

Measuring and Moderating Opinion Polarization in Online Social Networks

A Thesis

submitted to the designated
by the General Assembly of Special Composition
of the Department of Computer Science and Engineering
Examination Committee

by

Antonis Matakos

in partial fulfillment of the requirements for the degree of

MASTER OF SCIENCE IN COMPUTER SCIENCE

WITH SPECIALIZATION

IN SOFTWARE

University of Ioannina

June 2017

DEDICATION

Dedicated to my family.

ACKNOWLEDGEMENTS

First and foremost, I would like to thank my advisor Panayiotis Tsaparas and my co-advisor Evimaria Terzi for their guidance and advice in completing this Thesis. Working beside them has been a great experience that not only helped me develop my skills, but most importantly, their passion and work ethic has been a big influence for me. I am also grateful to Evaggelia Pitoura for her precious insight and comments.

I would also like to thank my colleagues for creating a pleasant work environment, and for their invaluable help and advice. It has been a privilege to work among them.

Finally, I would like to express my gratitude towards my family and friends, for their great help and support through all these years. Without them, this thesis would not have been possible.

TABLE OF CONTENTS

List of Figures	iii
List of Tables	v
List of Algorithms	vi
Abstract	vii
Εκτεταμένη Περίληψη	viii
1 Introduction	1
1.1 Scope	1
1.2 Roadmap	4
2 Related work	5
2.1 Filter bubbles and echo chambers	5
2.2 Quantifying and reducing polarization	6
2.3 Opinion Mining	7
2.4 Influence and Opinion Maximization	7
3 The Polarization Index	8
4 Moderating Internal Opinions	13
4.1 Problem Definition	13
4.2 Problem complexity	14
4.3 Algorithms	15
5 Moderating Expressed Opinions	17
5.1 Problem definition	17

5.2	Problem complexity	18
5.3	Algorithms	18
6	Experiments	21
6.1	Datasets	21
6.2	Evaluation of the Polarization Index	24
6.3	Heuristic Algorithms for Opinion Moderation	26
6.4	Evaluation of algorithms for MODERATEINTERNAL	27
6.5	Evaluation of the algorithms for MODERATEEXPRESSED	33
6.6	Scalability	39
6.7	Case study	40
7	Conclusions	42
	Bibliography	44
A	Proof of Theorem 1	48

LIST OF FIGURES

- 3.1 Three examples of graphs for the polarization index. 10
- 3.2 Comparison of Network disagreement index and Polarization index values 11

- 6.1 Performance of the algorithms for the MODERATEINTERNAL problem on
Karate 28
- 6.2 Performance of the algorithms for the MODERATEINTERNAL problem on
Books 29
- 6.3 Performance of the algorithms for the MODERATEINTERNAL problem on
Blogs 29
- 6.4 Performance of the algorithms for the MODERATEINTERNAL problem on
Elections 30
- 6.5 Performance of the algorithms for the MODERATEINTERNAL problem on
Hashtags P 30
- 6.6 Performance of the algorithms for the MODERATEINTERNAL problem on
Hashtags NP 31
- 6.7 Comparison with optimal 31
- 6.8 Selected nodes by GreedyInt on *Karate* 32
- 6.9 Selected nodes by GreedyInt on *Books* 32
- 6.10 Performance of the algorithms for the MODERATEEXPRESSED problem on
Karate 34
- 6.11 Performance of the algorithms for the MODERATEEXPRESSED problem on
Books 35
- 6.12 Performance of the algorithms for the MODERATEEXPRESSED problem on
Blogs 35
- 6.13 Performance of the algorithms for the MODERATEEXPRESSED problem on
Elections 36

6.14 Performance of the algorithms for the MODERATEEXPRESSED problem on <i>Hashtags P</i>	36
6.15 Performance of the algorithms for the MODERATEEXPRESSED problem on <i>Hashtags NP</i>	37
6.16 Performance of the algorithms for the MODERATEEXPRESSED problem on Elections top-1%	37
6.17 Selected nodes by GreedyExt on <i>Karate</i>	38
6.18 Selected nodes by GreedyExt on <i>Books</i>	38
6.19 Comparison of the running times (in secs) for the Sherman-Morrison and iterative implementation of GreedyExt, for varying size of n	40

LIST OF TABLES

- 6.1 Dataset Statistics 23
- 6.2 Dataset polarization index and randomization values 25
- 6.3 Comparison of polarization metrics on all Datasets 25
- 6.4 Running times (secs) of all algorithms for $k = 0.1n$ in the *Elections*
dataset. 39
- 6.5 Characteristics of the first ten nodes selected by GreedyExt in *Elections*
dataset. 41

LIST OF ALGORITHMS

ABSTRACT

Antonis Matakos, M.Sc. in Computer Science, Department of Computer Science and Engineering, University of Ioannina, Greece, June 2017.

Measuring and Moderating Opinion Polarization in Online Social Networks.

Advisor: Panayiotis Tsaparas, Associate Professor.

The polarization of society over controversial social issues has been the subject of study in social sciences for decades [22, 33]. The widespread usage of online social networks and social media and the tendency of people to connect and interact with like-minded individuals has only intensified the phenomenon of polarization [4]. In this thesis, we consider the problem of measuring and reducing polarization of opinions in a social network. Using a standard opinion formation model [15], we define the *polarization index*, which, given a network and the opinions of the individuals in the network, it quantifies the polarization observed in the network. Our measure captures the tendency of opinions to concentrate in network communities, creating an echo-chamber. Given a numerical measure of polarization in the network, we consider the problem of reducing polarization by convincing individuals (e.g., through education, exposure to diverse viewpoints, or incentives) to adopt a more neutral stand towards controversial issues. We formally define the `MODERATEINTERNAL` and `MODERATEEXPRESSED` problems, and we prove that both our problems are NP-hard. By exploiting the linear-algebraic characteristics of the opinion formation model we design polynomial-time algorithms for both problems, and efficient heuristics. We conduct our experiments on real-world datasets, with data from twitter representing well-known controversies. We demonstrate the validity of our metric, by comparing against other metrics, and observing the obtained value in both polarized and non-polarized settings. We also showcase the efficiency and the effectiveness of our algorithms and heuristics in practice.

ΕΚΤΕΤΑΜΕΝΗ ΠΕΡΙΛΗΨΗ

Αντώνης Ματάκος, Μ.Δ.Ε. στην Πληροφορική, Τμήμα Μηχανικών Η/Υ και Πληροφορικής, Πανεπιστήμιο Ιωαννίνων, Ιούνιος 2017.

Τίτλος Διατριβής.

Επιβλέπων: Παναγιώτης Τσαπάρας, Επίκουρος Καθηγητής.

Η πόλωση της κοινωνίας πάνω σε αμφιλεγόμενα κοινωνικά ζητήματα έχει υπάρξει ως αντικείμενο έρευνας στις κοινωνικές επιστήμες εδώ και δεκαετίες [22, 33]. Η διαδεδομένη χρήση των online κοινωνικών δικτύων και social media αναμενόταν να κάνει τους ανθρώπους πιο ανοιχτούς σε διαφορετικές ιδέες, νοοτροπίες και απόψεις, και φαινόταν σαν ένα βήμα προς τον εκδημοκρατισμό και την διαποικίληση των κοινωνιών. Όμως, η εύκολη πρόσβαση σε άφθονη πληροφορία και η τάση των ανθρώπων να συνδέονται και να αλληλεπιδρούν με ομοϊδεάτες, έχει οδηγήσει στο αντίθετο αποτέλεσμα. Σαν συνέπεια, αντί να γεφυρωθεί το χάσμα μεταξύ των ανθρώπων, το φαινόμενο της πόλωσης έχει γίνει πιο έντονο [4]. Η πόλωση κατακερματίζει την κοινωνία σε ομάδες, με αποτέλεσμα την κατάπτωση του δημόσιου διαλόγου και της αλληλοκατανόησης μεταξύ των διάφορων πλευρών, καταστρέφοντας την ομαλή και δημοκρατική λειτουργία των κοινωνιών. Συνεπώς, η δημιουργία μηχανισμών για την μείωση της πόλωσης είναι ζήτημα υψίστης σημασίας.

Σε αυτή τη διατριβή, θεωρούμε το πρόβλημα της μέτρησης και της μείωσης της πόλωσης των απόψεων σε ένα κοινωνικό δίκτυο. Χρησιμοποιούμε ένα καθιερωμένο μοντέλο σχηματισμού απόψεων [15], το οποίο υποθέτει ότι κάθε χρήστης έχει μια εσωτερική άποψη που είναι σταθερή, και μια εκφραζόμενη άποψη που εξαρτάται από την εσωτερική του άποψη και τις εκφραζόμενες απόψεις του κοινωνικού δικτύου. Με βάση αυτό το μοντέλο ορίζουμε το *polarization index*, το οποίο, δεδομένου ενός δικτύου και εσωτερικών απόψεων των ατόμων στο δίκτυο, ποσοτικοποιεί την παρατηρούμενη πόλωση στο δίκτυο. Η μετρική μας συλλαμβάνει την τάση των

απόψεων να συγκεντρώνονται σε δικτυακές κοινότητες, δημιουργώντας θαλάμους αντήχησης(echo-chambers).

Δοθείσας μιας αριθμητικής μετρικής της πόλωσης στο δίκτυο, θεωρούμε το αλγοριθμικό πρόβλημα της μείωσης της πόλωσης πείθοντας άτομα(για παράδειγμα, μέσω μόρφωσης, έκθεσης σε ποικίλες απόψεις ή δίνοντας κίνητρα) να υιοθετήσουν μια πιο ουδέτερη στάση προς αμφιλεγόμενα ζητήματα. Ορίζουμε τυπικά τα MODERATEINTERNAL, όπου μετριάζουμε τις εσωτερικές απόψεις των χρηστών, και MODERATEEXPRESSED προβλήματα, όπου μετριάζουμε τις εκφραζόμενες απόψεις. Από υπολογιστικής άποψης, αποδεικνύουμε ότι και τα δύο προβλήματα είναι NP-δύσκολα. Εκμεταλλευόμενοι τα αλγεβρικά χαρακτηριστικά του μοντέλου σχηματισμού απόψεων, σχεδιάζουμε πολυωνυμικού χρόνου αλγορίθμους και για τα δύο προβλήματα, καθώς και αποτελεσματικούς ευριστικούς. Διεξάγουμε τα πειράματά μας σε δεδομένα απο τον πραγματικό κόσμο, με δεδομένα απο το twitter που αφορούν γνωστές διαμάχες. Επιδεικνύουμε την εγκυρότητα της μετρικής, συγκρίνοντας με άλλες μετρικές, και παρατηρώντας την τιμή που μας δίνει για πολωμένα και μη-πολωμένα περιβάλλοντα. Τέλος, αποδεικνύουμε την αποδοτικότητα και την αποτελεσματικότητα των αλγορίθμων μας στην πράξη.

CHAPTER 1

INTRODUCTION

1.1 Scope

1.2 Roadmap

1.1 Scope

In the past decades, online social networks and social media have emerged as the primary vehicle for the public discourse. Today, discussions take place primarily on Facebook and Twitter, where information and viewpoints are exchanged, and opinions are shaped. In this new world, users have easy access to information, but also to a public podium and a broad audience for their opinions. A consequence of this paradigm shift is the displacement of traditional media as the primary news outlet, as the networked public sphere provides anyone with an outlet to speak, inquire, and investigate.

Empowering ordinary users to express and share their opinions online seems like a step towards making individuals more open to different ideas, cultures, and viewpoints, and thus making societies overall more democratic and diverse. Nevertheless, it has been observed that the easy and uninhibited access to information and expression often leads to the opposite effect. Users tend to favor content that agrees with their existing worldview, and create connections with like-minded individuals, creating “*echo-chambers*”. Without any kind of moderation, current social-media platforms gravitate towards a state in which net-citizens are constantly reinforcing their existing

opinions. Additionally, social networks mostly provide personalized content, that the users are most likely to find agreeable. The end effect is that users get entrenched in a confine of comforting information — a *“filter bubble”*. Filter bubbles have been linked to increased partisanship and amplified ideological segregation. In such cases, instead of smoothing the differences, online social networks reinforce them, thus leading to increased *polarization* [4, 5].

Online polarization has been observed over a variety of issues and topics, ranging from frivolous (the dress controversy¹) to decisive and consequential (the increasing divide in US politics²). Polarization separates individuals into sides that have little or no communication with and understanding of each other, and has a corrosive and detrimental effect to the functioning of communities, societies, and democracies. It is thus of critical importance to devise mechanisms for reducing polarization. This is typically achieved by raising awareness and educating individuals about the different sides of an issue, with the goal of moderating extreme opinions and reaching a common ground. This is an arduous and costly process that may span a generation to yield results.

In this thesis, we take an algorithmic approach to the problem of measuring and reducing polarization. In order to measure polarization, we consider a popular opinion formation model [15]. In this model, opinions are modeled as real numbers ranging from -1 to 1, depending on the viewpoint of the user. Each user u has an internal opinion s_u that is given as input and it is fixed, and an expressed opinion z_u that depends on their own internal opinion and the expressed opinions in their social network. Using a random walk interpretation of the opinion formation model, we can interpret z_u as the expected opinion that node u will reach when taking a random walk in the social network. High value of z_u implies that the user is surrounded mostly by single-minded individuals with extreme opinions, while low value implies that the social network of u adopts moderate and diverse opinions. We view the absolute value $|z_u|$ as a measure of the degree of the polarization of user u . Given the vector of expressed opinions \mathbf{z} for the whole network, the length of the opinion vector $\|\mathbf{z}\|^2$ captures the degree of polarization in the network. We refer to $\|\mathbf{z}\|^2$ as the *polarization index* $\pi(\mathbf{z})$ of the network.

Given this numeric measure of polarization, we are interested in algorithms for

¹https://en.wikipedia.org/wiki/The_dress

²<http://www.people-press.org/2014/06/12/political-polarization-in-the-american-public/>

reducing polarization in the network. We assume that we can reduce polarization by convincing people (through education or other means) to adopt a more moderate opinion. Given a budget value k , we want to find the best set of k individuals in the network, such that convincing them to moderate their opinions (in our model, set their opinion value to zero) will minimize the polarization index of the network. We consider two variants of this problem: the `MODERATEINTERNAL` problem, and the `MODERATEEXPRESSED` problem. In the `MODERATEINTERNAL` problem we moderate the internal opinion of the users, that is, for each user u in the selected set we set $s_u = 0$. This is the case where through education we expose users to the viewpoint of the other side, and lead them to adopt a moderate viewpoint. In the `MODERATEEXPRESSED` problem we moderate the expressed opinions, that is, for each user u in the selected set we set $z_u = 0$. This is a case where we give incentives to users to adopt a moderate public opinion, and propagate a balanced viewpoint.

From the computational point of view, we prove that both problems are NP-hard. We propose algorithms that exploit the properties of the opinion formation model so as to efficiently construct the solution set, as well as efficient heuristics. We experiment on real datasets and we demonstrate the effectiveness of our algorithms in decreasing polarization.

In summary, in this thesis we make the following contributions:

- We define a novel polarization index for quantifying polarization in a network, based on the opinions of users under a popular opinion formation model [15]. Our measure takes into account both the existing opinions of the users, and the network structure. To the best of our knowledge we are the first to use this model to measure polarization.
- We define two novel problems, `MODERATEINTERNAL` and `MODERATEEXPRESSED` for reducing polarization in a network. We show that both problems are NP-hard, and propose efficient algorithms for solving them. Our algorithms exploit a linear-algebraic view of the opinion-formation model we adopt.
- We experiment on real data, including a Twitter network from 2016 US Elections. We demonstrate that our polarization index is successful in capturing polarization, by comparing the value that it takes in different opinion settings, and against other state-of-the-art metrics. We also showcase and that our algorithms are effective in reducing polarization, and are efficient in practice.

1.2 Roadmap

The thesis is structured as follows. In Chapter 2 we review related work on measuring and reducing polarization, opinion formation models, and opinion mining. In Chapter 3 we define the polarization index, and provide the intuition behind it. In Chapter 4 we define and study the MODERATEINTERNAL problem, and in Chapter 5 the MODERATEEXPRESSED problem. Chapter 6 presents the experimental evaluation of our metric and algorithms, and Chapter 7 concludes the thesis. Appendix A contains the full proof for the NP-hardness of the MODERATEINTERNAL problem.

CHAPTER 2

RELATED WORK

2.1 Filter bubbles and echo chambers

2.2 Quantifying and reducing polarization

2.3 Opinion Mining

2.4 Influence and Opinion Maximization

Although, to the best of our knowledge, we are the first to introduce and study the `MODERATEEXPRESSED` and the `MODERATEINTERNAL` problems, our work is related to recent work on polarization and opinion maximization.

2.1 Filter bubbles and echo chambers

While social media have the potential to expose individuals to more diverse viewpoints, they can also limit exposure to attitude-challenging information, which leads to a radicalization of attitudes and false perceptions about events. This has led to theories about the effects of “echo-chambers” [4, 5, 18], where users are only exposed to information by like-minded individuals, and “filter bubbles” [4, 31], where algorithms only present personalized content that agrees with the user’s attitude. Recent lines of work [18] have investigated the strength of echo chambers and filter bubbles, and found that opinion-challenging information reduces the likelihood of a news story’s

exposure.

2.2 Quantifying and reducing polarization

The phenomenon of polarization has been the subject of study in social sciences for decades [22, 33]. There has been a lot of work on measures for quantifying the polarization observed in online social networks and social media [2, 10, 16, 20, 3, 11] and model its emergence [11, 34]. The main characteristic of those works is that the measures proposed are based on the structural characteristics of the underlying graph and they do not consider the existing opinions, or an opinion formation model, when quantifying polarization. Vicario et al. [13] study polarization while incorporating opinion dynamics, assuming a variation of the Bounded Confidence Model (BCM). This variation of the model, has the limitation that it can only converge in states where opinions form clusters of a single value. The closest to our definition is the notion of tension for measuring polarization [6, 11], which focuses on pairwise disagreements of opinions over the edges of the network. The metric does not consider the overall distribution of opinions, and it will fail to detect the emergence of echo-chambers in the network, where like-minded individuals only interact with each other.

Given the negative effects of polarization and fragmentation on the well-being and the well-functioning of societies there has been work that focuses on methods for decreasing polarization. Such studies focus on proposing mechanisms that will expose online social-media users to content that is not necessarily aligned with their prior beliefs. The work in this direction can be split into work that focuses on (a) *how* to present information to users and (b) *who* to approach with the new information. In terms of (a) there has been work focusing on user studies as well as the design of the appropriate interfaces that predispose users positively towards diverse ideas presented to them [29, 28, 35]. Clearly, our work is complementary to the above as we focus on the algorithmic aspects of decreasing polarity.

In terms of who to approach, the recent work by Garimella et al. [17] considers the introduction of edges that will reduce the observed polarization in a social network. Although this work is related to ours, it focuses on graph-theoretic measures of polarization and does not take into consideration the opinions of individuals neither does it explicitly consider an opinion-formation model. Furthermore, it considers the

addition of links, rather than the moderation of opinions. Therefore, both our model and our problem are different from that in [17].

2.3 Opinion Mining

In our work we assume that we are given user opinions as input, and we focus on using these opinions to measure and moderate polarization. In our experiments we assume that opinions can be inferred by the actions of the users (membership in known communities, or following specific accounts). In networks where we have information about the attributes of the users, or the content they contribute, it may be possible to obtain more fine-grained and nuanced opinion values by applying opinion mining and sentiment analysis techniques [26]. Opinion mining deals with the inference of the semantics of a given text. Concept-based techniques have come to prominence recently [8, 7], along with new neural network approaches, using deep neural networks [32, 9]. Our framework for measuring and moderating polarization can be extended to include an opinion mining algorithm as the first step of the pipeline.

2.4 Influence and Opinion Maximization

At a high level, our work is also related to the line of work on influence and opinion maximization [23, 19]. In these works the goal is to select a set of individuals that will adopt a product or an opinion so as to maximize the overall adoption in the network. The closest to ours is the work of Gionis et al. [19], where the goal is to find a set of individuals who will change their opinion (internal or expressed), such that the sum of expressed opinions is maximized. Both our work and theirs assume the same opinion-formation model. However, our goal is different; rather than maximizing the positive expressed opinion we aim at minimizing the polarization index. The difference in the objectives results in differences in the problem properties and the algorithmic techniques that need to be developed.

CHAPTER 3

THE POLARIZATION INDEX

In this chapter, we define the polarization index we will use in the thesis, and we provide the necessary background for understanding and analyzing our metric.

Throughout the thesis, we consider a social graph $G = (V, E)$ with n nodes and m edges. Each edge (i, j) is associated with a weight $w_{ij} \geq 0$, which expresses the strength of the connection between i and j , and the influence they exert to each other.

We adopt the opinion-formation model of Friedkin and Johnsen [15], which assumes that every person i in the network has a persistent *internal opinion* s_i , and an *expressed opinion* z_i which depends both on their internal opinion s_i and the expressed opinions of their neighbours. More precisely, the expressed opinion of node i is computed as the weighted average of their internal opinion and the expressed opinions of the neighbours of i , $N(i)$, in G :

$$z_i = \frac{w_{ii}s_i + \sum_{j \in N(i)} w_{ij}z_j}{w_{ii} + \sum_{j \in N(i)} w_{ij}}, \quad (3.1)$$

where w_{ii} denotes the importance that node i places on their own opinion. It has been shown that if every person i iteratively updates their expressed opinion, then the expressed opinions converge to a unique opinion vector \mathbf{z} .

In our setting, opinions can be both positive and negative. Thus, we assume that they take values in the interval $[-1, 1]$, where -1 reflects a negative opinion, and 1 a positive one, while 0 corresponds to a neutral position.

In the absence of any polarization all users would express a neutral opinion, i.e., $z_i = 0$ for all $i \in V$. The absolute value of the opinion of a user $|z_i|$ captures how

extreme the user opinion is. We quantify polarization in the network by measuring how far we are from the state of complete neutrality. We measure this by looking at the length of the vector \mathbf{z} under the L_2^2 norm. To make the value of our metric independent of the size of the network, we normalize by the number of nodes in the graph. More formally, we have the following definition.

Definition 3.1 (Polarization Index). Given a network $G = (V, E)$ and the opinion vector \mathbf{z} defined over the nodes of the network, we define the **polarization index** $\pi(\mathbf{z})$ as: $\pi(\mathbf{z}) = \frac{\|\mathbf{z}\|^2}{n}$.

We now give some additional background and intuition behind our metric. First, an equivalent way of obtaining the *expressed opinion* vector \mathbf{z} from the *internal opinion* vector \mathbf{s} is the following [6]: if \mathbf{L} is the Laplacian matrix of graph $G = (V, E)$, and \mathbf{I} is the identity matrix, then $\mathbf{z} = (\mathbf{L} + \mathbf{I})^{-1} \mathbf{s}$. We will refer to the matrix $\mathbf{Q} = (\mathbf{L} + \mathbf{I})^{-1}$ as the *fundamental matrix*.

Second, there is a direct connection between the opinion formation model, and random walks with absorbing nodes, as it is shown in [19]. More specifically, given the graph $G = (V, E)$, with n vertices, and weights w_{ij} for the edges $(i, j) \in E$, we construct the *augmented graph* $H = (V \cup X, E \cup R)$ as follows. For each vertex $v_i \in V$, we add a new vertex x_i in X . We also add a *directed* edge (v_i, x_i) in R , with weight w_{ii} . The node x_i corresponds to the internal opinion of node v_i .

Now consider a random walk on graph H that starts from a vertex $v \in V$. The nodes in X are *absorbing*. That is, when reaching these nodes, the random walk terminates. For each absorbing node x_i we can compute the probability $P(x_i | v_j)$, that the random walk that started from v_j terminates at node x_i . It was shown in [19] that $Q(j, i) = P(x_i | v_j)$, that is, the j -th row of the matrix \mathbf{Q} is a probability distribution over all nodes in X . Therefore, we have that

$$z_j = \sum_{i=1}^n P(x_i | v_j) s_i.$$

We can think of the probability $P(x_i | v_j)$ as the probability that node v_j adopts the opinion of node v_i . This probability depends on the structure of the graph: the more the paths that connect v_j with node v_i , the higher the probability $P(x_i | v_j)$. The probability is also affected by the weights w_{ij} and w_{ii} since they determine the probability that a specific edge is followed. For example, high w_{ii} weight means that the user is more likely to be absorbed in his own opinion node than follow a path

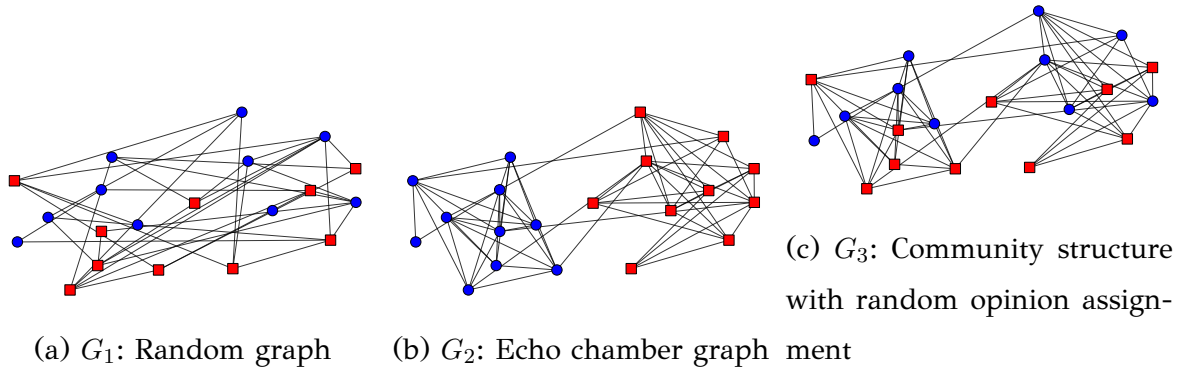


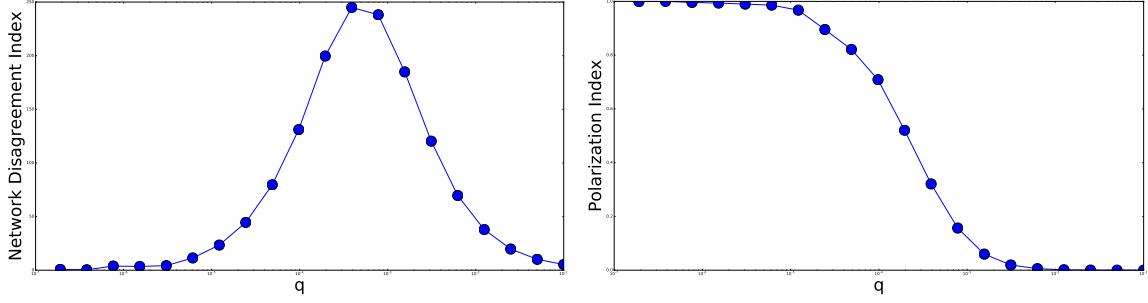
Figure 3.1: Three examples of graphs for the polarization index.

in the graph to some other node. The expressed opinion of node z_j is the expected value of the internal opinion of the node at the point of absorption.

The implications of this connection are the following. For a specific node v_j , the value $|z_j|$ is minimized if node v_j has equal probability to reach positive and negative opinions, that is, it has a balanced view of the opinions in the network. On the other hand, if the user is trapped in a filter-bubble of like-minded friends, all with extreme opinions, the value of $|z_j|$ will be high. The polarization index becomes high if we have echo chambers in the network, that is, we have communities in the graph, that are homogeneous with respect to their internal opinions.

To illustrate this point, consider the three graphs shown in Figure 3.1. The graph G_1 in Figure 3.1(a) consists of 20 nodes, 10 with opinion -1, and 10 with opinion +1, that are randomly interconnected. The graphs in Figure 3.1(b) and Figure 3.1(c) are the same and they consist of two densely connected subgraphs of size 10 that are sparsely interconnected. In Figure 3.1(b), the opinions are aligned with the communities in the graph: the nodes in the left community have opinion -1 (blue round nodes), while the nodes in the right community have opinion +1 (red round nodes). In Figure 3.1(c), the opinions are randomly assigned.

We compute the polarization index for all three graphs. In the first graph G_1 , edges are created at random, and thus the graph has no community structure, and opinions mix randomly in the network. Therefore, each node has more or less equal probability to adopt a positive or negative opinion, resulting in a low polarization index of 0.03. In the second graph G_2 , there is a clear echo chamber effect: positive nodes speak mostly with positive nodes, and negative nodes speak mostly with negative nodes. As a result the polarization index is high, 0.30. In the third graph G_3 , although there



(a) Network disagreement index for exponentially increasing values of q (b) Polarization index for exponentially increasing values of q

Figure 3.2: Comparison of Network disagreement index and Polarization index values

is a clear community structure in the graph, the opinions are equally distributed in the two communities. Therefore, although the nodes tend to communicate mostly with the nodes within their community (the probability of adopting the opinion of a node in a different community is small), both opinions are equally represented in each community, resulting in a small polarization index of 0.03.

We also demonstrate the superiority of our metric in capturing polarization, compared to its closest analogue, the *network disagreement index* defined in [14], in the following example. We consider a graph that consists of 1000 nodes, with two single-opinion clusters of 500 nodes, one with opinion -1, and the other with opinion +1. For each pair of nodes in the same cluster, we assign an edge between them with probability 0.1, while we vary the inter-cluster edge probability q , and observe the effect on the π and the network disagreement index. The result is illustrated in Figure 3.2. The largest value of inter-cluster edge probability is 0.1 (and rightmost value in the plot), which represents the case where there are no clusters, and in each subsequent computation the probability is halved. The network becomes more polarized as the number of edges between the cluster decreases, since there is less communication between the two sides. The example shows that while the network disagreement index increases in value as the number of edges between clusters decreases, its value starts to decrease after some point. Our metric, on the other hand, increases in value monotonically as the number of inter-cluster edges decreases, which accurately captures the increasing polarization in the network.

Given this measure of polarization in the network, our next goal is to minimize it by convincing a small set of users to adopt more neutral positions. We consider two possible ways to achieve this. The first is to convince users, via education and

exposure to different viewpoints, to change their internal opinions. The second is by giving incentives to users to express and propagate a neutral opinion. In both cases we say that we *moderate* the opinions of the users. Depending on whether we moderate the internal or the expressed opinions of users we define the `MODERATEINTERNAL` and the `MODERATEEXPRESSED` problems respectively. We define and study these two problems in the following chapters.

CHAPTER 4

MODERATING INTERNAL OPINIONS

4.1 Problem Definition

4.2 Problem complexity

4.3 Algorithms

In this chapter we define the MODERATEINTERNAL problem, we analyze its complexity, and we design efficient and effective algorithms for solving it.

4.1 Problem Definition

When moderating internal opinions, we seek a small set of nodes, T_s , whose *internal* opinions would be set to zero, such that the polarization index is minimized. We use $\pi(\mathbf{z} \mid T_s)$ to denote the polarization index after setting the internal opinions of the nodes in T_s to zero. The formal problem definition is the following.

Problem 1 (MODERATEINTERNAL). *Given a graph $G = (V, E)$, a vector of internal opinions \mathbf{s} , and an integer k , identify a set T_s of k nodes such that changing the internal opinions of the nodes in T_s to 0, minimizes the polarization index $\pi(\mathbf{z} \mid T_s)$.*

4.2 Problem complexity

We prove the following Theorem for the hardness of the the MODERATEINTERNAL problem.

Theorem 4.1. *The MODERATEINTERNAL problem is NP-hard.*

Proof. We only give some intuition for the proof. The full proof appears in the Appendix A.

Our proof uses a reduction from the m -SUBSETSUM problem, where given a set of N positive integer numbers v_1, \dots, v_N , a value m , and a target value b , we ask if there is a set of numbers B of size m , such that $\sum_{v_i \in B} v_i = b$.

Given an instance of the m -SUBSETSUM problem, we construct an instance of MODERATEINTERNAL as follows. The graph is a star with $N + 1$ nodes: we have a central node u_0 , and a spoke node u_i for each integer v_i . For the center of the star (node u_0) we have that $w_{00} = t$, for an appropriately selected value of t (we will discuss this below), and $s_0 = -1$. The weight of the edge (u_0, u_i) from the center to node u_i is $w_{0i} = v_i$, and the weight of node u_i to its internal opinion is also $w_{ii} = v_i$. The opinion of all spoke nodes is $s_i = 1$. We set $k = N - m$, and we ask for a set of nodes T_s , $|T_s| = k$, such that, when setting $s_i = 0$ for $u_i \in T_s$, $\pi(\mathbf{z} | T_s) = \|\mathbf{z}\|^2$ is minimized.

Assume that we have selected the set T_s , $|T_s| = k$. We can prove that

$$\pi(\mathbf{z} | T_s) = \frac{N + 4}{4} z_0^2 + \frac{N - k}{2} z_0 + \frac{N - k}{4}.$$

Therefore, $\pi(\mathbf{z} | T_s)$ is determined by the expressed opinion of the center node z_0 . Let $R = V \setminus T_s \cup \{u_0\}$ denote the set of spoke nodes whose opinion was *not* set to 0. Using the equations for the expressed opinions of the opinion formation model we can show the following for the value of z_0 (details in the Appendix):

$$z_0 = \frac{\sum_{u_i \in R} v_i - 2t}{W + 2t}.$$

For the sake of the argument, assume for a moment that we achieve the minimum $\pi(\mathbf{z} | T_s)$ for $z_0 = 0$. Then clearly, we need to select a set of nodes in T_s , such that for the nodes in R we have $\sum_{u_i \in R} v_i = 2t$. Setting $t = b/2$ we can prove that we minimize $\pi(\mathbf{z} | T_s)$ if and only if there is a set of nodes R such that $\sum_{u_i \in R} v_i = b$, which proves the reduction. However, the value $z_0 = 0$ does not minimize $\pi(\mathbf{z} | T_s)$. In the full proof, we determine the optimal value of z_0 , and the value of t that achieves this

optimal when there is a set of nodes R such that $\sum_{u_i \in R} v_i = b$, and thus complete the reduction. \square

Furthermore, we observe that $\pi(\mathbf{z} \mid T_s)$ is not monotone with respect to T_s . That is, it is not necessarily true that the more nodes we make neutral, the lower the polarization. This can be seen by considering a simple graph consisting of two nodes, u and v , with internal opinions -1 and 1 and $w_{uu} = w_{vu} = w_{vv} = 1$. In this case $\pi(\mathbf{z}) = 2/9$. If we change the internal opinion of the negative node to neutral, then $\pi(\mathbf{z}) = 5/9$. Thus, making a node neutral, causes the polarization index to increase. This observation implies that designing an algorithm for MODERATEINTERNAL is challenging.

4.3 Algorithms

In this section, we present our algorithms for MODERATEINTERNAL. For the following, we assume that the matrix \mathbf{Q} has been pre-computed, and it is given as input to the algorithm.

The BOMP algorithm: The *Binary Orthogonal Matching Pursuit* (BOMP) algorithm is inspired by the connection of the MODERATEINTERNAL problem to the problem of sparse approximation [30]. First, we establish this connection and then we describe the BOMP algorithm.

As we have already discussed in Chapter 3, *expressed opinion* vector \mathbf{z} can be computed as $\mathbf{z} = \mathbf{Q}\mathbf{s}$, where $\mathbf{Q} = (\mathbf{L} + \mathbf{I})^{-1}$. Note that we have that $\mathbf{Q}\mathbf{s} = \mathbf{Q}\mathbf{S}\mathbf{1}$, where \mathbf{S} is the diagonal matrix with $\mathbf{S}_{ii} = s_i$, and $\mathbf{1}$ is the vector of all ones. For the rest of the discussion we will use $\mathbf{R} = \mathbf{Q}\mathbf{S}$.

Now, let \mathbf{s}' denote the vector \mathbf{s} after we set k of its entries to zero – these entries will correspond to users whose internal opinions become neutral. Our goal is to find the vector \mathbf{s}' that minimizes $\|\mathbf{Q}\mathbf{s}'\|^2$. Note that $\mathbf{Q}\mathbf{s}' = \mathbf{R}\mathbf{1} - \mathbf{R}\mathbf{x}$, where \mathbf{x} is a vector with 1's at the positions of the selected nodes, and zeros everywhere else. Since $\mathbf{R}\mathbf{1} = \mathbf{z}$, the original expressed opinion vector, our problem can be stated as follows: Find the best *binary* vector \mathbf{x} with k non-zero entries (i.e., $\|\mathbf{x}\|_0 = k$) such that $\|\mathbf{z} - \mathbf{R}\mathbf{x}\|^2$ is minimized. This is the definition of the sparse approximation problem [30], where we restrict the solution to binary vectors.

Inspired by [24] we will approximate the solution to this problem using a variation of a known algorithm from signal processing [27, 12] called *nonnegative orthogonal*

matching pursuit (NNOMP). The NNOMP algorithm is designed to find a sparse vector \mathbf{x} (with no more than k non-zero entries) with *non-negative* yet *real* coefficients that when multiplied to a matrix \mathbf{R} is minimizes $\|\mathbf{z} - \mathbf{R}\mathbf{x}\|^2$ for a target vector \mathbf{z} .

In our problem, the vector \mathbf{x} is a binary vector and thus \mathbf{x} essentially selects a subset of columns from \mathbf{R} and uses their sum to approximate the target vector \mathbf{z} . Our algorithm, *Binary Orthogonal Matching Pursuit* (BOMP), is a variant of NNOMP and it proceeds in iterations. At iteration t , BOMP starts with a vector \mathbf{x}^{t-1} with $(t-1)$ entries of value 1. These entries correspond to the columns of the matrix \mathbf{R} that have been selected up to this iteration. Let $\hat{\mathbf{z}}^{t-1} = \mathbf{R}\mathbf{x}^{t-1}$ denote the approximation of the target vector \mathbf{z} constructed so far. The algorithm selects the column from \mathbf{R} (not selected so far) that has the largest dot-product with the residual $\mathbf{z} - \hat{\mathbf{z}}^{t-1}$ of the target vector. The set of selected indexes is augmented with this new index to produce vector \mathbf{x}^t . The algorithm terminates when we have selected k columns. The set of columns define the set T_s of nodes whose internal opinions will be set to zero.

The computational complexity of the BOMP algorithm is $O(kn^2)$. In each of the k iterations, the algorithm computes the dot-product of every candidate index to be added to set of selected indices with the residual vector. This step is the most computationally expensive, requiring time $O(n^2)$. All the other steps require at most $O(n)$ time, resulting in $O(kn^2)$ complexity in total.

The GreedyInt algorithm: We also consider a greedy algorithm for the problem. The algorithm builds the selected set of nodes T_s iteratively. It starts with an empty set T_s^0 , and at each step t it adds to the existing solution T_s^{t-1} the node v which, when added to the solution, $T_s^t = T_s^{t-1} \cup \{v\}$, it causes the largest decrease $\pi(\mathbf{z} | T_s^{t-1}) - \pi(\mathbf{z} | T_s^t)$ in the objective function. We will denote this algorithm as GreedyInt.

The GreedyInt algorithm can be implemented efficiently by exploiting the observation that the effect of neutralizing a node in the graph in the expressed opinion \mathbf{z} can be computed by the subtraction of the corresponding column of the matrix \mathbf{R} from \mathbf{z} . Therefore, for each candidate node v we need time $O(n)$ to compute $\pi(\mathbf{z} | T_s^{t-1} \cup \{v\})$, resulting in complexity $O(n^2)$ for each iteration, and $O(kn^2)$ for the algorithm in total.

CHAPTER 5

MODERATING EXPRESSED OPINIONS

5.1 Problem definition

5.2 Problem complexity

5.3 Algorithms

In this Chapter, we define the MODERATEEXPRESSED problem, we analyze its complexity, and we design an efficient algorithm for solving it.

5.1 Problem definition

When moderating expressed opinions, we seek a small set of nodes T_z to set their *expressed* opinions to zero, such that the polarization index is minimized. We use $\pi(\mathbf{z} \mid T_z)$ to denote the polarization index after setting the expressed opinions of the nodes in T_z to zero. The formal problem definition is the following.

Problem 2 (MODERATEEXPRESSED). *Given a graph $G = (V, E)$ a vector of internal opinions \mathbf{s} , and an integer k , identify a set T_z of k nodes such that fixing the expressed opinions of the nodes in T to 0, minimizes the polarization index $\pi(\mathbf{z} \mid T_z)$.*

5.2 Problem complexity

We prove the following Theorem for the hardness of the MODERATEEXPRESSED problem.

Theorem 5.1. *The MODERATEEXPRESSED problem is NP-hard.*

Proof. The proof of the theorem follows closely the proof of hardness in [19], so we only provide the correspondence between the two proofs. The proof exploits the equivalence between the opinion formation model and absorbing random walks, shown in [19].

Similar to the proof in [19] our proof uses a reduction from the VERTEX COVER ON REGULAR GRAPHS problem (VCRG) [14]. We show that there exists a set of nodes Y for a regular graph G_{VC} in the VCRG problem such that $|Y| \leq K$ and Y is a vertex cover if and only if there exists a set T_s for a graph G in the MODERATEEXPRESSED problem, such that $|T_z| \leq k$ and $\pi(\mathbf{z} | T_z) \leq \theta$, for $\theta = \frac{n}{2(d+1)^2}$. In our construction we set $G = G_{VC}$, and we initialize the vector \mathbf{s} , such that $s_i = 1$ for all $i \in V$. The proof then proceeds in the same way as in [19]. We can show that we can achieve a value $\pi(\mathbf{z} | T_z)$ less than θ if and only if the nodes T_z that we select in G (to make absorbing) define a vertex cover in G_{VC} . \square

Using a similar example as the one we used in Section 4.2 we can show that $\pi(\mathbf{z} | T_z)$ is also not monotone with respect to T_z , implying again that it is not straightforward to design an algorithm for solving MODERATEEXPRESSED.

5.3 Algorithms

Our algorithm for MODERATEEXPRESSED is a greedy algorithm, which we call GreedyExt. GreedyExt is an iterative algorithm which starts with an empty set T_z^0 . At each step t the algorithm adds to the existing solution T_z^{t-1} the node v_i , which, when setting $z_i = 0$, it causes the largest decrease $\pi(\mathbf{z} | T_z^{t-1}) - \pi(\mathbf{z} | T_z^t)$ in the objective function.

A naive implementation of the GreedyExt algorithm is computationally expensive. At each step of the algorithm we need to check n nodes, and for each node compute the new opinion vector after setting the expressed opinion of the node to zero. The most straightforward way to do this is by multiplying \mathbf{s} directly with \mathbf{Q} , in $O(n^2)$ time; recall that $\mathbf{Q} = (\mathbf{L} + \mathbf{I})^{-1}$, and thus it is a dense matrix. Alternatively, one can

iteratively apply Equation (3.1) and achieve the same computation in time $O(mI)$, where I the number of iterations it will take until convergence and m the number of edges of G . In our experiments, this computation converges in about a hundred iterations. Thus if we implement GreedyExt using the first method its running time becomes $O(kn^3)$, while with the second method the running time is $O(knmI)$. Our experiments with large graphs show that both these computations are impractical when dealing with medium-size datasets.

In order to improve the overall running time of GreedyExt, we exploit the *Sherman-Morrison* formula [21], a special case of the *Woodbury matrix identity*, to speed up the computation of the updated polarization index after adding a new node to the solution set. The identity states that the inverse of a matrix after adding a *rank-1* correction matrix to it can be computed by doing a *rank-1* correction to the inverse of the original matrix. Formally, given an invertible matrix \mathbf{A} and vectors \mathbf{u} and \mathbf{v}^T , the Sherman-Morrison formula states that:

$$(\mathbf{A} + \mathbf{u}\mathbf{v}^T)^{-1} = \mathbf{A}^{-1} - \frac{\mathbf{A}^{-1}\mathbf{u}\mathbf{v}^T\mathbf{A}^{-1}}{1 + \mathbf{v}^T\mathbf{A}^{-1}\mathbf{u}} \quad (5.1)$$

Consider now the case where we want to add node v_i to the solution set T_z . Let $\pi(\mathbf{z}')$ denote the new polarization index after setting $z_i = 0$. We can express the polarization index as $\pi(\mathbf{z}') = \|\mathbf{Q}'\mathbf{s}'\|^2$. In this equation, $\mathbf{Q}' = (\mathbf{L}' + \mathbf{I})^{-1}$, where \mathbf{L}' is the updated Laplacian with a row of zeros at the i -th index, and \mathbf{s}' is the updated internal opinion vector, with $s_i = 0$ at the i -th entry. To understand the update process of \mathbf{Q} and \mathbf{s} , note that in the random walk interpretation, setting $z_i = 0$ is equivalent to removing all outgoing edges from node v_i and keeping only the edge to the node x_i , while setting $s_i = 0$.

We can express the Laplacian matrix \mathbf{L}' as a *rank-1* correction of the Laplacian \mathbf{L} , that is, $\mathbf{L}' = \mathbf{L} + \mathbf{u}\mathbf{v}^T$, where \mathbf{u} is the unit vector with 1 at the i -th entry, and \mathbf{v}^T is the negative i -th row of the Laplacian. Following the Sherman-Morrison formula (Equation 5.1) we have that $\mathbf{Q}' = \mathbf{Q} - \mathbf{B}$, where

$$\mathbf{B} = \frac{\mathbf{Q}\mathbf{u}\mathbf{v}^T\mathbf{Q}}{1 + \mathbf{v}^T\mathbf{Q}\mathbf{u}},$$

We can also write $\mathbf{s}' = \mathbf{s} - \bar{\mathbf{s}}$, where $\bar{\mathbf{s}}$ is a vector with $\bar{s}_i = s_i$, and zero in all other

entries. Thus, we have:

$$\begin{aligned}
\|\mathbf{Q}'\mathbf{s}'\|^2 &= \|(\mathbf{Q} - \mathbf{B})(\mathbf{s} - \bar{\mathbf{s}})\|^2 \\
&= \|\mathbf{Q}\mathbf{s} - \mathbf{B}\mathbf{s} - \mathbf{Q}\bar{\mathbf{s}} + \mathbf{B}\bar{\mathbf{s}}\|^2 \\
&= \left\| \mathbf{z} - \frac{\mathbf{Q}\mathbf{u}\mathbf{v}^T\mathbf{z}}{1 + \mathbf{v}^T\mathbf{Q}\mathbf{u}} - \mathbf{Q}\bar{\mathbf{s}} + \frac{\mathbf{Q}\mathbf{u}\mathbf{v}^T\mathbf{Q}\bar{\mathbf{s}}}{1 + \mathbf{v}^T\mathbf{Q}\mathbf{u}} \right\|^2.
\end{aligned} \tag{5.2}$$

In order to efficiently compute the quantity in Equation (5.2), we perform the operations in such an order, so that we never need to compute any $n \times n$ matrix. As a result we can compute Equation (5.2) in time $O(n)$, which is better than the $O(mI)$ complexity of the power-iteration, given that $m = nd$, where d is the average degree of the graph.

First, we compute the vector $\mathbf{w} = \frac{\mathbf{Q}\mathbf{u}}{1 + \mathbf{v}^T\mathbf{Q}\mathbf{u}}$. This can be computed in linear time. Since \mathbf{u} is the unit vector with 1 in one entry, and zero everywhere else, $\mathbf{Q}\mathbf{u}$ can be obtained in $O(1)$ via column selection. Given the vector $\mathbf{Q}\mathbf{u}$ we can compute $\mathbf{v}^T\mathbf{Q}\mathbf{u}$ in $O(n)$ time, and then obtain \mathbf{w} by scaling $\mathbf{Q}\mathbf{u}$, bringing the total computational cost of \mathbf{w} to $O(n)$.

Given the vector \mathbf{w} we can now compute $\mathbf{w}\mathbf{v}^T\mathbf{z}$ (the second term in Equation (5.2)) in linear time, by first computing the dot-product $\mathbf{v}^T\mathbf{z}$, and then scaling the vector \mathbf{w} with the result. Also, we can compute the vector $\mathbf{Q}\bar{\mathbf{s}}$ in $O(n)$ time, by first selecting the column of \mathbf{Q} and then scaling it by s_i . The term $\mathbf{w}\mathbf{v}^T\mathbf{Q}\bar{\mathbf{s}}$ (the last term in Equation (5.2)) can be computed as before in linear time. All other computations are computations on vectors, resulting in $O(n)$ total cost for the computation of Equation (5.2).

We repeat the above procedure n times to find the best candidate node. For the selected node, we compute the updated matrix $\mathbf{Q}' = \mathbf{Q} - \mathbf{B}$ using the Sherman-Morrison formula in $O(n^2)$. This brings the total computational cost of GreedyExt to $O(kn^2)$.

CHAPTER 6

EXPERIMENTS

6.1 Datasets

6.2 Evaluation of the Polarization Index

6.3 Heuristic Algorithms for Opinion Moderation

6.4 Evaluation of algorithms for ModerateInternal

6.5 Evaluation of the algorithms for ModerateExpressed

6.6 Scalability

6.7 Case study

In this Chapter, we present an experimental evaluation of the polarization index, and of our algorithms for both problems. The goals of our experiments are to validate the polarization index, study the properties of the proposed algorithms, and evaluate their performance and scalability.

6.1 Datasets

We consider five datasets representing different types of social networks. We use networks that are partitioned into opposing communities, and there is ground-truth data about the community membership of the nodes. Thus, we can naturally assign

internal opinions -1 and 1 to the nodes depending on their community membership. We consider the following datasets:

*Karate*¹: This dataset represents a social network of friendships between 34 members of a karate club at a US university in the 1970s. The social network is partitioned into two distinct equal-sized communities that correspond to two fractions built around two rival instructors.

*Books*²: This is a network of books about US politics published around the time of the 2004 presidential election and sold by the online bookseller Amazon.com. Edges between books represent frequent co-purchasing. Books are classified as *Liberal*, *Conservative*, and *Neutral*. There are in total 43 liberal books, 49 conservative and 13 neutral . We handled the neutral nodes by assigning to them internal opinion zero.

*Blogs*³: A directed network of hyperlinks between weblogs on US politics, recorded in 2005 by Adamic and Glance [1]. Blogs are classified as either *Liberal* or *Conservative*. We converted the social graph into an undirected one and only kept the largest connected component. The resulting dataset contains two communities with 636 and 586 nodes each, and 19,089 edges.

Elections: This dataset is the network between the Twitter followers of Hillary Clinton and Donald Trump collected in the period 15/12/2016-15/01/2017 – around the time of the 2016 presidential elections. Members of this network are assigned an internal opinion of 1 or -1 based on which one of the two candidates they follow. Followers of both candidates are assigned a neutral opinion. Since the dataset is prohibitively large (20M followers), we only considered the network formed by the first 50,000 users, according to their user id. We took the largest connected component and iteratively pruned nodes to guarantee that every node has degree greater than 1. The resulting network had a disproportionately large number of Clinton followers so we subsampled her followers to ensure that the ratio of followers for each side reflected the one in the full dataset. In the resulting network there are 7,715 Hillary Clinton followers, 8,336 Donald Trump followers, and 2,216 Neutral followers, for a total of 18,267 users with 204,040 connections between them. As before, we treat the network as undirected.

Hashtags: Using the followers of Clinton and Trump that we collected, we also created

¹<https://networkdata.ics.uci.edu/data.php?id=105>

²<https://networkdata.ics.uci.edu/data.php?id=8>

³<https://networkdata.ics.uci.edu/data.php?id=102>

Table 6.1: Dataset Statistics

Dataset	Nodes	Edges	Avg Degree	Diameter	Positive	Negative	Neutral
<i>Karate</i>	34	78	4.58	5	17	17	0
<i>Books</i>	105	441	8.4	7	43	49	13
<i>Blogs</i>	1,222	16,717	27.36	8	636	586	0
<i>Elections</i>	18,267	204,040	22.33	8	7,715	8,336	2,216
<i>Hashtags P</i>	18,890	269,696	28.55	7	12,281	6,612	0
<i>Hashtags NP</i>	18,890	269,696	28.55	7	12,408	4,102	2,383

“topical” networks based on the hashags that they tweeted, where we assign the opinions according to the specific hashtag that the users tweeted. We considered two pairs of hashtags: The #maga and #imwithher hashtags, which we expect to be polarized, and the #halloween and #walkingdead hashtags for which we do not expect to have polarization. We selected these hashtags since they are among the most popular in the dataset. We sampled users that have tweeted at least one hashtag from both pairs, and we created the follow network between them. Again, we kept the largest connected component and iteratively pruned nodes to guarantee that every node has degree greater than 1. The resulting network has 18,890 nodes and 269,696 edges. Using the graph we created, we consider two possible settings for the opinions: In the first, we assign opinion -1 to the users that have tweeted the hashtag #maga, 1 if they have tweeted #imwithher, and 0 if they have tweeted both. We will refer to this dataset as *Hashtags P*. In the second, we assign opinion -1 to the users that have tweeted the hashtag #halloween, 1 if they have tweeted #walkingdead, and 0 if they have tweeted both. We will refer to this dataset as *Hashtags NP*. These two different settings allow us to study the behavior of our metric in a polarized (*Hashtags P*), and non-polarized network (*Hashtags NP*).

Table 6.1 summarizes the statistics of our datasets. For all datasets we treat the graphs as undirected. When applying the opinion formation model, we set all edge weights, and all opinion weights to be 1. In order to handle the cases where there is an imbalance of opinions, we normalize the opinion values by subtracting the mean opinion and dividing by the difference of the maximum and the minimum. In this way, the mean opinion value becomes zero, which we consider to be the moderate stance.

6.2 Evaluation of the Polarization Index

In this section we evaluate the metric in its ability to identify polarization. We compute the value of the metric for the five different datasets. Table 6.2 shows the values we obtain. In order to understand if the values are indicative of polarization, we perform a *randomization* test, where we randomly assign the internal opinions on the graph. A randomized assignment of opinions that is independent of the network structure does not create any opinion clusters, and thus it corresponds to a non-polarized state. We compare the value of the π with that of the random assignment, in order to understand the significance of the polarization index value.

We create 100 random assignments of the opinion values, and we report the average and standard deviation of the polarization index values we obtain for these cases. We observe that the polarization index values are significantly higher than those in the randomized datasets in all networks except for the *Elections* and *Hashtags NP* datasets, where the π values are small, and close to that of the random assignment. In the case of the *Hashtags NP* dataset we obtain essentially the same π value.

It is interesting to contrast the π values for the *Hashtags P* and *Hashtags NP* datasets. In the first case, the polarization is much higher, which agrees with our intuition that these hashtags are adopted by different communities that do not interact with each other. In the second case the polarization index is much lower, and close to that of the random assignment. This suggests that the distribution of the opinions in the second case cuts across the natural communities that appear in the graph. Although users are organized in weakly connected communities, positive and negative opinions appear in both, and as a result there is no polarization and echo-chamber effect. This experiment highlights the importance of taking the opinions into account for measuring polarization; looking only at the network structure it is not possible to differentiate between these two cases.

Finally, we perform a comparison between our metric, and the two closest state-of-the-art metrics for polarization. The first one, the *network disagreement index*, was defined by Dandekar et al. [11], and calculates the sum of pairwise distances of opinions across all edges in the Graph. The metric was defined to characterize an opinion formation *process* as polarizing but its value is also indicative of polarization in the network. The second one, is the *controversy score* defined by Garimella et al. [16], used to characterize whether a particular conversation graph based on a *topic* is

Table 6.2: Dataset polarization index and randomization values

Dataset	π	Mean π for random assignments	Std. deviation
<i>Karate</i>	0.089	0.022	0.00499
<i>Books</i>	0.107	0.007	0.00172
<i>Blogs</i>	0.029	0.012	0.00027
<i>Elections</i>	0.012	0.011	0.00007
<i>Hashtags P</i>	0.028	0.005	0.00004
<i>Hashtags NP</i>	0.0049	0.0044	0.00005

Table 6.3: Comparison of polarization metrics on all Datasets

Dataset	π	Network Disagreement Index	Controversy Score
<i>Karate</i>	0.089	7.238	0.71
<i>Books</i>	0.107	16.371	0.51
<i>Blogs</i>	0.029	231.15	0.39
<i>Elections</i>	0.012	2067.67	0.12
<i>Hashtags P</i>	0.028	3430.02	-
<i>Hashtags NP</i>	0.0049	1337.07	-

polarized. Although not exactly analogous, we consider both of these metrics similar to our own, and hence perform a comparison in the values of the metrics, for all datasets. The values are reported in the Table 6.3

We observe that the network disagreement index and the controversy score values are relative to the π values. There is a clear correlation in the values when looking across datasets. In the *Karate* dataset, the controversy score gave a high value, since it is based on the network structure, and the *Karate* dataset exhibits a highly clustered structure. It is also worth mentioning that the controversy score was too computationally intense to carry out on the *Hashtags* dataset. However, the metric would be unable to discern between the *Hashtags P* and *Hashtags NP* datasets since it does not explicitly consider opinions.

6.3 Heuristic Algorithms for Opinion Moderation

In addition to the algorithms we described in Chapters 4 and 5, we also consider a few more scalable heuristics. Our reasoning in the design of the heuristics is that in order to moderate the overall expressed opinion we need to convert to neutral the opinions of individuals that express extreme opinions, individuals belonging to extreme neighborhoods, or individuals that are influential in the network. The following algorithms implement this reasoning.

ExtremeExpressed: This heuristic works iteratively and at each step it selects to neutralize the node v with the highest expressed opinion $|z_v|$. Since it requires $O(n)$ time to find the most extreme node, the complexity of the algorithm is determined by the time required to compute the updated \mathbf{z} vector after neutralizing a node. As we have shown in Section 4.3, we can efficiently calculate the new \mathbf{z} vector by subtracting the column of \mathbf{Q} corresponding to the neutralized node from the current \mathbf{z} vector. Therefore, the algorithm has complexity $O(kn)$. In the case of the MODERATEEXPRESSED problem, the fastest way to compute the new \mathbf{z} is by iteratively updating the z_i values as defined in Equation (3.1), until convergence. The updates are implemented using efficient matrix-vector multiplication. This takes time $O(mI)$, where I is the number of iterations required for convergence, and m is the number of edges in the graph, leading to complexity $O(kmI)$ for the algorithm. In practice, we have found that the algorithm converges in less than 100 iterations.

ExtremeNeighbors: In this heuristic we select the next node to neutralize based on how extreme the neighborhood of the node is. The intuition is that neutralizing this node will have an effect on many extreme nodes. The algorithm at each step changes the opinion of the node v whose neighbors have the highest absolute sum of expressed opinions, that is, $v = \arg \max_{i \in V} |\sum_{j \in N(i)} \mathbf{z}_j|$. For every node we need to check its neighbors, which takes $O(m)$ time, and then update \mathbf{z} , accordingly. Therefore, if we use the efficient update of \mathbf{z} as above, the complexity of the algorithm when solving MODERATEINTERNAL is $O(k(n + m))$. Using the iterative method to compute \mathbf{z} , the complexity of the algorithm when solving MODERATEEXPRESSED is $O(k(n + m)I)$.

Pagerank: The idea behind this heuristic is that in order to moderate the overall opinion, it is a good idea to neutralize the nodes that are central in the network. This will result in maximum spread of a balanced viewpoint. We use PageRank [25] to measure the centrality of a node. The algorithm selects the nodes in decreasing order

of their PageRank value. The algorithm’s complexity is $O((m + n)I + n \log n)$, where $O((n + m)I)$ is the time the PageRank values, and $O(n \log n)$ is the time required to sort the nodes.

6.4 Evaluation of algorithms for ModerateInternal

We first evaluate our algorithms with respect to the value they achieve for the objective function π . We evaluate on all five networks. For the *Hashtags* network we evaluate for both the hashtags #maga and #imwithher and the #halloween and #walkingdead hashtags to set the opinions. Figures 6.1 6.2 6.3 6.4 6.5 6.6 show the value of $\pi(\mathbf{z} \mid T_s)$ for different sizes of the solution set $|T_s| = k$, for all datasets. For the smaller datasets *Karate* and *Books* we let k range over the full size of the dataset. This is impractical for the larger *Blogs*, *Elections* and *Hashtags* datasets, hence we consider T_s up to 10% of the dataset; we plot the value of π in increments of 1%.

As expected, the GreedyInt algorithm achieves the best performance in all datasets. The performance of GreedyInt is consistently matched by BOMP and ExtremeExpressed. However, in the *Hashtags P* and *Hashtags NP* BOMP loses its effectiveness after some point, and surprisingly, in the case of *Hashtags NP* after $k = 4\%$ even starts to increase polarization. This also coincides with ExtremeExpressed also losing somewhat its effectiveness, and indicates that the change in performance is due to a particular feature of the dataset. Nevertheless, this is worth investigating further. The Pagerank and ExtremeNeighbors algorithms are significantly worse, and in the big datasets achieve only a minimal reduction in π . While we expected the BOMP algorithm to be competitive with GreedyInt the performance of ExtremeExpressed was a surprise. We also compare the BOMP and GreedyInt algorithms against the optimal for k up to 50% of the graph, for the smallest dataset *Karate*, where this computation is possible. We observe that the GreedyInt algorithm behaves optimally, while BOMP achieves performance very close to optimal for $k \leq 6$, and coincides with it for $k > 6$.

Our results indicate in order to minimize polarization in the MODERATEINTERNAL problem, the best strategy is to moderate the nodes with the most extreme opinions. The Pagerank and ExtremeNeighbors algorithms that take into account how well a given node is connected to the network do not perform well.

We further investigate this observation by visualizing the nodes selected by GreedyInt

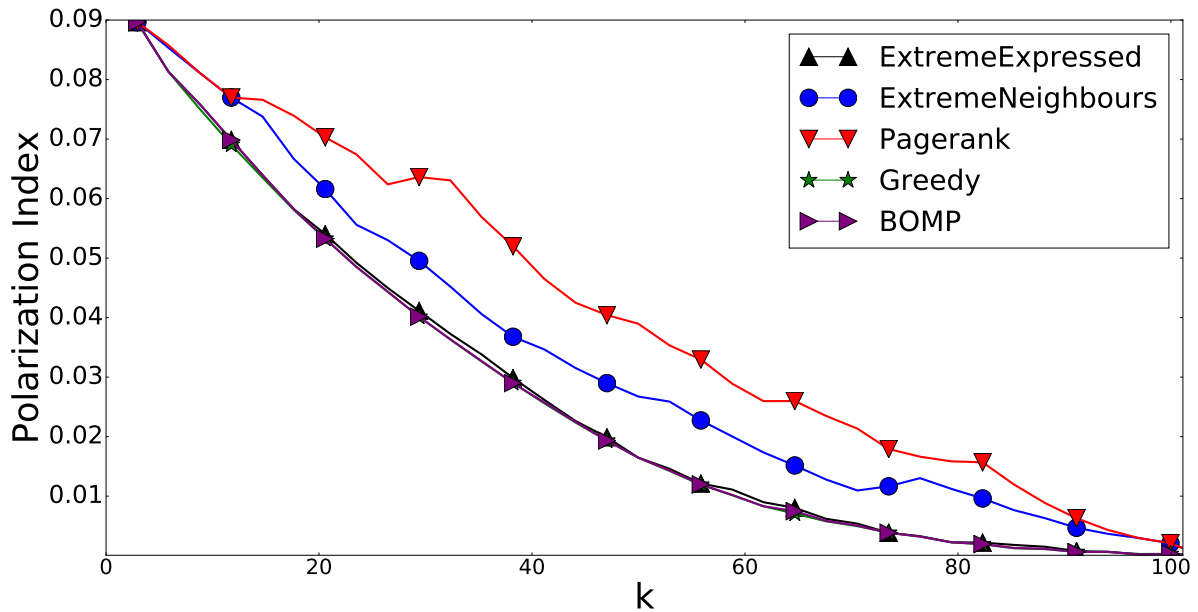


Figure 6.1: Performance of the algorithms for the MODERATEINTERNAL problem on Karate

in Figures 6.8 and 6.9, for the two smaller datasets, *Karate* and *Books*. In the visualization, we assign different color and shape to the nodes of the different communities. The nodes are numbered according to their selection order by GreedyInt. The first ten nodes are colored in orange-red and have larger size.

The visualization further confirms the behavior of GreedyInt: the nodes that are selected first are nodes on the outskirts of the network. This means that the impact on \mathbf{z} is bigger when moderating fringe nodes with extreme opinions, instead of central nodes. The broader implication of this is that in the case of the MODERATEINTERNAL problem the best we can do for moderating polarization is to change the opinions one user at-a-time, rather than “diffusing” moderation in the network. This in part due to the fact that the internal opinion is only one of the contributing factors to the expressed opinion of an individual, and thus its change has a limited effect.

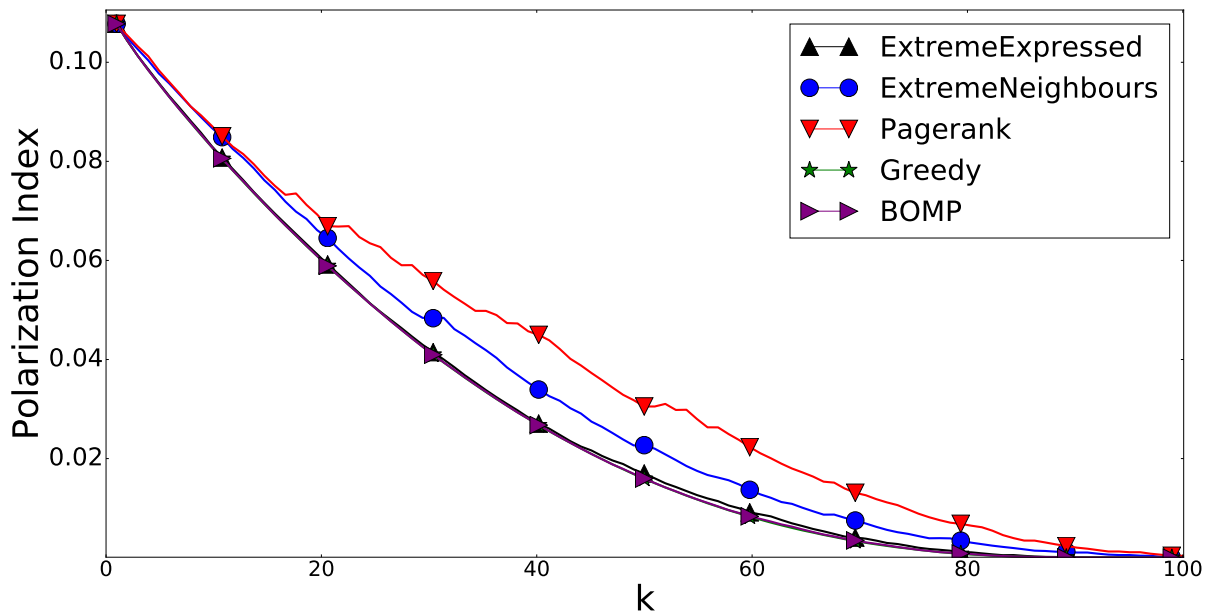


Figure 6.2: Performance of the algorithms for the MODERATEINTERNAL problem on Books

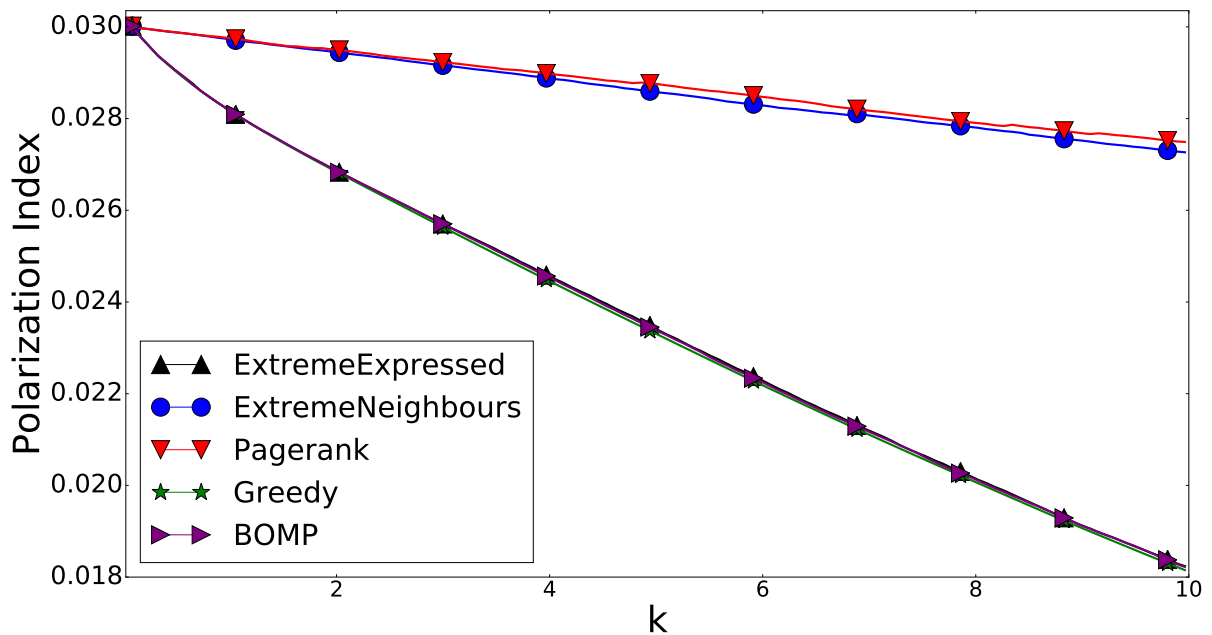


Figure 6.3: Performance of the algorithms for the MODERATEINTERNAL problem on Blogs

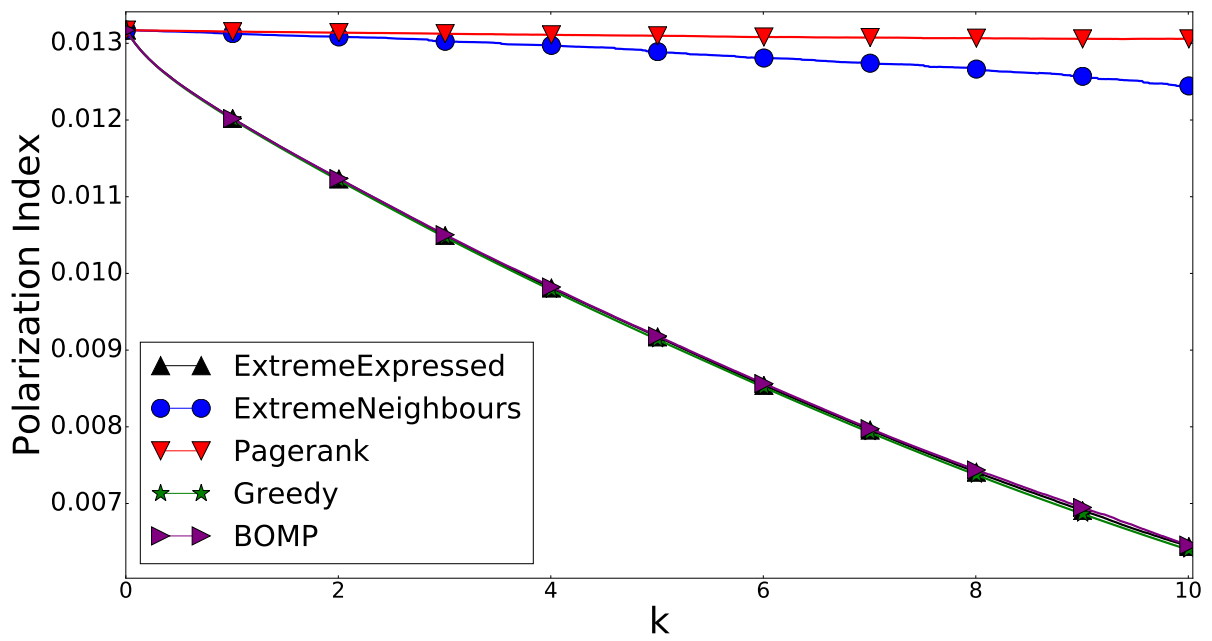


Figure 6.4: Performance of the algorithms for the MODERATEINTERNAL problem on Elections

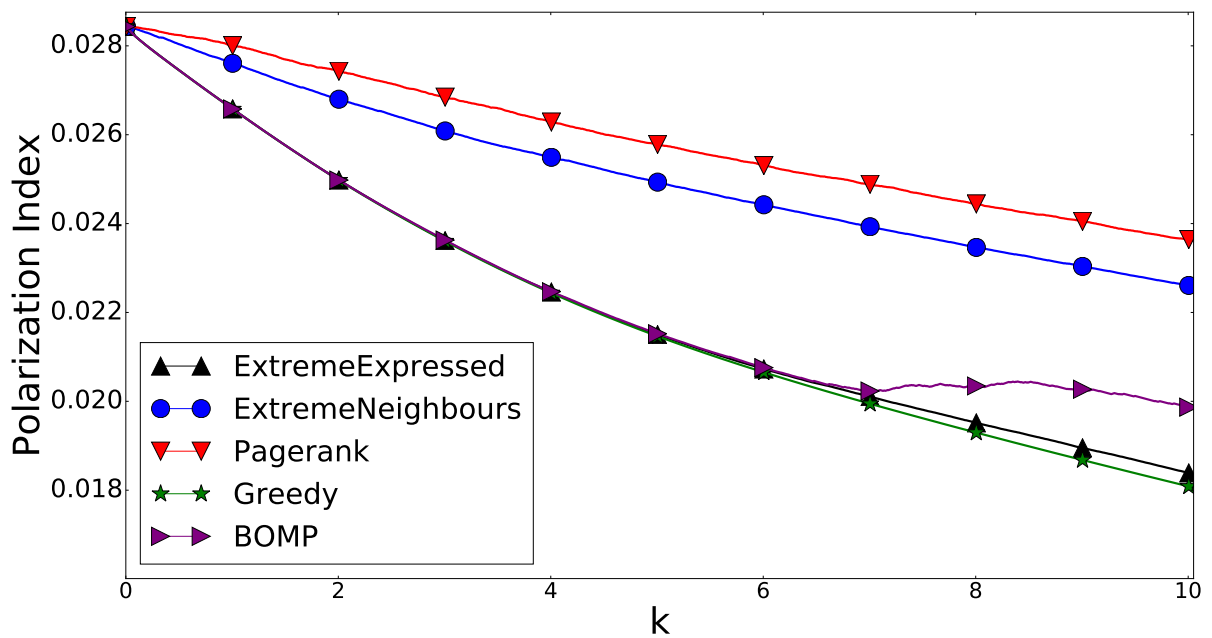


Figure 6.5: Performance of the algorithms for the MODERATEINTERNAL problem on Hashtags P

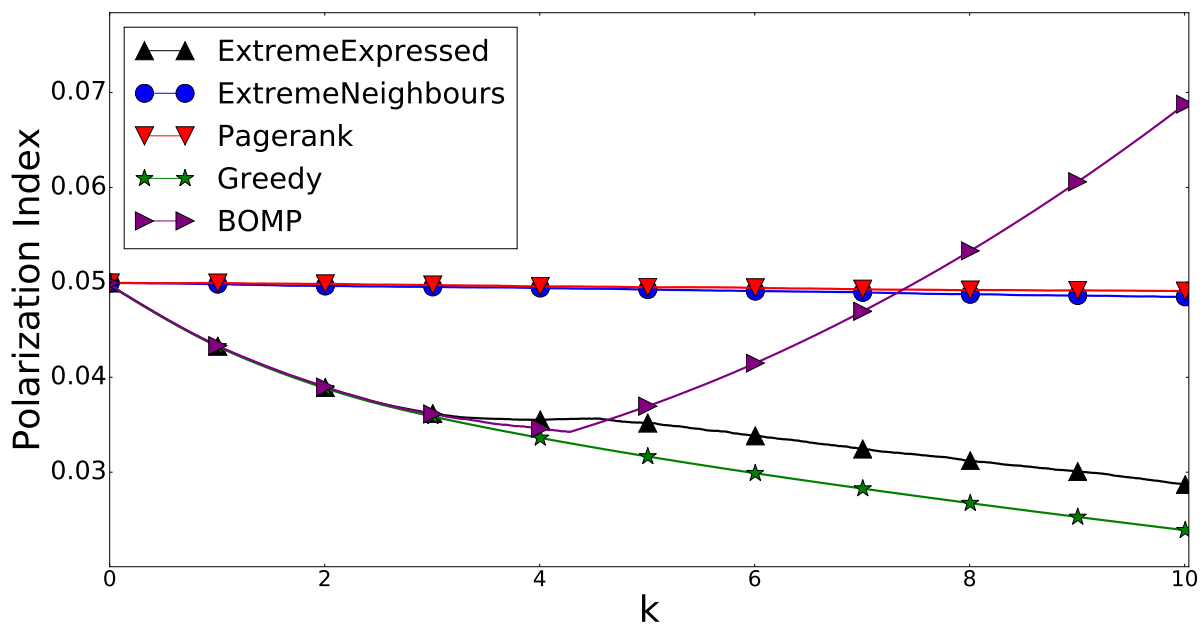


Figure 6.6: Performance of the algorithms for the MODERATEINTERNAL problem on *Hashtags NP*

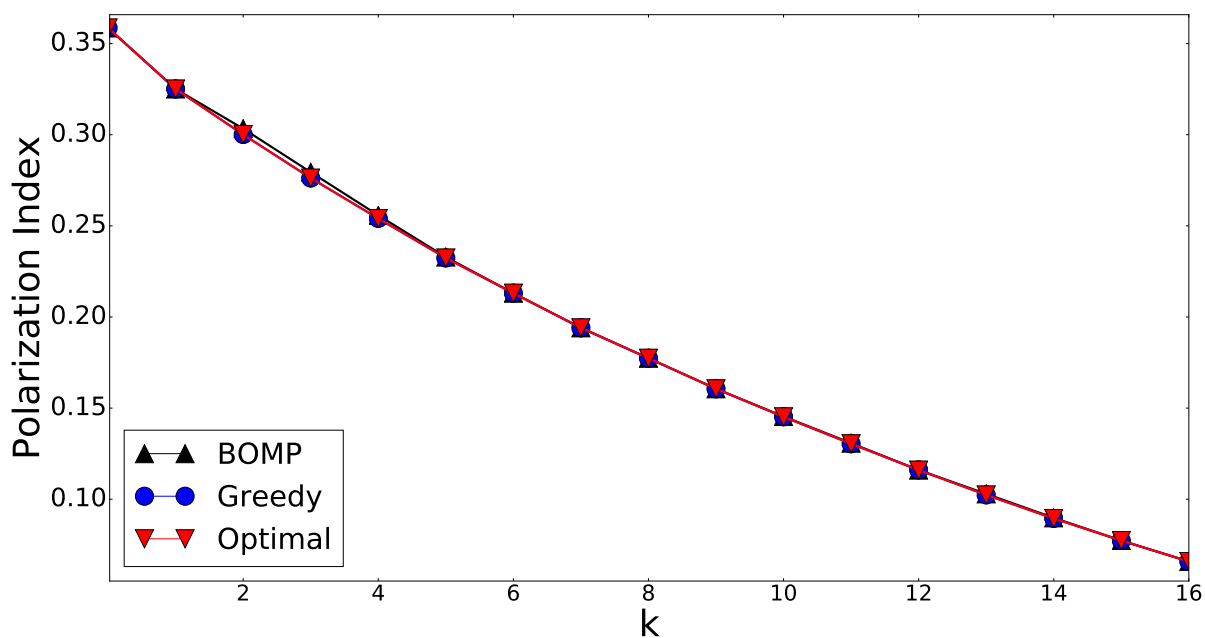


Figure 6.7: Comparison with optimal

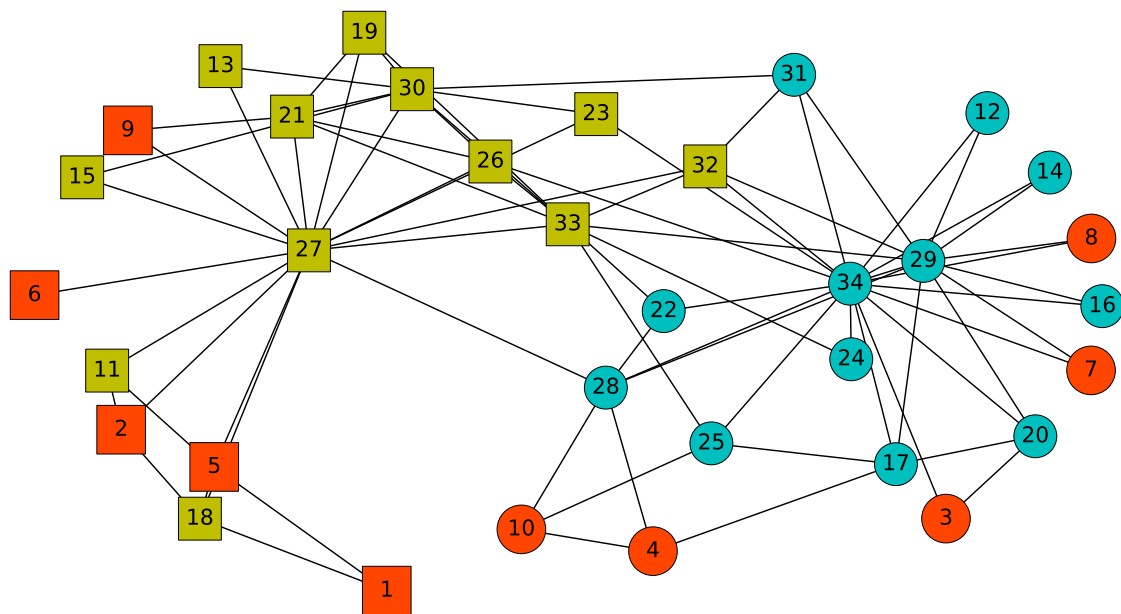


Figure 6.8: Selected nodes by GreedyInt on *Karate*

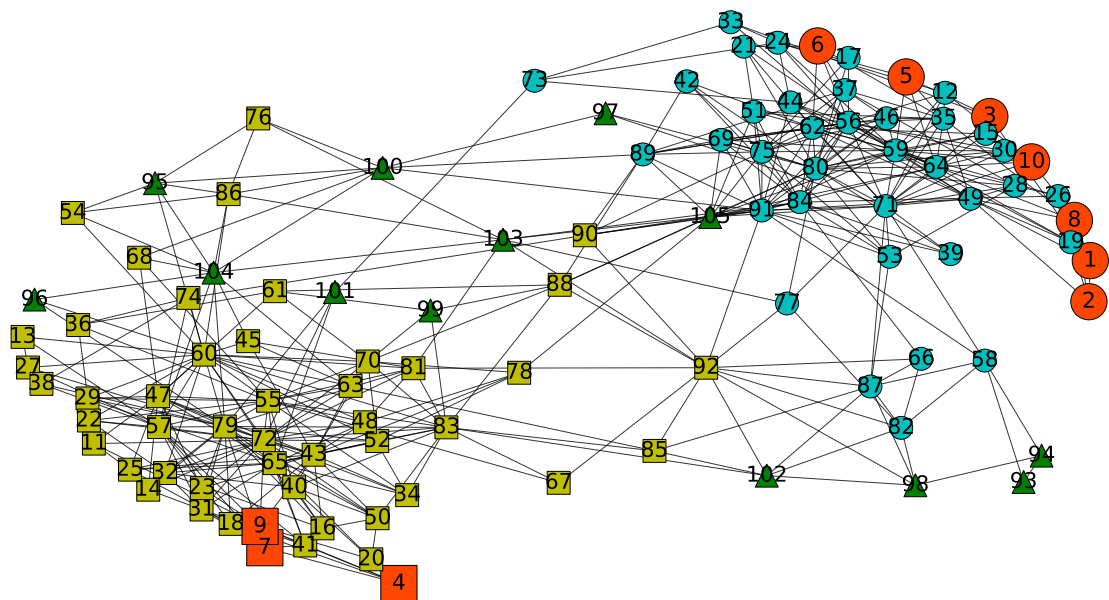


Figure 6.9: Selected nodes by GreedyInt on *Books*

6.5 Evaluation of the algorithms for ModerateExpressed

For the evaluation of the MODERATEEXPRESSED problem we follow the same methodology as for MODERATEINTERNAL. Figures 6.10 6.11 6.12 6.13 6.14 6.15 show the $\pi(\mathbf{z} | T_s)$ as a function of the size of T_s for all datasets.

As expected, the GreedyExt is again the best-performing algorithm. However, the performance of the other algorithms changes depending on the dataset. For the *Karate*, *Books*, *Blogs* and *Hashtags P* datasets, ExtremeNeighbors and Pagerank achieve performance close to that of GreedyExt, especially for smaller values of k , while ExtremeExpressed is clearly the worst performer. As the size of the solution increases, Pagerank and ExtremeNeighbors seem to lose their effectiveness, while ExtremeExpressed catches up with them. These results indicate that when moderating expressed opinions, it is a good strategy to select nodes that are relatively central and express an extreme opinion. After selecting a sufficient number of influential nodes, the gains of moderating central nodes is diminished, there is more benefit in neutralizing extreme nodes which were not affected by the influential ones, essentially, adopting the approach of moderating one node at the time.

However, we observe a very different picture in the *Elections* and *Hashtags NP* datasets, where ExtremeExpressed is almost as good as GreedyExt, and Pagerank and ExtremeNeighbors perform poorly. Note that, according to the randomization test, the *Elections* and *Hashtags NP* datasets are not very polarized. Therefore, there is sufficient mixing of opinions and it is not possible to moderate a large number of nodes by neutralizing an influential node. The one-node-at-the-time approach works better. In order to further investigate this claim we “zoom in” in the performance of the algorithms in the *Elections* dataset for k up to the top-1% of the nodes. Now, Pagerank and ExtremeNeighbors appear competitive for small k , but their performance fades in comparison to ExtremeExpressed as k increases. This is in stark contrast to the *Hashtags P* dataset, which is of similar size with *Elections* and it is very polarized, where the algorithms that change influential nodes achieve a good performance.

A final observation is that the reduction in $\pi(\mathbf{z})$ is significantly higher for the MODERATEEXPRESSED problem than for MODERATEINTERNAL for the same dataset, for the same k . This is expected, as the moderation of the expressed opinion has a much larger effect in the opinion of the individual, and the opinions in her social network, than the moderation of the internal opinion.

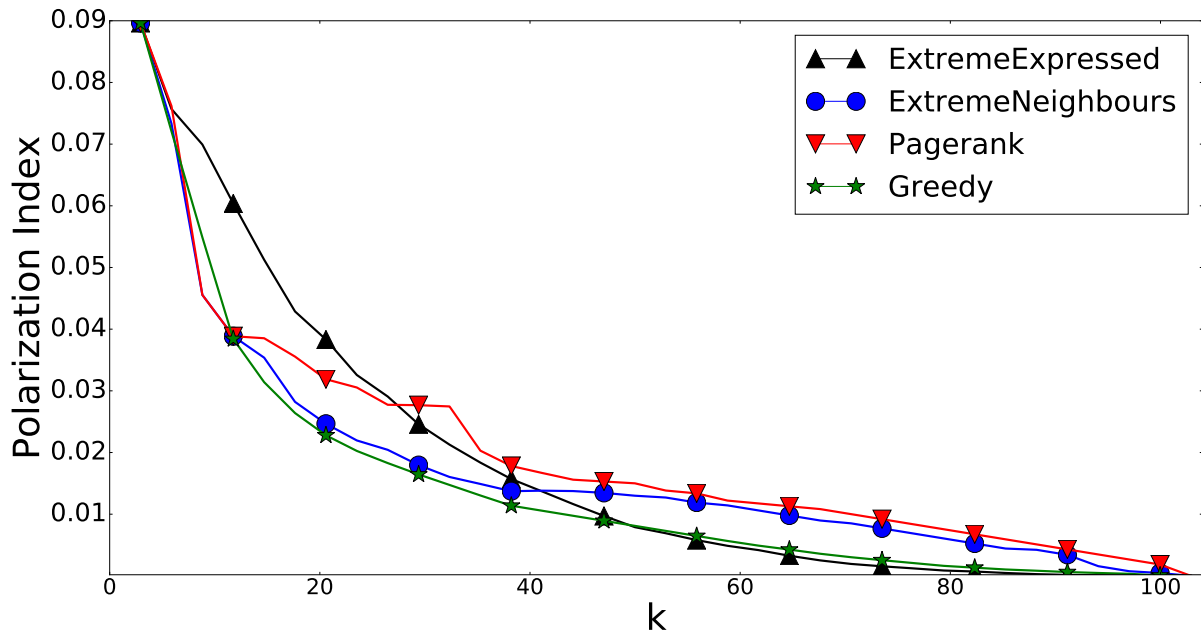


Figure 6.10: Performance of the algorithms for the `MODERATEEXPRESSED` problem on Karate

In Figures 6.17 and 6.17, we visualize again the selected nodes by GreedyExt for the *Karate* and *Books* datasets. The selection is different from the one we obtained for the `MODERATEINTERNAL` problem (Figures 6.9 and 6.8), and highlights the different nature of the two problems. In the solution of `MODERATEEXPRESSED` the nodes selected are more central in the graph. It is obvious that changing the expressed opinion of a node has a bigger impact on the opinions of the neighbors of that node. As a result, GreedyExt tries to pick nodes that are both central and extreme. The first selection of GreedyExt is the node that it is ranked first for both Pagerank and ExtremeExpressed. This combination is essential in achieving high reduction of π . As the selection process continues, the selections of GreedyExt alternate between central and fringe nodes as the algorithm is trying to “cover” different parts of the graph, and moderate opinions of nodes that are not easily reached by the central nodes.

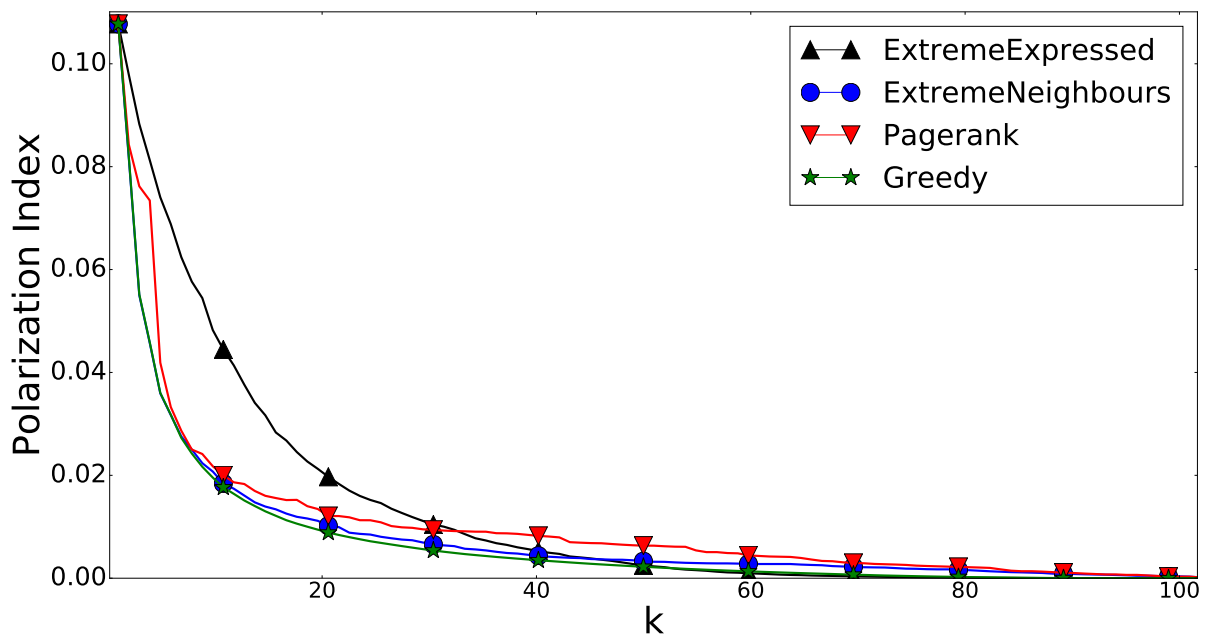


Figure 6.11: Performance of the algorithms for the MODERATEEXPRESSED problem on Books

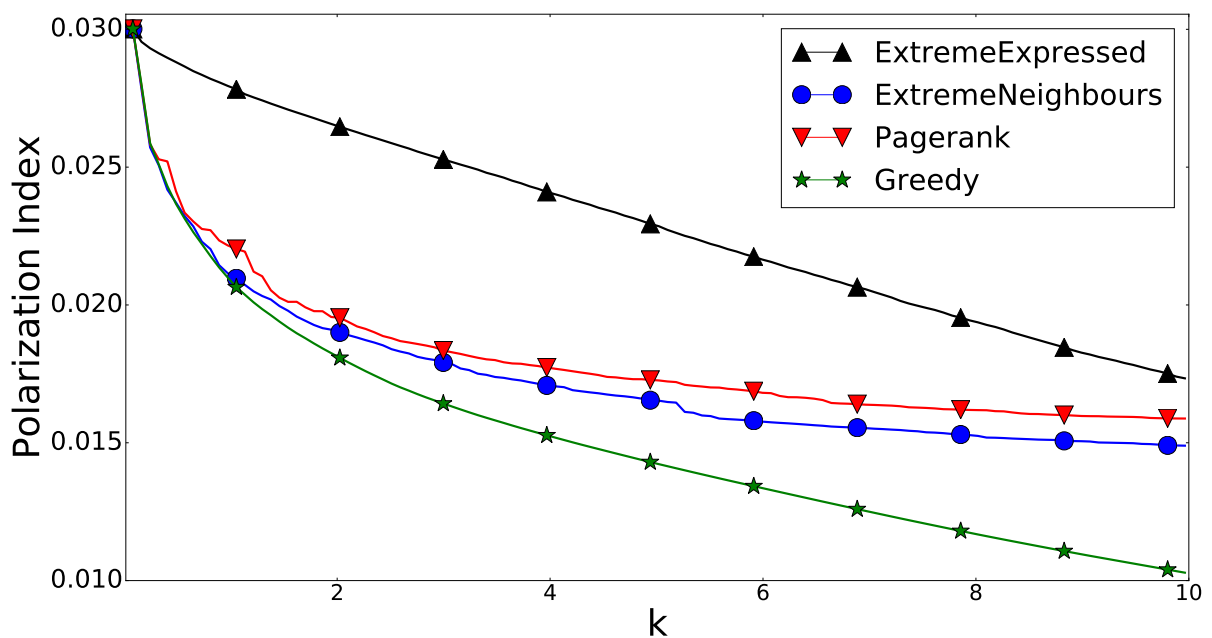


Figure 6.12: Performance of the algorithms for the MODERATEEXPRESSED problem on Blogs

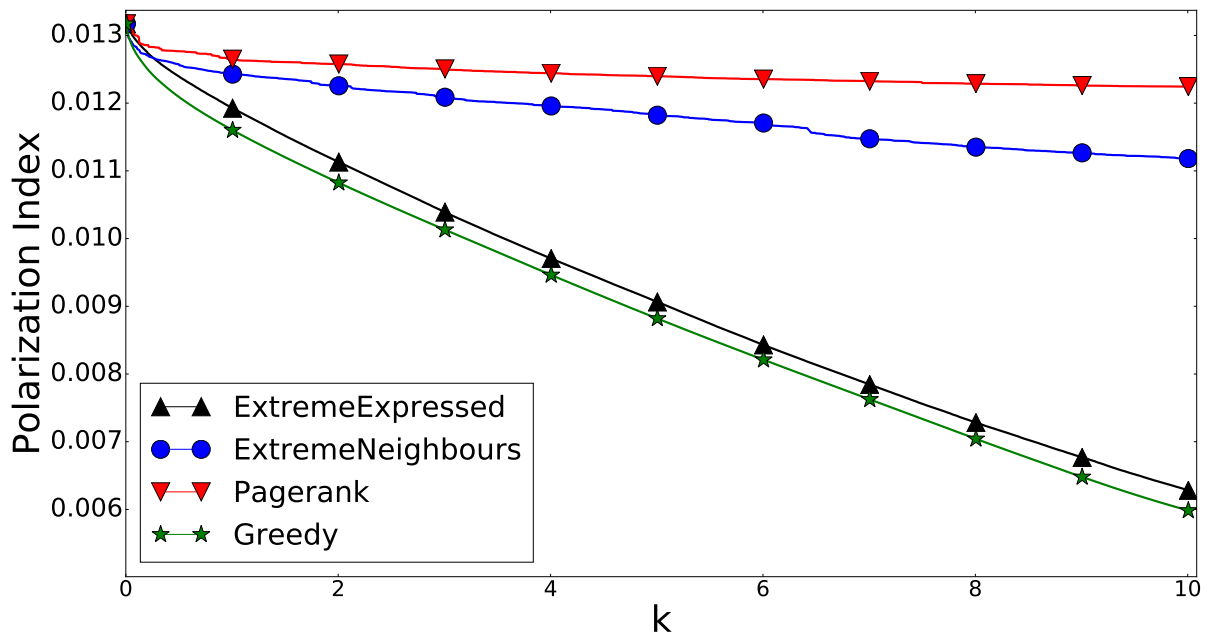


Figure 6.13: Performance of the algorithms for the MODERATEEXPRESSED problem on Elections

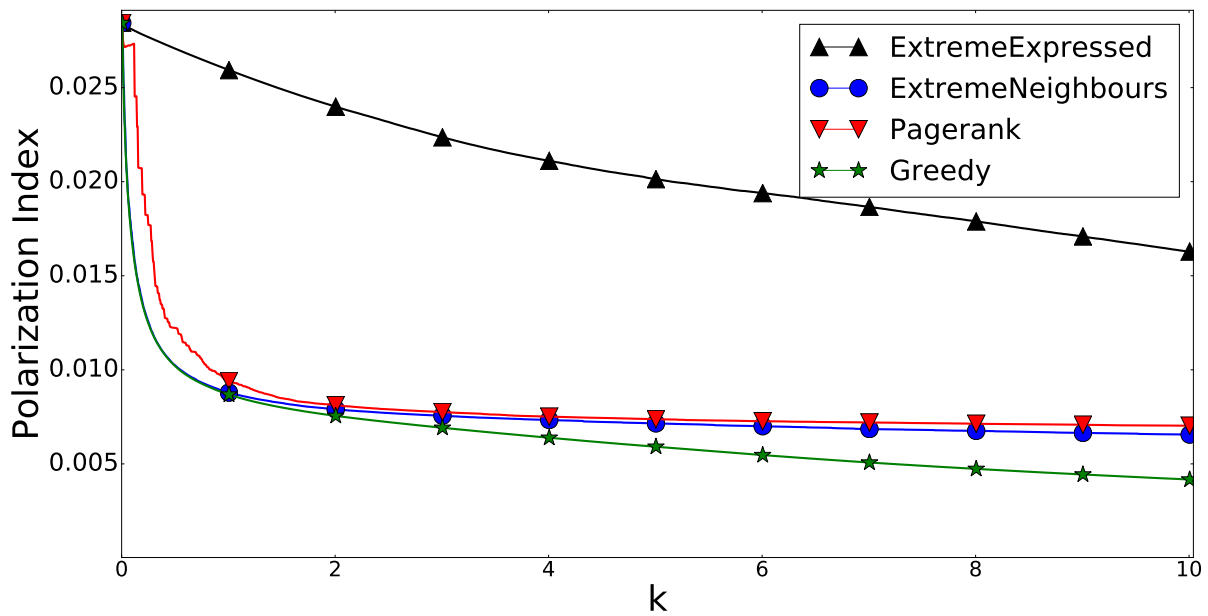


Figure 6.14: Performance of the algorithms for the MODERATEEXPRESSED problem on Hashtags P

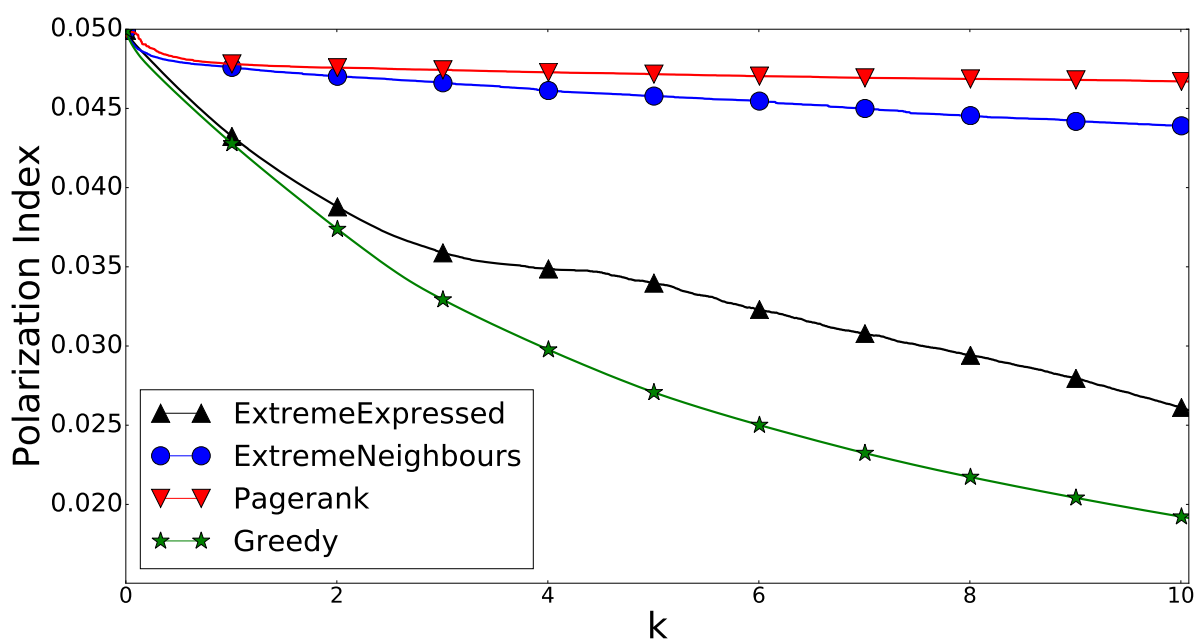


Figure 6.15: Performance of the algorithms for the MODERATEEXPRESSED problem on *Hashtags NP*

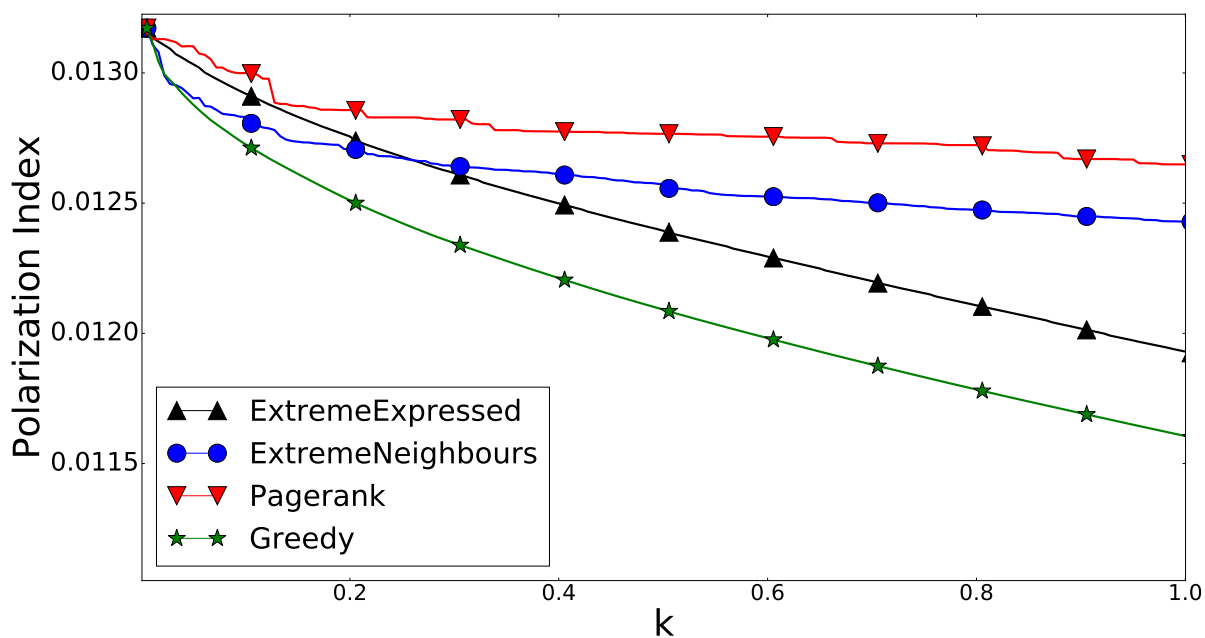


Figure 6.16: Performance of the algorithms for the MODERATEEXPRESSED problem on Elections top-1%

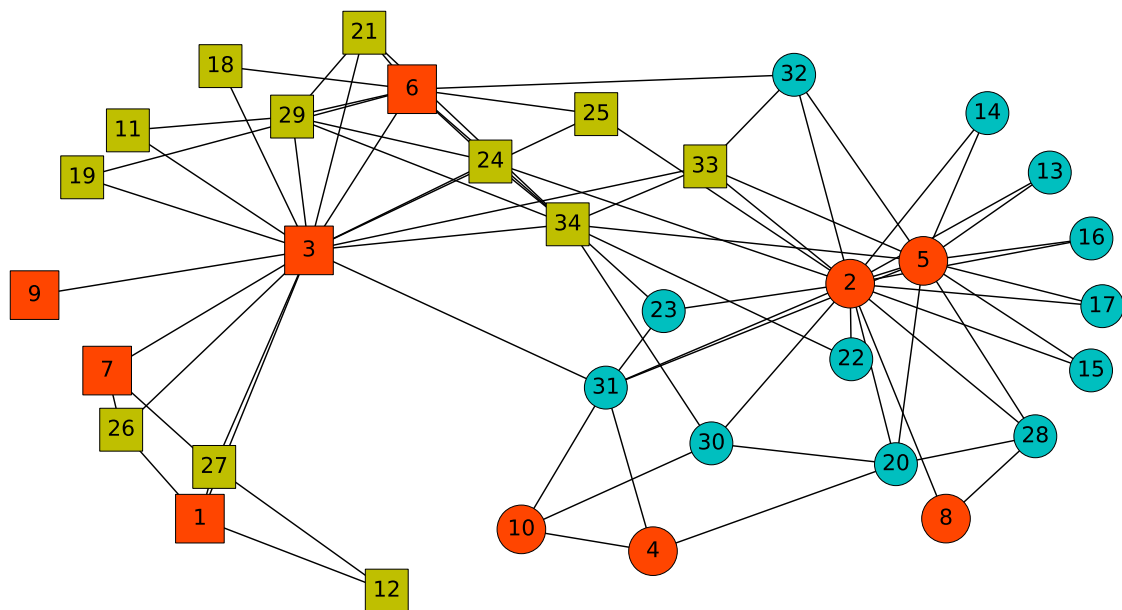


Figure 6.17: Selected nodes by GreedyExt on *Karate*

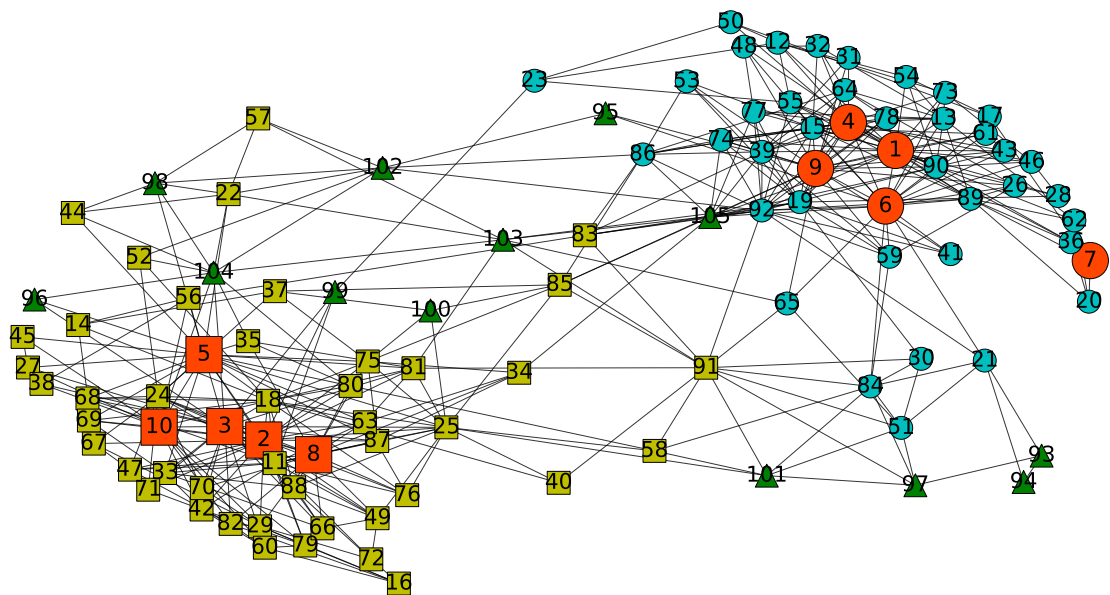


Figure 6.18: Selected nodes by GreedyExt on *Books*

6.6 Scalability

We now evaluate the scalability of our algorithms. Table 6.4 shows the running time for all algorithms on the *Elections* dataset for MODERATEINTERNAL and MODERATEEXPRESSED and $k = 0.1n$. All experiments were conducted on a machine with an Intel Core i7-4790 CPU and 16GB RAM. The algorithms are implemented in Python using the networkx and numpy libraries.

As expected from the theoretical analysis, for the MODERATEINTERNAL problem ExtremeExpressed, ExtremeNeighbors and Pagerank far outperform the GreedyInt and BOMP in terms of running time. Given that ExtremeExpressed is matching the performance of GreedyInt and BOMP, this indicates that this is an effective heuristic for very large datasets.

Finally, we study the effect of Sherman-Morrison formula in the computation of the update of the \mathbf{z} vector when GreedyExt considers a candidate node. We consider two implementations of GreedyExt: one that uses the Sherman-Morrison formula (Section 5.3), and one that computes the \mathbf{z} vector using Equation (3.1) iteratively. We construct different samples from the *Elections* dataset, of size $2.8K$, $6K$, $9.9K$, $13.8K$ and $18.2K$. Figure 6.19 shows the comparison of the two implementations for one update of \mathbf{z} . The x -axis is the size of the graph and the y -axis is the running time (in secs). The plot is in log-log scale. Clearly, the Sherman-Morrison implementation is one order magnitude faster, making the algorithm scalable for larger datasets.

Our algorithms use the fundamental matrix \mathbf{Q} , and thus require quadratic amount of memory and time. They are applicable to medium-to-large networks, such as an ego-network, or the network induced by the users of specific hashtags, or a subset

Table 6.4: Running times (secs) of all algorithms for $k = 0.1n$ in the *Elections* dataset.

MODERATEINTERNAL		MODERATEEXPRESSED	
Algorithm	Running time (secs)	Algorithm	Running time (secs)
BOMP	2725		
GreedyInt	2930	GreedyExt	16326
ExtremeExpressed	6	ExtremeExpressed	87
ExtremeNeighbors	106	ExtremeNeighbors	121
Pagerank	7	Pagerank	17

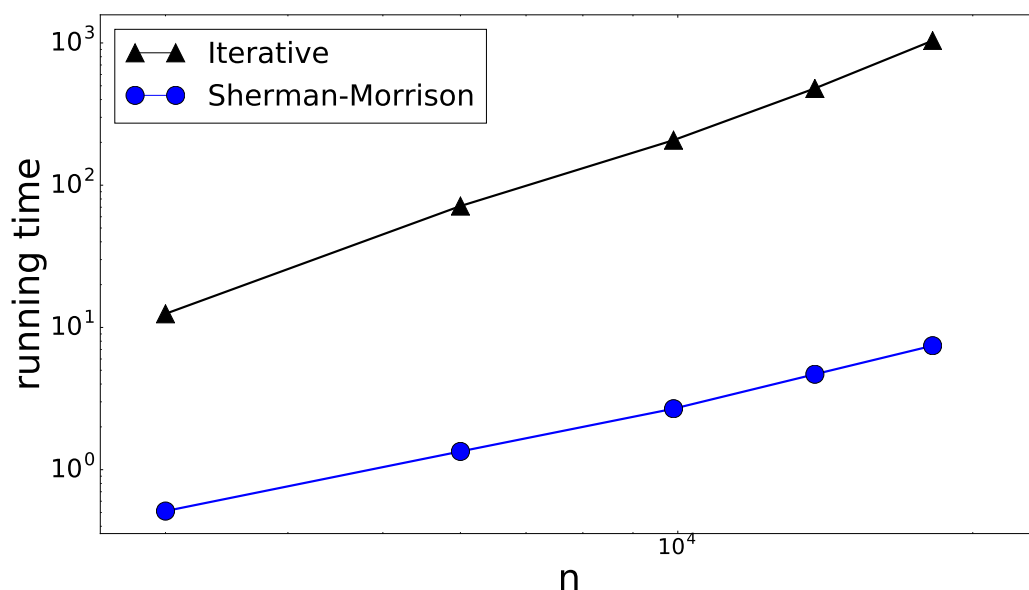


Figure 6.19: Comparison of the running times (in secs) for the Sherman-Morrison and iterative implementation of GreedyExt, for varying size of n .

of Twitter followers, but they cannot be used for massive networks of millions of nodes. In such cases, we can use the iterative computation of the \mathbf{z} vector. This computation is very similar to the computation of PageRank, and lends itself to a distributed implementation. Using existing distributed computation techniques, we can compute the polarization index for very large networks. Our algorithms can combine the efficient heuristics we described to reduce the number of candidate nodes to be considered (e.g., consider only the top nodes in terms of z_i , or PageRank value).

6.7 Case study

We conclude by taking a closer look at the characteristics of individuals that were selected by GreedyExt in the *Elections* dataset. For this, we pick the first 10 nodes selected by GreedyExt and we rank them according to three other measures: the extremity of their expressed opinion, measured by the node's $|z_i|$, their centrality, measured by their pagerank score and their degree. In the last three columns of Table 6.5 we show the rank of these first 10 nodes in the three rankings. In the same table we report the internal opinion of the node (-1 for Trump and $+1$ for Clinton) and their original expressed opinions z_i .

Table 6.5: Characteristics of the first ten nodes selected by GreedyExt in *Elections* dataset.

	Opinion	z_i	$ z_i $ rank	Degree rank	Pagerank rank
1	positive	0.045	10683	16	11
2	neutral	-0.049	10211	34	21
3	positive	0.03	11704	9	8
4	negative	-0.35	32	12473	4861
5	negative	-0.03	12582	52	37
6	positive	0.02	13634	2	1
7	negative	-0.04	10844	114	59
8	negative	-0.06	8366	257	157
9	negative	-0.07	7486	601	170
10	positive	0.04	11269	23	19

We observe that GreedyExt mainly selects central nodes, but also selects a node with very extreme expressed opinion high in the list (4-th pick). Nine out of the top-10 nodes are clearly very central in the network as they are ranked high both by their degree and their PageRank scores. This is in line with the previous observation that GreedyExt initially tries to diffuse as much neutrality as possible and then tries to cover individuals that were not reached. We also note in the top-10 selected nodes we have five Trump followers, four Clinton followers, and a neutral user. Note that this matches relatively closely the proportions of Trump, Clinton and Neutral followers in the full dataset, which agrees with our previous observations on *Karate* and *Books*, and indicates that it is a good strategy to take a balanced approach to moderating opinions.

CHAPTER 7

CONCLUSIONS

In this thesis we considered the problem of polarization in online social networks. Using a popular opinion formation model, we proposed the *polarization index*, a novel measure for quantifying the degree of polarization in the network that takes into account both the network structure and the existing opinions of users. We then considered the problem of identifying a small set of individuals, such that, if we convince them to adopt a moderate opinion, this will minimize the polarization index. We defined two variants of the problem, and showed that both variants are NP-hard. We proposed efficient algorithms by exploiting the mathematical properties of the opinion formation model. Experiments with real data demonstrate the validity of our model, and the effectiveness of our algorithms in reducing polarization. Our experiments also highlight the properties and the differences of the two problems we considered.

In our work we assumed that the opinions are given as input for the computation of the polarization index. An interesting future direction for our work is to use opinion mining techniques to derive the opinions of the users in the social network. Such techniques can be used as the first step in our pipeline. Alternatively, we could integrate ideas from opinion and sentiment mining into the computation of a polarization metric, or in the moderation algorithms.

Furthermore, our approach to moderation is to set the opinions of the users to zero. An alternative approach would be to set the user opinions to values other than zero, so as to minimize polarization. In the case of the internal opinions, there are interesting

connections of this problem with the algebraic properties of the fundamental matrix Q that are worth exploring in future work.

BIBLIOGRAPHY

- [1] L. A. Adamic and N. Glance. The political blogosphere and the 2004 u.s. election: Divided they blog. In *International Workshop on Link Discovery*, LinkKDD, 2005.
- [2] L. Akoglu. Quantifying political polarity based on bipartite opinion networks. In *International Conference on Weblogs and Social Media, ICWSM*, 2014.
- [3] V. Amelkin, A. K. Singh, and P. Bogdanov. A distance measure for the analysis of polar opinion dynamics in social networks. *CoRR*, abs/1510.05058, 2015.
- [4] E. Bakshy, S. Messing, and L. Adamic. Exposure to ideologically diverse news and opinion on facebook. *Science*, 2015.
- [5] A. Bessi, F. Zollo, M. D. Vicario, M. Puliga, A. Scala, G. Caldarelli, B. Uzzi, and W. Quattrociocchi. Users polarization on facebook and youtube. *PLoS ONE*, 11(8), 2016.
- [6] D. Bindel, J. M. Kleinberg, and S. Oren. How bad is forming your own opinion? *Games and Economic Behavior*, 92:248–265, 2015.
- [7] E. Cambria, D. Olsher, and D. Rajagopal. Senticnet 3: A common and common-sense knowledge base for cognition-driven sentiment analysis. In *AAAI Conference on Artificial Intelligence*, 2014.
- [8] E. Cambria, S. Poria, F. Bisio, R. Bajpai, and I. Chaturvedi. *The CLSA Model: A Novel Framework for Concept-Level Sentiment Analysis*, pages 3–22. Springer International Publishing, Cham, 2015.
- [9] T. Chen, R. Xu, Y. He, Y. Xia, and X. Wang. Learning user and product distributed representations using a sequence model for sentiment analysis. *IEEE Comp. Int. Mag.*, 11(3):34–44, 2016.

- [10] M. Conover, J. Ratkiewicz, M. R. Francisco, B. Gonçalves, F. Menczer, and A. Flammini. Political polarization on twitter. In *International Conference on Weblogs and Social Media ICWSM*, 2011.
- [11] P. Dandekar, A. Goel, and D. T. Lee. Biased assimilation, homophily, and the dynamics of polarization. *Proceedings of the National Academy of Sciences*, 110(15):5791–5796, 2013.
- [12] G. Davis, S. Mallat, and Z. Zhang. Adaptive time-frequency decompositions with matching pursuits. *Optical Engineering*, 33, 1994.
- [13] M. Del Vicario, A. Scala, G. Caldarelli, H. E. Stanley, and W. Quattrociocchi. Modeling confirmation bias and polarization. *Scientific Reports*, 7:40391, Jan. 2017.
- [14] U. Feige. Vertex cover is hardest to approximate on regular graphs. *Technical report MCS03-15 of the Weizmann Institute*, 2003.
- [15] N. E. Friedkin and E. Johnsen. Social influence and opinions. *Journal of Mathematical Sociology*, 1990.
- [16] K. Garimella, G. D. F. Morales, A. Gionis, and M. Mathioudakis. Quantifying controversy in social media. In *ACM International Conference on Web Search and Data Mining, WSDM*, pages 33–42, 2016.
- [17] V. R. K. Garimella, G. D. F. Morales, A. Gionis, and M. Mathioudakis. Reducing controversy by connecting opposing views. In *ACM WISDOM International Conference on Web Search and Data Mining*, 2017.
- [18] R. K. Garrett. Echo chambers online?: Politically motivated selective exposure among internet news users¹. *Journal of Computer-Mediated Communication*, 14(2):265–285, 2009.
- [19] A. Gionis, E. Terzi, and P. Tsaparas. Opinion maximization in social networks. In *SIAM International Conference on Data Mining*, pages 387–395, 2013.
- [20] P. H. C. Guerra, W. M. Jr., C. Cardie, and R. Kleinberg. A measure of polarization on social media networks based on community boundaries. In *International Conference on Weblogs and Social Media, ICWSM*, 2013.

- [21] W. W. Hager. Updating the inverse of a matrix. *SIAM Rev.*, 31(2):224–239, June 1989.
- [22] D. J. Isenberg. Group polarization: A critical review and meta-analysis. *Journal of personality and social psychology*, 50(6), 1986.
- [23] D. Kempe, J. Kleinberg, and E. Tardos. Maximizing the spread of influence through a social network. In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 137–146, 2003.
- [24] T. Lappas, M. Crovella, and E. Terzi. Selecting a characteristic set of reviews. In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 832–840, 2012.
- [25] P. Lawrence, B. Sergey, R. Motwani, and T. Winograd. The pagerank citation ranking: Bringing order to the web. Technical report, Stanford University, 1998.
- [26] B. Liu. *Sentiment analysis and opinion mining*, 2012.
- [27] S. Mallat. *A Wavelet Tour of Signal Processing, Third Edition: The Sparse Way*. Academic Press, 3rd edition, 2008.
- [28] S. A. Munson, S. Y. Lee, and P. Resnick. Encouraging reading of diverse political viewpoints with a browser widget. In *International Conference on Weblogs and Social Media, ICWSM*, 2013.
- [29] S. A. Munson and P. Resnick. Presenting diverse political opinions: how and how much. In *International Conference on Human Factors in Computing Systems, CHI*, pages 1457–1466, 2010.
- [30] B. K. Natarajan. Sparse approximate solutions to linear systems. *SIAM J. Comput.*, 24(2):227–234, 1995.
- [31] E. Pariser. *The Filter Bubble: What the Internet Is Hiding from You*. Penguin Group , The, 2011.
- [32] S. Poria, E. Cambria, and A. Gelbukh. Aspect extraction for opinion mining with a deep convolutional neural network. *Know.-Based Syst.*, 108(C):42–49, Sept. 2016.

- [33] C. R. Sunstein. The law of group polarization. *Journal of political philosophy*, 10(2):175–195, 2002.
- [34] M. D. Vicario, A. Scala, G. Caldarelli, H. E. Stanley, and W. Quattrociocchi. Modeling confirmation bias and polarization. *CoRR*, abs/1607.00022, 2016.
- [35] V. Vydiswaran, C. Zhai, D. Roth, and P. Pirolli. Overcoming bias to learn about controversial topics. *Journal of the Association for Information Science and Technology*, 2015.

APPENDIX A

PROOF OF THEOREM 1

Theorem 1. *The MODERATEINTERNAL problem is NP-hard.*

Proof. Our proof uses a reduction from the m -SUBSETSUM problem, where given a set of N positive integer numbers v_1, \dots, v_N , a value m , and a target value b , we ask if there is a set of numbers B of size m , such that $\sum_{v_i \in B} v_i = b$.

Given an instance of the m -SUBSETSUM problem, we construct an instance of MODERATEINTERNAL as follows. The graph is a star with $N + 1$ nodes: we have a central node u_0 , and a spoke node u_i for each integer v_i . For the center of the star (node u_0) we have that $w_{00} = t$, for an appropriately selected value of t (we will discuss this below), and $s_0 = -1$. The weight of the edge (u_0, u_i) from the center to node u_i is $w_{0i} = v_i$, and the weight of node u_i to its internal opinion is also $w_{ii} = v_i$. The opinion of all spoke nodes is $s_i = 1$. We set $k = N - m$, and we ask for a set of nodes T_s , $|T_s| = k$, such that, when setting $s_i = 0$ for $u_i \in T_s$ $\pi(\mathbf{z} | T_s) = \|\mathbf{z}\|^2$ is minimized.

The intuition of the proof is that the expressed opinion of the center node z_0 determines $\pi(\mathbf{z})$. The value of z_0 is determined by the weight t of the internal opinion of u_0 , and the weights of the edges of nodes whose opinion is not set to zero. If we select t appropriately, we can guarantee that $\|\mathbf{z}\|^2$ is minimized when the nodes whose opinion is not set to zero sums to the value b .

Formally, assume that we have selected the set T_s , $|T_s| = k$. Assume that $u_0 \notin T_s$. Also let $R = V \setminus T_s \cup \{u_0\}$ denote the set of spoke nodes whose opinion was *not* set to 0. According to the opinion formation model, the equations for the expressed

opinions of the spoke nodes are as follows. For every node $u_i \in R$, $z_i = \frac{z_0}{2} + \frac{1}{2}$. while for every node $u_i \in T_s$, $z_i = \frac{z_0}{2}$.

We can thus write:

$$\begin{aligned}\pi(\mathbf{z} | T_s) &= \|\mathbf{z}\|^2 = z_0^2 + k\frac{1}{4}z_0^2 + (N-k)\frac{1}{4}(z_0^2 + 2z_0 + 1) \\ &= \frac{N+4}{4}z_0^2 + \frac{N-k}{2}z_0 + \frac{N-k}{4}.\end{aligned}$$

Recall that we want to minimize $\pi(\mathbf{z} | T_s)$. To find the value of z_0 that minimizes $\pi(\mathbf{z} | T_s)$, we take the derivative of the expression above, we set it zero, and solve for z_0 . We get that the value of z_0 that minimizes $\pi(\mathbf{z})$ is:

$$z_0^* = \frac{k-N}{N+4}.$$

It follows that the minimum value of $\pi(\mathbf{z} | T_s)$ is

$$\pi^* = \frac{(N-k)(k+4)}{4(N+4)}.$$

We now set the value of t such that if the set of numbers in R sums to the value of b , then z_0 achieves the z_0^* value. First we compute the value of z_0 as a function of t . In the following we set $W = \sum_{i=1}^N v_i$. We have that:

$$\begin{aligned}z_0 &= \sum_{i=1}^N \frac{v_i z_i}{W+t} - \frac{t}{W+t} = \sum_{u_i \in T_s} \frac{v_i z_0}{2(W+t)} + \sum_{u_i \in R} \frac{v_i(z_0+1)}{2(W+t)} - \frac{t}{W+t} \\ &= \frac{\sum_{i=1}^N v_i}{2(W+t)} z_0 + \frac{\sum_{u_i \in R} v_i}{2(W+t)} - \frac{t}{W+t} = \frac{W}{2(W+t)} z_0 + \frac{\sum_{u_i \in R} v_i - 2t}{2(W+t)}\end{aligned}$$

Solving for z_0 we get :

$$z_0 = \frac{\sum_{u_i \in R} v_i - 2t}{W + 2t}.$$

We want the minimum to be achieved when $\sum_{u_i \in R} v_i = b$. Setting $z_0 = z_0^*$ we get:

$$\frac{b-2t}{W+2t} = \frac{K-N}{N+4}$$

Solving for t we get:

$$t = \frac{(N+4)b + (N-k)W}{2(k+4)}.$$

Now, we want to prove the following. There is a set B of m numbers such that $\sum_{v_i \in B} v_i = b$, if and only if there is a set of nodes T_s of size $k = N - m$ such that when setting their internal opinion to zero, $\pi(\mathbf{z} | T_s) < \pi^* + \epsilon$ for some appropriate value of ϵ .

The forward direction is easy. If there exists this set B , then there is a set T_s such that when setting their opinions to zero, for the set R we have that

$$z_0 = \frac{\sum_{u_i \in R} v_i - 2t}{W + 2t} = \frac{b - 2t}{W + 2t} = \frac{k - N}{N + 4},$$

and therefore $\pi(\mathbf{z} | T_s) = \pi^*$.

For the backwards direction, if no such set of numbers exists, then it is not possible to find a set of nodes T_s such the nodes in R give $z_0 = \frac{K-N}{N+4}$ that minimizes $\pi(\mathbf{z} | T_s)$. Therefore, there must be an ϵ such that $\pi(\mathbf{z} | T_s) \geq \pi^* + \epsilon$.

To set ϵ note that for any $z_0 \neq z_0^*$

$$|z_0 - z_0^*| = \left| \frac{\sum_{u_i \in R} v_i - b}{W + 2t} \right| \geq \frac{1}{W + 2t} = \frac{k + 4}{(N + 4)(W + b)},$$

where the inequality follows from the fact that the values v_1, \dots, v_N, b are integers and their difference is at least one. Now, let \mathbf{z}^* be the vector with z_0^* that achieves the minimum value π^* . For any other \mathbf{z} we have

$$\begin{aligned} \pi(\mathbf{z}) - \pi^* &= \frac{N + 4}{4} (z_0^2 - (z_0^*)^2) + \frac{N - k}{2} (z_0 - z_0^*) \\ &= (z_0 - z_0^*) \left(\frac{N + 4}{4} z_0 + \frac{N + 4}{4} z_0^* - \frac{2(N + 4)k - N}{4} \right) \\ &= (z_0 - z_0^*) \left(\frac{N + 4}{4} z_0 - \frac{N + 4}{4} z_0^* \right) = \frac{N + 4}{4} (z_0 - z_0^*)^2 \\ &\geq \frac{N + 4}{4} \left(\frac{1}{W + 2t} \right)^2 = \frac{(k + 4)^2}{4(N + 4)(W + b)^2}. \end{aligned}$$

So it suffices to set $\epsilon < \frac{(k+4)^2}{4(N+4)(W+b)^2}$.

Finally, in our computations so far we have assumed that our set T_s does not contain node u_0 . This is not a restrictive assumption. Consider a solution T_s , where $u_0 \in T_s$, and $s_0 = 0$. Then, since s_0 is the only negative opinion value in our instance, it follows that $z_0 \geq 0$, and for any node $u_i \in R$ we have that $z_i = \frac{1}{2}z_0 + \frac{1}{2} \geq \frac{1}{2}$. There are $N + 1 - k$ nodes in R . Therefore,

$$\pi(\mathbf{z} | T_s) \geq \frac{N + 1 - k}{2}.$$

Note that $\pi^* = (N - k)(k + 4)/4(N + 4) \leq (N - k)/4$, since $k \leq N$. Therefore, $\pi(\mathbf{z}) \geq 2\pi^* + 1/4$. Selecting $\epsilon < \pi^* + \frac{1}{4}$ guarantees that $\pi(\mathbf{z}|T_s) > \pi^* + \epsilon$. Thus, if there is a set T_s such that $\pi(\mathbf{z}|T_s)$ is minimized, it cannot contain u_0 . \square

AUTHOR'S PUBLICATIONS

Antonis Matakos, Panayiotis Tsaparas:

Temporal mechanisms of polarization in online reviews. ASONAM 2016: 529-532

Antonis Matakos, Evimaria Terzi, Panayiotis Tsaparas:

Measuring and Moderating Opinion Polarization in Social Networks (to appear in ECML PKDD 2017)

SHORT BIOGRAPHY

Antonis Matakos was born in Serres, Greece in 1992. He received his BSc degree from the Department of Computer Science & Engineering of University of Ioannina in 2015. In 2015 he became an MSc student at the same institution under the supervision of Panayiotis Tsaparas. His academic interests are in the area of Algorithm Data Mining, and its applications on Social Networks and Computational Sociology.