks [6, 11, 31]. A majority of existing feature selection algorithms for social media data exclusively focus on positive interactions among linked instances [26, 27, 44, 45]. However, in many real-world social networks, instances may also be negatively interconnected. Recent work shows that negative links have an added value over positive links, and the leverage of negative links could improve various learning tasks such as community detection [22, 28], recommendation [41] and link prediction [7, 14, 18, 24, 37]. To take advantage of negative links, I study a novel problem of unsupervised feature selection in signed social networks and propose a novel framework SignedFS. In particular, I provide a principled way to model positive and negative links for user preference learning. Then I embed the user preference learning into feature selection. Also, I revisit the homophily effect and balance theory in signed social networks and incorporate signed graph regularization into the feature selection framework to capture the first-order proximity and the second-order proximity in signed social networks. Experiments on real-world signed social networks demonstrate the effectiveness of our proposed framework. Further experiments are conducted to understand the impacts of negative links for feature selection.

4.1 Problem Statement

To formally define the problem unsupervised feature selection on signed social networks, I first present the notations.

Let $\mathcal{U} = \{u_1, u_2, \dots, u_n\}$ be the set of n users in a signed network \mathcal{G} . \mathcal{G} can be decomposed into a positive component \mathcal{G}_p and a negative component \mathcal{G}_n in which $\mathbf{A}^p \in \mathbb{R}^{n \times n}$ is the corresponding adjacency matrix for the positive component \mathcal{G}_p such that $\mathbf{A}_{ij}^p = 1$ if u_i has a positive link to u_j , and $\mathbf{A}_{ij}^p = 0$ otherwise. Similarly, $\mathbf{A}^n \in \mathbb{R}^{n \times n}$ denotes the adjacency matrix of \mathcal{G}_n where $\mathbf{A}_{ij}^n = 1$ if u_i has a negative link to u_j , and $\mathbf{A}_{ij}^n = 0$ otherwise. Let $\mathcal{F} = (f_1, f_2, \dots, f_d)$ be a set of d features and $\mathbf{X} \in \mathbb{R}^{n \times d}$ denotes the content information of all n instances. With these notations, the problem of unsupervised feature selection in signed social networks can be formally stated as follows:

Given: the feature set \mathcal{F} , content matrix \mathbf{X} and adjacency matrix \mathbf{A} for a signed network \mathcal{G} with positive links \mathbf{A}^p and negative links \mathbf{A}^n , **Select:** A subset of most relevant features $\mathcal{S} \in \mathcal{F}$ by exploiting both content information \mathbf{X} and signed network information \mathbf{A}^p and \mathbf{A}^n .

4.2 The Proposed Framework - SignedFS

In this section, I illustrate the proposed unsupervised feature selection in signed social networks in details. The workflow of the proposed framework SignedFS is shown in Figure 4.1. As can be observed from the figure, it consists of three components: first, I show how to learn user preference representation from both positive and negative links (Section 4.2.1); second, I show how to embed the user preference representation into feature selection when we are lack of label information (Section 4.2.2); third, I show how to employ the first-order and the second-order proximity in signed

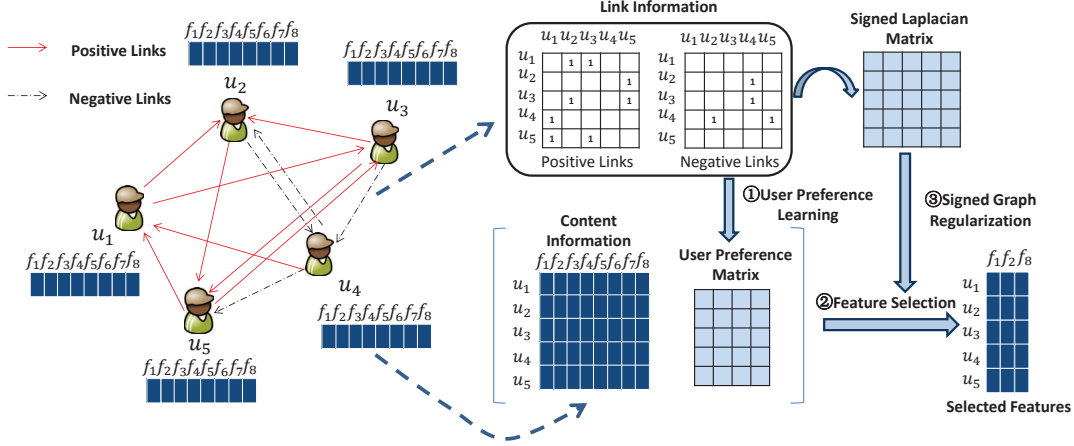


Figure 4.1: Illustration of the Proposed SignedFS Framework.

social networks to make user preference representation to be consistent with signed network structure via a signed graph regularization (Section 4.2.3).

4.2.1 Modeling Positive and Negative Links

In social media, a user establish relations with others due to a variety of hidden factors. These factors are often referred as user preferences including hobbies, geographical location, religion, etc. It has been widely studied in previous research that both positive and negative links are relevant to user preference [38, 43]. Considering the fact that negative links possess unique characteristics compared with positive links, I attempt to model positive and negative links independently to learn the user preference representation (phase 1 in Figure 4.1). Let $\mathbf{U} = [\mathbf{U}_{1*}; \mathbf{U}_{2*}; \dots; \mathbf{U}_{n*}] \in \mathbb{R}^{n \times c}$ be the user preference representation where \mathbf{U}_{i*} denotes user preference of u_i . It should be noticed that in real-world signed social networks, a user only has a small portion of links with others, resulting in a sparse and low rank network structure. Therefore, I employ low-rank matrix factorization method to learn user preference representation. Specifically, to capture the properties of positive and negative links independently, I collectively factorize \mathbf{A}^p and \mathbf{A}^n into a unified low rank representa-

tion \mathbf{U} via the following optimization problem:

$$\min_{\mathbf{U}, \mathbf{V}^p, \mathbf{V}^n} \beta_1 \|\mathbf{O}^p \odot (\mathbf{A}^p - \mathbf{U}\mathbf{V}^p\mathbf{U}')\|_F^2 + \beta_2 \|\mathbf{O}^n \odot (\mathbf{A}^n - \mathbf{U}\mathbf{V}^n\mathbf{U}')\|_F^2, \quad (4.1)$$

where β_1 and β_2 balances the contribution of positive links and negative links in learning user preference representation. \mathbf{O}^p and \mathbf{O}^n are defined as follows:

$$\mathbf{O}_{ij}^p = \begin{cases} 1, & \text{if } \mathbf{A}_{ij}^p = 1 \\ 0, & \text{otherwise} \end{cases}, \quad (4.2)$$

$$\mathbf{O}_{ij}^n = \begin{cases} 1, & \text{if } \mathbf{A}_{ij}^n = 1 \\ 0, & \text{otherwise} \end{cases}. \quad (4.3)$$

In the above formulation, I approximate the positive link from u_i to u_j with $\mathbf{U}_i\mathbf{V}^p\mathbf{U}'_j$ where $\mathbf{V}^p \in \mathbb{R}^{c \times c}$ captures the correlations among user preference representation for positive links. \odot is Hadamard product (element-wise product) where $(\mathbf{X} \odot \mathbf{Y})_{ij} = \mathbf{X}_{ij} \times \mathbf{Y}_{ij}$ for any two matrices \mathbf{X} and \mathbf{Y} of the same size. The Hadamard product operator is imposed since I only use existing positive links to learn user preference representation. Similarly, I approximate negative links with $\mathbf{U}\mathbf{V}^n\mathbf{U}'$. Since negative links are also related to user preference representation, I factorize \mathbf{A}^n into the same low-rank space \mathbf{U} . The correlation matrix \mathbf{V}^n is used to capture the unique properties of negative links.

4.2.2 Modeling Feature Information

After I model user preference representation, I now introduce how to employ them to guide feature selection in the content space (phase 2 in Figure 4.1). In social media platforms, labels are costly and labor intensive to obtain. Without label information, it would be difficult to assess feature relevance. Fortunately, since user preference representations encode latent factors of users, they are correlated with

the features, at least with some relevant features. Therefore, I leverage the user preference representations \mathbf{U} to take the role of class labels to guide feature selection via a multivariate linear regression model with a $\ell_{2,1}$ -norm sparse regularization term:

$$\min_{\mathbf{W}} \|\mathbf{X}\mathbf{W} - \mathbf{U}\|_F^2 + \alpha \|\mathbf{W}\|_{2,1}, \quad (4.4)$$

where $\mathbf{W} \in \mathbb{R}^{d \times c}$ is a feature weight matrix and each row of \mathbf{W} , i.e., \mathbf{W}_{i*} measures the importance of the i -th feature. The $\ell_{2,1}$ -norm regularization term is imposed on \mathbf{W} to achieve a joint feature sparsity across k different dimensions of user preference representation. α controls the sparsity of the model.

However, signed social networks may contain a lot of noisy links. For example, illegitimate users such as spammers and bots will generate a large amount of fake links to imitate normal users. In addition to that, network structure may also not be complete, mainly because of the imperfect data collection and data crawling process, or the network itself is partially observed. Therefore, I propose to embed the latent representation learning into feature selection to make the feature selection results more robust to noisy and incomplete positive and negative links, resulting in the following optimization framework:

$$\begin{aligned} \min_{\mathbf{W}, \mathbf{U}, \mathbf{V}^p, \mathbf{V}^n} & \|\mathbf{X}\mathbf{W} - \mathbf{U}\|_F^2 + \alpha \|\mathbf{W}\|_{2,1} + \frac{\beta_1}{2} \|\mathbf{O}^p \odot (\mathbf{A}^p - \mathbf{U}\mathbf{V}^p\mathbf{U}')\|_F^2 \\ & + \frac{\beta_2}{2} \|\mathbf{O}^n \odot (\mathbf{A}^n - \mathbf{U}\mathbf{V}^n\mathbf{U}')\|_F^2, \end{aligned} \quad (4.5)$$

where the parameter α controls the sparsity of the model.

4.2.3 Signed Graph Regularization

In Section 4.1, I revisit the homophily effect and balance theory by verifying the existence of the first-order and second-order proximity in signed social networks. In this subsection, I introduce how to model the first-order and the second-order

proximity for unsupervised feature selection in signed social networks (phase 3 in Figure 4.1).

I first construct a user proximity matrix by employing both the first-order and the second-order proximity. Given the adjacency matrix of a signed network \mathbf{A} where $\mathbf{A}_{ij} = 1$, $\mathbf{A}_{ij} = -1$ and $\mathbf{A}_{ij} = 0$ denote positive, negative and missing links from u_i to u_j . The first-order proximity matrix \mathbf{P}_1 is defined as $\mathbf{P}_1 = \mathbf{A}$, where $\mathbf{P}_{1ij} = 1$ indicates that u_j is a friend of u_i and preferences of the two users are similar while $\mathbf{P}_{1ij} = -1$ indicates that u_j is a foe of u_i and preferences of the two users are dissimilar. The second-order proximity matrix is defined as $\mathbf{P}_2 = \mathbf{O} \odot \mathbf{A}^2$, where \mathbf{O} is defined as follows:

$$\mathbf{O}_{ij} = \begin{cases} 0, & \text{if } \mathbf{P}_{1ij} \neq 0 \\ 1, & \text{otherwise} \end{cases}. \quad (4.6)$$

where $\mathbf{P}_{2ij} > 0$ and $\mathbf{P}_{2ij} < 0$ denote similarity and dissimilarity between u_i and u_j . In the above formulation, I capture the second-order proximity from u_i to u_k with $(\mathbf{A}^2)_{ik} = \sum_{j=1}^n a_{ij}a_{jk}$. To show that \mathbf{A}^2 can capture the second-order proximity, the proof is as follows: (1) to verify that \mathbf{A}^2 can capture the proximity between friend of my friend and me, I should prove that if both u_i and u_k have a friend u_j , u_i and u_j should be similar with each other in the second-order proximity matrix. In other word, if $\text{sign}(\mathbf{A}_{ij}) = 1$ and $\text{sign}(\mathbf{A}_{jk}) = 1$, I should prove that $\text{sign}(\mathbf{A}_{ik}) = 1$ which might seem obvious in the above formulation; (2) to verify that \mathbf{A}^2 can capture the proximity between enemy of my enemy and me, I should prove that if both u_i and u_k have an enemy u_j , u_i and u_j should be similar with each other in the second-order proximity matrix. That is if $\text{sign}(\mathbf{A}_{ij}) = -1$ and $\text{sign}(\mathbf{A}_{jk}) = -1$, I should prove that $\text{sign}(\mathbf{A}_{ik}) = 1$, which is also true in the above formulation. Though the second-order proximity (balance theory) may not be always hold in signed networks [42], in an aggregate sense, the second-order proximity from network structure should be

maximally preserved. Thus $(\mathbf{A}^2)_{ik}$ can capture the second-order proximity from u_i to u_k . The Hadamard product operator is imposed to avoid the confliction between the first-order proximity and the second-order proximity. User proximity matrix can be constructed by $\mathbf{P} = \mathbf{P}_1 + \theta\mathbf{P}_2$, where $\mathbf{P}_{ij} > 0$, $\mathbf{P}_{ij} < 0$ denote similarity and dissimilarity between u_i and u_j , respectively. The parameter θ controls the weight of the first-order and the second-order proximities in the model. In this paper, I empirically set the weight $\theta = 0.1$.

To integrate user proximity in feature selection, the basic idea is to make preference of two user \mathbf{U}_{i*} and \mathbf{U}_{j*} as close as possible if u_i and u_j are similar ($\mathbf{P}_{ij} > 0$) while as far as possible if u_i and u_j are dissimilar ($\mathbf{P}_{ij} < 0$). Since the signed Laplacian matrix aims to separate pairs with negative links rather than to force pairs with positive links closer [22], user proximity could be mathematically formulated by the signed graph regularization:

$$\begin{aligned} & \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n |\mathbf{P}_{ij}| \times \|\mathbf{U}_{i*} - \text{sgn}(\mathbf{P}_{ij})\mathbf{U}_{j*}\|^2 \\ & = \text{tr}(\mathbf{U}'\mathbf{L}\mathbf{U}), \end{aligned} \tag{4.7}$$

where $\text{sgn}(\mathbf{P}_{ij})$ denotes the sign of \mathbf{P}_{ij} . $\mathbf{L} = \mathbf{D} - \mathbf{P}$ is a signed Laplacian matrix [22] constructed from \mathbf{P} and the signed degree matrix $\mathbf{D} \in \mathbb{R}^{n \times n}$ is a diagonal matrix with $\mathbf{D}_{ii} = \sum_{j=1}^n |\mathbf{P}_{ij}|$.

With the modeling of user proximity by signed graph regularization, the final objective function of the proposed SignedFS framework is formulated as follows:

$$\begin{aligned} & \min_{\mathbf{W}, \mathbf{U}, \mathbf{V}^p, \mathbf{V}^n} \|\mathbf{X}\mathbf{W} - \mathbf{U}\|_F^2 + \alpha \|\mathbf{W}\|_{2,1} + \frac{\gamma}{2} \text{tr}(\mathbf{U}'\mathbf{L}\mathbf{U}) \\ & + \frac{\beta_1}{2} \|\mathbf{O}^p \odot (\mathbf{A}^p - \mathbf{U}\mathbf{V}^p\mathbf{U}')\|_F^2 \\ & + \frac{\beta_2}{2} \|\mathbf{O}^n \odot (\mathbf{A}^n - \mathbf{U}\mathbf{V}^n\mathbf{U}')\|_F^2, \end{aligned} \tag{4.8}$$

where γ is a regularization parameter for the modeling of user proximity in signed social networks.

4.3 Optimization

In this section, I introduce the alternating optimization algorithm for solving the optimization problem of the proposed SignedFS framework with time complexity analysis.

In Eq.(4.8), the coupling between $\mathbf{U}, \mathbf{V}^p, \mathbf{V}^n$ and \mathbf{W} makes it difficult to find the global optimal solutions for all four variables simultaneously. Therefore, I propose to employ an alternating optimization scheme to solve it which has been widely adopted for a variety of real-world problems [3].

First, I fix \mathbf{U}, \mathbf{V}^p and \mathbf{V}^n and update \mathbf{W} . Specifically, when \mathbf{U}, \mathbf{V}^p and \mathbf{V}^n are fixed, the objective function is convex w.r.t. the feature weight matrix \mathbf{W} . I take the partial derivative of objective function w.r.t. \mathbf{W} and set it to be zero:

$$2\mathbf{X}'(\mathbf{X}\mathbf{W} - \mathbf{U}) + 2\alpha\mathbf{H}\mathbf{W} = 0, \quad (4.9)$$

where $\mathbf{H} \in \mathbb{R}^{d \times d}$ is a diagonal matrix with its i -th diagonal element as:

$$\mathbf{H}_{ii} = \frac{1}{2\|\mathbf{W}_{i*}\|_2}. \quad (4.10)$$

It can be noticed that $\mathbf{X}'\mathbf{X}$ is a positive definite matrix and $\alpha\mathbf{H}$ is a diagonal matrix with positive entries which is positive definite as well. Therefore, their summation should also be positive definite. Hence, \mathbf{W} has a closed form solution, which is:

$$\mathbf{W} = (\mathbf{X}'\mathbf{X} + \alpha\mathbf{H})^{-1}\mathbf{X}'\mathbf{U}. \quad (4.11)$$

¹In practice, $\|\mathbf{W}_{i*}\|_2$ could be close to zero. Thus, I regularize $\mathbf{H}_{ii} = \frac{1}{2\|\mathbf{W}_{i*}\|_2 + \epsilon}$, where ϵ is a very small constant.

By substituting the above solution of \mathbf{W} into Eq.(4.8), we have:

$$\begin{aligned}
& \min_{\mathbf{U}, \mathbf{V}^p, \mathbf{V}^n} \mathcal{J}(\mathbf{U}, \mathbf{V}^p, \mathbf{V}^n) \\
& = \text{tr}(\mathbf{U}'\mathbf{U}) - \text{tr}(\mathbf{U}'\mathbf{X}\mathbf{M}^{-1}\mathbf{X}'\mathbf{U}) + \frac{\beta_1}{2} \|\mathbf{O}^p \odot (\mathbf{A}^p - \mathbf{U}\mathbf{V}^p\mathbf{U}')\|_F^2 \\
& + \frac{\beta_2}{2} \|\mathbf{O}^n \odot (\mathbf{A}^n - \mathbf{U}\mathbf{V}^n\mathbf{U}')\|_F^2 + \frac{\gamma}{2} \text{tr}(\mathbf{U}'\mathbf{L}\mathbf{U}) \tag{4.12} \\
& = \text{tr}(\mathbf{U}'(\mathbf{I}_n - \mathbf{X}\mathbf{M}^{-1}\mathbf{X}')\mathbf{U}) + \frac{\beta_1}{2} \|\mathbf{O}^p \odot (\mathbf{A}^p - \mathbf{U}\mathbf{V}^p\mathbf{U}')\|_F^2 \\
& + \frac{\beta_2}{2} \|\mathbf{O}^n \odot (\mathbf{A}^n - \mathbf{U}\mathbf{V}^n\mathbf{U}')\|_F^2 + \frac{\gamma}{2} \text{tr}(\mathbf{U}'\mathbf{L}\mathbf{U}),
\end{aligned}$$

where $\mathbf{M} = \mathbf{X}'\mathbf{X} + \alpha\mathbf{H}$.

Similarly, we fix other variables to update \mathbf{U} , \mathbf{V}^p and \mathbf{V}^n iteratively. Since their closed form solutions are hard to obtain, we employ gradient descent to update them. In particular, the partial derivative of the objective function w.r.t. \mathbf{U} , \mathbf{V}^p and \mathbf{V}^n can be calculated as follows:

$$\begin{aligned}
\frac{\partial \mathcal{J}}{\partial \mathbf{U}} & = (\mathbf{I}_n - \mathbf{X}\mathbf{M}^{-1}\mathbf{X}')\mathbf{U} + (\mathbf{I}_n - \mathbf{X}\mathbf{M}^{-1}\mathbf{X}')'\mathbf{U} \\
& + \beta_1(-(\mathbf{O}^p \odot \mathbf{O}^p \odot \mathbf{A}^p)\mathbf{U}\mathbf{V}^{p'} - (\mathbf{O}^p \odot \mathbf{O}^p \odot \mathbf{A}^p)'\mathbf{U}\mathbf{V}^p \\
& + (\mathbf{O}^p \odot \mathbf{O}^p \odot \mathbf{U}\mathbf{V}^p\mathbf{U}')\mathbf{U}\mathbf{V}^{p'} \\
& + (\mathbf{O}^p \odot \mathbf{O}^p \odot \mathbf{U}\mathbf{V}^p\mathbf{U}')'\mathbf{U}\mathbf{V}^p) \\
& + \beta_2(-(\mathbf{O}^n \odot \mathbf{O}^n \odot \mathbf{A}^n)\mathbf{U}\mathbf{V}^{n'} - (\mathbf{O}^n \odot \mathbf{O}^n \odot \mathbf{A}^n)'\mathbf{U}\mathbf{V}^n \\
& + (\mathbf{O}^n \odot \mathbf{O}^n \odot \mathbf{U}\mathbf{V}^n\mathbf{U}')\mathbf{U}\mathbf{V}^{n'} \\
& + (\mathbf{O}^n \odot \mathbf{O}^n \odot \mathbf{U}\mathbf{V}^n\mathbf{U}')'\mathbf{U}\mathbf{V}^n) + \gamma\mathbf{L}\mathbf{U}, \tag{4.13}
\end{aligned}$$

$$\begin{aligned}
\frac{\partial \mathcal{J}}{\partial \mathbf{V}^p} & = \beta_1(\mathbf{U}'(\mathbf{O}^p \odot \mathbf{O}^p \odot \mathbf{U}\mathbf{V}^p\mathbf{U}')\mathbf{U} \\
& - \mathbf{U}'(\mathbf{O}^p \odot \mathbf{O}^p \odot \mathbf{A}^p)\mathbf{U}), \tag{4.14}
\end{aligned}$$

$$\begin{aligned}
\frac{\partial \mathcal{J}}{\partial \mathbf{V}^n} & = \beta_2(\mathbf{U}'(\mathbf{O}^n \odot \mathbf{O}^n \odot \mathbf{U}\mathbf{V}^n\mathbf{U}')\mathbf{U} \\
& - \mathbf{U}'(\mathbf{O}^n \odot \mathbf{O}^n \odot \mathbf{A}^n)\mathbf{U}). \tag{4.15}
\end{aligned}$$

Algorithm 1: SignedFS Algorithm

Input : $\{\mathbf{X}, \mathbf{A}^p, \mathbf{A}^n, c, \alpha, \beta_1, \beta_2, \gamma\}$

Output: ranking of features in a descending order

- 1 Initialize \mathbf{U} , \mathbf{V}^p and \mathbf{V}^n randomly;
 - 2 Initialize \mathbf{H} as an identity matrix;
 - 3 $\mathbf{A} = \mathbf{A}^p - \mathbf{A}^n$, $\mathbf{P}_1 = \mathbf{A}$, $\mathbf{P}_2 = \mathbf{O} \odot \mathbf{A}^2$, $\mathbf{P} = \mathbf{P}_1 + \theta\mathbf{P}_2$;
 - 4 $\mathbf{L} = \mathbf{D} - \mathbf{P}$;
 - 5 **while** *not converge* **do**
 - 6 Set $\mathbf{M} = \mathbf{X}'\mathbf{X} + \alpha\mathbf{H}$;
 - 7 Calculate $\frac{\partial \mathcal{J}}{\partial \mathbf{U}}$, $\frac{\partial \mathcal{J}}{\partial \mathbf{V}^p}$ and $\frac{\partial \mathcal{J}}{\partial \mathbf{V}^n}$;
 - 8 Update $\mathbf{U} \leftarrow \mathbf{U} - \lambda_u \frac{\partial \mathcal{J}}{\partial \mathbf{U}}$;
 - 9 Update $\mathbf{V}^p \leftarrow \mathbf{V}^p - \lambda_p \frac{\partial \mathcal{J}}{\partial \mathbf{V}^p}$;
 - 10 Update $\mathbf{V}^n \leftarrow \mathbf{V}^n - \lambda_n \frac{\partial \mathcal{J}}{\partial \mathbf{V}^n}$;
 - 11 Update $\mathbf{W} \leftarrow \mathbf{M}^{-1}\mathbf{X}'\mathbf{U}$;
 - 12 Update \mathbf{H} through Eq.(4.10);
 - 13 **end**
 - 14 Rank features according to the values of $\|\mathbf{W}_{i*}\|_2$ in a descending order;
-

With these equations, the detailed algorithm of the proposed SignedFS framework is illustrated in Algorithm 1. At first, we initialize \mathbf{U} , \mathbf{V}^p , \mathbf{V}^n , \mathbf{H} and calculate user proximity matrix and signed Laplacian matrix. From line 5 to 13, we update \mathbf{U} , \mathbf{V}^p , \mathbf{V}^n and \mathbf{W} alternatively until achieving convergence. In each iteration, we first calculate \mathbf{M} , the computation cost of \mathbf{M} is $\mathcal{O}(nd^2)$. After obtain \mathbf{M} , we fix \mathbf{W} and update \mathbf{U} , \mathbf{V}^p and \mathbf{V}^n with gradient descent method. λ_u , λ_p , λ_n is the step size for the update \mathbf{U} , \mathbf{V}^p and \mathbf{V}^n . These step sizes can be determined by line search according to Armijo rule [2]. The computation cost of updating \mathbf{U} , \mathbf{V}^p

and \mathbf{V}^n are $\mathcal{O}(nd^2) + \mathcal{O}(n^2d) + \mathcal{O}(n^2c) + \mathcal{O}(nc^2)$, $\mathcal{O}(nc^2) + \mathcal{O}(n^2c) + \mathcal{O}(n^3)$ and $\mathcal{O}(nc^2) + \mathcal{O}(n^2c) + \mathcal{O}(n^3)$, respectively. Then we employ Eq.(4.11) to update \mathbf{W} , the computational cost of updating \mathbf{W} is $\mathcal{O}(nd^2) + \mathcal{O}(dn^2) + \mathcal{O}(d^3) + \mathcal{O}(d^2c) + \mathcal{O}(ncd)$. After we obtain the local optimal solution of \mathbf{W} , we rank the features in a descending order according to the values of $\|\mathbf{W}_{i*}\|_2$.

4.4 Experiments

In this section, I conduct experiments to evaluate the effectiveness of the proposed SignedFS framework. Details of two real-world datasets used in experiments can be found in Section 2.1. I begin by introducing the experimental settings. After that I present the comparison results between SignedFS and the state-of-the-art unsupervised feature selection methods. Finally, I discuss the impact of negative links and the effects of parameters of SignedFS.

4.4.1 Experimental Setting

Following is a commonly accepted way to assess unsupervised feature selection, I evaluate the proposed SignedFS in terms of clustering performance. To be specific, after I obtain the selected features, I employ K-means clustering based on the selected features. Since K-means may converge in local minimal, I repeat it 20 times and report the average clustering results. Two clustering evaluation metrics, clustering accuracy (ACC) and normalized mutual information(NMI) are used. The higher the ACC and NMI values are, the better the selected features are.

SignedFS is compared with the following state-of-the art unsupervised feature selection algorithms.

- Laplacian Score [16] selects features based on their ability to preserve data manifold structure.

- SPEC [52] evaluates features by spectral regression.
- NDFS [29] selects features by a joint nonnegative spectral analysis and $\ell_{2,1}$ -norm regularization.
- LUFS [45] utilizes social dimension extracted from links to guide feature selection.
- NetFS [27] embeds latent representation extracted from links into feature selection.

Among these baseline methods, Laplacian Score, SPEC and NDFS are traditional unsupervised feature selection methods which only use feature information \mathbf{X} . LUFS and NetFS are unsupervised feature selection algorithms for unsigned networks which only use positive links.

To fairly compare unsupervised feature selection methods, I set the parameters for all methods by a grid search strategy from the range of $\{0.001, 0.01, \dots, 100, 1000\}$. Afterwards, I compare the best clustering results of different feature selection methods.

4.4.2 Quality of Selected Features by SignedFS

In this subsection, I compare the quality of features selected by SignedFS and other above mentioned baseline algorithms. The number of selected features are varied among $\{400, 800, \dots, 1800, 2000\}$. In SignedFS, I have four regularization parameters α , β_1 , β_2 and γ . I empirically set these parameters as $\{\alpha = 1, \beta_1 = 10, \beta_2 = 1000, \gamma = 1000\}$ in Epinions and $\{\alpha = 1, \beta_1 = 1, \beta_2 = 100, \gamma = 1000\}$ in Wiki-rfa. More discussions about these parameters are given in Section 4.4.4. The comparison results of various feature selection algorithms on Epinions and Wiki-rfa datasets are shown in Table 4.1 and Table 4.2. I make the following observations from these two tables:

Table 4.1: Clustering Performance of Different Feature Selection Algorithms in Epinions

Accuracy									
	400	600	800	1000	1200	1400	1600	1800	2000
LapScore	11.48	11.34	10.95	11.79	12.54	11.61	11.29	11.19	12.79
SPEC	21.1	16.93	17.73	17.96	17.91	18.73	18.75	18.57	17.38
NDFS	12.18	11.29	11.92	12.16	12.32	12.14	11.92	13.19	11.78
LUFS	16.23	17.02	18.47	17.44	17.54	19.10	19.29	17.63	18.54
NetFS	18.59	19.62	19.21	18.80	18.43	18.77	17.82	19.76	19.98
SignedFS	23.24	21.76	21.69	22.11	21.27	21.88	20.04	20.64	21.20
NMI									
	400	600	800	1000	1200	1400	1600	1800	2000
LapScore	0.0274	0.0272	0.0268	0.0368	0.0268	0.0273	0.0275	0.0263	0.0267
SPEC	0.0166	0.0175	0.0174	0.0241	0.0250	0.0253	0.0264	0.0262	0.0260
NDFS	0.0149	0.0147	0.0146	0.0146	0.0146	0.0146	0.0146	0.0146	0.0146
LUFS	0.0161	0.0160	0.0176	0.0182	0.0186	0.0189	0.0191	0.0172	0.0199
NetFS	0.0180	0.0190	0.0228	0.0175	0.0179	0.0111	0.0156	0.0147	0.0208
SignedFS	0.0382	0.0368	0.0384	0.0387	0.0372	0.0384	0.0400	0.0386	0.0379

- SignedFS consistently outperforms traditional feature selection algorithms LapScore, SPEC and NDFS on both datasets with significant clustering performance gain in most cases. I also perform pairwise Wilcoxon signed-rank test between SignedFS and these three traditional unsupervised feature selection methods, it shows SignedFS is significantly better (p -value=0.05). The superiority of SignedFS can be attributed to the utilization additional link information while traditional methods are mainly based on the data i.i.d. assumption.
- SignedFS also obtains better clustering performance than the other two feature selection methods LUFS and NetFS on linked data. A major reason is that LUFS and NetFS only exploit positive links while SignedFS leverages both

Table 4.2: Clustering Performance of Different Feature Selection Algorithm in Wiki-rfa

Accuracy									
	400	600	800	1000	1200	1400	1600	1800	2000
LapScore	70.92	70.94	70.93	70.31	70.52	70.89	70.92	71.13	71.37
SPEC	71.76	72.11	72.02	71.76	71.76	71.90	71.76	71.83	71.56
NDFS	72.94	72.73	72.94	72.75	72.78	72.94	72.94	72.94	72.94
LUFS	75.55	75.55	73.79	74.11	74.14	73.24	73.21	73.28	73.89
NetFS	72.81	72.91	72.94	72.73	72.68	72.70	72.97	72.97	72.97
SignedFS	79.10	79.52	78.59	78.15	78.18	78.63	79.27	80.94	80.63
NMI									
	400	600	800	1000	1200	1400	1600	1800	2000
LapScore	0.0026	0.0026	0.0026	0.0026	0.0026	0.0026	0.0026	0.0026	0.0026
SPEC	0.0070	0.0105	0.0083	0.0045	0.0029	0.0030	0.0026	0.0026	0.0026
NDFS	0.0026	0.0030	0.0028	0.0025	0.0026	0.0026	0.0027	0.0026	0.0026
LUFS	0.0014	0.0014	0.0007	0.0004	0.0005	0.0007	0.0007	0.0008	0.0011
NetFS	0.0044	0.0038	0.0036	0.0036	0.0038	0.0038	0.0037	0.0037	0.0037
SignedFS	0.0154	0.0157	0.0149	0.0147	0.0140	0.0156	0.0181	0.0337	0.0334

positive links and negative links into a coherent model to obtain better features. It indicates the potential of using negative links for feature selection. I will further discuss how negative links affect the performance of feature selection in Section 4.4.3.

- We can see that when we gradually increase the number of selected features from 400 to 2000, the clustering performance in terms of ACC and NMI does not vary a lot. In particular, when a small number of features are selected, SignedFS already gives us very good performance. For example, when 400 features are selected in the Epinions dataset, the clustering performance is already very high. A small number of selected features is very appealing in practice as

it significantly reduces the memory storage costs and computational costs for further learning tasks.

4.4.3 Impact of Negative Links

In Table 4.1 and Table 4.2, I have already shown that compared with the methods which only leverage positive links, SignedFS shows effectiveness in improving clustering performance. In this subsection, I further investigate how the negative links help select relevant features. As can be shown in the objective function of SignedFS in Eq.(4.8), I have two terms involving the negative links, the first term $\|\mathbf{O}^n \odot (\mathbf{A}^n - \mathbf{U}\mathbf{V}^n\mathbf{U}')\|_F^2$ models the negative links for user preference representation learning while the second term $tr(\mathbf{U}'\mathbf{L}\mathbf{U})$ leverages negative links to capture the first-order and the second-order proximity in signed social networks. To study how these two terms affect the performance of feature selection, I define the following variants to eliminate the effects of negative links from SignedFS framework.

- SignedFS\I: I eliminate the term which uses the negative links in learning user preference representation ($\|\mathbf{O}^n \odot (\mathbf{A}^n - \mathbf{U}\mathbf{V}^n\mathbf{U}')\|_F^2$) by setting $\beta_2 = 0$.
- SignedFS\II: I eliminate the term that leverages negative links in modeling user proximity for signed social networks ($tr(\mathbf{U}'\mathbf{L}\mathbf{U})$) by setting $\gamma = 0$.
- SignedFS\I,II: I eliminate both terms mentioned above by setting $\beta_2 = 0$ and $\gamma = 0$.

I compare these three variants of SignedFS with the original SignedFS framework, and the performance comparison results are shown in Figure 4.2. Due to space limit, I only show the results on the Wiki-rfa dataset as we have the similar observations on the Epinions. From the figure, we can see that SignedFS\I, SignedFS\II has significantly lower clustering performance than the SignedFS framework. It demonstrates

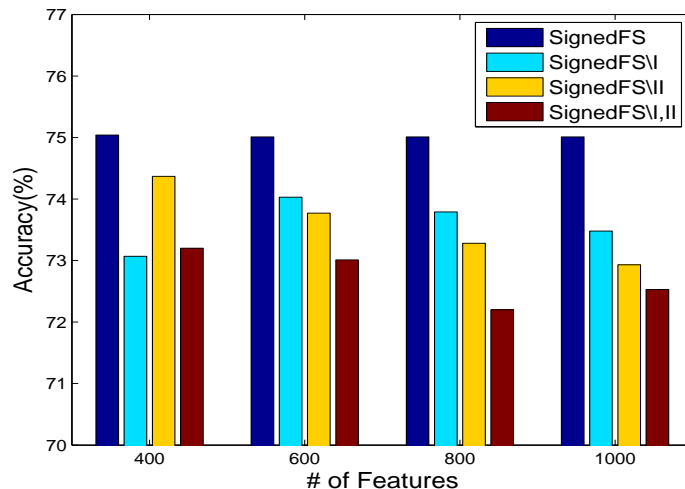


Figure 4.2: The Impact of Negative Links for SignedFS on Wiki-rfa

the effectiveness of leveraging negative links in modeling user preference representation and user proximity for unsupervised feature selection. When both terms are eliminated, SignedFS\I,II obtains even lower clustering performance when the number of selected features is varied from 400 to 1000. It further validates the potential of using negative links for feature selection.

4.4.4 Parameter Analysis

The proposed SignedFS has four important parameters. α controls the sparsity of the model. β_1 and β_2 balances the contribution of positive and negative links in learning user preference representation. γ controls the modeling of user proximity for feature selection. I study the effect of each parameter by fixing the others to investigate how it affects the performance of feature selection. Since I make the similar observation on both datasets, I only report the experimental result of ACC on Wiki-rfa dataset to save space.

First, I fix $\{\beta_1 = 10, \beta_2 = 10, \gamma = 100\}$ and vary α as $\{0.01, 0.1, 1, 10, 100\}$. As shown in Figure 4.3 (a), the clustering performance first increases and then reaches

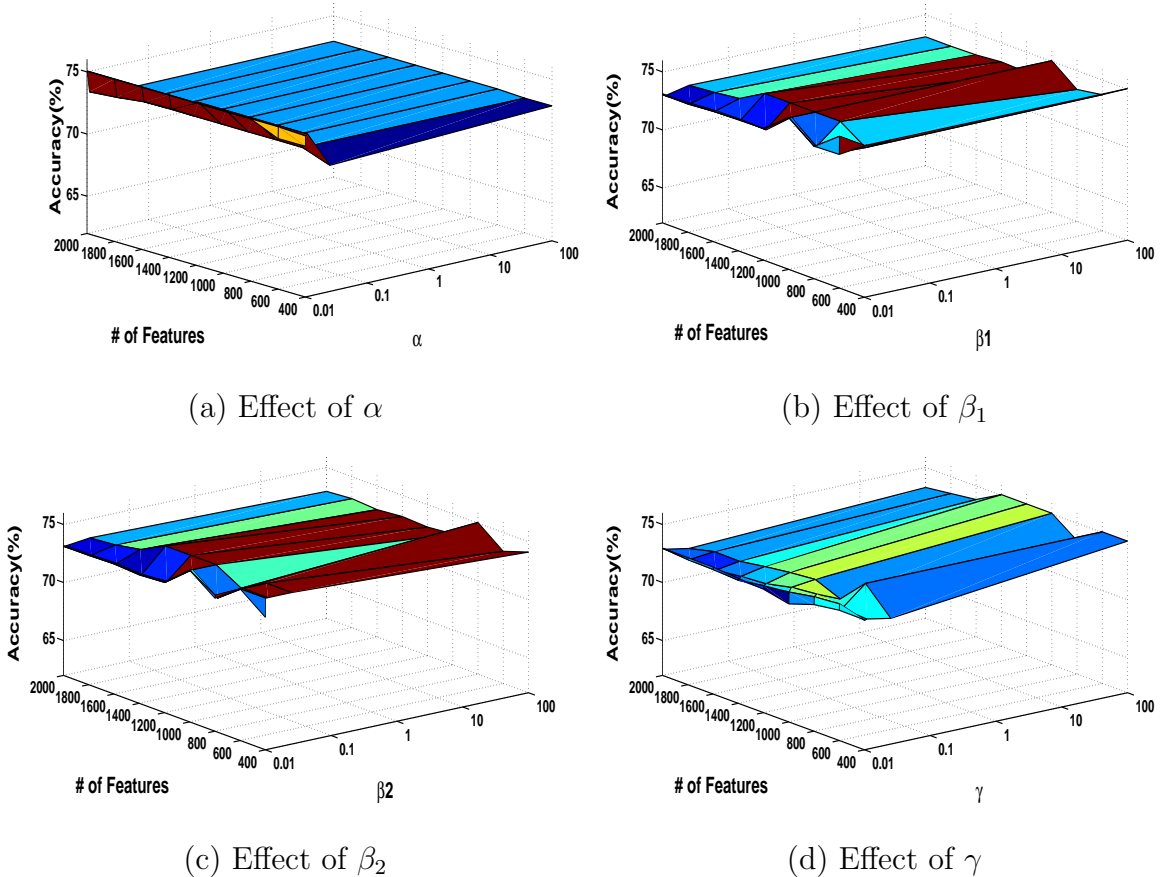


Figure 4.3: Parameter Analysis of SignedFS on Wikipedia

the peak values when $\alpha = 0.1$. If I continuously increase the value of α , the clustering performance decreases. Therefore, I could empirically set the values of α among the range of 0.01 to 1. Second, to investigate how β_1 affects the clustering performance, I vary β_1 as $\{0.01, 0.1, 1, 10, 100\}$ by fixing $\{\alpha = 1, \beta_2 = 10, \gamma = 100\}$. The result is presented in Figure 4.3 (b). Similarly, the clustering performance first increases, reaches its maximal value when $\beta_1 = 10$ and then degrades. Next, to study the impact of β_2 , I set $\{\alpha = 1, \beta_1 = 10, \gamma = 100\}$, and vary β_2 as $\{0.01, 0.1, 1, 10, 100\}$. The result is presented in Figure 4.3 (c). The performance variation w.r.t. β_2 has a similar trend as the variation of β_1 , which suggests that negative links are equally important as positive links in finding relevant features. Finally, I fix $\{\alpha = 1, \beta_1 = 10, \beta_2 = 10\}$ and

vary γ as $\{0.01, 0.1, 1, 10, 100\}$ to investigate the effect of γ . As depicted in Figure 4.3 (d), with the increase of γ , the clustering performance gradually increases and then keeps stable. The clustering performance is relatively more sensitive to the number of selected features than these regularization parameters, which is still an open problem in unsupervised feature selection.

UNSUPERVISED SENTIMENT ANALYSIS IN SIGNED SOCIAL NETWORKS

Huge volumes of opinion-rich data is user-generated in social media at an unprecedented rate, easing the analysis of individual and public sentiments. Sentiment analysis has shown to be useful in probing and understanding emotions, expressions and attitudes in the text. However, the distinct characteristics of social media data present challenges to traditional sentiment analysis [4, 20, 21, 39, 46]. First, social media data is often noisy, incomplete and fast-evolved which necessitates the design of a sophisticated learning model. Second, sentiment labels are hard to collect which further exacerbates the problem by not being able to discriminate sentiment polarities. Meanwhile, opportunities are also unequivocally presented. Social media contains rich sources of sentiment signals in textual terms and user interactions [1, 15], which could be helpful in sentiment analysis. While there are some attempts to leverage implicit sentiment signals in positive user interactions [20, 36, 40, 46], little attention is paid to signed social networks with both positive and negative links. The availability of signed social networks motivates us to investigate if negative links also contain useful sentiment signals [25].

In this chapter, with the preliminary analysis on negative links in Section 2.2, I study the problem of sentiment analysis with signed social networks under an unsupervised scenario. In essence, I aim to answer the following questions: (1) how to employ sentiment signals revealed by negative links in Section 2.1 for sentiment analysis in signed social networks? (2) how to explicitly model positive and negative interactions among users for sentiment analysis in an unsupervised way? To answer the questions, I propose an unsupervised sentiment analysis framework - SignedSenti.

5.1 Problem Statement

To formally define the problem unsupervised feature selection on signed social networks, I first present the notations.

Let $\mathcal{T} = \{t_1, t_2, \dots, t_m\}$ be a set of m text posts and $\mathcal{F} = \{f_1, f_2, \dots, f_d\}$ be a set of d textual terms. As shown in Figure 5.1, the matrix representation of \mathcal{T} is $\mathbf{X} \in \mathbb{R}^{m \times d}$. Each text post may be a review or a comment for a product or an article, respectively. Assume these m text posts are describing a set of l items $\mathcal{O} = \{o_1, \dots, o_l\}$ (e.g., $\{o_1, \dots, o_4\}$ in Figure 5.1). Their relations are encoded in a text-item relation matrix $\mathbf{O} \in \{0, 1\}^{m \times l}$ where $\mathbf{O}_{i,j} = 1$ if text post t_i is about item o_j , otherwise $\mathbf{O}_{i,j} = 0$. Also, we assume that these m text posts are generated by n distinct social media users $\mathcal{U} = \{u_1, u_2, \dots, u_n\}$. Matrix $\mathbf{T} \in \{0, 1\}^{n \times m}$ shows the authorship between users and text posts such that $\mathbf{T}_{i,j} = 1$ if text post t_j is posted by user u_i , $\mathbf{T}_{i,j} = 0$ otherwise. In addition to positive user interactions, social media users can also be negatively connected, I use $\mathbf{A} \in \mathbb{R}^{n \times n}$ to denote the signed adjacency matrix where $\mathbf{A}_{ij} = 1$, $\mathbf{A}_{ij} = -1$ and $\mathbf{A}_{ij} = 0$ represent positive, negative and missing links from user u_i to u_j , respectively. The relations among posts \mathcal{T} , items \mathcal{O} and users \mathcal{U} are shown in the middle of Figure 5.1; while an illustration of matrices \mathbf{O} , \mathbf{T} and \mathbf{A} are demonstrated at the bottom of Figure 5.1.

With above notations and preliminary validation of signed link based partial order assumption in Section 2.2, I now define the problem of unsupervised sentiment analysis with signed social networks can be as follows:

Given: a set of social media posts \mathcal{T} , a set of items \mathcal{O} , a set of social media users \mathcal{U} , and available relations including the user-text relation \mathbf{T} , user-user relation \mathbf{A} (either positive or negative) and text-item relation \mathbf{O} ;

Infer: the sentiment polarities of all posts in \mathcal{T} .

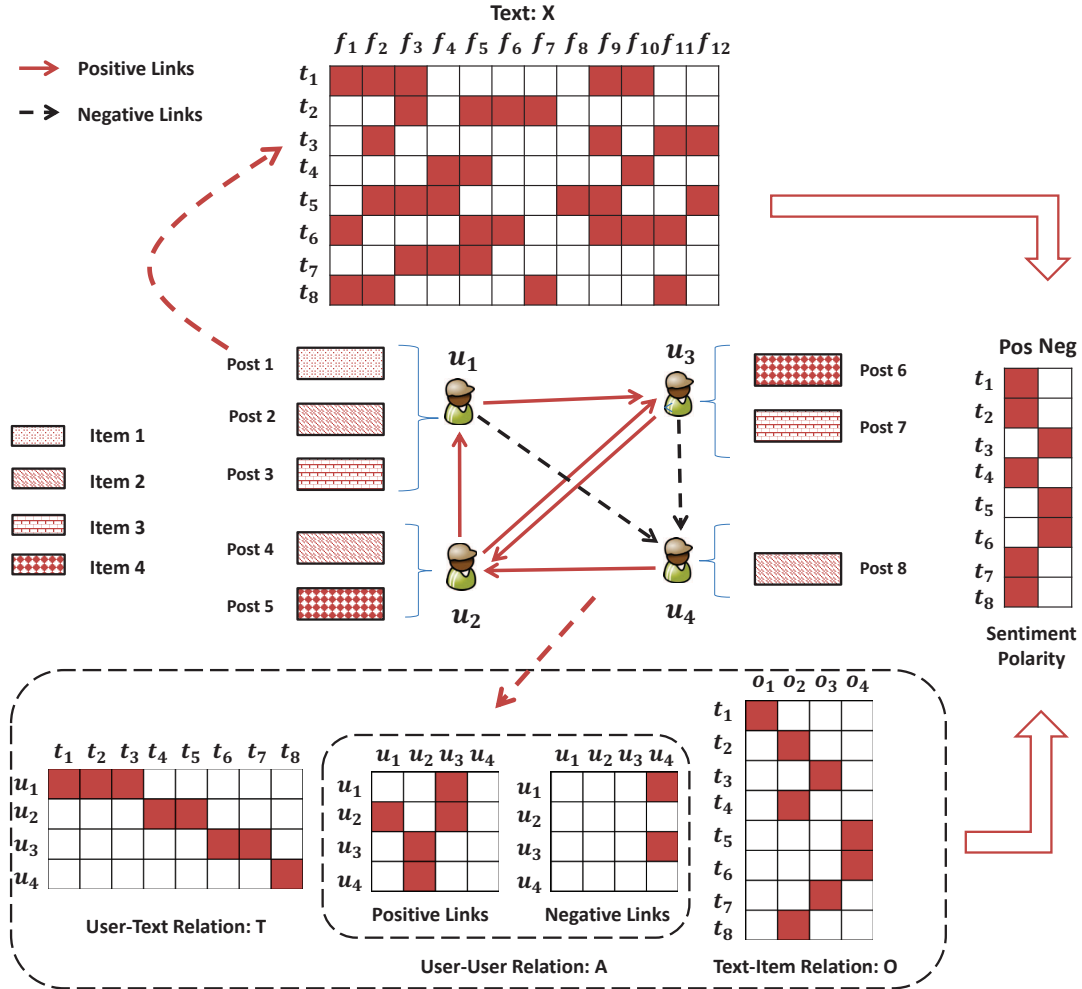


Figure 5.1: An Illustration of Unsupervised Sentiment Analysis With Signed Social Networks.

5.2 The Proposed Framework-SignedSenti

In this section, I discuss how to model both positive and negative user interactions in understanding and predicting sentiment polarities in an unsupervised scenario.

5.2.1 Basic Model for Unsupervised Sentiment Analysis

Unsupervised sentiment analysis is naturally a clustering problem. Specifically, I would like to cluster text posts into k different sentiment groups. Let $\mathbf{U} \in \mathbb{R}^{m \times k}$ be the text-sentiment cluster matrix such that $\mathbf{U}_{ij} = 1$ if text post t_i belongs to class

c_j , and $\mathbf{U}_{ij} = 0$ otherwise. In essence, it can be modeled by solving the following nonnegative matrix factorization problem:

$$\begin{aligned} \min_{\mathbf{U}, \mathbf{V}} \|\mathbf{X} - \mathbf{UV}'\|_F^2 + \gamma(\|\mathbf{U}\|_F^2 + \|\mathbf{V}\|_F^2) \\ \text{s.t. } \mathbf{U} \geq 0, \mathbf{V} \geq 0, \mathbf{U} \in \{0, 1\}^{m \times k}, \mathbf{U}'\mathbf{1} = \mathbf{1}, \end{aligned} \quad (5.1)$$

where $\mathbf{V} \in \mathbb{R}^{d \times k}$ is a term-sentiment matrix, and each row of \mathbf{V} shows the distribution of each term in these k sentiment groups. $\gamma(\|\mathbf{U}\|_F^2 + \|\mathbf{V}\|_F^2)$ is introduced to avoid overfitting.

5.2.2 Sentiment Signals from Textual Terms

It has been widely studied in literature [48] that the overall sentiment of a text post is strongly correlated with sentiment of terms in the post. In other words, some terms may contain strong sentiment signals in identifying sentiment polarities. For example, the words of “wonderful” and “appealing” in a text post may express positive emotions while the words of “terrible” and “disappointed” could express negative emotions. The rich sentiment signals in terms help to bridge the gap between the difficulties in obtaining sentiment labels and the necessity of label supervision in sentiment analysis. To leverage sentiment signals in rich textual information, I employ a widely used sentiment lexicon SentiWordNet [12] to obtain sentiment polarities of terms. SentiWordNet contains positive, negative and objective scores between 0 and 1 for all synsets in WordNet. In WordNet, there are a total of 117,659 words and phrases. Let $\mathbf{P} \in \mathbb{R}^{d \times k}$ be a term-sentiment indication matrix which encodes sentiment signals of words. Since our task is polarity sentiment analysis, I set $k = 2$ and let \mathbf{P}_{i1} denote the positive score of term f_i while \mathbf{P}_{i2} represents the negative score of term f_i . To take advantage of the textual sentiment signal, I force the above term-sentiment matrix \mathbf{V} in the base model to be consistent with the term-sentiment

indication matrix \mathbf{P} by minimizing:

$$\min_{\mathbf{V}} \|\mathbf{V} - \mathbf{P}\|_F^2. \quad (5.2)$$

It should be noted that the number of sentiment signals, i.e., k should be adapted according to the needs whether to perform binary or multi-class sentiment polarity analysis.

5.2.3 Exploiting Positive and Negative Interactions

The signed link based partial order assumption suggests that for each text post, its sentiment is more similar to posts in its positive linked set than posts in its negative linked set. In other words, it indicates that friends are more likely to reveal similar sentiments than foes on the same item. As $\mathbf{U} \in \mathbb{R}^{m \times k}$ denotes the sentiment polarity hard assignment matrix, I use $\|\mathbf{U}_{i*} - \mathbf{U}_{j*}\|_2^2$ to represent the sentiment similarity between two text posts t_i and t_j . To model the signed link based partial order assumption, there are two cases that I need to discuss. For each text post t_i , (1) if another text post t_j in its positive linked set is more closer to the text post t_k in its negative linked set, i.e., $\|\mathbf{U}_{i*} - \mathbf{U}_{j*}\|_2^2 - \|\mathbf{U}_{i*} - \mathbf{U}_{k*}\|_2^2 < 0$, I do not need to penalize it; (2) if its negative linked set is more closer to its positive linked set, i.e., $\|\mathbf{U}_{i*} - \mathbf{U}_{j*}\|_2^2 - \|\mathbf{U}_{i*} - \mathbf{U}_{k*}\|_2^2 > 0$, I should add a penalty to pull the sentiment of t_i be more closer to t_j than to t_k . Mathematically, it can be formulated by solving the following objective function:

$$\min \sum_{(i,j,k) \in \Omega} \max(0, \|\mathbf{U}_{i*} - \mathbf{U}_{j*}\|_2^2 - \|\mathbf{U}_{i*} - \mathbf{U}_{k*}\|_2^2), \quad (5.3)$$

where Ω denotes all triplets that satisfies the signed link based partial order assumption, i.e., $\Omega = \{(i, j, k) | i \in \mathcal{T}, j \in \mathcal{P}(t_i), k \in \mathcal{N}(t_i)\}$. The above penalty term can be

further reformulated as:

$$\begin{aligned}
& \sum_{(i,j,k) \in \Omega} \max(0, \|\mathbf{U}_{i*} - \mathbf{U}_{j*}\|_2^2 - \|\mathbf{U}_{i*} - \mathbf{U}_{k*}\|_2^2) \\
& = \sum_{(i,j,k) \in \Omega} w_{ij}^k \text{tr}(\mathbf{M}_{ij}^k \mathbf{U} \mathbf{U}'),
\end{aligned} \tag{5.4}$$

where \mathbf{M} is a sparse matrix with all entries equal to zero except that $\mathbf{M}_{ij} = \mathbf{M}_{ji} = \mathbf{M}_{kk} = -1$ and $\mathbf{M}_{ik} = \mathbf{M}_{ki} = \mathbf{M}_{jj} = 1$. \mathbf{M}_{ij}^k is the matrix \mathbf{M} with elements associated with triplet (i, j, k) and w_{ij}^k is defined as follows:

$$w_{ij}^k = \begin{cases} 1 & \text{if } \text{tr}(\mathbf{M}_{ij}^k \mathbf{U} \mathbf{U}') > 0 \\ 0 & \text{otherwise} \end{cases}. \tag{5.5}$$

5.2.4 Objective Function of SignedSenti

With the model components of sentiment signals from terms and the signed link based partial order assumption, the final objective function of unsupervised sentiment analysis with signed social network can be formulated as follows:

$$\begin{aligned}
& \min_{\mathbf{U}, \mathbf{V}} \|\mathbf{X} - \mathbf{U} \mathbf{V}'\|_F^2 + \alpha \sum_{(i,j,k) \in \Omega} w_{ij}^k \text{tr}(\mathbf{M}_{ij}^k \mathbf{U} \mathbf{U}') \\
& + \beta \|\mathbf{V} - \mathbf{P}\|_F^2 + \gamma (\|\mathbf{U}\|_F^2 + \|\mathbf{V}\|_F^2) \\
& \text{s.t } \mathbf{U} \geq 0, \mathbf{V} \geq 0, \mathbf{U} \in \{0, 1\}^{m \times k}, \mathbf{U}' \mathbf{1} = \mathbf{1}.
\end{aligned} \tag{5.6}$$

Parameters α and β control the contribution of sentiment signals from terms and signed social networks, respectively.

The problem in Eq. (5.6) is difficult to solve due to the discrete constraint on \mathbf{U} . To tackle this issue, I relax the objective function by reformulating it as an orthogonal constraint. After the relaxation, Eq.(5.6) can be rewritten as:

$$\begin{aligned}
& \min_{\mathbf{U}, \mathbf{V}} \|\mathbf{X} - \mathbf{U} \mathbf{V}'\|_F^2 + \alpha \sum_{(i,j,k) \in \Omega} w_{ij}^k \text{tr}(\mathbf{M}_{ij}^k \mathbf{U} \mathbf{U}') \\
& + \beta \|\mathbf{V} - \mathbf{P}\|_F^2 + \gamma (\|\mathbf{U}\|_F^2 + \|\mathbf{V}\|_F^2) \\
& \text{s.t } \mathbf{U} \geq 0, \mathbf{V} \geq 0, \mathbf{U}' \mathbf{U} = \mathbf{I}.
\end{aligned} \tag{5.7}$$

5.3 Optimization Algorithm for SignedSenti

The objective function of the proposed SignedSenti framework is not convex w.r.t. both \mathbf{U} and \mathbf{V} simultaneously. Hence, I introduce an alternating algorithm to solving its optimization problem.

Update \mathbf{U} : First, I fix \mathbf{V} to update \mathbf{U} . Specifically, when \mathbf{V} is fixed, the objective function is convex w.r.t. the text-sentiment matrix \mathbf{U} . Thus, \mathbf{U} can be obtained by solving the following optimization problem:

$$\begin{aligned} \min_{\mathbf{U}} \mathcal{J}(\mathbf{U}) &= \|\mathbf{X} - \mathbf{U}\mathbf{V}'\|_F^2 + \alpha \sum_{(i,j,k) \in \Omega} w_{ij}^k \text{tr}(\mathbf{M}_{ij}^k \mathbf{U}\mathbf{U}') + \gamma \|\mathbf{U}\|_F^2 \\ \text{s.t. } \mathbf{U} &\geq 0, \mathbf{U}'\mathbf{U} = \mathbf{I}. \end{aligned} \quad (5.8)$$

The Lagrangian of Eq. (5.8) is:

$$\begin{aligned} \min_{\mathbf{U}} \mathcal{L}(\mathbf{U}) &= \|\mathbf{X} - \mathbf{U}\mathbf{V}'\|_F^2 + \alpha \sum_{(i,j,k) \in \Omega} w_{ij}^k \text{tr}(\mathbf{M}_{ij}^k \mathbf{U}\mathbf{U}') \\ &+ \gamma \|\mathbf{U}\|_F^2 + \text{tr}(\Gamma_u(\mathbf{U}'\mathbf{U} - \mathbf{I})) - \text{tr}(\Lambda_u \mathbf{U}'). \end{aligned} \quad (5.9)$$

where Γ_u and Λ_u are the Lagrange multipliers for constraints $\mathbf{U}'\mathbf{U} = \mathbf{I}$ and $\mathbf{U} \geq 0$, respectively. To compute \mathbf{U} , I take the partial derivative of Eq. (5.9) w.r.t. \mathbf{U} and set it to be zero:

$$\Lambda_u = 2(\mathbf{U}\mathbf{V}'\mathbf{V} - \mathbf{X}\mathbf{V} + \gamma\mathbf{U} + \mathbf{U}\Gamma_u) + \alpha \sum_{(i,j,k) \in \Omega} w_{ij}^k (\mathbf{M}_{ij}^k \mathbf{U} + \mathbf{M}_{ij}^{k'} \mathbf{U}). \quad (5.10)$$

With the KKT complementary condition for the nonnegativity constraint of \mathbf{U} , i.e., $(\Lambda_u)_{ij} \mathbf{U}_{ij} = 0$, I have:

$$\begin{aligned} (\mathbf{U}\mathbf{V}'\mathbf{V} - \mathbf{X}\mathbf{V} + \gamma\mathbf{U} + \frac{\alpha}{2} \sum_{(i,j,k) \in \Omega} w_{ij}^k (\mathbf{M}_{ij}^k \mathbf{U} + \mathbf{M}_{ij}^{k'} \mathbf{U}) \\ + \mathbf{U}\Gamma_u)_{ij} \mathbf{U}_{ij} = 0, \text{ where} \end{aligned} \quad (5.11)$$

$$\Gamma_u = -\frac{\alpha}{2} \sum_{(i,j,k) \in \Omega} w_{ij}^k (\mathbf{U}'(\mathbf{M}_{ij}^k \mathbf{U} + \mathbf{M}_{ij}^{k'} \mathbf{U})) - \mathbf{V}'\mathbf{V} + \mathbf{U}'\mathbf{X}\mathbf{V} - \gamma\mathbf{I}. \quad (5.12)$$

It leads to the following update rule for \mathbf{U} :

$$\mathbf{U}_{ij} \leftarrow \mathbf{U}_{ij} \sqrt{\frac{\mathbf{B}_{ij}}{\mathbf{E}_{ij}}}, \text{ where} \quad (5.13)$$

$$\mathbf{B} = 2\mathbf{X}\mathbf{V} + \alpha \sum_{(i,j,k) \in \Omega} w_{ij}^k (\mathbf{M}_{ij}^k \mathbf{U} + \mathbf{M}_{ij}^{k'} \mathbf{U})^- + 2\mathbf{U}\Gamma_u^-, \quad (5.14)$$

$$\mathbf{E} = 2(\mathbf{U}\mathbf{V}'\mathbf{V} + \gamma\mathbf{U}) + \alpha \sum_{(i,j,k) \in \Omega} w_{ij}^k (\mathbf{M}_{ij}^k \mathbf{U} + \mathbf{M}_{ij}^{k'} \mathbf{U})^+ + 2\mathbf{U}\Gamma_u^+. \quad (5.15)$$

Update \mathbf{V} : Likewise, I fix \mathbf{U} to update \mathbf{V} . When \mathbf{U} is fixed, the objective function is convex w.r.t. the term-sentiment matrix \mathbf{V} . Hence, \mathbf{V} can be obtained by solving:

$$\begin{aligned} \min_{\mathbf{V}} \mathcal{J}(\mathbf{V}) &= \|\mathbf{X} - \mathbf{U}\mathbf{V}'\|_F^2 + \beta \|\mathbf{V} - \mathbf{P}\|_F^2 + \gamma \|\mathbf{V}\|_F^2 \\ \text{s.t. } \mathbf{V} &\geq 0. \end{aligned} \quad (5.16)$$

The Lagrangian of Eq. (5.16) is:

$$\mathcal{L}(\mathbf{V}) = \|\mathbf{X} - \mathbf{U}\mathbf{V}'\|_F^2 + \beta \|\mathbf{V} - \mathbf{P}\|_F^2 + \gamma \|\mathbf{V}\|_F^2 - \text{tr}(\Lambda_v \mathbf{V}'), \quad (5.17)$$

where Λ_v is the Lagrange multipliers for the constraints $\mathbf{V} \geq 0$. I take the partial derivative of Eq. (5.17) w.r.t. \mathbf{V} and set it to be zero:

$$\Lambda_v = 2(\mathbf{V}\mathbf{U}'\mathbf{U} - \mathbf{X}'\mathbf{U} + \beta(\mathbf{V} - \mathbf{P}) + \gamma\mathbf{V}). \quad (5.18)$$

Similarly, with the KKT complementary condition for the nonnegativity constraint of \mathbf{V} , i.e., $(\Lambda_v)_{ij} \mathbf{V}_{ij} = 0$, we have:

$$2(\mathbf{V}\mathbf{U}'\mathbf{U} - \mathbf{X}'\mathbf{U} + \beta(\mathbf{V} - \mathbf{P}) + \gamma\mathbf{V})_{ij} \mathbf{V}_{ij} = 0, \quad (5.19)$$

which leads to the following update rule for \mathbf{V} :

$$\mathbf{V}_{ij} \leftarrow \mathbf{V}_{ij} \sqrt{\frac{\mathbf{X}'\mathbf{U} + \beta\mathbf{P}}{\mathbf{V}\mathbf{U}'\mathbf{U} + (\beta + \gamma)\mathbf{V}}}. \quad (5.20)$$

With these update rules, the detailed algorithm of the proposed SignedSenti framework is illustrated in Algorithm 2. At the very beginning, we initialize \mathbf{U} , \mathbf{V} randomly

and calculate \mathbf{M} from \mathbf{T} , \mathbf{A} and \mathbf{O} . From line 3 to 7, we update \mathbf{U} and \mathbf{V} iteratively until converge. To update \mathbf{U} , we need to calculate w_{ij}^k and Γ_u at first. According to Eq.(5.5) and Eq.(5.12), the computation cost of obtaining w_{ij}^k and Γ_u are $\mathcal{O}(m^2k)$ and $\mathcal{O}(k^2d + kmd + m^2k + k^2m)$ respectively. With w_{ij}^k and Γ_u , we employ Eq.(5.13) to update \mathbf{U} , the computational cost of updating \mathbf{U} is $\mathcal{O}(kmd + m^2k + k^2m)$. The total cost of computing \mathbf{V} according to Eq.(5.20) is $\mathcal{O}(kmd)$. After we obtain \mathbf{U} , sentiment polarities of text texts can be obtained by performing K-Means on \mathbf{U} .

Algorithm 2: SignedSenti Algorithm

Input : $\{\mathbf{X}, \mathbf{T}, \mathbf{A}, \mathbf{O}, \mathbf{P}, k, \alpha, \beta, \gamma\}$

Output: sentiment polarity for each text post.

- 1 Initialize \mathbf{U} , \mathbf{V} randomly;
 - 2 Compute \mathbf{M} based on \mathbf{T} , \mathbf{A} and \mathbf{O} ;
 - 3 **while** *not converge* **do**
 - 4 Calculate w_{ij}^k according to Eq.(5.5) ;
 - 5 Compute Γ_u according to Eq.(5.12) ;
 - 6 Update \mathbf{U} according to Eq.(5.13);
 - 7 Update \mathbf{V} according to Eq.(5.20);
 - 8 **end**
 - 9 Employing \mathbf{U} to predict sentiment polarity of text posts.
-

5.4 Experiments

In this section, I conduct experiments to evaluate the effectiveness of the proposed SignedSenti framework. I begin by introducing the experimental settings. After that, I present the comparison results between SignedSenti and the state-of-the-art unsupervised sentiment analysis methods. Finally, I discuss the sensitivity of parameters.

5.4.1 Experimental Setting

Following a common way to assess the performance of unsupervised sentiment analysis, I take clustering accuracy as the evaluation metric. Higher clustering accuracy often indicates better performance. SignedSenti is compared with the following baseline methods:

- **SentiStrength** [47]: SentiStrength is a lexicon-based unsupervised method that extracts sentiment strength from informal English with pre-defined sentiment lexicon.
- **MPQA** [50]: It predicts sentiment polarity of text posts according to a manually labeled sentiment lexicon MPQA.
- **SentiWordNet** [12]: It determines sentiment scores of text posts via a widely used sentiment lexicon SentiWordNet.
- **K-Means**: As one of the most representative clustering methods, it partitions the text posts into k sentiment polarities on the original textual terms.
- **NMF** [34]: Nonnegative matrix factorization is a popular method in text mining. It is also a variant of the proposed SignedSenti model by setting $\alpha = \beta = 0$.
- **SignedSenti-T**: It is a variant of the proposed SignedSenti that only employs the textual information for sentiment analysis. Specifically, I set $\alpha = 0$.
- **SignedSenti-L**: It is a variant of the proposed SignedSenti that does not explicitly leverage sentiment signals from textual terms. In particular, I set $\beta = 0$.

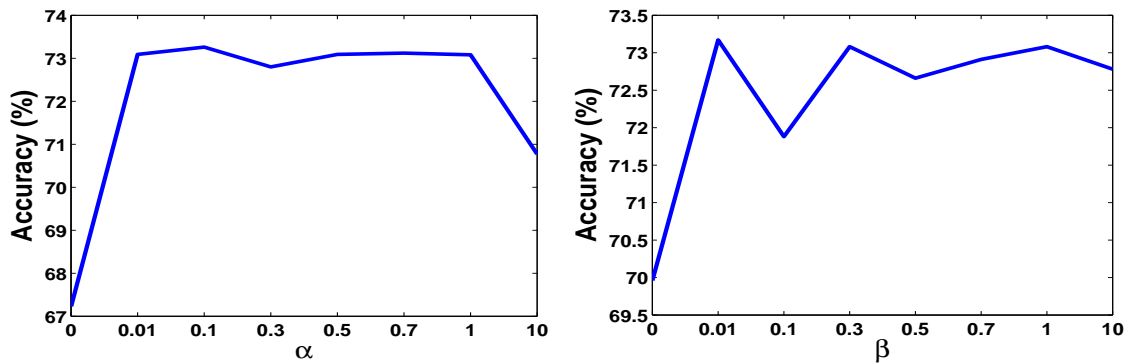
5.4.2 Sentiment Polarity Prediction Performance

In this subsection, I compare SignedSenti with other baseline algorithms shown in Section 5.4.1. Noticed that in SigendSenti, we have three regularization parameters α , β , γ . I empirically set these parameters as $\{\alpha = 1, \beta = 0.5, \gamma = 0.7\}$ in Epinions and $\{\alpha = 1, \beta = 1, \gamma = 0.1\}$ in Slashdot. More discussions about the effectiveness of these parameters will be presented later. The comparison results of various unsupervised sentiment analysis algorithms on Epinions and Slashdot datasets are shown in Table 5.1. I make the following observations:

- SignedSenti consistently outperforms other baseline methods on both datasets with significant performance gain. I also perform pairwise Wilcoxon signed-rank test [8] between SignedSenti and these baseline methods, it shows SignedSenti is significantly better with a significance level of 0.05. The superiority of the proposed SignedSenti can be attributed to the utilization of external sources, including textual sentiment signals and positive (negative) user interactions.
- In general, traditional lexicon-based unsupervised methods such as SentiStrength, MPQA and SentiWordNet do not perform well in the unsupervised case. These observations show the necessity to build a sophisticated learning model to automatically predict the sentiment polarities of text posts.
- SignedSenti also obtains better performance than traditional document clustering methods K-Means and NMF. The reason is that social media texts are often noisy and incomplete, hence without the guide of any sentiment signals or user interactions, it is difficult to discriminate the sentiment polarities of different text posts.
- The clustering accuracy of SignedSenti is higher than its variant SignedSenti-

Table 5.1: Sentiment Polarity Prediction Accuracy.

Method	Epinions	Slashdot
SentiStrength	0.521	0.628
MPQA	0.662	0.684
SentiWordNet	0.645	0.586
K-Means	0.644	0.677
NMF	0.637	0.648
SignedSenti-T	0.649	0.672
SignedSenti-L	0.714	0.700
SignedSenti	0.723	0.731



(a) Effect of α (b) Effect of β
Figure 5.2: Parameter Analysis of SignedSenti on Slashdot.

T. SignedSenti-T only leverages sentiment signals from terms and does not explicitly consider user interactions. Its inferiority to SignedSenti indicates that in addition to textual sentiment signals, positive and negative links also contain implicit rich sentiment signals that can boost the sentiment polarity prediction.

5.4.3 Parameter Analysis

The proposed SignedSenti has two important parameters α and β which controls the contribution of implicit sentiment signals from positive (negative) user interactions and textual terms respectively. I study the effect of each parameter by fixing the other to investigate how it affects the clustering performance. I only report the experimental result on Slashdot as we have similar observations on Epinions. In particular, I first fix $\{\beta = 1, \gamma = 0.1\}$ and vary α as $\{0, 0.01, 0.1, 0.3, 0.5, 0.7, 1, 10\}$. As shown in Figure 5.2(a), when α increase from 0 to 0.01 the performance increases dramatically which further validates the effectiveness of leveraging implicit sentiment signals in positive and negative interactions. If I continuously increase α , the performance is relatively stable in fairly large ranges $[0.01, 1]$, then it decreases when $\alpha > 1$. Similarly, to investigate how β affects the performance, I vary β as $\{0, 0.01, 0.1, 0.3, 0.5, 0.7, 1, 10\}$ by fixing $\{\alpha = 1, \gamma = 0.1\}$. The result is presented in Figure 5.2(b). Likewise, the performance increases significantly at the very beginning due to the increase of β from 0 to 0.01. After that, with the increase of β , the performance fluctuates in ranges of 71.5 and 73.5. To summary, the clustering performance is rather stable when I tune these two parameters in a wide range, which is very appealing in practice.

CONCLUSION AND FUTURE WORK

As the most distinct characteristic of social media, various types of social relations are essential for social data mining. A majority of existing methods for social data mining only consider positive interactions among connected instances while negative links are also prevailing in real-world social networks such as distrust relations in Epinions and foe links in Slashdot. Even though negative links have some added value over positive links, it is difficult to directly employ them for learning tasks because of its distinct characteristics. In this thesis, I propose two novel unsupervised learning tasks to derive actionable patterns and gain insights from signed social networks - (1) unsupervised feature selection in signed social network; (2) unsupervised sentiment analysis with signed social network;

For unsupervised feature selection in signed social network, I propose a principled framework SignedFS. It first models both positive and negative links for a unified user preference representation. Then it embeds the user preference learning into feature selection. In addition, I model user proximity in signed social networks by signed graph regularization. Also, I conduct the experiment on two real-world datasets, Epinions and Wiki-rfa. The results show that SignedFS significantly improve the clustering performance and further experiments show that negatives links play an important role in the feature selection process.

For unsupervised sentiment analysis with signed social network, I proposed a principled framework SignedSenti. Methodologically, I propose to incorporate the signed social relations and sentimental signals from terms into a unified framework when we are lack of sentiment labels. I also conduct experiments on two real-world

signed social networks Epinions and Slashdot. The results show that the proposed SignedSenti has significantly better performance than state-of-the-art methods.

Future work can be focused on two aspects. First, in addition to social media data, relations between instances of other kinds of networks such as gene networks and citation networks also exhibit some implicit negative interactions. I would like to investigate how to employ negative links in such networks to solve some problem. Second, as shown in [42], for some social media sites without explicit negative links such as Facebook and Twitter, negative links can be predicted from explicit positive links. Therefore, I would like to adapt the SignedFS and SignedSenti framework to signed social networks with explicit positive links and predicted negative links.

REFERENCES

- [1] R. P. Abelson. Whatever became of consistency theory? *Personality and Social Psychology Bulletin*, 1983.
- [2] D. P. Bertsekas. Nonlinear programming. 1999.
- [3] J. C. Bezdek and R. J. Hathaway. Some notes on alternating optimization. In *Proceedings of the AFSS International Conference on Fuzzy Systems*, pages 288–300. 2002.
- [4] J. Bollen, H. Mao, and A. Pepe. Modeling public mood and emotion: Twitter sentiment and socio-economic phenomena. *Proceedings of the International AAAI Conference on Web and Social Media*, 11:450–453, 2011.
- [5] D. Cai, C. Zhang, and X. He. Unsupervised feature selection for multi-cluster data. In *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 333–342. ACM, 2010.
- [6] K. Cheng, J. Li, and H. Liu. Featureminer: A tool for interactive feature selection. In *Proceedings of the 25th ACM International Conference on Conference on Information and Knowledge Management*. ACM, 2016.
- [7] K.-Y. Chiang, N. Natarajan, A. Tewari, and I. S. Dhillon. Exploiting longer cycles for link prediction in signed networks. In *Proceedings of the 20th ACM international conference on Information and knowledge management*, pages 1157–1162. ACM, 2011.
- [8] J. Demšar. Statistical comparisons of classifiers over multiple data sets. *Journal of Machine learning research*, 7(Jan):1–30, 2006.
- [9] P. Doreian. Network autocorrelation models: Problems and prospects. *Spatial Statistics: Past, Present, Future*. Ann Arbor, Michigan Document Services, 1989.
- [10] R. O. Duda, P. E. Hart, and D. G. Stork. *Pattern classification*. John Wiley & Sons, 2012.
- [11] J. G. Dy and C. E. Brodley. Feature selection for unsupervised learning. *The Journal of Machine Learning Research*, 5:845–889, 2004.
- [12] A. Esuli and F. Sebastiani. Sentiwordnet: A publicly available lexical resource for opinion mining. In *Proceedings of LREC*, volume 6, pages 417–422. Citeseer, 2006.
- [13] Q. Gu and J. Han. Towards feature selection in network. In *Proceedings of the 20th ACM international conference on Information and knowledge management*, pages 1175–1184. ACM, 2011.
- [14] R. Guha, R. Kumar, P. Raghavan, and A. Tomkins. Propagation of trust and distrust. In *Proceedings of the 13th international conference on World Wide Web*, pages 403–412. ACM, 2004.

- [15] E. Hatfield, J. T. Cacioppo, and R. L. Rapson. *Emotional contagion*. Cambridge university press, 1994.
- [16] X. He, D. Cai, and P. Niyogi. Laplacian score for feature selection. In *Advances in neural information processing systems*, pages 507–514, 2005.
- [17] F. Heider. Attitudes and cognitive organization. *The Journal of psychology*, 21(1):107–112, 1946.
- [18] C.-J. Hsieh, K.-Y. Chiang, and I. S. Dhillon. Low rank modeling of signed networks. In *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 507–515. ACM, 2012.
- [19] X. Hu, J. Tang, H. Gao, and H. Liu. Unsupervised sentiment analysis with emotional signals. In *Proceedings of the 22nd international conference on World Wide Web*, pages 607–618. ACM, 2013.
- [20] X. Hu, L. Tang, J. Tang, and H. Liu. Exploiting social relations for sentiment analysis in microblogging. In *Proceedings of the sixth ACM international conference on Web search and data mining*, pages 537–546. ACM, 2013.
- [21] S. D. Kamvar and J. Harris. We feel fine and searching the emotional web. In *Proceedings of the fourth ACM international conference on Web search and data mining*, pages 117–126. ACM, 2011.
- [22] J. Kunegis, S. Schmidt, A. Lommatzsch, J. Lerner, E. W. De Luca, and S. Albayrak. Spectral analysis of signed graphs for clustering, prediction and visualization. In *Proceedings of the SIAM International Conference on Data Mining*. SIAM, 2010.
- [23] T. La Fond and J. Neville. Randomization tests for distinguishing social influence and homophily effects. In *Proceedings of the 19th international conference on World wide web*, pages 601–610. ACM, 2010.
- [24] J. Leskovec, D. Huttenlocher, and J. Kleinberg. Predicting positive and negative links in online social networks. In *Proceedings of the 19th international conference on World wide web*, pages 641–650. ACM, 2010.
- [25] J. Leskovec, D. Huttenlocher, and J. Kleinberg. Signed networks in social media. In *Proceedings of the SIGCHI conference on human factors in computing systems*, pages 1361–1370. ACM, 2010.
- [26] J. Li, K. Cheng, S. Wang, F. Morstatter, T. Robert, J. Tang, and H. Liu. Feature selection: A data perspective. *arXiv:1601.07996*, 2016.
- [27] J. Li, X. Hu, L. Wu, and H. Liu. Robust unsupervised feature selection on networked data. In *Proceedings of the 2016 SIAM International Conference on Data Mining*, 2016.

- [28] Y. Li, J. Liu, and C. Liu. A comparative analysis of evolutionary and memetic algorithms for community detection from signed social networks. *Soft Computing*, 18(2):329–348, 2014.
- [29] Z. Li, Y. Yang, J. Liu, X. Zhou, and H. Lu. Unsupervised feature selection using nonnegative spectral analysis. In *Proceedings of the AAAI conference on Artificial Intelligence*, 2012.
- [30] B. Liu. Sentiment analysis and opinion mining. *Synthesis lectures on human language technologies*, 5(1):1–167, 2012.
- [31] H. Liu and L. Yu. Toward integrating feature selection algorithms for classification and clustering. *Knowledge and Data Engineering, IEEE Transactions on*, 17(4):491–502, 2005.
- [32] M. McPherson, L. Smith-Lovin, and J. M. Cook. Birds of a feather: Homophily in social networks. *Annual review of sociology*, pages 415–444, 2001.
- [33] B. O’Connor, R. Balasubramanyan, B. R. Routledge, and N. A. Smith. From tweets to polls: Linking text sentiment to public opinion time series. In *Proceedings of the International AAAI Conference on Web and Social Media*, pages 122–129, 2010.
- [34] V. P. Pauca, F. Shahnaz, M. W. Berry, and R. J. Plemmons. Text mining using non-negative matrix factorizations. In *Proceedings of the SIAM International Conference on Data Mining*, volume 4, pages 452–456. SIAM, 2004.
- [35] C. R. Shalizi and A. C. Thomas. Homophily and contagion are generically confounded in observational social network studies. *Sociological methods & research*, 40(2):211–239, 2011.
- [36] M. Speriosu, N. Sudan, S. Upadhyay, and J. Baldridge. Twitter polarity classification with label propagation over lexical links and the follower graph. In *Proceedings of the First workshop on Unsupervised Learning in NLP*, pages 53–63. Association for Computational Linguistics, 2011.
- [37] P. Symeonidis, E. Tiakas, and Y. Manolopoulos. A unified framework for link and rating prediction in multi-modal social networks. *International Journal of Social Network Mining*, 1(3-4):225–253, 2013.
- [38] M. Szell, R. Lambiotte, and S. Thurner. Multirelational organization of large-scale social networks in an online world. *Proceedings of the National Academy of Sciences*, 107(31):13636–13641, 2010.
- [39] M. Taboada, J. Brooke, M. Tofiloski, K. Voll, and M. Stede. Lexicon-based methods for sentiment analysis. *Computational linguistics*, 37(2):267–307, 2011.
- [40] C. Tan, L. Lee, J. Tang, L. Jiang, M. Zhou, and P. Li. User-level sentiment analysis incorporating social networks. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1397–1405. ACM, 2011.

- [41] J. Tang, C. Aggarwal, and H. Liu. Recommendations in signed social networks. In *Proceedings of the 25th International Conference on World Wide Web*, pages 31–40. International World Wide Web Conferences Steering Committee, 2016.
- [42] J. Tang, S. Chang, C. Aggarwal, and H. Liu. Negative link prediction in social media. In *Proceedings of the Eighth ACM International Conference on Web Search and Data Mining*, pages 87–96. ACM, 2015.
- [43] J. Tang, X. Hu, and H. Liu. Is distrust the negation of trust?: the value of distrust in social media. In *Proceedings of the 25th ACM conference on Hypertext and social media*, pages 148–157. ACM, 2014.
- [44] J. Tang and H. Liu. Feature selection with linked data in social media. In *Proceedings of the SIAM International Conference on Data Mining*, pages 118–128. SIAM, 2012.
- [45] J. Tang and H. Liu. Unsupervised feature selection for linked social media data. In *Proceedings of the ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 904–912. ACM, 2012.
- [46] J. Tang, C. Nobata, A. Dong, Y. Chang, and H. Liu. Propagation-based sentiment analysis for microblogging data. In *Society for Industrial and Applied Mathematics Publications*, 2015.
- [47] M. Thelwall, K. Buckley, G. Paltoglou, D. Cai, and A. Kappas. Sentiment strength detection in short informal text. *Journal of the American Society for Information Science and Technology*, 61(12):2544–2558, 2010.
- [48] H. Wang, Y. Lu, and C. Zhai. Latent aspect rating analysis on review text data: a rating regression approach. In *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 783–792, 2010.
- [49] Y. Wang, S. Wang, J. Tang, H. Liu, and B. Li. Unsupervised sentiment analysis for social media images. In *Proceedings of the Twenty-Fourth International Joint Conference on Artificial Intelligence, IJCAI 2015, Buenos Aires, Argentina*, pages 2378–2379, 2015.
- [50] J. Wiebe, T. Wilson, and C. Cardie. Annotating expressions of opinions and emotions in language. *Language resources and evaluation*, 39(2-3):165–210, 2005.
- [51] T. Wilson, J. Wiebe, and P. Hoffmann. Recognizing contextual polarity in phrase-level sentiment analysis. In *Proceedings of the conference on human language technology and empirical methods in natural language processing*, pages 347–354. Association for Computational Linguistics, 2005.
- [52] Z. Zhao and H. Liu. Spectral feature selection for supervised and unsupervised learning. In *Proceedings of the 24th international conference on Machine learning*, pages 1151–1157. ACM, 2007.