

# Hidden Opinions

Shaden Shabayek\*

November 18, 2020

## Abstract

This paper widens the scope of analysis of opinion dynamic models by introducing a novel heuristic: individuals choose to express their opinion or hide it, as a function of their local popularity. Intuitively, individuals who hide their opinion could be interpreted as individuals who have a low popularity such that even if they speak-up (or *tweet*) they will not be heard. Local popularity captures the idea that immediacy causes higher influence. Locally popular individuals express their opinion and can interact with like-minded or ideologically-opposed peers, namely expression entails debates and discussions. In the presence of hidden opinions, I show that the interactions between locally popular individuals and the magnitude of their influence explains whether consensus or polarization prevails. The primary mechanism at play is that the influence structure allows for consensus of opinion locally but communication between ideologically opposed expressers lead to global disagreement. The main contribution of this paper is to provide a unifying theoretical framework to assess different long-run opinion patterns with a focus on the topology of the network. I provide a measure of polarization and I run simulations to show the extent to which the topology of the network affects long-run opinion patterns.

**Keywords:** Naive learning, repulsive influence, opinion polarization.

**JEL Classification numbers:** D83, D91, Z1.

---

\*Phd candidate, Université Paris1 Panthéon-Sorbonne, Paris School of Economics, 48 boulevard Jourdan 75014 Paris, France. Contact information: shaden.shabayek@psemail.eu.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Related work in social psychology</b>	<b>11</b>
2.1	Hidden Profiles . . . . .	11
2.2	Dynamic social Impact theory . . . . .	12
<b>3</b>	<b>The model</b>	<b>13</b>
3.1	Set-up . . . . .	13
3.2	Expression heuristic . . . . .	14
3.3	Micro-foundation . . . . .	16
<b>4</b>	<b>Opinion Dynamics</b>	<b>18</b>
4.1	Long-run opinions of expressers . . . . .	18
4.1.1	Non-monotonic opinion updating . . . . .	20
4.1.2	Moderate long-run opinions . . . . .	21
4.2	The process of interpersonal influence . . . . .	24
4.3	Patterns of long-run opinions: consensus and bi-polarization . . . . .	28
4.3.1	Consensus . . . . .	30
4.3.2	Bi-polarization . . . . .	31
<b>5</b>	<b>Network topology and opinion patterns</b>	<b>32</b>
5.1	Example . . . . .	34
5.2	Agregate statistics . . . . .	35
<b>6</b>	<b>Conclusion and the way forward</b>	<b>42</b>
<b>7</b>	<b>Appendix</b>	<b>45</b>
7.1	Proof of proposition 1 . . . . .	45
7.2	Proof of proposition 2 . . . . .	46
7.3	Proof of lemma 1 . . . . .	49
7.4	Proof of theorem 1 . . . . .	51
7.5	Proof of proposition 3 . . . . .	52
7.6	Proof of lemma 2 . . . . .	54
7.7	Network statistics . . . . .	55
7.8	Toy Networks . . . . .	56

# 1 Introduction

By taking a simple walk on social media, one can effortlessly notice that not all voices are equal. Social media have become a new networked public sphere where influencers and politicians express themselves and interact among each other and supposedly with the rest of the population. In particular, opinion polarization is well documented empirically, yet its drivers remain unclear. So have social media connected or disconnected our societies? My paper investigates this question by adding to the study of opinion dynamics the following *expression heuristic*: individuals choose to express their opinion or hide it based on how popular they are within their social network.

Individuals who hide their opinion could be interpreted as individuals who have a low popularity such that even if they speak-up they will not be heard or considered. One can think of an individual who has very few followers on Twitter. Alternatively, a second interpretation could be that hiding one's opinion is less costly than expressing it. The cost of expression could be the time spent arguing with more eloquent and persuasive peers or the cost of social isolation when one's opinion drifts from the average group viewpoint or even the psychological cost of online bullying. To account for popularity I use a simple local centrality measure defined as the number of *direct* friends of a given individual divided by their average number of friends. Intuitively, the influence of an individual over their friends is higher when those same individuals have a small number of influence sources. I use a **local measure** rather than a global one because immediacy of interaction is associated with higher social influence. Both of these **two novel ingredients** are grounded in social psychology literatures, which I review in section 2. **not clear why 2**

The *expression heuristic* departs from early models of opinion formation, e.g. French (1956) [20], DeGroot (1971) [10], where individuals are typically consensus-seeking, in the sense that all individuals express their opinion regardless of their network position and adapt their opinion towards the opinions within their social circle. In DeGroot Models, regardless of the specific topology of the **network, as long as it is (strongly) connected,** individuals end up reaching consensus. Furthermore, **aperiodic, footnote** there is an active line of research that models disagreement in social contexts (See Flache et al. (2017) [11] for a survey). Disagreement is modeled either by introducing repulsive or negative

influence when an individual interacts with dissimilar others or by modeling individuals who only take into account opinions of others when they are close enough. Nevertheless, the role played by the specific topology of the network is not yet at the center stage of this literature. The results are often based on random pair interaction<sup>1</sup> or a specific topology is given by assumption.<sup>2</sup> In this paper, I build on both strands by introducing the expression heuristic to account for the impact of the structure of the social network on the dynamics of opinions, when influence can be positive or negative. In particular, I show that in the presence of hidden opinions, the study of the interactions between locally *popular* individuals who interact with like-minded or ideologically-opposed peers can explain whether consensus or polarization prevails. Since influence is stronger locally, clusters can form. But some members within a given cluster who are popular enough and interact with ideologically-opposed peers can fall into disagreement, which causes opinions to polarize across clusters.

Formally, a group of individuals are connected through an exogenous network of interpersonal relationships. The neighborhood of each individual is their direct set of friends. Each individual is allocated an initial opinion which lies between -1 and 1. This initial opinion represents the stance or attitude concerning a given issue.<sup>3</sup> At each time period, individuals observe the expressed opinions from the previous period and update their opinion at the current period.<sup>4</sup> An individual chooses to express (hereafter an *expresser*) when their local popularity is above a given threshold of expression. Otherwise, they choose to hide their opinion (hereafter a *consensual* individual). Choosing to express or hide as a function of the network position, determines how individuals update their opinion while interacting with peers at each time period. Individuals who choose to express undergo an attractive effect (positive influence or assimilation) when interacting at a given period

---

<sup>1</sup>See for example Deffuant et al. (2000) [9] or Grow (2017) [15] where individuals are randomly paired and engage in dyadic interaction.

<sup>2</sup>For example, Krueger et al. (2017) [21] in a framework of a voter model assume that the social network has a double-clique topology consisting of two complete graphs connected by some cross links with each other. Axelrod (1997) [3] in his seminal paper assumes that individuals are distributed over a 10 by 10 grid.

<sup>3</sup>I do not assume that there is any correlation in initial opinions. In particular two neighbors can have initially very different view points about the issue to be discussed. One way to explain this, could be that there is some common ground or an attribute which led individuals to form a link, e.g. colleagues at work, same political orientation, etc. But the opinion about the issue  $\theta$  to be discussed is not related to that attribute. Network formation is out of the scope of this paper and I look at a given network as a snapshot at a given date.

<sup>4</sup>This is like scrolling down your feed on Twitter or your timeline on Facebook.

model or  
result?  
we have the  
impression that its a result  
  
how do expressers  
update? not clear  
  
clear model  
clarify it here!!  
  
untill here

with like-minded peers who also choose to express. But when they interact with expressing peers that are ideologically-opposed, they undergo a repulsive effect (negative influence or distancing). In particular, expressers only interact with peers who also choose to express. Intuitively, expression allows for a debate or a discussion to take place. Individuals who choose to hide, update their opinions à la DeGroot [10]. That is, at each time period, their updated opinion is the average of opinions expressed within their social circle at the previous period. In other words, expressers only pay attention to expressing neighbors, while consensual individuals pay attention to all their neighbors.

This paper makes several contributions. First, I provide a novel unifying framework which explains how consensus or polarization of opinions can prevail in the long-run. Second, I characterize the process of interpersonal influence when introducing the expression heuristic. This characterization is challenging because, unlike linear models with assimilative influence, my model is non-linear. It is non-linear in the sense that the influence structure or the *hearing* matrix<sup>5</sup> can vary across periods in the presence of expressers who repulse or attract each other. Many expressers can be interlinked directly or indirectly. Hence two individuals can be initially like-minded and influence each other positively. But in subsequent periods they can become ideologically-opposed and influence each other negatively, if one of them is repulsed by another expressing neighbor. Intuitively, this situation occurs when two individuals are somehow like-minded initially. But one of them starts adopting an extreme point of view in an unreasonable fashion under the influence of a third expressing peer, so her friend starts bringing to the table arguments which support the opposing view. In other words, the weights in the hearing matrix can depend on the opinion itself, as in Hegselmann and Krause (2002) [16] and opinions of expressers can get updated monotonously or non-monotonously.<sup>6</sup> I overcome this difficulty by first characterizing the evolution of opinions of expressers and I show that repulsion to the upper or lower bound of the opinion interval occurs at a much faster rate than attraction. Reaching an agreement can require a long lasting debate, while conflict can escalate quite rapidly.

<sup>5</sup>I give later in the paper a formal definition of the hearing matrix. This matrix provides information on who listens to whom or who pays attention to whom and the magnitude of this attention.

<sup>6</sup>The non-monotonous case occurs when two expressers are initially like-minded but one of them becomes extreme too fast, so her opinion difference with her initially like-minded neighbors starts growing.

Proposition 1 shows that when a given expresser has no neighbors who also express then, this individual remains stubborn throughout all periods of interaction and their long-run opinion is simply their initial opinion. It also characterizes the evolution of opinions of expressers when they only have like-minded neighbors. Each expresser of the group undergoes only the attractive effect. In particular, the group can have members who are ideologically-opposed that are not neighbors. But since they are not directly linked, they do not push each other to become extreme. Hence the whole group can reach consensus by gathering different viewpoints. This result sheds light on the design of media tools or the formation of discussion groups for initiatives related to participatory democracy.<sup>7</sup>

Proposition 2 characterizes the long-run opinions of a given group of expressers when a pair of initially ideologically-opposed neighbors belong to the group. I show that generically, the long-run opinions of the group of connected expressers reach the upper or lower bound of the opinion interval  $[-1, 1]$ . That is, all the members of the connected group of expressers become extreme and adopt either opinion 1 or  $-1$ . One special case where moderate opinions<sup>8</sup> of expressers can survive in the long-run occurs when an expresser has like-minded neighbors. But they are indirectly connected to at least two ideologically-opposed expressing individuals. This case depicts a political left-right spectrum where parties on the far right and far left are ideologically-opposed and interact often together but between both parties, many moderate parties survive.

I study the process of interpersonal influence with both types of individuals, starting from the time period by which repulsion opportunities in the course of a discussion are exhausted. I show that opinions of the group of individuals converge in the long-run. In particular, I show that the opinions of consensual individuals vanish in the long-run and remain hidden forever. Their long-run opinions simply become convex combinations of opinions of expressers to whom they are connected.

Furthermore I study of the emergence of two specific long-run patterns: consensus and bipolarization of opinions. I focus on these two patterns because they match two situations of interest for our societies. First, consensus of opinion is consistent with a society which fully agrees on a

---

<sup>7</sup>For example, think of the Citizens' convention for climate (*convention citoyenne pour le climat*) in France, where 150 citizens were randomly selected to work in smaller groups on different propositions.

<sup>8</sup>Strictly smaller than 1 and strictly higher than  $-1$ .

position concerning one issue, or in other words a situation with full cooperation. There are different types of consensus. For example, consensus can obtain when there is a unique opinion leader and the rest of the society is formed of consensual individuals.<sup>9</sup> Second, bi-polarization corresponds to situations where society is divided into two large groups; such that there is strong agreement within each group and disagreement across groups. I choose to focus on the case of bi-polarization because it takes a very large group in society to over-turn a policy or to elect a president.

In DeGroot like models, consensus prevails provided that the *hearing* matrix is (strongly) connected. It means that if every individual is directly or indirectly connected by a path to any other individual in the network, then long-run opinions form consensus. However, in the framework of my model, consensus prevails for very specific initial opinion distributions and network structures. I show in proposition 3 that consensus prevails if and only if there is a unique expresser, or there **not clear** exists multiple sets of connected expressers<sup>10</sup> such that first within each set, each member has like-minded neighbors. Second, the average initial opinions across two groups of connected expressers should be relatively close. As explained above, the long-run opinions of consensual individuals depend on the long-run opinions of expressers to whom they are connected. Since those expressers all have similar viewpoints, consensual individuals adopt similar viewpoints. The *type* of consensus obtained will depend first on how many expressers there are in the society. Their number is directly related to the structure of the network. Second, it depends on the expressed opinions within the neighborhood of each expresser. This relies on the initial opinion distribution and the local structure of the network.

Unlike consensus, long-run opinion bi-polarization is challenging to characterize. As mentioned above, this difficulty lies in the fact that my model is non-linear due to indirect influence between expressers that can make a pair that is initially like-minded to become ideologically-opposed. Hence, mapping back a long-run bi-polarized opinion vector to an exact set of network structures and initial opinion distributions is tedious<sup>11</sup>. In addition, consensual individuals are influenced by expressers

---

<sup>9</sup>Of course this situation may not be optimal, since many members of the society can have *insightful* viewpoints or pieces of information that get lost or are not taken into account or not shared with others.

<sup>10</sup>It means that there exists a path connecting any two members of the set and no path that connects an expresser from this set to an individual who does not belong to the set.

<sup>11</sup>Hegselmann and Krause (2002) [16] have a non-linear model and they discuss this point extensively: *though elementary, the model is nonlinear in that the structure of the model changes with the states of the model given by*

to whom they are linked directly and indirectly. Put simply, they receive influence of different magnitudes from many expressers depending on how far or close they are from those expressers in the network. That is, one needs to account for all the possible paths of different length connecting consensual individuals to expressers who may be ideologically-opposed. Straightforward sufficient conditions for long-run bi-polarized opinions can be formalized for network structures and initial opinion distributions, where the average path length between a consensual player and an expresser is very small. For example, consider a multi-star network, that is there are multiple *stars* where each has many followers (leaves) and the stars (center) are connected directly or indirectly among themselves. Then for this network structure, it is sufficient to account for the degree of each expresser and the opinion distribution within each expresser’s neighborhood. These networks are very similar to the common hub and spoke structure on Twitter<sup>12</sup>.

Furthermore, I provide compelling necessary conditions for long-run bi-polarization of opinions. I show that if a long-run opinion vector is bi-polarized then necessarily, first consensual individuals who remain moderate in the long-run do not exist, second both extreme influence groups influence an equal share of the society. The first condition means that there does not exist consensual individuals who receive an equal amount of influence from two ideologically-opposed extreme opinion groups. Those individuals occupy a very particular location in the network, because they are not locally popular enough to express and they are at an equal distance from two ideologically-opposed groups of expressers. Those individuals could be interpreted as neutral TV hosts, or non-biased journalists or intermediaries in general.<sup>13</sup>

**Simulations.** I explore the model through simulations.<sup>14</sup> The objective is to relate the topology of the network to the long-run opinion patterns. Polarization of long-run opinions is measured

---

*the opinions of the agents (see Section 2). Not only that helpful mathematical tools like Markov chains are no longer applicable, it turns out, moreover, that rigorous analytical results are difficult to obtain. For that reason we carry out the analysis of the above nonlinear model to a large extent by simulations on the computer. Indeed, though we are proud to present analytical results, too, we want to emphasize the importance of careful computer simulations for social dynamics in general, and opinion dynamics in particular, whenever nonlinearities are involved.*

<sup>12</sup>See for example the report [Mapping Twitter Topic Networks: From Polarized Crowds to community clusters](#) by the Pew Research Center (2014)

<sup>13</sup>Such individuals typically do not exist in modern communication networks such as Twitter where politicians can *mention* each other and interact directly. While in past decades, debates between politicians, open to the public were usually moderated by TV hosts.

<sup>14</sup>The Matlab code for the simulations can be found on the GitHub repository [shadenshabayek/Hidden-Opinions](#).



simply by looking at the variance of final opinions. I generate the initial opinions uniformly at random in the interval  $[-1, 1]$ . I study the evolution of opinions of a large set of individuals in scale-free networks. The degree distribution within scale-free networks follows a power law. Due to this inequality in the degree distribution, expressers and consensual individuals co-exist. The objective of the simulations is to make obvious the impact of the interaction structure - location in the network and initial opinion - of expressers on the long-run opinion pattern of the whole group of individuals. Scale-free networks are ment to model Twitter-like networks where there are *stars* or *influencers* with many followers and where stars have the opportunity to interact together. First, I provide two examples where I show how the density of connections among expressers affects the average level of polarization, for two specific network topologies of scale-free networks. Second, I provide aggregate statistics on different network topologies by varying the number of initial hubs and the number of connections a newly added node has. Doing so allows me to relate different scale-free network topologies to the average level of polarization. Those statistics are divided in two groups: (i) (inter-related) measures of degree inequality (ii) and measures of connectivity. In particular, I show that average polarization level can be relatively low even if expressers are densely connected among each other. This happens precisely because for those same network topologies, consensual individuals are connected to many influence sources.

**Related literature.** This paper is related to models that study the dynamics of opinions in communication networks. Those models are challenging because their ultimate goal is to understand human behavior within a complex system formed by society. Naturally, the study of opinion dynamics is a multi-disciplinary topic. Different fields such as economics (learning in networks, for literature surveys see Golub and Sadler (2017) [14] and Acemoglu and Ozdaglar (2011) [2]), sociology (the community cleavage problem, see for example the survey by Flache et al. (2017) [11] or the article by Friedkin (2015) [13]), statistical physics and computer science (community detection, Malliaros and Vazirgiannis (2013) [24] survey the literature) have tackled this problem from different angles. The literature is too large to survey here but I have selected a few papers that are closest to my work.

My paper is broadly related to the literature on naive learning often referred to as the Degroot Model and considered as a credible alternative to the Bayesian branch, given the complexity of networks (Banerjee et al. (2018) [4]). In Degroot Models, agents update their opinions at each step by taking the average of the opinions of their neighbors. My paper extends and departs from DeGroot models or models of assimilative influence along two dimensions. First, not all agents share their piece of information. Second, agents who do share their piece of information can be subject to repulsive influence. Precisely, those two features together contribute not only to explaining disagreement (repulsive effect) but also the rise of bi-polarization (group size).

Furthermore, this paper is closely related to a strand of the literature on Naive Learning which introduces *stubborn* agents or agents that remain attached to their initial opinion to a certain extent, in order to model disagreement (See Friedkin and Johnsen (1990) [12] and Friedkin(2015) [13], Acemoglu et al (2013) [1]) . Two papers that are closest to mine are Yildiz et al. (2013) [34] and Sadler (2019) [27]. Both papers introduce *stubborn* agents in a voter model set-up where opinions are discrete and can take only two values either 0 or 1. My paper extends these two papers by considering stubborn agents as a special case that can occur when an expresser does not update their opinion because they have no expressing neighbors. Furthermore, I allow for opinions to take continuous values between  $-1$  and  $1$ . Hence, my model can also explain the survival of moderate opinions in the long-run.

why  $-1$  and  $1$

My paper also fits in the family of bounded confidence models (See Hegselmann and Krauss (2002) [16], Jager and Amblard (2005) [17]). The key ingredient of those models is to consider the difference between the opinions of individuals when opinion updating is taking place. In particular, Hegselmann and Krauss (2002) [16] consider a model where agents update their opinions by taking an average over the opinions of neighbors whose opinion difference falls within a confidence interval.

Finally, explaining polarization has been tackled by a handful of recent papers, yet there is no consensus in the literature about its main drivers. Bolletta and Pin (2019) [7] introduce a network formation model and argue that under certain conditions when agents optimally choose their links, the network can become disconnected and consensus of opinions cannot be reached. Banisch and Olbrich (2019) [5] explain the emergence of polarization by introducing reinforcement learning,

where agents optimally adopt one viewpoint when they get positive feedback from peers. But their focus is not on the network structure itself and how it could be one of the drivers. In particular they fix a random geometric network to account for the structure of interactions.

**Outline.** The remainder of the paper is organized in three **main** blocks. First, section 2 reviews relevant literatures in social psychology that lay the ground for the main behavioral assumptions of my model. Second, I present the model in section 3 and I characterize the overall process of interpersonal influence in section 4. Third, through simulations I relate the topology of networks to long-run opinion patterns in section 5. Section 6 concludes and mostly discusses the way forward.

## 2 Related work in social psychology

The two main novel ingredients of the model are the expression heuristic and the local popularity measure. The former ingredient is grounded in a well known literature in social psychology called Hidden Profiles<sup>15</sup> and the latter is borrowed from the Dynamic Impact Social Theory. In what follows I briefly review both literatures.

### 2.1 Hidden Profiles

Stasser and Titus (1985) [31] document how individuals in social contexts, do not always share their own opinion or the information they hold. The starting point of their research is to challenge the common belief that a group of individuals should be able to take a better decision than each individual on their own by pooling the members' knowledge and expertise. Namely, group discussion or communication is believed to have a corrective function because members can each have incomplete information but together they can gather the different pieces of the puzzle. The authors ran an experiment in which they simulate a political set-up where a group has to elect one of three candidates: *Best*, *Okay* and *Ohum*. In a first protocol, they distributed a different subset of desirable traits of *Best* and a different subset of *Okay*'s undesirable traits over the members of the group, such that from each one's individual perspective *Okay* appeared more positive than

---

<sup>15</sup>See the survey *Hidden Profiles: a brief history* by Stasser and Titus (2003) [32]

*Best*. Before discussion *Best* received 25% of endorsement. Since the whole group had complete (but dispersed) information about *Best* they could exchange it and come to the conclusion that *Best* was actually the best candidate. Yet after group discussion, surprisingly the percentage of endorsement for *Best* remained at 24%. This finding suggests that unique information held by some members of the group about candidates were not being shared. In a later study, Stasser, Taylor and Hanna (1989) [30] showed that unique pieces of information are less likely to be mentioned during group discussion. One explanation is that social status, expertise or popularity can be a driver for expression of opinion. In fact, Larson et al. (1996) [19] suggest that repeating a unique piece of information, leading to the formation of group opinion during a discussion, is more likely by higher status members (experts, leaders, etc.) rather than lower status members. They ran an experiment with residents, interns and 3rd-year medical students and they show that residents were more likely to repeat (unique) information when compared to interns and students.

Using the findings of the above literature, I introduce an expression heuristic to a dynamic opinion formation model. An individual chooses to express their opinion or hide it based on a popularity measure that is meant to capture different hierarchical and expertise levels.

## 2.2 Dynamic social Impact theory

Latané’s Dynamic Social Impact Theory(1981) [22] suggests that social influence has three determinants: strength, immediacy and the number of influence sources. Strength refers to social status, level of expertise or persuasiveness, while immediacy refers to closeness in space, time or the possibility of direct contact. The theory bridges the influence processes at an individual level using these three determinants, to outcomes at the level of a social system. The main statement of the theory is that *total impact of a group of people on an individual is a multiplicative function of their strength, immediacy, and number*. Latané, Nowak and Liu (1994) [23] use this theory to study through simulations the dynamics of attitude change in groups and societies. Rather than studying attitude distribution as the usual percentage frequencies of different attitude choices, they study the distribution of attitudes in *space*. Using *immediacy*, they are able to explain phenomena such as attitude clustering because individuals are more influenced by nearby individuals.

To incorporate *immediacy* in my model, I use a local popularity measure as opposed to global measures (e.g. eigenvector centrality). Individuals choose to express when their local popularity is above a given threshold. To incorporate *strength*, the local popularity of a given individual is inversely related to the average level of neighborhood connectivity. Meaning that an individual will have a higher impact over their neighbors if those same neighbors are not exposed to many other influence sources.

### 3 The model

#### 3.1 Set-up

A group of individuals  $N = \{1, \dots, n\}$  is embedded in a connected and symmetric network  $G$  of interpersonal relationships, with typical entries  $g_{ij} = g_{ji} \in \{0, 1\}$ . Each node represents an individual. The set of friends, colleagues or acquaintances of individual  $i \in N$  is denoted by  $N_i = \{j \in N : g_{ij} = 1\}$  and  $d_i = |N_i|$  is the cardinality of  $N_i$ . For all  $i \in N$ , I assume that  $g_{ii} = 1$ . A chain of friends of friends of length  $l$  between two individuals  $i \neq j \in N$ , hereafter called *path*, is defined as follows: there exists a sequence of individuals  $i = k_0, k_1, \dots, k_l = j \in N$  such that  $g_{ik_1} \times g_{k_1k_2} \times \dots \times g_{k_lj} > 0$ . A set of individuals  $\mathcal{C} \subset N$  is called an *essential class*<sup>16</sup> if there does not exist a path starting at an individual  $i \in \mathcal{C}$  and ending at an individual  $j \in N \setminus \mathcal{C}$ . Finally, the social status or relative expertise is represented by a local centrality measure which I call *local popularity*.

**Definition 1 (Local popularity)** Let  $i \in N$  be an individual and  $N_i$  the set of their direct friends. The local popularity of  $i \in N$  is:

$$\delta_i = \frac{d_i - 1}{\frac{1}{d_i - 1} \sum_{j \neq i \in N_i} (d_j - 1)} \quad (1)$$

---

<sup>16</sup>For a formal definition see the subsection *Classification of indices* p.12 in Seneta (1981) [28].

This centrality measure is inversely related to the average degree of friends. In other words, the popularity or influence of an individual is relatively low when their direct friends have many influence sources. It is a measure of relative neighborhood connectivity, where the numerator is the **connectivity** of a given individual and the denominator is the neighborhood connectivity. The network statistic  $\delta_i$  of individual  $i \in N$  can also be interpreted as their strength of persuasiveness, which could be high if they are linked with many friends who have relatively small social circles and hence exposed to very little social influence. Finally, this local popularity measure does not require the knowledge of the complete structure of the network of interpersonal relationship.<sup>17</sup> Furthermore I assume that each individual knows their degree and the degree of their direct friends.

**Assumption 1** *Each individual  $i \in N$  knows  $d_i$  and  $d_j$  for all  $j \in N_i$ .*

## 3.2 Expression heuristic

Each individual is endowed with an exogenous initial opinion  $\alpha_{i,0} \in [-1, 1]$  which represents their attitude, belief or stance about an issue<sup>18</sup>  $\theta \in [-1, 1]$ . Individuals exchange opinions about the issue over periods  $t \geq 0$ . First, each individual  $i \in N$  observes opinions in their neighborhood  $N$  at the previous period, denoted by  $\alpha_{j,t-1} \in [-1, 1]$  for all  $j \in N_i$ . Second, each individual chooses whether **to hide their opinion by repeating the average of** observed opinions or to express an updated opinion. To make this choice they rely on their local popularity parameter (1) and use the following heuristic:

$$\begin{cases} \delta_i < \delta^* & \text{play hide} \\ \delta_i \geq \delta^* & \text{play express} \end{cases}$$

---

<sup>17</sup>Informally, an individual can have an idea of their local impact on their direct friends. Since it's fair to assume that they know their number of friends and can have an idea of the number of friends of their friends, for example the number of friends on Facebook or the followers on Twitter of friends.

<sup>18</sup>The *issue* is not modeled as a state of nature intentionally. Doing so, I cover a whole range of issues which are not necessarily factual knowledge (e.g. the quality of a product or weight of an ox could be considered as factual knowledge which could be verified by weighing the ox). The readers could have in mind the following interpretation: there is a build-up of individual opinions about a policy which will impact the whole society and for which we have to wait a long period of time before we know whether it is successful or not. Examples: opinions about vaccines or about risk reducing measures for the coronavirus pandemic.

where  $\delta^*$  denotes the expression threshold. In other words, when local popularity is higher than the *expression threshold*  $\delta^*$ , players choose to express their opinion. Otherwise, they opt for repeating opinions they have heard within their social circle. In the remainder of the paper, the set of individuals who choose to express will be labelled  $E = \{i \in N, \text{ s.t. } \delta_i \geq \delta^*\}$ . Similarly, the set of individuals who choose to hide their opinions will be labelled  $C = \{i \in N, \text{ s.t. } \delta_i < \delta^*\}$ .

**Hide & Express.** When an individual chooses to hide, they update à la DeGroot; that is they update their opinion by taking the average of opinions of neighbors expressed at the previous period:

$$\alpha_{i,t} = \bar{\alpha}_{i,t-1} = \frac{1}{d_i} \sum_{j \in N} g_{ji} \alpha_{j,t-1} \quad (2)$$

When an individual chooses to express, their opinion update will depend on their neighbors, who also choose to express. To account for the difference in opinions of two neighbors who interact, I compare their opinion difference to a threshold  $\tau \in (0, 1)$ .

**Definition 2** *Two individuals  $i \neq j \in N$  who choose to express are : (i) like-minded at period  $t \geq 1$  if  $|\alpha_{i,t} - \alpha_{j,t}| < \tau$ , (ii) ideologically-opposed at period  $t \geq 1$  if  $|\alpha_{i,t} - \alpha_{j,t}| \geq \tau$ . Moreover,  $\underline{N}_{i,t} = \{j \in N_i \cap E, |\alpha_{i,t} - \alpha_{j,t}| < \tau\}$  and  $\bar{N}_{i,t} = \{j \in N_i \cap E, |\alpha_{i,t} - \alpha_{j,t}| \geq \tau\}$ .*

Expression allows for a debate or discussion and opinion update of expressers follows the law of motion described below. It incorporates an attractive and a repulsive effect among direct neighbors who express their opinions. For  $\forall i \in E$  and  $\mu \in (0, 1/2)$ :

$$\alpha_{i,t} = \alpha_{i,t-1} + \mu \Delta_{\alpha_{i,t-1}} \text{ s.t. } \alpha_{i,t} \in [-1, 1] \quad (3)$$

Where

$$\Delta_{\alpha_{i,t-1}} = \sum_{j \in \underline{N}_{i,t-1}} (\alpha_{j,t-1} - \alpha_{i,t-1}) + \sum_{j \in \bar{N}_{i,t-1}} (\alpha_{i,t-1} - \alpha_{j,t-1})$$

**Example 1** *Consider a network  $G$  with only two connected individuals 1 and 2 who choose to express. Suppose that  $\tau = 0.5$  and initial opinions are  $\alpha_{1,0} = -0.7$  and  $\alpha_{2,0} = 0.7$ . In period  $t = 1$ ,*

$\alpha_{1,1} = \alpha_{1,0} + \mu(\alpha_{1,0} - \alpha_{2,0}) = -0.7(1 + \mu) - \mu 0.7 < \alpha_{1,0} = -0.7$  and  $\alpha_{2,1} = \alpha_{2,0} + \mu(\alpha_{2,0} - \alpha_{1,0}) = 0.7(1 + \mu) + \mu 0.7 > \alpha_{2,0} = 0.7$ . The updated opinion of individual 1 becomes more negative or pushed-down towards  $-1$ , while the updated opinion of individual 2 more positive or pushed-up towards  $1$ . Individuals 1 and 2 repulse each other.

Finally, notice that a positive weight of  $\mu$  is assigned to the opinion of like-minded expressing neighbors. While a negative weight of  $-\mu$  is assigned to the opinion of ideologically-opposed expressing neighbors. Moreover, the weights  $\mu$  or  $-\mu$  assigned to the opinion of expressing neighbors can change across periods. This happens because expressing neighbors receive influence from their own expressing neighbors (if any). Hence, an initially like-minded expressing neighbor of individual  $i$  can become in subsequent periods ideologically-opposed, if their opinion difference with  $i$  becomes larger than  $\tau$ . Equivalently, the number of elements in the sets  $\overline{N}_{i,t}$  and  $\underline{N}_{i,t}$  of like-minded or ideologically-opposed expressing neighbors of individual  $i \in N$  at a given period  $t \geq 1$  can vary across two periods. For instance, at period  $t - 1 \geq 0$  individual  $j$  can belong to  $\underline{N}_{i,t-1}$  but at period  $t$  they can move to the set  $\overline{N}_{i,t}$  as they can be influenced when updating her opinion by other expressing neighbors.

### 3.3 Micro-foundation

The choice to hide and express as a function of local popularity can be micro-founded by a simple game where individuals face a social cost of expression. Intuitively this cost is a psychological cost of disagreement with peers or, the cost of a long lasting debate because of strong opposing views or, the perceived cost of social isolation when one's viewpoint is faraway from peers. The payoff of individual  $i \in N$  given action  $a_{i,t} \in \{\text{express}, \text{hide}\} = \{\alpha_{i,t}, \overline{\alpha}_{i,t-1}\}$  is the following:

$$\pi_i(a_{i,t}, a_{-i,t-1}) = -f(\delta_i)(\alpha_{i,t} - a_{i,t})^2 - (1 - f(\delta_i))(\overline{\alpha}_{i,t-1} - a_{i,t})^2 \quad (4)$$

The function  $f$  has support  $(0, 1)$  and  $f' > 0$ . The first term of the payoff function (4) is the cost borne by each player when they doesn't express their opinion, in other words it is the cost of being consensual. The second term is the cost of expression, which is a cost borne when a player expresses



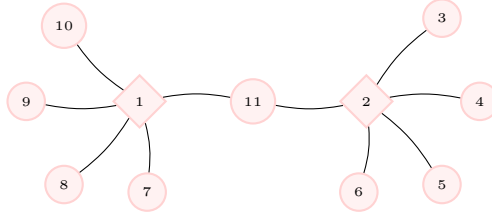


Figure 1: Network  $G_1$ , diamond shaped nodes correspond to individuals who choose to *express*.

their own opinion rather than her neighborhood average opinion. Both terms are weighted by an increasing function of the local popularity parameter  $\delta_i$ . The higher  $\delta_i$  is, the higher the influence of player  $i$  within their own neighborhood and the lower the cost of expression or social disagreement. At each time period, individuals observe the actions (opinions) of neighbors at the previous period and myopically best-respond at the current period. Given the payoff function (4), any  $i \in N$  such that  $\delta_i \geq \delta^* = f^{-1}(\frac{1}{2})$  best-responds with *express*, while any  $i \in N$  such that  $\delta_i < \delta^* = f^{-1}(\frac{1}{2})$  best-responds with *hide*. Hence the local centrality parameter of each individual is a sufficient statistic for the choice of the optimal binary action at each stage-game and individuals can be classified into two groups<sup>19</sup>:  $E = \{i \in N, \delta_i \geq \delta^*\}$  and  $C = \{i \in N, \delta_i < \delta^*\}$ .

**Assumption 2** For the remainder of the paper, the expression threshold is normalized to  $\delta^* = 1$ .

Given the local popularity measure  $\delta_i$  of an individual  $i$ , the threshold  $\delta^* = 1$  allows all individuals to play *express* when they are all equally popular. Namely, in any  $d$ -regular network where each individual has degree  $d$ , any player  $i \in N$  expresses because  $\delta_i = \frac{d}{\frac{1}{d}d \times d} = \delta^*$ . Intuitively,  $d$ -regular networks have perfect assortativity meaning that there is no difference in the level of popularity, expertise or leadership.<sup>20</sup>

**Example 2** Consider the network  $G_1$  in figure 1. Individuals 1 and 2 both have local popularity  $\delta_1 = \delta_2 = 5 / ((1/5)(1 + 1 + 1 + 1 + 2)) = 25/6 > \delta^*$ , therefore they play *express* and  $E = \{1, 2\}$ .

<sup>19</sup>Notice that here the choice to play *express* or *hide* is independent of the opinion of a player. A relevant extension would be to make the choice to *express* or *hide* depend on the opinion of a player with relation to the opinions of neighbors. For example, even if a player has a high enough local popularity enabling them to *express*, they may refrain from *expressing* if their opinion is too extreme relative to their friends.

<sup>20</sup>For example, think of a homogenous group of Phd students. Within the group all members have similar levels of expertise and they can all engage in a discussion. But adding to the group a faculty member will result in unequal levels of expertise, hence different group discussion dynamics.

While individuals  $j \in \{3, \dots, 9, 10\}$  have a local popularity  $\delta_j = 1/5 < \delta^*$  and individual 11 has  $\delta_{11} = 4/10 < \delta^*$ , hence they play hide and  $C = \{3, \dots, 10, 11\}$ .

## 4 Opinion Dynamics

Given the above model, I am interested in studying the long-run opinions, for a network structure  $G$  and a vector of initial opinions  $\alpha_0$ . Expressers only pay attention to other individuals who also express, while consensual individuals pay attention to all their neighbors.<sup>21</sup> In other words, influence can flow from expressers to consensual individuals but not vice versa. I start by considering the evolution of opinions of expressers. Then I characterize in section 4.2 the overall dynamics of opinions and show that long-run opinions always converge.

### 4.1 Long-run opinions of expressers

Individuals who choose to express update their opinions at each time period according to the law of motion (3). Consequently, their opinion update will directly depend on the opinions of neighbors who choose to *express* (if any). The opinions of the latter will depend on the opinions of their own neighbors (if any) who choose to *express* and so on. Recall that  $E$  is the set of individuals who choose to express because their local popularity is higher than the expression threshold  $\delta^*$ . In order to account for the indirect effect of the opinions of expressers on other individuals who also express, I give a formal definition of a connected set of expressers. Then I characterize long-run opinions of expressers within a given connected set of expressers.

**Definition 3 (Connected set of expressers)** *Let  $G$  be a given network structure and let  $\mathcal{E}$  be a set of individuals such that (i)  $\forall i \in \mathcal{E}, \delta_i \geq \delta^*$ , (ii)  $\mathcal{E} \subseteq E$  and (iii) any pair of individuals  $i \neq j \in \mathcal{E} \subseteq E$  are connected in network  $G$  by a path of length  $l_{ij}$  of expressers belonging to the set  $\mathcal{E}$ , that is  $\exists g_{ik_1} \times g_{k_1k_2} \dots \times g_{k_lj} > 0$  for  $k_1, \dots, k_l \in \mathcal{E} \subseteq E$ . When expressers  $k_1, \dots, k_l \in \mathcal{E} \subseteq E$*

---

<sup>21</sup>Intuitively expression leaves room for discussions or a debates. One can also think of politicians who create content on Twitter and who cite their peers and peers respond and cite back, while followers will generally share the content created by influencers.

all have like-minded neighbors, I say that individual  $i$  is linked to individual  $j$  through a path of like-minded expressers and denote it by  $l_{ij}^+$ .

In a given network there could be multiple connected sets of expressers  $\mathcal{E}_1, \dots, \mathcal{E}_k$  such that  $E = \bigcup_{i=1}^k \mathcal{E}_i$ . Those sets could be singletons or they could contain more than one expresser. In the network  $G_1$  in figure 1, each of both individuals 1 and 2 form a set of connected expresser(s) on their own:  $\mathcal{E}_1 = \{1\} \subset E$  and  $\mathcal{E}_2 = \{2\} \subset E$ . Moreover,  $E = \mathcal{E}_1 \cup \mathcal{E}_2 = \{1, 2\}$ .

For a given set of connected expressers, if each pair of neighbors are like-minded, then their long-run opinions will be the average of their initial opinions. I formalize this idea in the following proposition.

**Proposition 1** *Let  $G$  be a network of interpersonal relationships,  $\alpha_0$  an initial opinion vector and consider  $\mathcal{E} \subseteq E$  a given set of connected expressers.*

(i) *(Stubborn) If  $|\mathcal{E}| = 1$  and  $\mathcal{E} = \{i\}$  then*

$$\forall t \geq 1, \alpha_{i,t} = \alpha_{i,0}$$

(ii) *(Like-minded) If  $|\mathcal{E}| = \kappa > 1$  and  $\forall i \neq j \in \mathcal{E}, \forall j \in N_i \cap \mathcal{E}, |\alpha_{i,0} - \alpha_{j,0}| < \tau$  then for  $\mu \in (0, 1/\kappa)$  and  $j_1, \dots, j_\kappa \in \mathcal{E}$ ,*

$$\exists t^* \geq 1, \forall t \geq t^*, \alpha_{i,t} = \frac{\alpha_{i,0} + \alpha_{j_1,0} + \dots + \alpha_{j_\kappa,0}}{\kappa}$$

**Proof.** See Appendix 7.1. ■

When a set of connected expressers contains at least one pair of ideologically-opposed neighbors, *generically* long-run opinions become extreme, that is they take the value 1 or  $-1$ . For very specific initial opinion distributions some expressers within the same set can hold moderate opinions (strictly between  $-1$  and  $1$ ) in the long-run. First, I discuss an example to explain when and why opinion

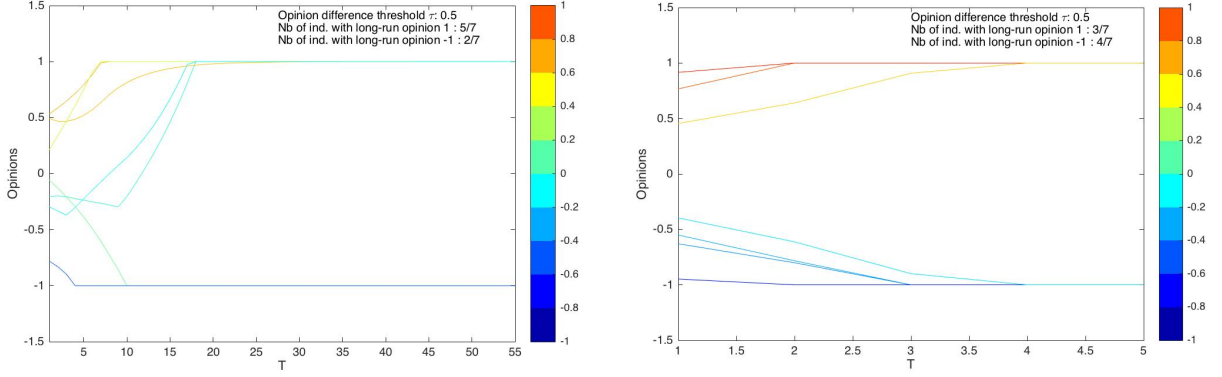


Figure 2: Non-monotonic and monotonic opinion updating in a circle

updating within a set of connected expressers can be non-monotonic. I do so, to show why the opinion of an individual who is not directly connected to an ideologically-opposed individual, can adopt an extreme opinion in the long-run. Second, I discuss why moderate opinions can survive in the long-run. Finally, I summarize the discussion in proposition 2. In all the subsequent examples, I assume that the opinion difference of two directly linked individuals at a given time period is  $\tau = 0.5$ .

#### 4.1.1 Non-monotonic opinion updating

The coexistence of pairs of directly linked individuals who are like-minded and pairs of directly linked individuals who are ideologically-opposed<sup>22</sup> within the same connected set of expressers can lead to either monotonic or non-monotonic opinion updating. Non-monotonic opinion updating can occur when an individual is indirectly connected to an ideologically-opposed individual through a chain of like-minded expressers.

To see this, consider a wheel  $G$  of seven individuals  $\{1, \dots, 7\}$  such that  $g_{12} = g_{23} = \dots = g_{67} = g_{71} = 1$  and the remaining entries of  $G$  are zero. All individuals express because their local popularity is higher or equal to  $\delta^*$  and together all seven form a unique connected set of expressers.

<sup>22</sup>For example, consider three expressers  $i, j$  and  $k$  such that  $g_{ij} = g_{ik} = 1$  and all three form together a connected set of expressers. Suppose that initial opinions are as follows:  $\alpha_{i,0} = 0.5$ ,  $\alpha_{j,0} = 0.6$ ,  $\alpha_{k,0} = -0.5$  and  $\tau = 0.5$ . Individuals  $i$  and  $j$  are neighbors that are initially like-minded. While, individuals  $i$  and  $k$  are neighbors who are initially ideologically-opposed.

Consider the following initial opinion vector:  $\alpha_0 = (-0.21, -0.06, 0.53, 0.49, 0.21, -0.78, -0.29)$ . Individuals 5 and 6 are initially ideologically-opposed with initial opinions  $\alpha_{5,0} = 0.21$  and  $\alpha_{6,0} = -0.78$ . The remaining individuals  $\{1, 2, 3, 4, 7\}$  have initially like-minded neighbors. The evolution of opinions is plotted in the left panel in figure 2, with time periods on the  $x$ -axis and opinions on the  $y$ -axis. The colormap on the east side of the figure is associated with the opinion interval  $[-1, 1]$  and the colors of the curves correspond to initial opinions.

Individuals 5 and 6 are pushed to the upper and lower bound of the opinion interval  $[-1, 1]$  after few periods of interaction, as they are ideologically opposed. The evolution of their opinion is monotonic. The opinion of individual 5 becomes more and more positive, while the opinion of individual 6 becomes more and more negative. However, individuals 1 and 7 update their opinions non-monotonically. Individual 7 and 6 are initially like-minded and the attractive effect is at play in the first few periods of interaction. Hence the opinion of individual 7 starts becoming more negative, because it converges towards the opinion of individual 6. But after a few periods, individual 6 becomes extreme *too fast* by reaching the lower bound  $-1$ . The opinion difference with their direct neighbor individual 7 becomes larger and larger, until a point where this difference becomes higher than the threshold  $\tau$  and individual 7 starts getting repulsed by the extreme opinion of individual 6. In other words, the repulsive effect takes over. Intuitively, this situation occurs when a given individual  $i$  is having a discussion with an individual  $j$  who is initially more or less like-minded, but individual  $i$  is more neutral than  $j$ . As the discussion goes on, individual  $j$  becomes too extreme in an unreasonable fashion that individual  $i$  starts defending the opposite view.

Finally, in the right panel of figure 2, I provide an example where opinion updating is monotonic. Initial opinions are given by  $\alpha_0 = (-0.55, 0.77, -0.63, -0.9478, 0.92, -0.4, 0.45)$ . In other words, each individual has an ideologically-opposed neighbor and opinions converge *monotonically* to the upper or lower bound.

#### 4.1.2 Moderate long-run opinions

There exists initial opinion distributions such that *moderate* opinions (i.e. with an opinion that is neither 1, nor  $-1$ ) survive in the long-run within a set of connected expressers containing at least

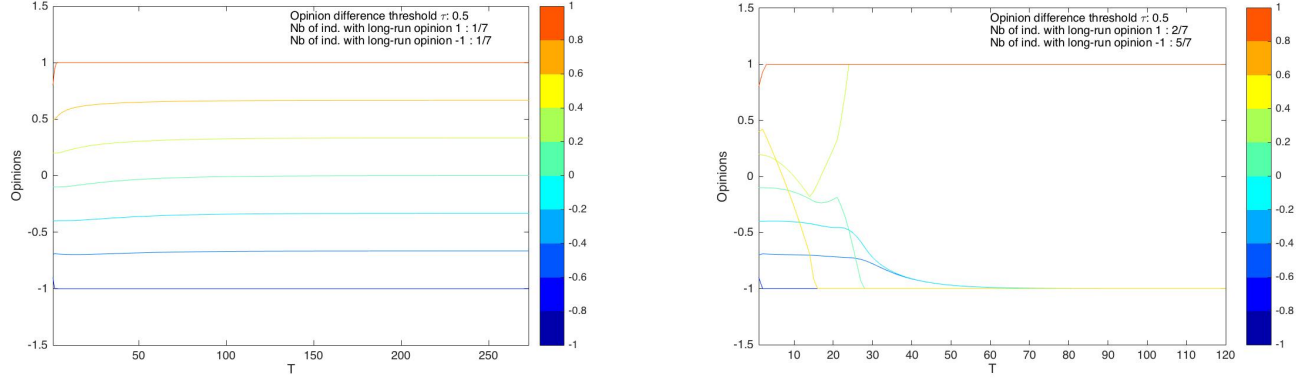


Figure 3: Moderate expressers

one ideologically-opposed pair. To see this, consider again a wheel  $G$  of 7 individuals such that  $g_{12} = g_{23} = \dots = g_{67} = g_{71} = 1$  and all the remaining entries of  $G$  are zero. Suppose that initial opinions are :  $\alpha_0 = (-0.9, -0.7, -0.4, -0.1, 0.2, 0.5, 0.8)$ . Individuals 1 and 7 are initially ideologically-opposed. Individuals  $\{2, \dots, 6\}$  all have neighbors who are initially like-minded. The evolution of opinions is plotted in the left panel of figure 3. Individuals  $i \in \{2, \dots, 5\}$  remain moderate in the long-run and never adopt an extreme opinion of 1 or  $-1$ . To understand why moderate opinions can persist in the long-run, first notice that individual 7 is repulsed by her direct ideologically-opposed neighbor individual 1. Second, the opinion difference of individuals 6 and 7 remains smaller than  $\tau$  even when 7 becomes extreme. Similarly, the opinion difference of individuals 1 and 2 remains smaller than  $\tau$  even when 1 becomes extreme. Intuitively, extreme individuals 1 and 7 influence in opposite *directions* the chain of like-minded individuals that separate them so that in the long-run each of these intermediate individuals have like-minded neighbors and remain moderate. This situation bears a resemblance to the left-right political spectrum in some countries, where the moderate parties survive in the long-run. To summarize, moderate opinions of expressers can persist in the long-run if such expressers only have like minded-neighbors and are linked to at least two ideologically-opposed expressers by a path of like-minded neighbors.

In the right panel of figure 3, all individuals hold extreme opinions in the long-run. Similarly to the setting in the left panel of figure 3, individuals  $i \in \{2, \dots, 6\}$  have initially like-minded neigh-

bors and are connected to individuals 1 and 7 directly or indirectly through a chain of like-minded neighbors. But the opinion difference of individuals 6 and 7 is initially larger than in the previous case (left panel of figure 3) and as individual 7 is pushed towards holding an extreme opinion, her opinion difference with individual 6 becomes larger and larger to a point where this difference becomes larger than  $\tau$ .

There is one last case to be discussed. Consider a wheel  $G$  of three individuals such that  $g_{12} = g_{23} = 0$  and all the remaining entries of the network  $G$  are zero. All three individuals express and all three form together a unique connected set of expressers. Suppose that, the initial opinion of individual 2 is  $\alpha_{i,0} = 0$ . Suppose that individuals 1 and 3 hold respectively the following initial opinions  $\alpha_{1,0} = -1$  and  $\alpha_{3,0} = 1$ . In this case, the long-run opinion of individual  $i$  is  $\alpha_{i,\infty} = 0$  because she is equally repulsed by both neighbors but in opposite directions.

I formalize the above discussion in the following proposition. Recall that  $\underline{N}_{i,t}$  and  $\overline{N}_{i,t}$  are the sets of neighbors of  $i \in N$  who are respectively like-minded and ideologically-opposed at period  $t \geq 0$ . Denote the set of expressers that have at least one initially ideologically-opposed neighbor by  $IO(\mathcal{E}) = \{i \neq j \in \mathcal{E}, i \in \overline{N}_{j,0} \neq \emptyset, j \in \overline{N}_{i,0} \neq \emptyset \text{ and } g_{ij} = 1\}$ .

**Proposition 2** *Let  $G$  be a network of interpersonal relationships,  $\alpha_0$  an initial opinion vector and consider  $\mathcal{E} \subseteq E$  a set of connected expressers. Suppose that  $IO(\mathcal{E}) \neq \emptyset$ .*

- (i) *If there exists at each  $t \geq 0$  paths of expressers with only like-minded neighbors connecting  $i \notin IO(\mathcal{E})$  to  $i_1, i_2 \dots i_k \in IO(\mathcal{E})$  then  $\alpha_{i,\infty} \in \text{conv}(\alpha_{i_1,\infty}, \dots, \alpha_{i_k,\infty})$ .*
- (ii) *Let  $i_1, i_2 \dots i_k \in IO(\mathcal{E})$ . If there exists  $i \in IO(\mathcal{E})$  such that (a)  $\alpha_{i,0} = 0$ , (b)  $\forall t \geq 0 \underline{N}_{i,t} = \emptyset$ , (c) and  $\sum_{j \in \overline{N}_{i,t}} \alpha_{j,t} = 0$  then  $\alpha_{i,\infty} = 0$ .*
- (iii) *otherwise if (i) and (ii) don't hold then  $\forall i \in \mathcal{E}, \alpha_{i,\infty} \in \{-1, 1\}$ .*

**Proof.** See Appendix 7.2. ■

To summarize, when the connected set of expressers is formed of only initially like-minded expressers, influence is positive and opinions get *attracted* to the average opinion of the group. However, when the set contains at least one ideologically-opposed pair of neighbors, the repulsive and attractive effect can both be at play. Finally, when a pair of expressers are initially like-minded they take a longer *time* to reach consensus but when the pair is ideologically-opposed they disagree at a much faster rate. I formalize this idea in the following lemma.

**Lemma 1** *Let  $i \neq j \in \mathcal{E} \subset E$  such that  $|\mathcal{E}| = 2$  and  $t_{\alpha_\infty} = \min\{t : |\alpha_t - \alpha_\infty| < \epsilon\}$ . If  $\alpha_\infty^a$  is the long-run opinion vector when  $|\alpha_{i,0} - \alpha_j, 0| < \tau$  and  $\alpha_\infty^r$  is the long-run opinion vector when  $|\alpha_{i,0} - \alpha_j, 0| \geq \tau$  then  $t_{\alpha_\infty^r} < t_{\alpha_\infty^a}$ .*

**Proof.** See Appendix 7.3. ■

## 4.2 The process of interpersonal influence

In order to study the overall dynamics of opinions and show that opinions converge in the long-run, I build a hearing matrix which takes into account who listens to whom. In other words, this hearing matrix takes into account the two opinion update rules depending on the type of individual: expresser (updates according to the law of motion (3)) or consensual (updates à la DeGroot). In particular I study the long-run behavior starting at the time period where pairs of ideologically-opposed expressing neighbors have repulsed each other towards the most extreme opinion. As I have shown in the previous section, given the parameter  $\mu$  in the law of motion (3), a pair of ideologically-opposed individuals repulse each other at a much faster rate than a pair of like-minded individuals who debate to reach a consensus. Formally let  $t^* \geq t$  be the time period by which the least ideologically-opposed pair of directly connected expressers, in the group of individuals  $N$ , have repulsed each other to reach opinions at the upper and lower bound of the opinion interval. That is for any period  $t$  beyond time period  $t^*$ , ideologically-opposed neighbors who express, are no longer updating their opinions and have long-run opinions that are either 1 or  $-1$ . Given a network  $G$  representing interpersonal relationships, denote by  $\tilde{G}$  the hearing matrix with typical entries  $\tilde{g}_{ij}$ .



**Consensual individuals.** For each individual  $i \in N$  such that  $\delta_i < \delta^*$ , the entries in the hearing matrix become  $\tilde{g}_{ij} = g_{ij}/d_i$ .

**Expressers.** For individuals who choose to express there are four cases to consider.

- (i) For all  $i \in N$  in a connected set of expressers  $\mathcal{E}$  such that  $|\mathcal{E}| = 1$  (stubborn), the entries of the hearing matrix  $\tilde{G}$  are:  $\tilde{g}_{ii} = 1$  and  $\tilde{g}_{ij} = 0$  for all  $j \in N_i$ .
- (ii) For all  $i \in N$  in a connected set of expressers  $\mathcal{E}$  such that  $|\mathcal{E}| = \kappa > 1$  with like-minded neighbors at period  $t^*$ , i.e.  $\forall i \neq j \in \mathcal{E}$  and  $j \in N_i \cap \mathcal{E}$ ,  $|\alpha_{i,t^*} - \alpha_{j,t^*}| < \tau$ , the entries of the hearing matrix  $\tilde{G}$  are:  $\tilde{g}_{ii} = 1 - |N_i \cap \mathcal{E}| \mu$  and  $\tilde{g}_{ij} = \mu$  for  $j \in N_i \cap \mathcal{E}$ .
- (iii) For all  $i \in N$  in a connected set of expressers  $\mathcal{E}$  such that  $|\mathcal{E}| = \kappa > 1$  with ideologically-opposed neighbors, i.e.  $\forall i \neq j \in \mathcal{E}$  and  $j \in N_i \cap \mathcal{E}$ ,  $|\alpha_{i,t^*} - \alpha_{j,t^*}| \geq \tau$ , the entries in the hearing matrix  $\tilde{G}$  are:  $\tilde{g}_{ii} = \tilde{g}_{jj} = 1$ ,  $\tilde{g}_{ik} = 0$  for all  $k \in N_i$  and  $\tilde{g}_{jk} = 0$  for all  $k \in N_j$ .

**Example 3** Consider network  $G$  in figure 1. The hearing matrix  $\tilde{G}$  has the following entries for expressers 1 and 2 who are both stubborn:  $\tilde{g}_{11} = \tilde{g}_{22} = 1$  and  $\tilde{g}_{1j} = \tilde{g}_{2j} = 0$  for all  $j \in N_1 \cup N_2$ . For the consensual individuals  $i \in \{3, 4, 5\}$ ,  $\tilde{g}_{ii} = \tilde{g}_{i2} = 1/2$ . For the consensual individuals  $i \in \{6, 7, 8\}$ ,  $\tilde{g}_{ii} = \tilde{g}_{i1} = 1/2$ . Finally for the consensual individual 9,  $\tilde{g}_{99} = \tilde{g}_{91} = \tilde{g}_{92} = 1/3$ . All the remaining entries are zero.

For a given network structure  $G$ , the process of interpersonal influence describing the evolution of opinions at period  $t \geq t^*$  is given by the following equation:

$$\alpha_{t+1} = \tilde{G}\alpha_t \tag{5}$$

By induction, the opinions at period  $t \geq t^*$  are given by  $\tilde{G}^t \alpha_{t^*}$  and the limit yields the long-run opinions. A few comments are in order. First, the entries of the hearing matrix  $\tilde{G}$  are all non-negative and all the diagonal entries are strictly positive. Moreover it has rows and columns that sum to one. Hence, the eigenvalues of  $\tilde{G}$  are all lower or equal to 1 and  $\lim_{t \rightarrow \infty} \tilde{G}^t$  exists. The

entry on the row  $i$  and column  $j$  of the matrix  $\lim_{t \rightarrow \infty} \tilde{G}^t$  is the weight (between 0 and 1) that the opinion of individual  $i$  at period  $t^*$  has in the final opinion of individual  $j$ . Second,  $\tilde{G}$  is a reducible<sup>23</sup> matrix, whenever there exists at least one individual who chooses to express. To see this, recall that consensual individuals account for the opinions of all their neighbors, while expressers only account for the opinions of neighbors who also express (when such neighbors exist). Hence, there always exists at least one path starting at a node that represents a consensual individual and that ends at a node representing an expresser. However, there does not exist any paths that start at a node representing an expresser and that end at a node representing a consensual player. Third, the multiplicity of the eigenvalue 1 is equal to the number of essential classes<sup>24</sup> in the hearing matrix  $\tilde{G}$ . To see this simply, consider a circle as a network structure with exactly  $k$  individuals, where each individual has two neighbors and where initial opinions are such that each individual has at least one neighbor who is ideologically-opposed. For this network structure, given the expression threshold  $\delta^*$ , all individuals choose to express. Since each individual has at least one ideologically-opposed neighbor, each individual reaches an extreme opinion of 1 or  $-1$  after few periods of interaction. In this setting, individuals no longer take into account the opinions of other expressers in the long-run and each individual forms an essential class on their own. Hence, the hearing matrix  $\tilde{G}$  is simply the identity matrix of size  $k$  and the multiplicity of the eigenvalue 1 is exactly  $k$ . Beyond this example, the only case where an essential class is not a singleton is the case where there is a group of individuals that form a connected set expressers (see definition 3) that are like-minded. In other words, there exists a path connecting each pair in this connected set of expressers at each time period of interaction, but no paths from any of those expressers to an individual outside this set. I summarize the above discussion in the following theorem and provide a formal proof which makes use of standard linear algebra results.

**Theorem 1** *Given  $\alpha_{t^*} \in [-1, 1]^n$  a vector of opinion at period  $t^*$  and a hearing matrix  $\tilde{G}$  associated*

---

<sup>23</sup>For a formal definition of an irreducible matrix see for example definition 1.6 in Seneta (1981) [28].

<sup>24</sup>The definition of an *essential class* can be found in section 3.1.

with the network structure  $G$ , the long-run opinions are :

$$\alpha_\infty = \left( \lim_{t \rightarrow \infty} \tilde{G}^t \right) \alpha_{t^*} = \mathcal{G} \mathbf{a}_{t^*} < \infty$$

where  $\mathcal{G}$  is the spectral projector<sup>25</sup> associated with the eigenvalue 1. Moreover, the algebraic multiplicity of the eigenvalue 1 is equal to the number of essential classes of the hearing matrix  $\tilde{G}$ .

**Proof.** See Appendix 7.4 ■

The columns corresponding to consensual individuals in the matrix  $\mathcal{G}$  are all zero, meaning that in the long-run the initial opinions of such individuals vanish. Their opinions remain hidden through out the periods of interaction. As for the columns corresponding to expressers, they have at least one strictly positive entry. In particular, the long-run opinions of consensual individuals are exactly convex combinations of initial opinions of expressers. In other words, the long-run opinion of consensual individuals is affected by the long-run opinions of all the expressers to whom they are connected to through a path of other consensual individuals. Hence, the total impact of the initial opinion of a given expresser  $i \in N$  over long-run opinions can be assessed by considering the total weight an expresser has in the long-run opinions of other individuals. To summarize, I consider the following statistic.

**Definition 4 (Spectral influence)** Given a network structure  $G$ , a hearing matrix  $\tilde{G}$  and its limit  $\mathcal{G}$ , the spectral influence of individual  $i \in N = \{1, \dots, n\}$  is:

$$s_i = \frac{1}{n} (\mathcal{G}' \mathbf{1})_i$$

---

<sup>25</sup>The hearing matrix at period  $t \geq t^*$  can be decomposed in the following way:  $\tilde{G}^t = U J^t U^{-1}$ , where  $U = [U_1 \ U_2]$  is the eigenspace, such that the column(s)  $U_1$  are the eigenvectors associated with the eigenvalue 1 and  $U_2$  the columns that correspond to the eigenvectors associated with the remaining eigenvalues strictly smaller than 1,  $U^{-1} = V = \begin{bmatrix} V_1 \\ V_2 \end{bmatrix}$  it's inverse,  $J = \begin{bmatrix} I_{p \times p} & \mathbf{0} \\ \mathbf{0} & \mathbf{K} \end{bmatrix}$ , where  $I_{p \times p}$  is an identity matrix of size  $p$  and  $p$  is the multiplicity of the eigenvalue 1 and  $\mathbf{K}$  is a diagonal matrix containing the remaining eigenvalues smaller than 1. The spectral projector associated with the eigenvalue 1 is  $\mathcal{G} = U_1 V_1$ .

**Example 4** Consider network  $G$  in figure 1 and suppose that the initial opinions of individuals 1 and 2 are  $\alpha_{0,1}$  and  $\alpha_{0,2}$ . Recall that only individuals 1 and 2 choose to express. Since both are not directly connected nor are they connected via a chain of expressers, each of them forms an essential class on their own. The spectral projector  $\mathcal{G}$  associated to the eigenvalue 1 of the hearing matrix  $\tilde{G}$  is a symmetric matrix of size 9 and is given by :  $\mathcal{G}_{11} = \mathcal{G}_{i1} = 1$  for  $i \in \{6, 7, 8\}$ . That is the long-run opinion of individuals 6, 7 and 8 is fully determined by the initial opinion of individual 1. As for individual 2,  $\mathcal{G}_{22} = \mathcal{G}_{j2} = 1$  for  $j \in \{3, 4, 5\}$ . The long-run opinion of individual 9 is equally determined by the initial opinions of expressers 1 and 2, that is  $\mathcal{G}_{91} = \mathcal{G}_{92} = 1/2$ . All the remaining entries of the matrix  $\mathcal{G}$  are zero. Hence, long-run opinions are  $a_{\infty,i} = \alpha_{0,1}$  for  $i \in \{1, 6, 7, 8\}$ ,  $a_{\infty,j} = \alpha_{0,2}$  for  $j \in \{2, 3, 4, 5\}$  and  $a_{\infty,9} = \frac{1}{2}\alpha_{0,1} + \frac{1}{2}\alpha_{0,2}$ . Expressers 1 and 2 have an identical spectral influence equal to  $s_1 = s_2 = 1/2$ .

### 4.3 Patterns of long-run opinions: consensus and bi-polarization

In this section, I characterize patterns of long-run opinions. I am interested in understanding the impact of the different parameters of the model, namely the structure of the network of interpersonal relationships  $G$  along with the threshold  $\tau$  of opinion differences and the distribution of initial opinions  $\alpha_0$ , on the pattern of long-run opinions  $\alpha_\infty \in [-1, 1]^n$  that I obtain. I focus on two patterns. First, I study consensus as a benchmark, where all the individuals in the long-run become like-minded.

**Definition 5 (Consensus)** Long-run opinions form consensus if  $\forall i \neq j \in N, |\alpha_{\infty,i} - \alpha_{\infty,j}| < \tau$ .

The above definition of consensus reflects the idea that the limiting opinions need not to be identical but the pairwise difference needs to be at most  $\tau$ . In other words, the matrix  $\mathcal{G}$  (see Theorem 1) in general will not have identical rows.<sup>26</sup> Second I consider bi-polarization<sup>27</sup> of long-run

<sup>26</sup>The spectral projector has identical rows if and only if the algebraic multiplicity of the unit eigenvalue is 1; meaning that computing the perron vector is sufficient to obtain the long-run opinions. But having more than two sets of connected expressers is translated into an algebraic multiplicity of the unit eigenvalue strictly higher than 1.

<sup>27</sup>I focus on opinion bi-polarization rather than the more general case of opinion polarization, where a polarized society is one that is divided into a small number (larger than two) of opposed groups. The special case of bi-polarization fits applications of the model where it would take a very large group - e.g. at least half of the population of interest - to over-turn a policy or to elect a president or to produce a divided public opinion. For example, one can think of Brexit, the election of Trump, the implementation of a Carbon tax in France or even at the beginning of 2020 divided views about risk reducing measures regarding the coronavirus.

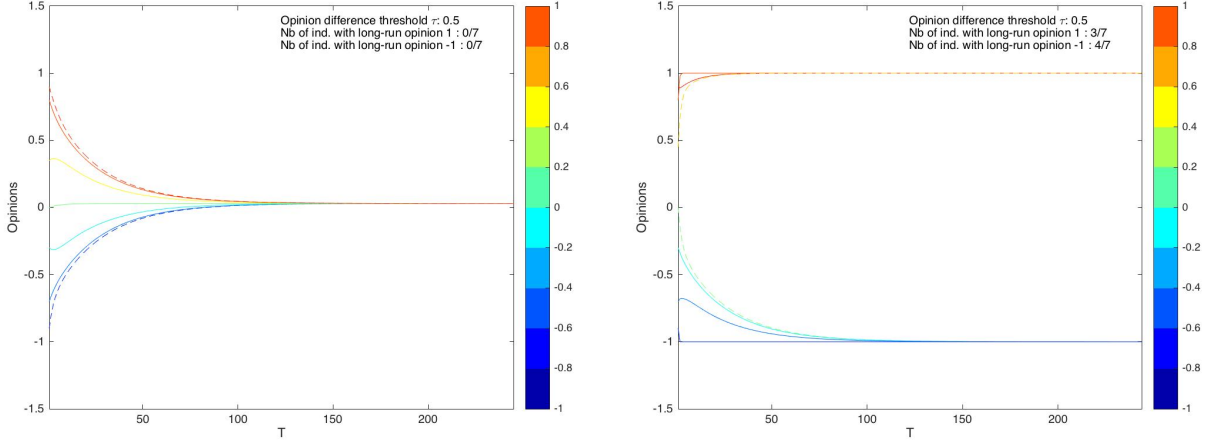


Figure 4: Consensus and bi-polarization

opinions, that is given a network structure and an initial opinion vector, the society of  $N$  individuals gets divided into two opinion groups of approximately equal size with strong agreement within each group and disagreement across both groups.

**Definition 6 (Bi-polarization)** *Long-run opinions are bi-polarized if (i) the society of  $N$  individuals is divided into two groups of size  $N_1$  and  $N_2$  such that  $|N_1 - N_2| < \epsilon$ , (ii) for  $k \in \{1, 2\}$ ,  $\forall i \neq j \in N_k$ ,  $|\alpha_{\infty, i} - \alpha_{\infty, j}| < \tau$ , (iii)  $|\bar{\alpha}_{\infty, N_1} - \bar{\alpha}_{\infty, N_2}| \geq \tau$  where  $\bar{\alpha}_{\infty, N_i}$  is the average long-run opinion of group  $N_i$  for  $i \in \{1, 2\}$ .*

In general, the higher the number of expressers induced by a given network structure  $G$  and the more the long-run opinion pattern depends on the distribution of initial opinions within each expresser's neighborhood. To see this, consider a situation where most individuals can *express*, this happens in regular network structures. When each individual has neighbors such that their initial opinion difference is lower than  $\tau$  then the group converges to consensus; even if there exists a pair of expressers that are not directly connected and that hold ideologically-opposed views. For example, let  $G$  be a line network with 7 individuals such that  $g_{12} = g_{23} = g_{34} = g_{45} = g_{56} = g_{67} = 1$  and the remaining entries are all zero. Consider the following initial opinion vector  $\alpha_0 = (-0.9, -0.7, -0.3, 0, 0.45, 0.8, 0.9)$ . For  $\tau = 0.5$ , individuals reach consensus even-though

individuals 2 and 6 are ideologically-opposed, as shown in the left panel of figure 4. However, for the same network structure if the initial opinion vector is  $\alpha_0 = (0, -0.3, -0.7, -0.9, 0.8, 0.9, 0.45)$  then this group of individuals doesn't reach consensus, as shown in the right panel of figure 4. This is because individuals 3 and 4 repulse each other towards long-run opinions  $a_{\infty,3} = -1$  and  $a_{\infty,4} = 1$ . In this case, long-run opinions become  $\alpha_\infty = (-1, -1, -1, -1, 1, 1, 1)$  which corresponds to a bi-polarized group of individuals.

#### 4.3.1 Consensus

Consensus of opinion depends on the number of expressers in the network, how they are connected (or not) to each other and the distribution of their initial opinions. Namely, consensus is reached whenever there is a unique expresser in the network or several expressers with like-minded neighbors who express and for any two disjoint sets of connected expressers, individuals are in average like-minded across both sets. In the former case, the initial opinion of the unique expresser fully determines the long-run opinions of the remaining consensual individuals. While in the latter case, long-run opinions of individuals can be different but that difference is at most  $\tau$ .

**Proposition 3** *Given a network structure  $G$ , an initial opinion vector  $\alpha_0 \in [-1, 1]^n$  and the set of expressers  $E \subset N$ , long-run opinions form a consensus if and only if exactly one of the following statements holds:*

- (i)  $|E| = 1$
- (ii) *There exists  $\kappa \geq 1$  sets of connected expressers s.t.  $E = \bigcup_{k=1}^{\kappa} \mathcal{E}_k$ ,  $\forall k_i \neq k_j \in \{1, \dots, \kappa\}$ ,  $|\bar{\alpha}_{0, \mathcal{E}_{k_i}} - \bar{\alpha}_{0, \mathcal{E}_{k_j}}| < \tau$  and for each set  $\mathcal{E}_{k_i} \subset E$ ,  $\forall i \neq j \in \mathcal{E}_{k_i}$  and  $j \in \mathcal{E}_{k_i} \cap N_i$ ,  $|\alpha_{0,i} - \alpha_{0,j}| < \tau$ .*

**Proof.** See Appendix 7.5 ■

**Corollary 1 (unique opinion leader)** *Let  $\alpha_0 \in [-1, 1]^n$  be a vector of initial opinions. If the network structure  $G$  is a star where the central node has a degree  $k > 1/\delta^*$  then the long-run opinions  $a_\infty$  form a consensus.*

These necessary and sufficient conditions map long-run opinions to the network structure and the initial distribution of opinions. To see this, recall that choosing to express or hide one’s opinion is determined by local popularity, which is a network statistic. Moreover, in the case where the structure of the network of interpersonal relationships allows the existence of more than one expresser, the distribution of initial opinions along with the structure of connections among expressers become crucial in reaching consensus. In particular, consensus can never be reached if there exists a pair of ideologically-opposed expressers who are directly linked.

### 4.3.2 Bi-polarization

Unlike consensus, characterizing the initial opinion distribution of expressers within a network isn’t sufficient to explain the emergence of a bi-modal long-run opinion distribution.<sup>28</sup> The characterization is challenging because the emergence of a bi-polarized society depends on where expressers are placed in the network, how they interact together and how many people they influence. If extreme opinion groups of expressers do form, the network structure together with the opinion difference threshold  $\tau$  determine to which extreme opinion group, consensual individuals will belong.

Recall that the long-run opinion of each consensual individual is exactly a convex combination of the opinions of expressers to whom they are directly or indirectly linked by a path of other consensual individuals. For some network structures and expressers’ opinion distributions, consensual individuals can receive influence from ideologically-opposed groups of expressers. In which case, consensual individuals remain moderate in the long-run. In figure 1, individual 11 holds a long-run opinion of *zero*, when individuals 1 and 2 hold respectively opinions 1 and  $-1$ . In other words, the existence of long-run moderate individuals can block the split of the population into two ideologically-opposed groups.

Intuitively, the network position of moderate consensual individuals allows them to receive influence from different opinion groups. They are initially linked to both influence groups and they don’t express an opinion themselves. Such individuals could be interpreted as intermediaries or neutral TV hosts or moderators, who act as buffers against opinion polarization.

---

<sup>28</sup>That is, a *central interval that is sparsely populated and left and right intervals that are densely populated*. Friedkin (2015) [13].

Nevertheless, the absence of moderate consensual individuals isn't sufficient for long-run opinions to be polarized. In particular, one needs to account for the size of both opinion groups. The split of expressers into two extreme opinion groups isn't sufficient to cause opinion bi-polarization. Both extreme opinion groups need to have influence over an equally large number of consensual individuals, so that the society becomes divided. In other words, consensual individuals to belong to one of both extreme opinion groups, they need to receive *enough* influence from the members of one group so that they hold similar extreme opinions in the long-run. Finally, for two extreme opinion groups to form, at least one pair of ideologically-opposed expressers need to interact together or there must exist at least two individuals initially at each end of the opinion spectrum. Recall that the influence of each expresser  $i \in E \subset N$  is summarized by their spectral influence  $s_i$  in definition 4.

**Lemma 2** *Suppose that  $E = E^+ \cup E^- \subset N$  such that  $|E^+| \geq 1$ ,  $|E^-| \geq 1$  and  $\forall i \in E^+$ ,  $\alpha_{i,\infty} = 1$  and  $\forall i \in E^-$ ,  $\alpha_{i,\infty} = -1$ . If long-run opinions are bi-polarized then*

- (i)  $|\sum_{i \in E^+} s_i - \sum_{i \in E^-} s_i| < \epsilon$
- (ii)  $\nexists k \in C$  s.t.  $1 - \tau/2 \leq \sum_{i \in E^+} \mathcal{G}_{ki} \leq \tau/2$

**Proof.** See Appendix 7.6. ■

## 5 Network topology and opinion patterns

In this section I explore the model through simulations. The objective of this section is to relate the topology of the network to long-run polarization of opinions. Polarization of opinions is measured simply by looking at the variance in final opinions. Clearly, the maximum level of polarization is 1 and the minimum level is 0. To focus on the network topology, I generate a large number of initial opinions (for example 1000) distributed uniformly at random over  $[-1, 1]$  for the same network structure with  $n$  individuals, then I compute the average level of polarization of the (1000) final opinion vectors obtained.



I focus on scale-free networks because it's a general family of networks where the topology allows for the co-existence of expressers and consensual individuals. I provide a number of (inter-related) network statistics that measure inequality in the degree distribution such as assortativity, the Gini coefficient of the degree distribution and per capita average degree. I also look at network statistics that measure connectivity such as neighborhood connectivity, average path length and whether expressers are densely or sparsely connected.<sup>29</sup> All the definitions of those network statistics can be found in appendix 7.7.

Formally, I study the evolution of opinions of a large set of individuals in scale-free networks. The degree distribution within scale-free networks follows a power law. Due to this inequality in the degree distribution, expressers and consensual individuals co-exist. Barabási and Albert (1999) [6] explain the emergence of such networks by the following two mechanisms: the network expands because new individuals (vertices) keep getting added to an existing network and they create links (attach) preferentially to individuals who are already well connected. The mechanism leading to the formation of links is out of the scope of this paper. The objective of the simulations is to make obvious the impact of the interaction structure - location in the network and initial opinion - of expressers on the long-run opinion pattern of the whole group of individuals. To that end, I generate scale-free networks with a large number of individuals and work with them as a snapshot at a given point in time where no new individuals are added and no new links are created. They are generated by providing an initial number of hubs  $h$  and a number of nodes  $n_c$  a newly added node connects to.

I start by providing two examples with high and low average polarization levels with 100 individuals,  $h = 5$  initial hubs and  $n_c = 1$ . Then I give aggregate statistics by varying the number of hubs  $h$  and the number of connections  $n_c$  of each newly added nodes.

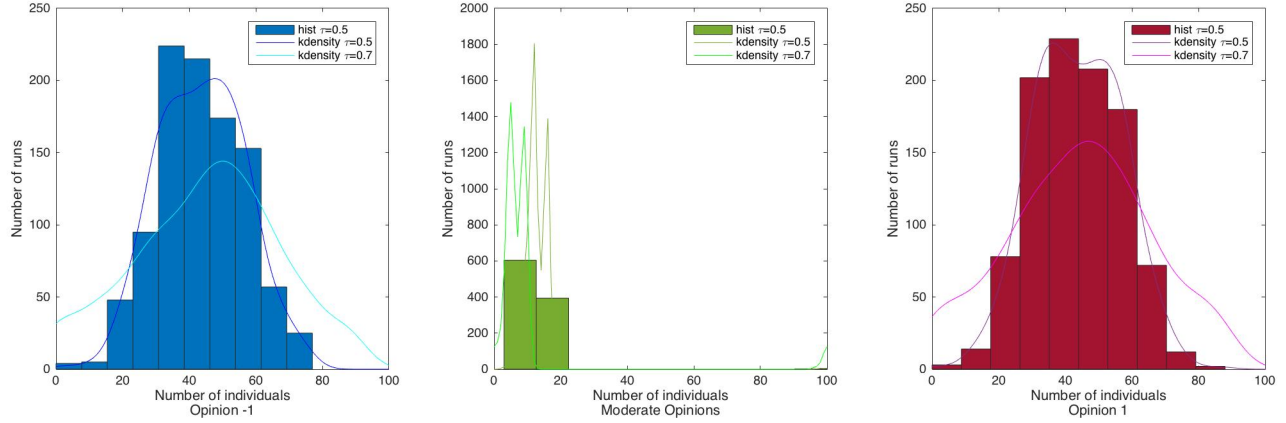


Figure 5: Kernel density over 1000 runs of groups of individuals with extreme opinion  $-1$  (blue),  $1$  (red) and moderate (green) in scale-free network  $G_1$

## 5.1 Example

I generate 1000 initial opinion vectors uniformly at random in the interval  $[-1, 1]$ . Then, I compute long-run opinions for several structures of scale-free networks with 100 individuals. I select two structures with the same number of hubs to illustrate the results of the model before moving to aggregate statistics. The first network structure  $G_1$  contains 14 expressers, represented by squares in figure 7. Out of those 14 expressers, 11 are linked through a path of expressers, highlighted in red. In other words, the network topology allows for interaction between the majority of expressers. There are a number of empirical papers that support the idea that ideologically-opposed individuals interact together often.<sup>30</sup> The first panel of figure 7 displays the initial opinions as colors given by the  $[-1, 1]$  colormap on the east side of the figure. The second panel displays the long-run opinions. The third panel shows the evolution of opinions, where dotted lines correspond to the evolution of opinions of consensual individuals and solid lines expressers. The second structure  $G_2$  contains 11 expressers and only two pairs of expressers are linked as shown in figure 8. I select those two

<sup>29</sup>Density or sparsity of connections among expressers means whether they are isolated or interact with other expressers. In a graph with  $k$  expressers than have no expressing neighbors, the subgraph of  $G$  restricted to expressers will have  $k$  components. One can think of some  $\delta^*$ -core and count the number of components to account for sparsity or density.

<sup>30</sup>See for example Conover et al. (2011) [8].

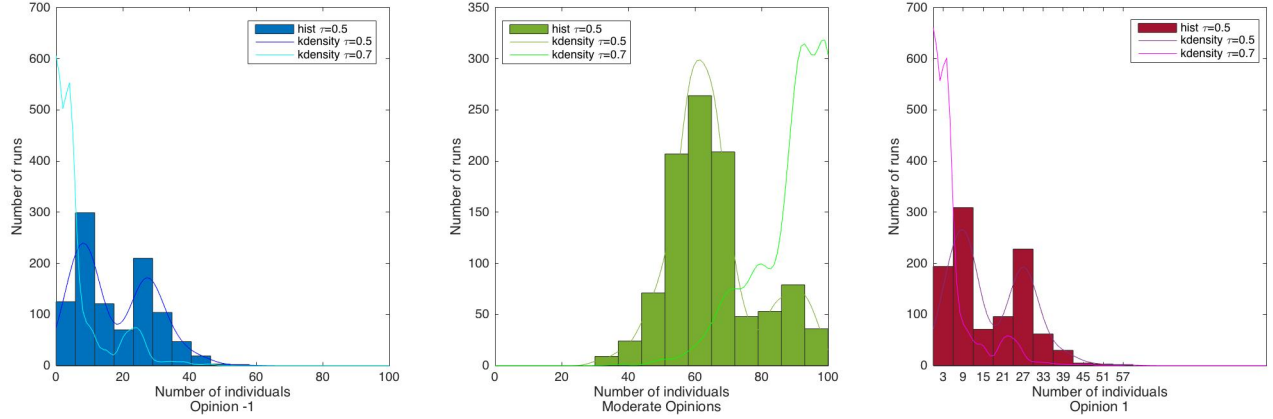


Figure 6: Kernel density over 1000 runs groups of individuals with extreme opinion  $-1$  (blue),  $1$  (red) and moderate (green) in scale-free network  $G_2$

network structures by generating for the same initial opinions a large number of scale-free networks with the same number of individuals and picking out two network structures with a high and low average level polarization. The average level of polarization in network  $G_1$  over the 1000 runs is 0.84, as opposed to only 0.43 for the network structure  $G_2$ . Figures 5 and 6 show the distribution over the 1000 runs of the size of the group of individuals holding extreme opinion  $-1$ , the group of individuals holding extreme opinion  $1$  and the group of individuals who hold a moderate opinion. For network  $G_1$  the average size of extreme groups is almost 50%, while in network  $G_2$  this same share is occupied by the group of moderate individuals.<sup>31</sup>

## 5.2 Agregate statistics

In order to study the impact of a specific topology of scale-free networks on the mean level of polarization, I vary the number of initial hubs  $h$  and the number of nodes  $n_c$  to which a newly node is added. Doing so allows me to look at network topologies with different levels of connectivity and a wide range of degree distributions. With  $n = 200$  individuals,  $h \in \{10, 20, 30, 40, 50, 60, 70\}$  and  $c_n \in \{1, 2, \dots, 8\}$  I obtain 56 scale-free network topologies<sup>32</sup> and I compute for each the mean

<sup>31</sup>To plot the distribution of the size of the opinion groups, for each of the 1000 runs I compute the number of individuals within each of the three groups, then I use this output to plot the histogram and the kernel density.

<sup>32</sup>Clearly with exactly the same number of hubs and the same  $n_c$  one obtains for different runs different network structures. But the variance in the estimated  $\alpha$  parameter is very small, meaning that they display similar degree

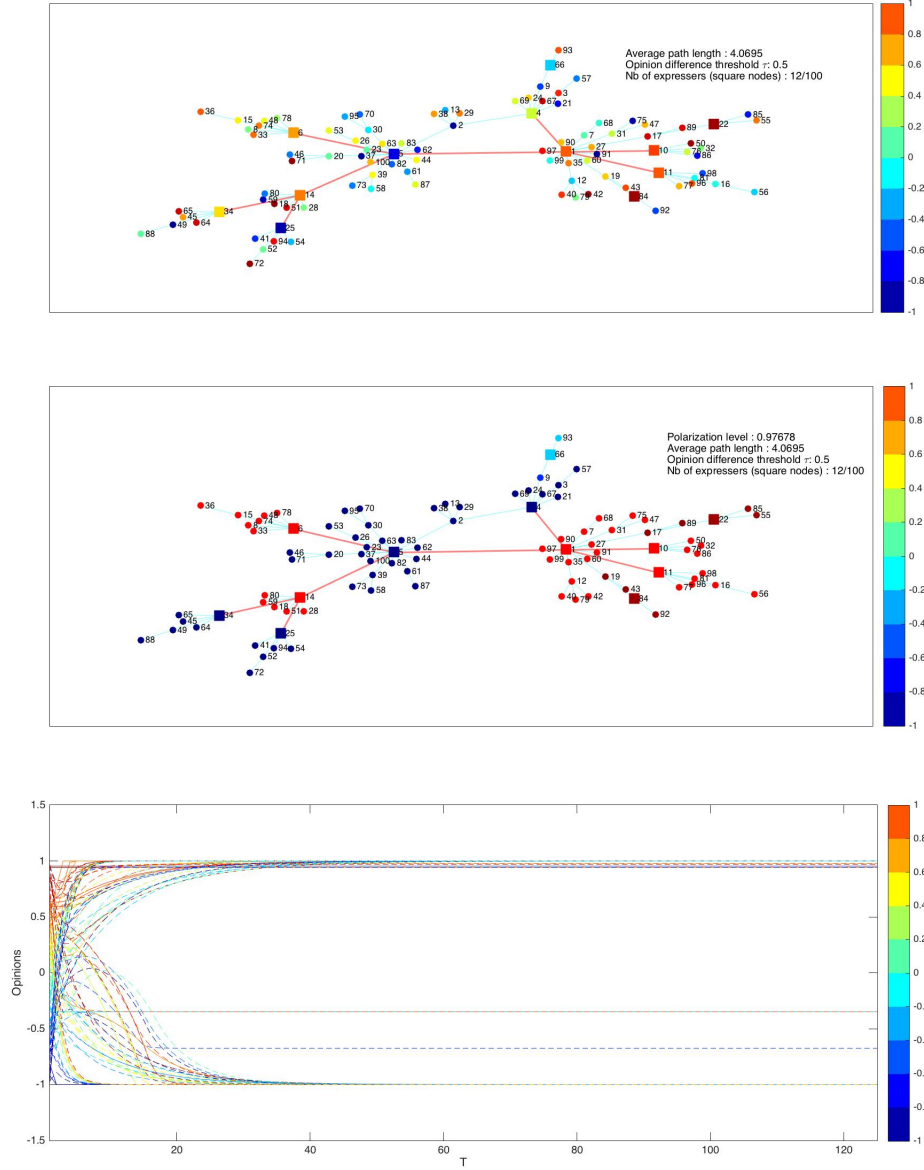


Figure 7: Initial opinions, final opinions and evolution of opinions in scale-free network  $G_1$ . Expressers are represented with a square marker. Links between expressers are highlighted in red.

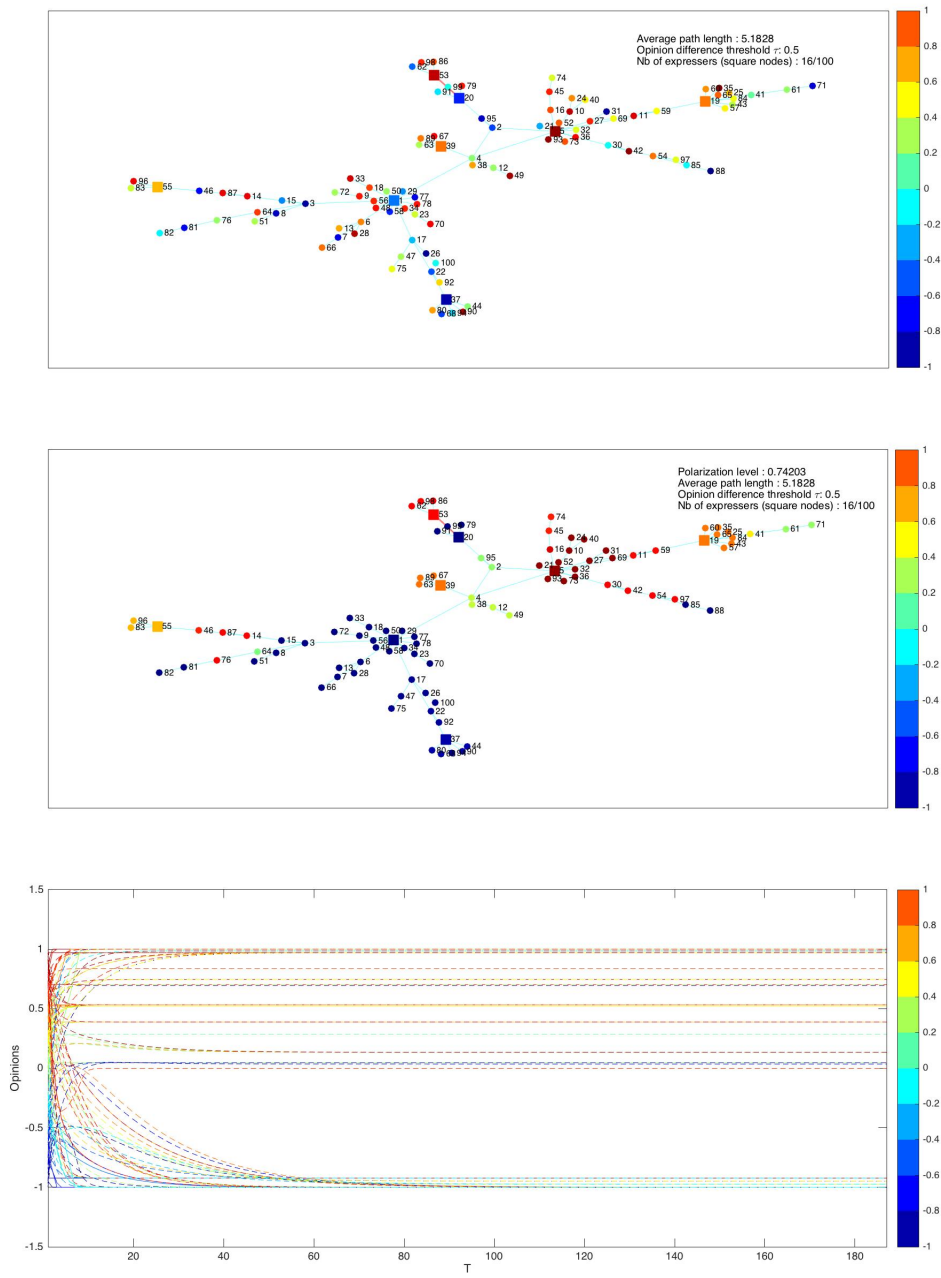


Figure 8: Initial and final opinions in the run with the highest level of polarization in scale-free network  $G_2$ . Links between expressers are highlighted in red.

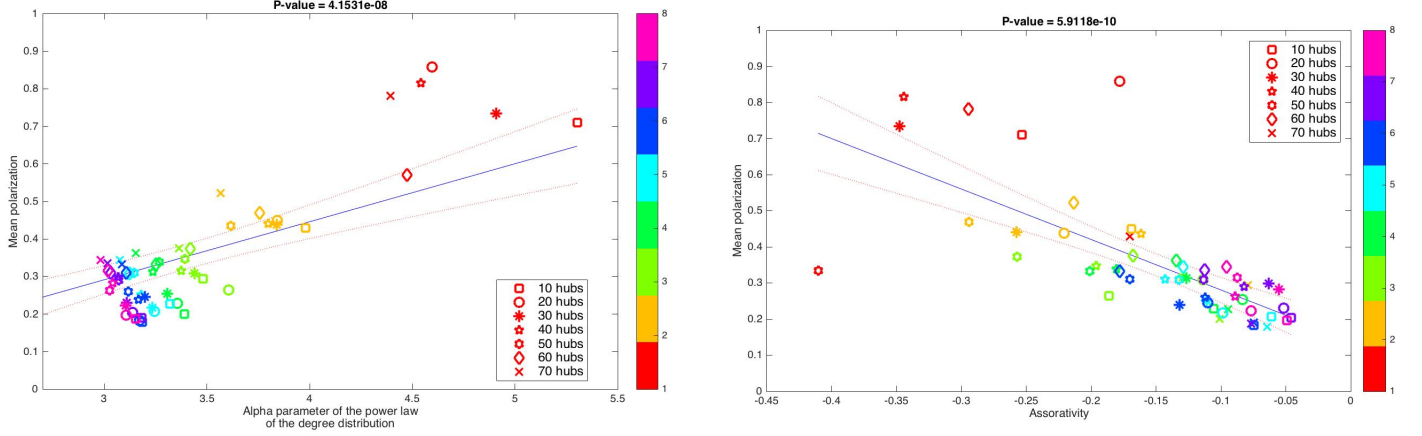


Figure 9: Markers correspond to different sizes of initial hubs  $h$  and colors correspond to the number  $n_c$  of connections a newly added node will have.

level of polarization over 1000 runs with the same 1000 initial opinion vectors.<sup>33</sup> Then I compare across those networks the average level of polarization and provide several network statistics. The network statistics could be divided in two groups: (i) (inter-related) measures of degree inequality (ii) measures of connectivity. Furthermore, I provide the  $\alpha$  parameter, which is the exponent of the power-law that fits best<sup>34</sup> the degree distribution of the considered scale-free network topology. In figures 9 to 13 the markers correspond to different sizes of initial hubs  $h$  and colors correspond to the number  $n_c$  of connections a newly added node will have. For all those figures, the  $y$ -axis corresponds to the level of mean polarization while the  $x$ -axis gives the value of the different network statistic that will be considered.

High levels of opinion polarization are expected to be observed for more skewed degree distribution where expressers have a very high degree. Unsurprisingly, figure 9 shows that the mean level of polarization is positively correlated with the  $\alpha$  parameter and negatively correlated with the assortativity coefficient of the degree distribution. In particular, in the right panel of figure 9 the two topologies with  $c_n = 1$  and respectively with  $h = 20$  and  $h = 30$  show a high level of

distributions and levels of connectivities.

<sup>33</sup>I do so to get information about the topology of the network, that are not related to a specific distribution of a given initial opinion vector.

<sup>34</sup>Based on a maximum likelihood estimator.

polarization and disassortative mixing.<sup>35</sup> Furthermore, the right panel of figure 11 and left panel of figure 12, indicate high levels of polarization for high values of the Gini coefficient of the degree distribution and very low network sparsity computed as the per capita mean degree.

The more expressers interact together the greater should be the possibilities of disagreement and consequently opinion polarization. Surprisingly, the left panel of figure 10 shows that the lowest levels of polarization are observed for topologies where the number of components in the subgraph of  $G$  restricted to expressers is smaller than 2. This means that many expressers are connected among each other, yet polarization is low. The right panel of both figures 10 and 11 bring an explanation. For those same network topologies, the average path length and the diameter are very small (respectively smaller than 3 and smaller than 5). Hence, those network topologies with a low levels of polarization display high connectivity. Thereby, expressers interact among each other and can fall into disagreement, but consensual individuals are exposed to more influence sources.<sup>36</sup> This explanation is also supported by figure 12 and the right panel of figure 13. In the right panel of figure 13 low levels of polarization are related to higher neighborhood connectivity. Yet for these high levels of neighborhood connectivity (pink, purple and blue markers), the number of expressers is dispersed over the whole range as shown in figure 12.

---

<sup>35</sup>Following Newman (2002) [26], a network is said to show assortative mixing if the nodes in the network that have many connections tend to be connected to other nodes with many connections.

<sup>36</sup>Intuitively, for those network topologies, expressers don't have a *fan base* which is solely devoted to each of them.

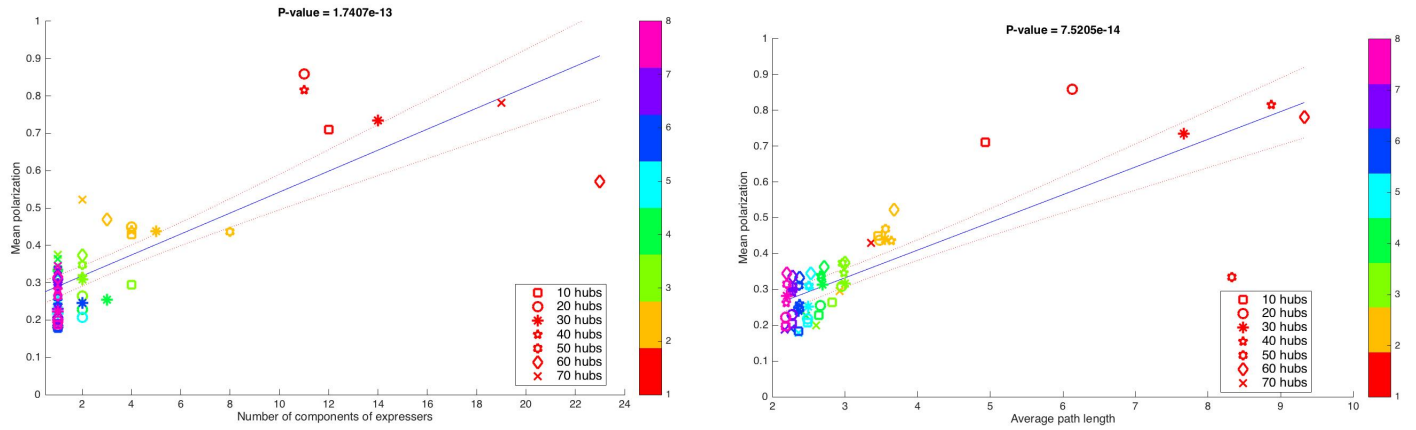


Figure 10

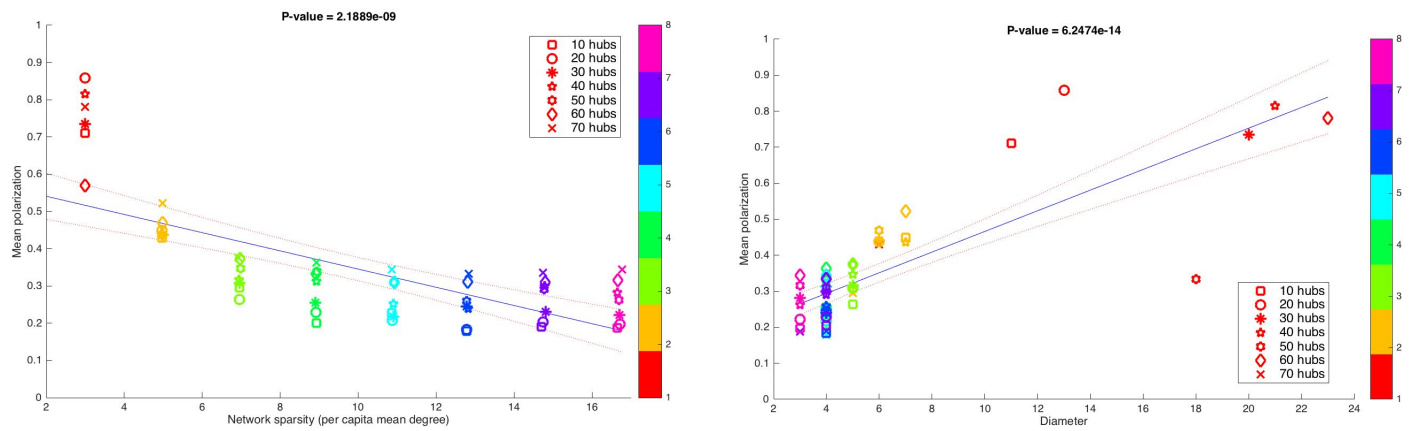


Figure 11



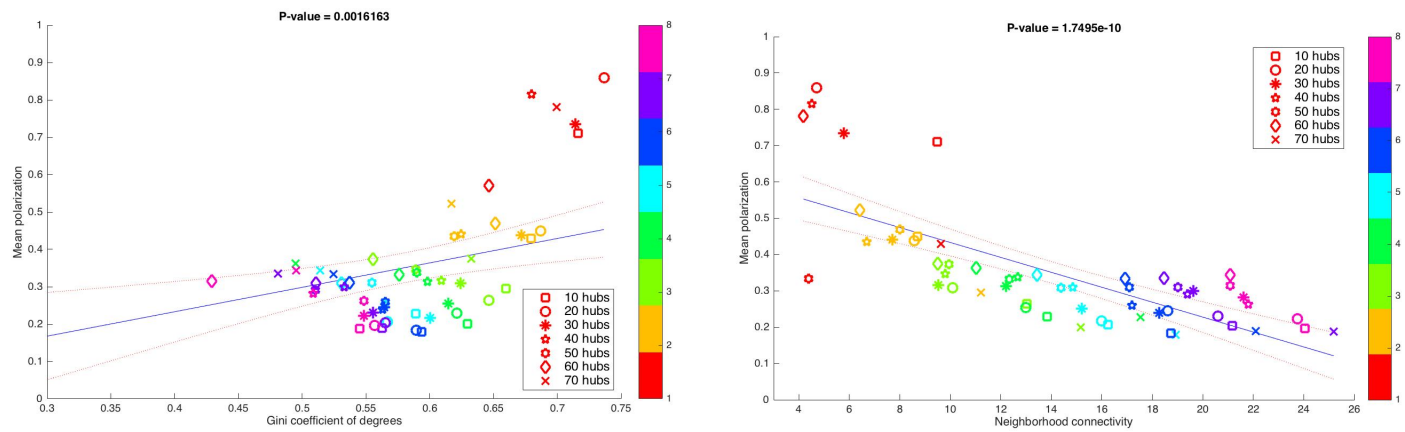


Figure 12

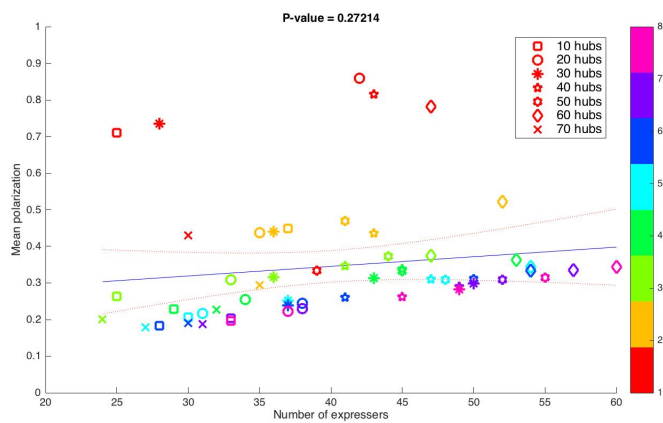


Figure 13

## 6 Conclusion and the way forward

This paper introduces several novel ingredients to classic opinion formation models with assimilative and repulsive influence. By linking opinion expression to the network topology, this paper brings the analysis of opinion dynamics a step closer to real life observed networks. It shows that opinion patterns can be explained by focusing on how influential individuals interact among each other and how they influence the (less popular) masses of users.

Introducing concepts grounded in behavioral sciences to model interactions between individuals and features of modern communication networks, is essential in understanding opinion dynamics. Namely this exercise is crucial for public policies aiming (or not!) to regulate platforms of social networking specialized in persuasive technology and that impact nearly half of the world population. A natural next step for this paper is to test the findings with real life data and consider the role of credibility and information verification. Namely, in a future research agenda I would like to study whether high centrality is related to higher credibility and whether this affects information verification by introducing a cost.

## References

- [1] Daron Acemoglu, Giacomo Como, Fabio Fagnani, and Asuman Ozdaglar. Opinion fluctuations and disagreement in social networks. *Mathematics of Operations Research*, 38(1), 2013.
- [2] Daron Acemoglu and Asuman Ozdaglar. Opinion dynamics and learning in social networks. *Dyn Games Appl*, 1:3–49, 2011.
- [3] Robert Axelrod. The dissemination of culture. a model with local convergence and global polarization. *The journal of conflict resolution*, 1977.
- [4] Abhijit Banerjee, Emily Breza, Arun g. Chandrasekhar, and Markus Mobius. Naive learning with uninformed agents. *mimeo*, 2018.
- [5] S. Banisch and E. Olbrich. Opinion polarization by learning from social feedback. *The Journal of Mathematical Sociology*, 43(2):76–103, 2019.
- [6] Albert-László Barabási and Réka Albert. Emergence of scaling in random networks. *Science*, 286(5439), 1999.
- [7] Ugo Bolletta and Paolo Pin. Polarization when people choose their peers. *mimeo*, 2019.
- [8] M.D. Conover, J. Ratkiewicz, M. Francisco, B. Gonçalves, A. Flammini, and F. Menczer. Political polarization on twitter. *Proceedings of the Fifth International AAAI conference on weblogs and social media*, 2011.
- [9] Guillaume Deffuant, David Neau, Frederic Amblard, and Gérard Weisbuch. Mixing beliefs among interacting agents. *Advances in complex systems*, 3:87–98, 2000.
- [10] Morris H. DeGroot. Reaching a consensus. *Journal of the American Statistical Association*, 69(345):118–121, 1974.
- [11] Andreas Flache, Michael Mas, Thomas Feliciani, Edmund Chattoe-Brown, Guillaume Deffuant, Sylvie Huet, and Jan Lorenz. Models of social influence: Towards the next frontiers. *Journal of Artificial Societies and Social Simulation*, 2017.
- [12] Noah Friedkin and Eugene Johnsen. Social influence and opinions. *Journal of mathematical sociology*, 15(3-4):193–205, 1990.
- [13] Noah E. Friedkin. The problem of social control and coordination of complex systems in sociology: a look at the community cleavage problem. *IEEE control systems magazine*, June, 2015.
- [14] Benjamin Golub and Evan Sadler. *The Oxford Handbook of the Economics of Networks*, chapter Learning in Social Networks. 2017.
- [15] A. Grow, A. Flache, and R. Wittek. Global diversity and local consensus in status beliefs: The role of network clustering and resistance to belief change. *Sociological Science*, 4:611–640, 2017.

- [16] Rainer Hegselmann and Ulrich Krause. Opinion dynamics and bounded confidence models, analysis, and simulation. *Journal of Artificial Societies and Social Simulation*, 5(3), 2002.
- [17] Wander Jager and Frédéric Amblard. Uniformity, bipolarization and pluriformity captured as generic stylized behavior with an agent-based simulation model of attitude change. *Computational and Mathematical Organization Theory*, 10:295–303, 2004.
- [18] J.Ding and Noah H. Rhee. On the equality of algebraic and geometric multiplicities of matrix eigenvalues. *Applied Mathematics letters*, 24, 2011.
- [19] James R. Larson Jr., Caryn Christensen, Ann S. Abbott, and Timothy M. Franz. Diagnosing groups: Charting the flow of information in medical decision-making teams. *Journal of Personality and social psychology*, 71(2):315–330, 1996.
- [20] J.R. P. French Jr. A formal theory of social power. *Psychol. Rev*, 63, 1956.
- [21] Tyll Krueger, Janusz Szubiński, and Tomasz Weron. Conformity, anticonformity and polarization of opinions: Insights from a mathematical model of opinion dynamics. *mimeo*, 2017.
- [22] Bibb Latané. The psychology of social impact. *American Psychologist*, 36(4):343–356, 1981.
- [23] Bibb Latané, Andrej Nowak, and James H. Liu. Measuring emergent social phenomena: dynamism, polarization, and clustering as order parameters of social systems. *Behavioral Science*, 39, 1994.
- [24] Fragkiskis Malliaros and Michalis Vazirgiannis. Clustering and community detection in directed networks: A survey. *Physics Reports*, 2013.
- [25] Carl D. Meyer. *Matrix Analysis and Applied linear algebra*. SIAM, 2000.
- [26] M.E.J. Newman. Assortative mixing in networks. *Phys. Rev. Lett*, 89, 2002.
- [27] Evan Sadler. Influence campaigns. *mimeo*, 2019.
- [28] E. Seneta. *Non-negative matrices and markov chains*. Springer Science and Business Media, 2006 edition, 1981.
- [29] Teresa M. Silva, Luís Silva, and Sara Fernandes. Convergence time to equilibrium distributions of autonomous and periodic non-autonomous graphs. *Linear Algebra and its Applications*, 488:199–215, 2016.
- [30] G. Stasser, L. Taylor, and C. Hanna. Information sampling in structured and unstructured discussions of three and six person groups. *Journal of Personality and social psychology*, 57:67–78, 1989.
- [31] Garold Stasser and William Titus. Pooling of unshared information in group decision making: Biased information sampling during discussion. *Journal of Personality and social psychology*, 48(6):1467–1478, 1985.
- [32] Garold Stasser and William Titus. Hidden profiles: a brief history. *Psychological Inquiry*, 14(3):304–313, 2003.
- [33] Lin Xiao and Stephen Boyd. Fast linear iterations for distributed averaging. *Systems and Control Letters*, 53, 2004.
- [34] E. Yildiz, A. Ozdaglar, D. Acemoglu, A. Saberi, and A. Scaglione. Binary opinion dynamics with stubborn agents. *ACM Trans. Econ. Comp.*, 2013.

## 7 Appendix

### 7.1 Proof of proposition 1

(i) When  $|\mathcal{E}| = 1$  it means that individual  $i \in \mathcal{E}$  has no direct neighbors who choose to express, hence individual  $i$  never updates their initial opinion and their long-run opinion is exactly their initial opinion  $\alpha_{0,i}$ .

(ii) Suppose without loss of generality that  $|\mathcal{E}| = \kappa$ . If  $|\mathcal{E}| = \kappa > 1$  and all individuals  $i \in \mathcal{E}$  are like-minded, that is  $\forall i \neq j \in \mathcal{E}, |\alpha_{0,i} - \alpha_{0,j}| < \tau$  then for  $\mu \in (0, 1/\kappa)$  the opinions get updated in the following way:

$$\begin{cases} \alpha_{1,1} = \alpha_{0,1} + \mu \sum_{j \neq 1 \in \mathcal{E}} g_{1j}(\alpha_{0,j} - \alpha_{0,1}) \\ \vdots \\ \alpha_{1,\kappa} = \alpha_{0,\kappa} + \mu \sum_{j \neq \kappa \in \mathcal{E}} g_{\kappa j}(\alpha_{0,j} - \alpha_{0,1}) \end{cases} \Leftrightarrow \begin{cases} \alpha_{1,1} = (1 - \mu d_1(\mathcal{E}))\alpha_{0,1} + \mu \sum_{j \neq 1 \in \mathcal{E}} g_{1j}\alpha_{0,j} \\ \vdots \\ \alpha_{1,\kappa} = (1 - \mu d_\kappa(\mathcal{E}))\alpha_{0,\kappa} + \mu \sum_{j \neq \kappa \in \mathcal{E}} g_{\kappa j}\alpha_{0,j} \end{cases}$$

where  $d_i(\mathcal{E}) = \sum_{j \in \mathcal{E}} g_{ij}$  corresponds to the number of expressers that are in the set of connected expressers  $\mathcal{E}$  and are also direct neighbors of individual  $i \in \mathcal{E}$ . Writing the above system in matrix notation and using induction we get the following relation :

$$\alpha_{t,\mathcal{E}} = M^t \alpha_{0,\mathcal{E}}$$

where  $\alpha_{t,\mathcal{E}} = (\alpha_{t,1}, \dots, \alpha_{t,\kappa})^T$ ,  $\alpha_{0,\mathcal{E}} = (\alpha_{0,1}, \dots, \alpha_{0,\kappa})^T$  and  $M$  an  $\kappa \times \kappa$  symmetric matrix with diagonal entries  $m_{ii} = 1 - d_i(\mathcal{E})\mu$  and off diagonal entries  $m_{ij} = \mu g_{ij}$ , for  $j \neq i \in \mathcal{E}$ . Hence,  $M$  is a symmetric matrix, with non-negative entries and whose columns and rows sum to one. In order to get the long-run opinions we need to compute  $\lim_{t \rightarrow \infty} \alpha_{t,\mathcal{E}} = \lim_{t \rightarrow \infty} M^t \alpha_{0,\mathcal{E}}$ .

**Claim 1**  $\lim_{t \rightarrow \infty} M^t$  exists.

This limit exists because all the eigenvalues of the matrix  $M$  are smaller or equal to 1. To see this, simply recall that by the Gershgorin Circle Theorem (1931), the eigenvalues of the square matrix  $M$  belong to the union of its Gershgorin disks.<sup>37</sup> write for each  $i \in \mathcal{E}$ ,  $D_i = \{x \in \mathbb{R} : |x - m_{ii}| \leq \sum_{j \neq i} |m_{ij}|\} = \{x \in \mathbb{R} : |x - (1 - d_i(\mathcal{E})\mu)| \leq d_i(\mathcal{E})\mu\}$ . Hence, the upper bound of the eigenvalues of  $M$  is given exactly by  $\max_{i \in \mathcal{E}} (1 - d_i(\mathcal{E})\mu) + d_i(\mathcal{E})\mu = 1$ . Now I will show that  $\lim_{t \rightarrow \infty} \alpha_{t,\mathcal{E}}$  is exactly the average of the initial opinions of individuals  $1, \dots, \kappa \in \mathcal{E}$ .

**Claim 2** Let  $\mathbf{1}_{p,q}$  be a matrix of ones of size  $p \times q$ .  $\lim_{t \rightarrow \infty} M^t = \frac{1}{\kappa} \mathbf{1}_{\kappa,1} \mathbf{1}_{1,\kappa}$ .

Intuitively, since at each time period every updated opinion of an expresser is a convex combination of the opinions of like-minded neighbors who also express, the long-run opinions converge to the average of initial opinions of the

<sup>37</sup>All the eigenvalues of  $M$  are real because  $M$  is a real symmetric matrix.

members of the connected set of expressers. Formally, I use theorem 1 in Xiao and Boyd (2004) [33], which states that  $\lim_{t \rightarrow \infty} M^t = \frac{1}{\kappa} \mathbf{1}_{\kappa,1} \mathbf{1}_{1,\kappa}$  if and only if (i) the vector  $\mathbf{1}$  is a left eigenvector of  $M$  associated with the eigenvalue one, (ii) the vector  $\mathbf{1}$  is a right eigenvector of  $M$  associated with the eigenvalue one, (iii) one is a simple eigenvalue of  $M$ . Conditions (i) and (ii) hold for the matrix  $M$  because it is symmetric and row stochastic. To see this, one can simply sum the entries over a given row  $i \in \mathcal{E}$ :  $m_{ii} + \sum_{j \neq i \in \mathcal{E}} m_{ij} = 1 - d_i(\mathcal{E})\mu + \sum_{j \neq i \in \mathcal{E}} g_{ij}\mu = 1 - d_i(\mathcal{E})\mu + d_i(\mathcal{E})\mu = 1$ . Since the matrix  $M$  is symmetric, it is also column stochastic and the vector one is a left and right eigenvector of the matrix  $M$  associated with the eigenvalue one. Finally, condition (iii) holds because the matrix  $M$  is irreducible with non-negative entries; because the set of individuals in  $\mathcal{E}$  is connected and they are all like-minded, in the sense of definition 3. Hence the eigenvalue 1 is simple (Perron-Frobenius Theorem).

## 7.2 Proof of proposition 2

**Warm-up for the proof and notations.** Recall that  $\underline{N}_{i,t-1} = \{j \in N_i \cap E \text{ s.t. } |\alpha_{i,t-1} - \alpha_{j,t-1}| < \tau\}$  is the set of expressing neighbors of individual  $i \in N$  such that their opinion difference is smaller than  $\tau$  (like-minded). Similarly, recall that  $\overline{N}_{i,t-1} = \{j \in N_i \cap E \text{ s.t. } |\alpha_{i,t-1} - \alpha_{j,t-1}| \geq \tau\}$ . Clearly it follows that  $\underline{N}_{i,t-1} \cup \overline{N}_{i,t-1} = N_i \cap E$ . Without loss of generality, suppose that individual  $i \in N$  chooses to express and belongs to the set of connected expressers  $\mathcal{E} \subseteq E$  such that  $|\mathcal{E}| > 1$  and there exists at least one pair of expressers in  $\mathcal{E}$  who are neighbors and initially ideologically-opposed. Each individual  $i \in \mathcal{E}$  updates their opinion at each time step according to the law of motion (3), given by :

$$\alpha_{i,t} = \alpha_{i,t-1} + \mu \sum_{j \in \underline{N}_{i,t-1}} (\alpha_{j,t-1} - \alpha_{i,t-1}) + \mu \sum_{j \in \overline{N}_{i,t-1}} (\alpha_{i,t-1} - \alpha_{j,t-1}) \quad \text{s.t } \alpha_{i,t} \in [-1, 1] \quad (6)$$

$$\Leftrightarrow \alpha_{i,t} = \alpha_{i,t-1}(1 + \mu(|\overline{N}_{i,t-1}| - |\underline{N}_{i,t-1}|)) + \mu \left( \sum_{j \in \underline{N}_{i,t-1}} \alpha_{j,t-1} - \sum_{j \in \overline{N}_{i,t-1}} \alpha_{j,t-1} \right) \quad \text{s.t } \alpha_{i,t} \in [-1, 1] \quad (7)$$

The size of the sets  $\overline{N}_{i,t-1}$  and  $\underline{N}_{i,t-1}$  can vary between two periods because a like-minded neighbor at period  $t-1$  can become ideologically-opposed at a subsequent period (see the examples in section 4.1). In other words, it's impossible to summarize the above system of equations by one time-invariant matrix because the entries or weights vary with time depending on the opinion differences between connected expressers and whether the opinion of a given individual has reached the upper or lower bound. Hence we need to account for the *flow* of influence through *chains* of expressers within a given set  $\mathcal{E}$ . To do so, it's convenient to write the opinion of a given individual  $i$  at period  $t$  belonging to the connected set of expressers  $\mathcal{E}$  as an **affine** combination of the opinions their neighbors  $j$  at period  $t-1$ :

$$\alpha_{i,t} = \sum_{j \in N_i \cap \mathcal{E}} a_{ij}(t-1) \alpha_{j,t-1} \quad (8)$$

where

$$a_{ij}(t-1) = \begin{cases} \mu & \text{if } i \neq j, |\alpha_{i,t-1} - \alpha_{j,t-1}| < \tau \text{ and } \alpha_{i,t-1} \in (-1, 1) \\ -\mu & \text{if } i \neq j, |\alpha_{i,t-1} - \alpha_{j,t-1}| \geq \tau \text{ and } \alpha_{i,t-1} \in (-1, 1) \\ 1 + \mu(|\overline{N}_{i,t-1}| - |\underline{N}_{i,t-1}|) & \text{if } i = j \text{ and } \alpha_{i,t-1} \in (-1, 1) \\ 0 & \text{if } i \neq j \text{ and } \alpha_{i,t-1} \in \{1 + \epsilon, -1 - \epsilon\} \\ 1 & \text{if } i = j \text{ and } \alpha_{i,t-1} \in \{1 + \epsilon, -1 - \epsilon\} \end{cases}$$

Let  $A(t-1)$  be the matrix with entries  $a_{ij}(t-1)$  for  $i, j \in \mathcal{E}$ . It follows that the opinions at period  $t$  of the expressers who belong to  $\mathcal{E}$  can be written as:

$$\alpha_{\mathcal{E},t} = A(t-1)A(t-2) \dots A(0)\alpha_{\mathcal{E},0} = B(t-1,0)\alpha_{\mathcal{E},0}$$

where  $B(t-1,0)$  is the matrix product of  $A(t-1)A(t-2) \dots A(0)$ . In other words it's a matrix that keeps track of the accumulated (positive and negative) weights between periods  $t-1$  and 0. In particular, the entry  $B_{ij}(t, t-1)$  reports the influence of  $j$  on  $i$ 's opinion between periods  $t$  and  $t-1$ . Recall that the set of expressers that have at least one ideologically-opposed neighbor is:

$$IO(\mathcal{E}) = \{i \neq j \in \mathcal{E}, i \in \overline{N}_{j,0} \neq \emptyset, j \in \overline{N}_{i,0} \neq \emptyset \text{ and } g_{ij} = 1\}$$

**Claim 3** *Let  $i_1, i_2 \dots i_k \in IO(\mathcal{E})$ . If there exists at each  $t \geq 0$  paths of expressers with only like-minded neighbors connecting  $i \notin IO(\mathcal{E})$  to  $i_1, i_2 \dots i_k$  then  $\alpha_{i,\infty} \in \text{conv}(\alpha_{i_1,\infty}, \dots, \alpha_{i_k,\infty})$ .*

**Idea of the proof.** Show that if  $i$  is connected indirectly to expressers with ideologically-opposed neighbors through a path of like-minded neighbors then  $\alpha_{i,t}$  can be written as a convex combination of the opinions of neighbors for all  $t \geq 1$ . In this case individual  $i$  holds in the long-run a moderate opinion  $\alpha_{i,\infty} \in (-1, 1)$ . Otherwise, their opinion keeps getting pushed to the upper or lower bound of the opinion interval and necessarily  $\alpha_{i,\infty} \in \{-1, 1\}$ .

**Proof.** Suppose that  $IO(\mathcal{E}) \neq \emptyset$ , that is there exists at least one pair of neighbors in  $\mathcal{E}$  that are initially ideologically-opposed. Let  $s$  be the time period such that for all  $i_k \in IO(\mathcal{E})$ ,  $a_{i_k j}(s) = 0$  for all  $j \neq i_k$ ,  $a_{i_k j}(s) = 1$  for  $i_k = j$ .

( $\Rightarrow$ ) Suppose that there exists at each  $t \geq 0$  paths of expressers with only like-minded neighbors connecting  $i \notin IO(\mathcal{E})$  to at least two individuals  $i_k \neq i_j \in IO(\mathcal{E})$ . That is there exists  $i, j_1, j_2 \dots, j_k \notin IO(\mathcal{E})$  such that for all  $t \geq s$ ,  $B_{i,i_k}(t, s) > 0$  and  $B_{i,i_j}(t, s) > 0$ . Since,

1. the opinion of a given individual  $j_k \in \{j_1, j_2, \dots\} \cup \{i\} \notin IO(\mathcal{E})$  at period  $t \geq s$  is a **convex** combination of the opinions of their neighbors at period  $t-1$  (because they all have like-minded neighbors):

$$\alpha_{j_k,t} = \sum_{k \in N_i \cap \mathcal{E}} a_{j_k k}(t-1)\alpha_{k,t-1}$$

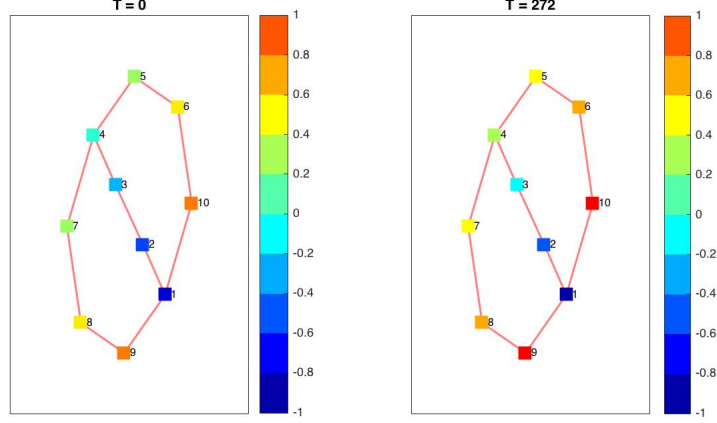


Figure 14

where  $a_{j_k k}$  takes the value  $\mu$  for any  $k \neq j_k$  and  $1 - |\underline{N}_{j_k, t}| \mu$  for  $k = j_k$ .

2. and at time period  $t \geq s$ ,  $\forall i_k \in IO(\mathcal{E})$ ,  $\alpha_{i_k, t} = \alpha_{i_k, s} \in \{-1, 1\}$ ,
3. it follows that  $\forall j_k \in \{j_1, j_2, \dots\} \cup \{i\} \notin IO(\mathcal{E})$ ,  $\alpha_{j_k, t} \in \text{conv}(\alpha_{i_1, s}, \dots, \alpha_{i_k, s})$ .
4. Moreover, since  $IO(\mathcal{E}) \neq \emptyset$  and  $\forall j_k \in \{j_1, j_2, \dots\} \cup \{i\} \notin IO(\mathcal{E})$  are connected (indirectly) to at least two individuals in  $IO(\mathcal{E})$  then  $\alpha_{j_k, t} \in \text{conv}(\alpha_{i_1, s}, \dots, \alpha_{i_k, s})$  and  $\alpha_{i, t} \notin \{-1, 1\}$ .
5.  $\forall t \geq 0$ ,  $\forall j_k \in \{j_1, j_2, \dots\} \cup \{i\} \notin IO(\mathcal{E})$  have like-minded neighbors at each period, hence the argument extends to the limit and  $\alpha_{j_k, t} = \lim_{t \rightarrow \infty} \alpha_{j_k, t} \in \text{conv}(\alpha_{i_1, s}, \dots, \alpha_{i_k, s}) = \text{conv}(\alpha_{i_1, \infty}, \dots, \alpha_{i_k, \infty})$ .

■

Before moving to the proof of the next claim, in figure 14, I provide an example that gives an intuition for the previous proof. I set  $\delta^*$  such that all individuals can express for the sake of the example and expressers are hence represented by a square. The colors of each node correspond to their opinion at period  $T$  indicated by the color map over  $[-1, 1]$  on the east side of each figure. The set of connected expressers is  $\mathcal{E} = E = N$  and  $IO(\mathcal{E}) = \{1, 9, 10\}$ . Individuals  $\{2, \dots, 8\}$  all have like-minded neighbors, as indicated by the colors in the left panel of figure 14 and they are all connected to the expressers in  $IO(\mathcal{E})$  by three paths of like-minded expressers. The right panel indicates the long-run opinions and we see that individuals  $\{2, \dots, 8\}$  still have like-minded neighbors. In particular, their long-run opinion is a convex combination of the opinions of the individuals in  $IO(\mathcal{E}) = \{1, 9, 10\}$ . Individual 4 is exactly at distance 3 from each of the three individuals in  $IO(\mathcal{E}) = \{1, 9, 10\}$  and their long-run opinion is  $\frac{1}{3}\alpha_{1, \infty} + \frac{1}{3}\alpha_{9, \infty} + \frac{1}{3}\alpha_{10, \infty} = \frac{1}{3}(-1) + \frac{1}{3}(1) + \frac{1}{3}(1) = \frac{1}{3}$ . Individuals in  $IO(\mathcal{E}) = \{1, 9, 10\}$  reach the upper and lower bound of the opinion period just after one period of interaction.

**Claim 4** Let  $i_1, i_2 \dots i_k \in IO(\mathcal{E})$ . If there exists  $i \in IO(\mathcal{E})$  such that (i)  $\alpha_{i, 0} = 0$ , (ii)  $\forall t \geq 0$   $\underline{N}_{i, t} = \emptyset$ , (iii) and



$\sum_{j \in \overline{N}_{i,t}} \alpha_{j,t} = 0$  then  $\alpha_{i,\infty} = 0$ .

**Idea of the proof.** Show that when a given individual  $i$  is initially neutral (opinion zero) with no like-minded friends, if she is repulsed by two neighbors who are themselves ideologically-opposed and the negative influence she receives at each period is exactly of the same magnitude. The simplest example is a circle with three individuals 1,2,3 where  $\alpha_{1,0} = 0$ ,  $\alpha_{2,0} = -0.8$  and  $\alpha_{3,0} = 0.8$ . Looking at expression (7), the opinions get updated in period  $t = 1$  as follows:

$$\begin{aligned}\alpha_{1,1} &= 0 + \mu(-(-0.8) - (0.8)) = 0 \\ \alpha_{2,1} &= -0.8 + \mu(-0 - (0.8)) = -0.8 - \mu 0.8 \\ \alpha_{3,1} &= 0.8 + \mu(-0 - (-0.8)) = 0.8 + \mu 0.8\end{aligned}$$

In subsequent periods, the opinion of individual 1 remains zero because  $\mu(-\alpha_{2,t} - \alpha_{3,t}) = 0$  for all  $t \geq 0$ . However, if  $\alpha_{2,0} \neq -\alpha_{3,0}$  then the opinion of individual 1 will be pushed to the upper or lower bound after a certain number of periods depending on which of  $|\alpha_{2,0}|$  and  $|\alpha_{3,0}|$  is larger.

**Proof.** Let  $i \in IO(\mathcal{E})$  such that  $\alpha_{i,0} = 0$  and  $\forall t \geq 0$ ,  $\overline{N}_{i,t} = \emptyset$ . Since individual  $i$  has no like-minded neighbors,  $a_{i,k}(t) = -\mu$  for all  $k \in \overline{N}_{i,0} = N_i$  and  $\forall t \geq 0$ . Their opinion gets updated as follows:

$$\alpha_{i,1} = 0 - \mu \sum_{j \in \overline{N}_{i,0}} \alpha_{j,0}$$

The opinion of  $i$  is zero at period 1 if and only if  $\sum_{j \in \overline{N}_{i,0}} \alpha_{j,0} = 0$ . Since, all the neighbors are in  $\overline{N}_{i,0}$  and  $\alpha_{i,0} = 0$  then even for an opinion threshold  $\tau > \epsilon$ ,  $\alpha_{j,0} \neq 0$ . Hence,  $\sum_{j \in \overline{N}_{i,0}} \alpha_{j,0} = 0$  if and only if the opinions  $\alpha_{j,0}$  for all  $j \in \overline{N}_{i,0}$  cancel out. By induction, the argument holds for all time periods.

Notice that if at any given period  $t$  this sum doesn't cancel out then the opinion of  $i$  will be different then zero in all subsequent period until it reaches the upper or lower bound. ■

### 7.3 Proof of lemma 1

Case 1:  $|\alpha_{i,0} - \alpha_{j,0}| < \tau$ . The law of motion 3 rewrites:

$$\begin{cases} \alpha_{i,t} = \alpha_{i,t-1} + \mu(\alpha_{j,t-1} - \alpha_{i,t-1}) \\ \alpha_{j,t} = \alpha_{j,t-1} + \mu(\alpha_{i,t-1} - \alpha_{j,t-1}) \end{cases} \Leftrightarrow \begin{cases} \alpha_{i,t} = (1 - \mu)\alpha_{i,t-1} + \mu\alpha_{j,t-1} \\ \alpha_{j,t} = (1 - \mu)\alpha_{j,t-1} + \mu\alpha_{i,t-1} \end{cases}$$

We can write the above system in matrix notation:

$$\begin{bmatrix} \alpha_{i,t} \\ \alpha_{j,t} \end{bmatrix} = \begin{bmatrix} 1-\mu & \mu \\ \mu & 1-\mu \end{bmatrix} \begin{bmatrix} \alpha_{i,t-1} \\ \alpha_{j,t-1} \end{bmatrix} \Leftrightarrow \begin{bmatrix} \alpha_{i,t} \\ \alpha_{j,t} \end{bmatrix} = \overbrace{\begin{bmatrix} 1-\mu & \mu \\ \mu & 1-\mu \end{bmatrix}^t}^{=M^t} \begin{bmatrix} \alpha_{i,0} \\ \alpha_{j,0} \end{bmatrix} \quad (\text{by induction})$$

Moreover, we can diagonalize the matrix  $M^t$  so that we can compute the limit easily:

$$M^t = \begin{bmatrix} 1-\mu & \mu \\ \mu & 1-\mu \end{bmatrix}^t = \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1^t & 0 \\ 0 & (1-2\mu)^t \end{bmatrix} \begin{bmatrix} 1/2 & 1/2 \\ -1/2 & 1/2 \end{bmatrix}$$

For  $\mu \in (0, 1/2)$ ,  $\lim_{t \rightarrow \infty} (1-2\mu)^t = 0$ . Notice that this is equivalent to upper bounding the distance between opinions at a given period  $t$  and the limiting opinions by the second highest eigenvalue.<sup>38</sup> It follows that when the opinions of  $i$  and  $j$  are close enough then they converge exactly to their average:

$$\alpha_\infty^a = \lim_{t \rightarrow \infty} \begin{bmatrix} \alpha_{i,t} \\ \alpha_{j,t} \end{bmatrix} = \begin{bmatrix} 1/2 & 1/2 \\ 1/2 & 1/2 \end{bmatrix} \begin{bmatrix} \alpha_{i,0} \\ \alpha_{j,0} \end{bmatrix} = \begin{bmatrix} \frac{\alpha_{i,0} + \alpha_{j,0}}{2} \\ \frac{\alpha_{i,0} + \alpha_{j,0}}{2} \end{bmatrix}$$

For  $\epsilon > 0$ , the time  $t_a$  it takes to reach convergence is:  $t_a \geq \frac{\log(\epsilon)}{\log(1-2\mu)}$ .

**Case 2:**  $|\alpha_{i,0} - \alpha_{j,0}| \geq \tau$ . The law of motion 3 rewrites:

$$\begin{cases} \alpha_{i,t} = \alpha_{i,t-1} + \mu(\alpha_{i,t-1} - \alpha_{j,t-1}) \\ \alpha_{j,t} = \alpha_{j,t-1} + \mu(\alpha_{j,t-1} - \alpha_{i,t-1}) \end{cases} \Leftrightarrow \begin{cases} \alpha_{i,t} = (1+\mu)\alpha_{i,t-1} - \mu\alpha_{j,t-1} \\ \alpha_{j,t} = (1+\mu)\alpha_{j,t-1} - \mu\alpha_{i,t-1} \end{cases}$$

We can write the above system in matrix notation:

$$\begin{bmatrix} \alpha_{i,t} \\ \alpha_{j,t} \end{bmatrix} = \begin{bmatrix} 1+\mu & -\mu \\ -\mu & 1+\mu \end{bmatrix} \begin{bmatrix} \alpha_{i,t-1} \\ \alpha_{j,t-1} \end{bmatrix} \Leftrightarrow \begin{bmatrix} \alpha_{i,t} \\ \alpha_{j,t} \end{bmatrix} = \overbrace{\begin{bmatrix} 1+\mu & -\mu \\ -\mu & 1+\mu \end{bmatrix}^t}^{=M^t} \begin{bmatrix} \alpha_{i,0} \\ \alpha_{j,0} \end{bmatrix} \quad (\text{by induction})$$

Moreover, we can diagonalize the matrix  $M^t$  :

$$M^t = \begin{bmatrix} 1+\mu & -\mu \\ -\mu & 1+\mu \end{bmatrix}^t = \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1^t & 0 \\ 0 & (1+2\mu)^t \end{bmatrix} \begin{bmatrix} 1/2 & 1/2 \\ -1/2 & 1/2 \end{bmatrix}$$

The limit opinions of  $i$  and  $j$  are:

$$\alpha_\infty^r \lim_{t \rightarrow \infty} \begin{bmatrix} \alpha_{i,t} \\ \alpha_{j,t} \end{bmatrix} = \lim_{t \rightarrow \infty} \frac{1}{2} \begin{bmatrix} 1 + (1+2\mu)^t & 1 - (1+2\mu)^t \\ 1 - (1+2\mu)^t & 1 + (1+2\mu)^t \end{bmatrix} \begin{bmatrix} \alpha_{i,0} \\ \alpha_{j,0} \end{bmatrix} = \lim_{t \rightarrow \infty} \frac{1}{2} \begin{bmatrix} \alpha_{i,0} + \alpha_{j,0} + (\alpha_{i,0} - \alpha_{j,0})(1+2\mu)^t \\ \alpha_{i,0} + \alpha_{j,0} + (\alpha_{j,0} - \alpha_{i,0})(1+2\mu)^t \end{bmatrix}$$

For any positive  $\mu$  this limit explodes. However, recall that opinions have an upper 1 and lower bound  $-1$ . It follows that when the opinions of  $i$  and  $j$  are faraway they diverge until they reach the upper and lower limit of opinions.

<sup>38</sup>For more details on this topic in linear algebra See Silva, Silva and Fernandes (2016) [29].

Moreover, there exists a time  $t$  for a given  $\mu > 0$  such that that we remain within the permitted bounds. To find this time  $t$  given  $\mu$ , we must solve:

$$\begin{cases} \frac{1}{2}(\alpha_{i,0} + \alpha_{j,0} + (\alpha_{i,0} - \alpha_{j,0})(1 + 2\mu)^t) = 1 & \text{if } \alpha_{i,0} > \alpha_{j,0} \\ \frac{1}{2}(\alpha_{i,0} + \alpha_{j,0} + (\alpha_{i,0} - \alpha_{j,0})(1 + 2\mu)^t) = -1 & \text{if } \alpha_{i,0} < \alpha_{j,0} \end{cases}$$

Given  $\mu$ , we get the following  $t_r$  (for integer values take the floor function):

$$t_r = \begin{cases} \frac{\log\left(\frac{2-\alpha_{i,0}-\alpha_{j,0}}{\alpha_{i,0}-\alpha_{j,0}}\right)}{\log(1+2\mu)} & \text{if } 1 \geq \alpha_{i,0} > \alpha_{j,0} \geq -1 \\ \frac{\log\left(\frac{-2-\alpha_{i,0}-\alpha_{j,0}}{\alpha_{i,0}-\alpha_{j,0}}\right)}{\log(1+2\mu)} & \text{if } -1 \leq \alpha_{i,0} < \alpha_{j,0} \leq 1 \end{cases}$$

For very small  $\epsilon$  and  $\mu \in (0, 1/2)$ , it takes a very large number of periods to reach full consensus while to reach 1 an  $-1$  the individuals take a finite number of time periods. In other words,  $t_r < t_a$  because we can always find a small enough  $\epsilon$  such that the inequality holds. Formally, we solve the inequality  $t_a > t_r$  for  $\epsilon > 0$ , for the case where  $\alpha_{i,0} > \alpha_{j,0}$  (similarly for the other case) and  $t_a$  at its lower bound:

$$\begin{aligned} \frac{\log(\epsilon)}{\log(1-2\mu)} &> \frac{\log\left(\frac{2-\alpha_{i,0}-\alpha_{j,0}}{\alpha_{i,0}-\alpha_{j,0}}\right)}{\log(1+2\mu)} \\ \Leftrightarrow \epsilon &< \exp\left(\frac{\log\left(\frac{2-\alpha_{i,0}-\alpha_{j,0}}{\alpha_{i,0}-\alpha_{j,0}}\right) \log(1-2\mu)}{\log(1+2\mu)}\right) \end{aligned}$$

## 7.4 Proof of theorem 1

Part (i): let  $\lambda$  be an eigenvalue of the matrix  $\tilde{G}$ . Recall that the algebraic multiplicity of  $\lambda$  is the number of times it is repeated as a root of the characteristic polynomial and the geometric multiplicity of  $\lambda$  is the maximum number of linearly independent eigenvectors associated with  $\lambda$ . An eigenvalue is semi-simple if its algebraic multiplicity is equal to its geometric multiplicity (definitions p.510, chapter 7, Meyer (2000) [25]). For  $\tilde{G} \in \mathbb{R}^{n \times n}$ ,  $\lim_{t \rightarrow \infty} \tilde{G}^t$  exists if and only if  $\rho(\tilde{G}) < 1$  (the spectral radius) or else  $\rho(\tilde{G}) = 1$  where  $\lambda = 1$  is the only eigenvalue on the unit circle and  $\lambda = 1$  is semi-simple (see *Limits of Powers* page 630, chapter 7, in Meyer (2000) [25]). Moreover, for every stochastic matrix, the spectral radius is 1 and it is semi-simple (p.696, Chapter 8 in Meyer (2000) [25] or see Corollary 2, page 2214, in Ding and Rhee (2011) [18]). Finally, since  $\tilde{G}$  has strictly positive diagonal entries, by the Gershgorin circle theorem  $\lambda = 1$  is the only eigenvalue on the unit circle.

Part (ii): when  $\lim_{t \rightarrow \infty} \tilde{G}^t$  exists, it is equal to the spectral projector associated with eigenvalue 1 (again see p.630, chapter 7, in Meyer (2000) [25]). To see this (the below details appear on p.629 of Meyer (2000) [25]), recall that  $\tilde{G}^t$  is row stochastic hence its Jordan form has the following structure :

$$J = \begin{bmatrix} I_{p \times p} & \mathbf{0} \\ \mathbf{0} & \mathbf{K} \end{bmatrix}$$

where  $I_{p \times p}$  is the identity matrix of size  $p$ , with  $p$  the algebraic multiplicity of the eigenvalue 1 and  $\mathbf{K}$  a diagonal matrix with entries corresponding to remaining eigenvalues which are strictly smaller than 1. Hence,  $\tilde{G}_\theta^t = PJ^tP^{-1}$ . Now write  $P = (P_1, P_2)$  where  $P_1$  are the columns that correspond to the eigenvectors associated with the eigenvalues 1 and  $P_2$  are the columns that correspond to the eigenvectors associated with the remaining eigenvalues which are strictly smaller than 1. Similarly  $P^{-1} = Q = (Q_1; Q_2)$  with  $Q_1$  the lines associated with the eigenvalues 1. Since  $K^t$  vanishes when  $t$  is large because all the diagonal entries are strictly smaller than one,  $\lim_{t \rightarrow \infty} \tilde{G}_\theta^t = P_1Q_1$  which is the spectral projector of the eigenvalue 1.

## 7.5 Proof of proposition 3

**Case (i):**  $|E| = 1$ . Suppose that  $E = \{i\}$  and that the initial opinion of this individual is  $\alpha_{0,i}$ . Then using proposition 1, this expresser remains stubborn forever because she does not have neighbors who also express. Hence they never updates their opinion and  $\alpha_{\infty,i} = \alpha_{0,i}$ . Moreover from theorem 1 the long-run opinion of consensual individuals is a convex combination of opinions of expressers. Since there is only one expresser then all consensual individuals have long-run opinion of  $\alpha_{0,i}$  and consensus prevails.

**Case (ii):**  $|E| > 1$ . Long-run opinions form consensus when  $\forall i \neq j \in N$ ,  $|\alpha_{\infty,i} - \alpha_{\infty,j}| < \tau$ .

1. Let  $E = \bigcup_{k=1}^{\kappa} \mathcal{E}_k$  be the set of expressers such that  $\kappa \geq 1$ .
2. For  $\kappa = 1$  there is a unique set of connected expressers  $E = \mathcal{E}_1$  and long-run opinions form a consensus if and only if  $\forall i \neq j \in \mathcal{E}_1$ , such that  $g_{ij} = 1$ ,  $|\alpha_{0,i} - \alpha_{0,j}| < \tau$ . Since each member within  $\mathcal{E}_1$  has like-minded neighbors, the members of  $\mathcal{E}_1$  converge to the average of their initial opinions as shown in proposition 1. Moreover, using theorem 1, the opinions of consensual individuals are convex combinations of the opinions of expressers. Since here there is only one set of connected expressers, the opinion of each consensual individual is exactly the average of opinions of the members of the set of connected expressers  $\mathcal{E}_1$ .
3. For  $\kappa > 1$ , without loss of generality suppose that the union of the two connected set of expressers  $\mathcal{E}_1$  and  $\mathcal{E}_2$  is equal to  $E$ . Long-run opinions form consensus if and only if (i)  $\forall i \neq j \in \mathcal{E}_k$  for  $k \in \{1, 2\}$ , such that  $g_{ij} = 1$ ,  $|\alpha_{0,i} - \alpha_{0,j}| < \tau$  (ii)  $|\bar{\alpha}_{0,\mathcal{E}_1} - \bar{\alpha}_{0,\mathcal{E}_2}| < \tau$  where  $\bar{\alpha}_{0,\mathcal{E}_k}$  is the average of initial opinions within the set  $\mathcal{E}_k$  for  $k \in \{1, 2\}$ . Since within each of both sets all members have like-minded neighbors within  $\mathcal{E}_1$  and  $\mathcal{E}_2$  opinions of expressers converge respectively to  $\bar{\alpha}_{0,\mathcal{E}_1}$  and  $\bar{\alpha}_{0,\mathcal{E}_2}$ . Moreover, consensus can prevail if and only if  $|\bar{\alpha}_{0,\mathcal{E}_1} - \bar{\alpha}_{0,\mathcal{E}_2}| < \tau$  because then any long-run opinion of a consensual individual  $i \in N \setminus E$  is at a distance

of at most  $\tau$  from any other long-run opinion of other individuals in the network. Formally,

$$\begin{aligned}
\alpha_{\infty_i} &= \sum_{j \in \mathcal{E}_1 \cup \mathcal{E}_2} \mathcal{G}_{ij} \alpha_{j,\infty} = \sum_{j_1 \in \mathcal{E}_1} \mathcal{G}_{ij_1} \bar{\alpha}_{0,\mathcal{E}_1} + \sum_{j_2 \in \mathcal{E}_2} \mathcal{G}_{ij_2} \bar{\alpha}_{0,\mathcal{E}_2} \\
&= \bar{\alpha}_{0,\mathcal{E}_1} \sum_{j_1 \in \mathcal{E}_1} \mathcal{G}_{ij_1} + \bar{\alpha}_{0,\mathcal{E}_2} \sum_{j_2 \in \mathcal{E}_2} \mathcal{G}_{ij_2} \\
&= \bar{\alpha}_{0,\mathcal{E}_1} \sum_{j_1 \in \mathcal{E}_1} \mathcal{G}_{ij_1} + \bar{\alpha}_{0,\mathcal{E}_2} (1 - \sum_{j_1 \in \mathcal{E}_1} \mathcal{G}_{ij_1}) \\
&= \bar{\alpha}_{0,\mathcal{E}_2} + (\bar{\alpha}_{0,\mathcal{E}_1} - \bar{\alpha}_{0,\mathcal{E}_2}) \sum_{j_1 \in \mathcal{E}_1} \mathcal{G}_{ij_1}
\end{aligned}$$

Hence for any expresser in  $\mathcal{E}_2$  the opinion difference with a given consensual individual  $i \in N \setminus E$  is at most  $\tau$  (similarly for any expresser in  $\mathcal{E}_1$ ) :

$$|\alpha_{\infty_i} - \bar{\alpha}_{0,\mathcal{E}_2}| = |(\bar{\alpha}_{0,\mathcal{E}_1} - \bar{\alpha}_{0,\mathcal{E}_2}) \sum_{j_1 \in \mathcal{E}_1} \mathcal{G}_{ij_1}| \leq |(\bar{\alpha}_{0,\mathcal{E}_1} - \bar{\alpha}_{0,\mathcal{E}_2})| < \tau$$

Furthermore, for two consensual individuals  $i \neq j \in N \setminus E$  their long-run is at most  $\tau$  because:

$$\begin{aligned}
|\alpha_{\infty_i} - \alpha_{\infty_j}| &= |(\bar{\alpha}_{0,\mathcal{E}_1} - \bar{\alpha}_{0,\mathcal{E}_2}) (\sum_{j_1 \in \mathcal{E}_1} \mathcal{G}_{ij_1} - \sum_{j_1 \in \mathcal{E}_1} \mathcal{G}_{jj_1})| \leq |(\bar{\alpha}_{0,\mathcal{E}_1} - \bar{\alpha}_{0,\mathcal{E}_2})| (1 - \sum_{j_1 \in \mathcal{E}_1} \mathcal{G}_{jj_1}) \\
&\leq |\bar{\alpha}_{0,\mathcal{E}_1} - \bar{\alpha}_{0,\mathcal{E}_2}| \\
&< \tau
\end{aligned}$$

The arguments easily extend for more than 2 sets of connected expressers. To see this, think of three sets of connected expressers  $\mathcal{E}_1$ ,  $\mathcal{E}_2$  and  $\mathcal{E}_3$  where the long-run opinions within each set are respectively  $\bar{\alpha}_{0,\mathcal{E}_1}$ ,  $\bar{\alpha}_{0,\mathcal{E}_2}$  and  $\bar{\alpha}_{0,\mathcal{E}_3}$ . For all  $i \in N \setminus E$ ,  $\alpha_{i,\infty} \in \text{conv}(\bar{\alpha}_{0,\mathcal{E}_1}, \bar{\alpha}_{0,\mathcal{E}_2}, \bar{\alpha}_{0,\mathcal{E}_3})$  and  $|\bar{\alpha}_{0,\mathcal{E}_2} - \bar{\alpha}_{0,\mathcal{E}_1}| < \tau$ ,  $|\bar{\alpha}_{0,\mathcal{E}_3} - \bar{\alpha}_{0,\mathcal{E}_1}| < \tau$ ,  $|\bar{\alpha}_{0,\mathcal{E}_2} - \bar{\alpha}_{0,\mathcal{E}_3}| < \tau$ , it follows that :

$$\begin{aligned}
|\alpha_{i,\infty} - \bar{\alpha}_{0,\mathcal{E}_1}| &= \left| \sum_{j_1 \in \mathcal{E}_1} \mathcal{G}_{ij_1} \bar{\alpha}_{0,\mathcal{E}_1} + \sum_{j_2 \in \mathcal{E}_2} \mathcal{G}_{ij_2} \bar{\alpha}_{0,\mathcal{E}_2} + \sum_{j_3 \in \mathcal{E}_3} \mathcal{G}_{ij_3} \bar{\alpha}_{0,\mathcal{E}_3} - \overbrace{\left( \sum_{j_1 \in \mathcal{E}_1} \mathcal{G}_{ij_1} + \sum_{j_2 \in \mathcal{E}_2} \mathcal{G}_{ij_2} + \sum_{j_3 \in \mathcal{E}_3} \mathcal{G}_{ij_3} \right)}^{=1} \bar{\alpha}_{0,\mathcal{E}_1} \right| \\
&= \left| \sum_{j_2 \in \mathcal{E}_2} \mathcal{G}_{ij_2} \bar{\alpha}_{0,\mathcal{E}_2} + \sum_{j_3 \in \mathcal{E}_3} \mathcal{G}_{ij_3} \bar{\alpha}_{0,\mathcal{E}_3} - \sum_{j_2 \in \mathcal{E}_2} \mathcal{G}_{ij_2} \bar{\alpha}_{0,\mathcal{E}_1} - \sum_{j_3 \in \mathcal{E}_3} \mathcal{G}_{ij_3} \bar{\alpha}_{0,\mathcal{E}_1} \right| \\
&= \left| \sum_{j_2 \in \mathcal{E}_2} \mathcal{G}_{ij_2} (\bar{\alpha}_{0,\mathcal{E}_2} - \bar{\alpha}_{0,\mathcal{E}_1}) + \sum_{j_3 \in \mathcal{E}_3} \mathcal{G}_{ij_3} (\bar{\alpha}_{0,\mathcal{E}_3} - \bar{\alpha}_{0,\mathcal{E}_1}) \right| \\
&< \sum_{j_2 \in \mathcal{E}_2} \mathcal{G}_{ij_2} \tau + \sum_{j_3 \in \mathcal{E}_3} \mathcal{G}_{ij_3} \tau \\
&\leq \tau
\end{aligned}$$

## 7.6 Proof of lemma 2

Recall that from theorem 1 that the limit of  $\tilde{G}^t$  when  $t$  is large, exists and is given by  $\mathcal{G}$ . Since  $\tilde{G}$  is row stochastic, we can show by induction that  $\tilde{G}^t$  is also row stochastic. The row stochasticity of  $\tilde{G}^t$  is a linear condition, hence continuous so it is preserved by the limits. It follows that :

$$\sum_{j \in N} \sum_{i \in N} \mathcal{G}_{ij} = |N| \quad (9)$$

Let  $\mathbf{a}_\infty$  be a bi-polarized long-run opinion vector (see definition 6) and suppose that

$$| \sum_{i \in E^+} (\mathcal{G}' \mathbf{1})_i - \sum_{i \in E^-} (\mathcal{G}' \mathbf{1})_i | > \epsilon \quad (10)$$

Equation (10) can be rewritten as:

$$\sum_{i \in E^+} (\mathcal{G}' \mathbf{1})_i - \sum_{i \in E^-} (\mathcal{G}' \mathbf{1})_i = \nu \text{ for } |\nu| > 1$$

Moreover,  $\sum_{i \in E^+} (\mathcal{G}' \mathbf{1})_i = \sum_{j \in N} \sum_{i \in E^+} \mathcal{G}_{ji}$  and for all  $j \in N$ ,  $\sum_{i \in E^+} \mathcal{G}_{ji} + \sum_{i \in E^-} \mathcal{G}_{ji} = 1$ . Hence:

$$\begin{aligned} \sum_{i \in E^+} (\mathcal{G}' \mathbf{1})_i - \sum_{i \in E^-} (\mathcal{G}' \mathbf{1})_i &= \nu \Leftrightarrow \sum_{j \in N} \sum_{i \in E^+} \mathcal{G}_{ji} - \sum_{j \in N} \sum_{i \in E^-} \mathcal{G}_{ji} = \nu \\ &\Leftrightarrow \sum_{j \in N} \left( \sum_{i \in E^+} \mathcal{G}_{ji} - \sum_{i \in E^-} \mathcal{G}_{ji} \right) = \nu \\ &\Leftrightarrow \sum_{j \in N} \left( 1 - 2 \sum_{i \in E^-} \mathcal{G}_{ji} \right) = \nu \\ &\Leftrightarrow \frac{|N| - \nu}{2} = \sum_{j \in N} \sum_{i \in E^-} \mathcal{G}_{ji} = \sum_{j \in N} \sum_{i \in N^-} \mathcal{G}_{ji} \end{aligned}$$

Since  $|\nu| > 1$  and there does not exist moderate individuals, the size of the group  $N^-$  is strictly smaller than  $|N|/2$ .

Let  $k \in N$  be a consensual individual. It follows from theorem 1 that:

$$\alpha_{k,\infty} = \sum_{i \in E} \mathcal{G}_{k,i} \alpha_{i,\infty} = \sum_{i \in E^+} \mathcal{G}_{k,i} - \sum_{i \in E^-} \mathcal{G}_{k,i}$$

Moreover, recall that for all  $k \in N$ ,  $\sum_{i \in E} \mathcal{G}_{k,i} = \sum_{i \in E^+} \mathcal{G}_{k,i} + \sum_{i \in E^-} \mathcal{G}_{k,i} = 1$ . Individual  $k$  doesn't hold the

extreme opinion of the members of  $E^+$ , if and only if, for  $i \in E^+$ :

$$\begin{aligned} |\alpha_{k,\infty} - \alpha_{i,\infty}| \geq \tau &\Leftrightarrow |(2 \sum_{i \in E^+} \mathcal{G}_{ki} - 1) - 1| \geq \tau \\ &\Leftrightarrow 1 - \sum_{i \in E^+} \mathcal{G}_{ki} \geq \frac{\tau}{2} \\ &\Leftrightarrow 1 - \frac{\tau}{2} \geq \sum_{i \in E^+} \mathcal{G}_{ki} \end{aligned}$$

Similarly, individual  $k$  doesn't hold the extreme opinion of the members of  $E^-$ , if and only if, for  $i \in E^-$ :

$$\begin{aligned} |\alpha_{k,\infty} - \alpha_{i,\infty}| \geq \tau &\Leftrightarrow |(2 \sum_{i \in E^+} \mathcal{G}_{ki} - 1) - (-1)| \geq \tau \\ &\Leftrightarrow \sum_{i \in E^+} \mathcal{G}_{ki} \geq \frac{\tau}{2} \end{aligned}$$

## 7.7 Network statistics

Consider a graph  $G$ , with  $N$  vertices and  $M$  edges.

- Assortativity or assortative mixing in Newman (2002) [26], assortativity of an observed network is given by :

$$r = \frac{M^{-1} \sum_i j_i k_i - [M^{-1} \sum_i \frac{1}{2}(j_i + k_i)]^2}{M^{-1} \sum_i \frac{1}{2}(j_i^2 + k_i^2) - [M^{-1} \sum_i \frac{1}{2}(j_i + k_i)]^2}$$

where  $j_i$  and  $k_i$  are the degrees of the vertices at the ends of the  $i$ th edge with  $i = 1, \dots, M$ .

- Neighborhood connectivity: average degree in the neighborhood of a given node  $i \in N$

$$\frac{1}{d_i} \sum_{j \in N_i} d_j$$

where  $d_i$  is the degree of node  $i$  and  $N_i$  is the neighborhood.

- Per capita average degree:

$$\frac{\sum_{i \in N} d_i}{N}$$

where  $d_i$  is the degree of node  $i$ .

- Average path length :

$$l = \frac{1}{N(N-1)} \sum_{i \neq j} d(v_i, v_j)$$

where  $d(v_i, v_j)$  is the shortest path between nodes  $v_1$  and  $v_2$  (computed here with Dijkstra).

## 7.8 Toy Networks

I consider a number of *toy* networks with the same number of individuals and compare the level of polarization for each. I study two families of *toy* networks. In the first family expressers are densely connected like in multi-star networks, regular networks and 5-regular Cayley trees. In the second family expressers are sparsely connected. A multi-star network is a network structure where (i) there are individuals (expressers or *stars*) with a very large number of followers (consensual), (ii) the stars are connected among each other directly or indirectly. This family of networks is ment to model Twitter-like communities where there are *stars* or *influencers* with many followers and where stars interact together. Intuitively, *stars* create content, while followers *share*. A 5-regular Cayley tree, is ment to model organizations or companies or even the political structure in representative democracies. In such a network all nodes have the same degree except for the leaves at the lowest level of the tree. Individuals at each level of the tree *express* while all the leaves choose to be consensual. Finally, the 3-regular network is ment to model networks where individuals have the same local popularity or expertise and consequently they all express. These networks can model discussions in relatively small groups which don't contain opinion leaders. Notice that with initial opinions distributed uniformly at random the probability, if you take two expressers that are neighbors, then the probability<sup>39</sup> that they are ideologically opposed is :

$$\mathbb{P}(|\alpha_{i,0} - \alpha_{j,0}| \geq \tau) = \frac{2 + \tau(2 - 3\tau)}{4}$$

Unsurprisingly, for  $\tau = 0.5$  this probability is approximatively  $1/2$  and this explains the high level of polarization in Toy networks where all the expressers are densely connected.

---

<sup>39</sup>The easiest way to compute this probability is to compute the area between  $y = x + \tau$  and  $y = x - \tau$ . Figure 15 shows this area. The are between the two red lines over the total area is probability that  $|x - y| < \tau = 0.5$ . The area between the two blue lines over the total area is the probability that  $|x - y| < \tau = 1$ . In general the area between two lines of the same color is  $1 - (2 + \tau(2 - 3\tau))/4$ .



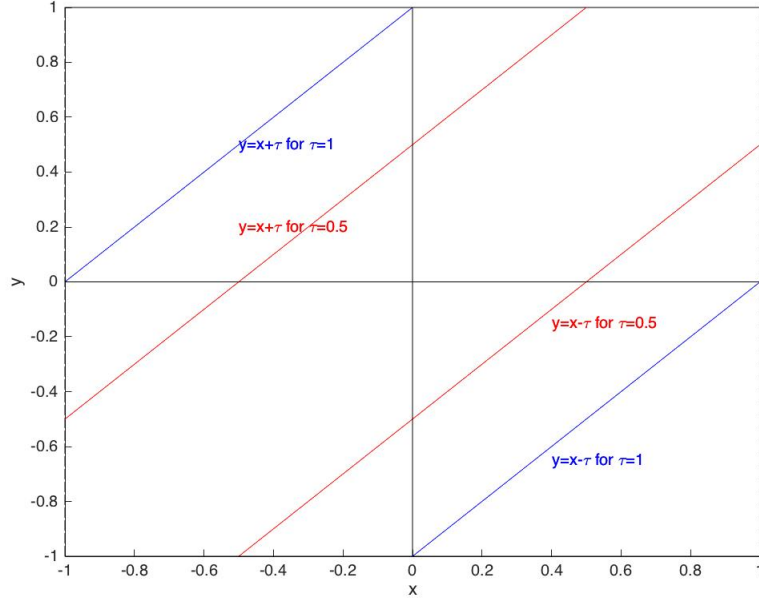


Figure 15: The area between the two red lines over the total area is probability that  $|x-y| < \tau = 0.5$ , the area between the two blue lines over the total area is the probability that  $|x-y| < \tau = 1$ , Area between two lines of the same color  $1 - (2 + \tau(2 - 3\tau))/4$

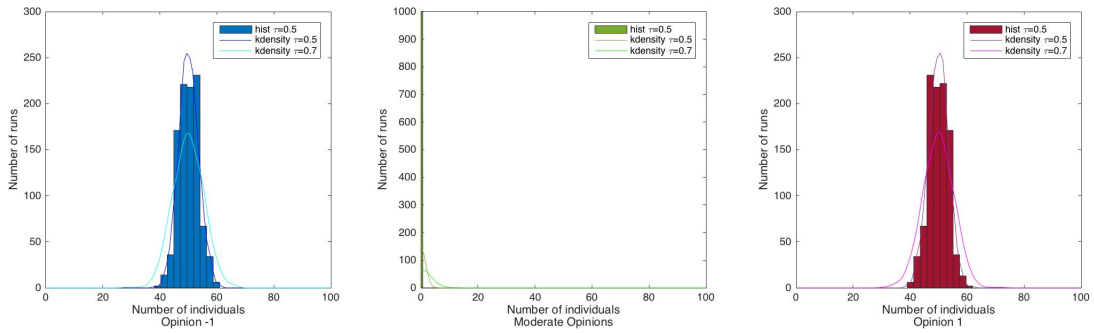


Figure 16: kernel density of group sizes over 1000 runs for a 3-regular network with 100 individuals

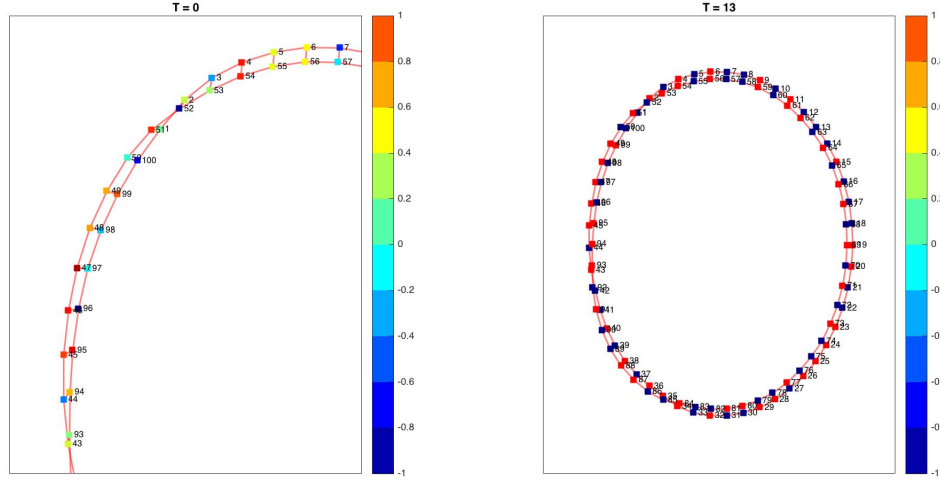


Figure 17: Initial and final opinions of one run out of 1000 in a 3-regular network with 100 individuals

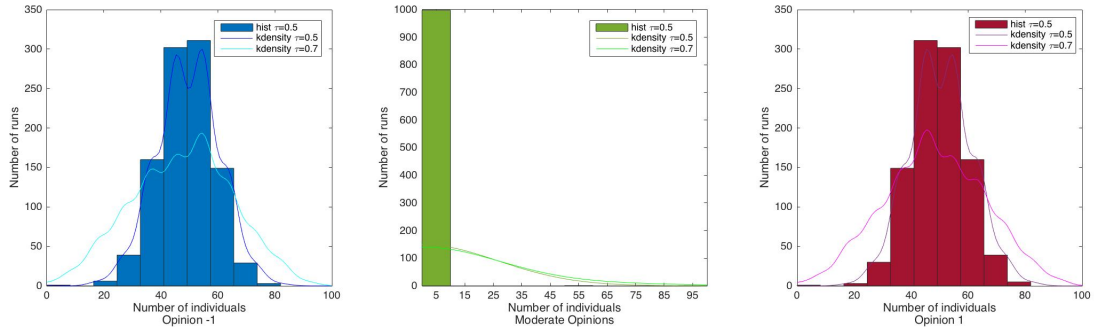


Figure 18: kernel density of group sizes over 1000 runs for a multi-star network with 100 individuals

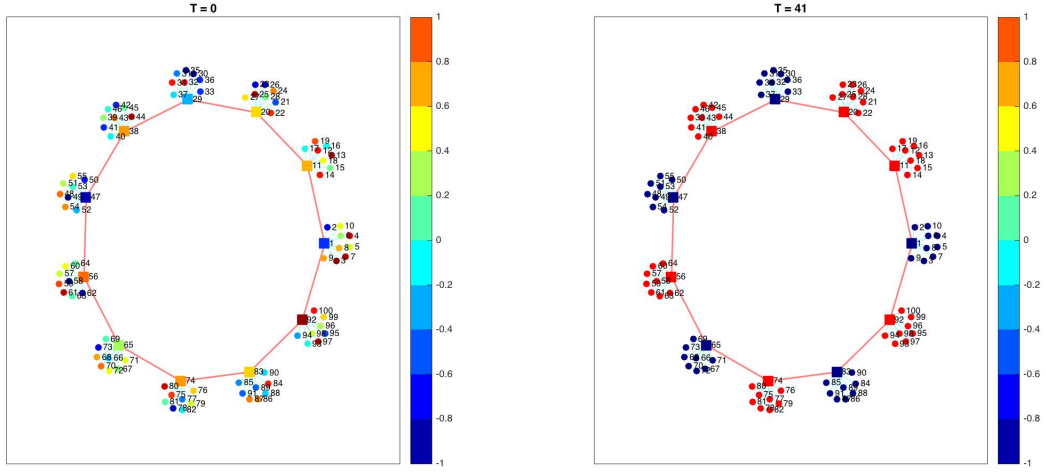


Figure 19: Initial and final opinions of one run out of 1000 in a multi-star network with 100 individuals showing bi-polarization

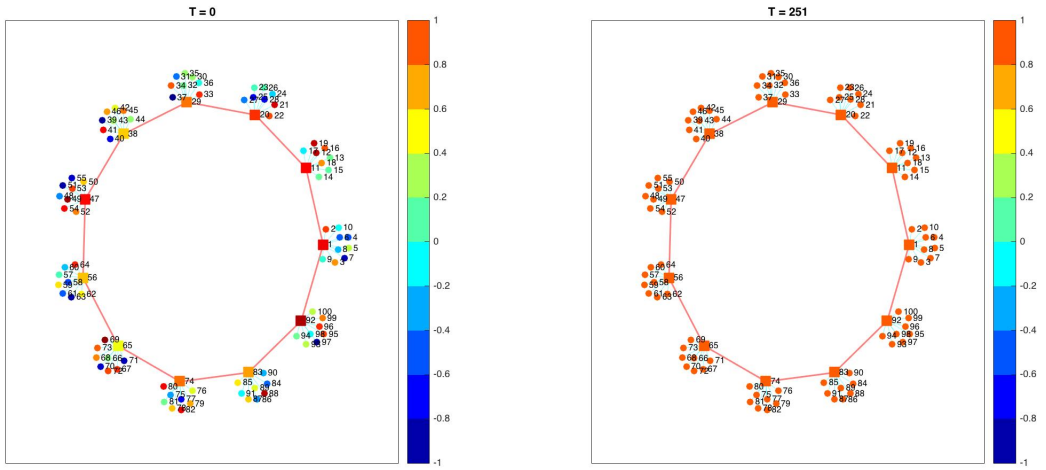


Figure 20: Initial and final opinions of the only run out of 1000 in a multi-star network with 100 individuals showing consensus

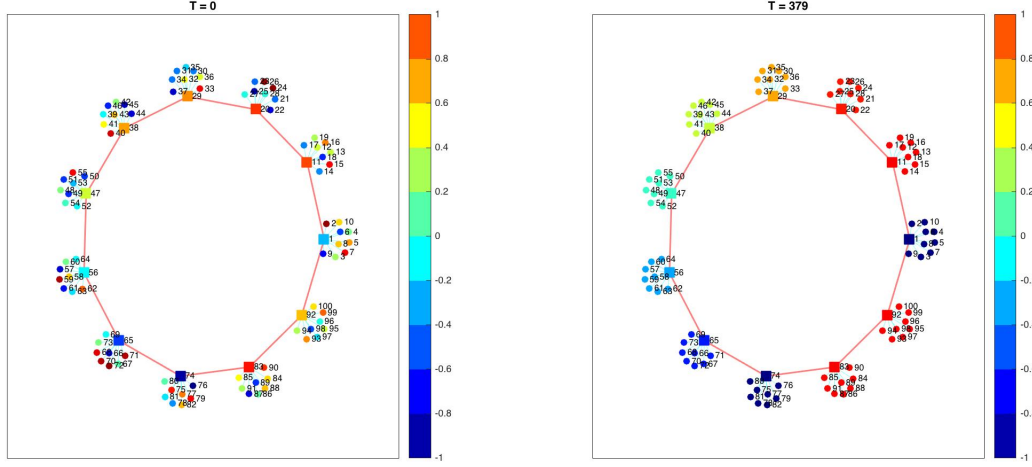


Figure 21: Initial and final opinions of the only run out of 1000 in a multi-star network with 100 individuals showing long-run moderate opinions

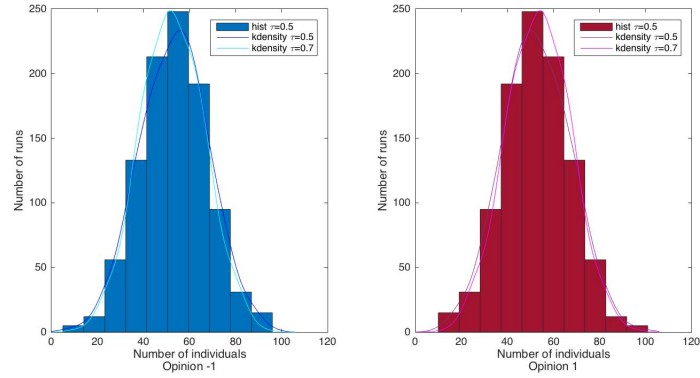


Figure 22: kernel density of group sizes over 1000 runs for a 5-regular Cayley tree with 106 individuals

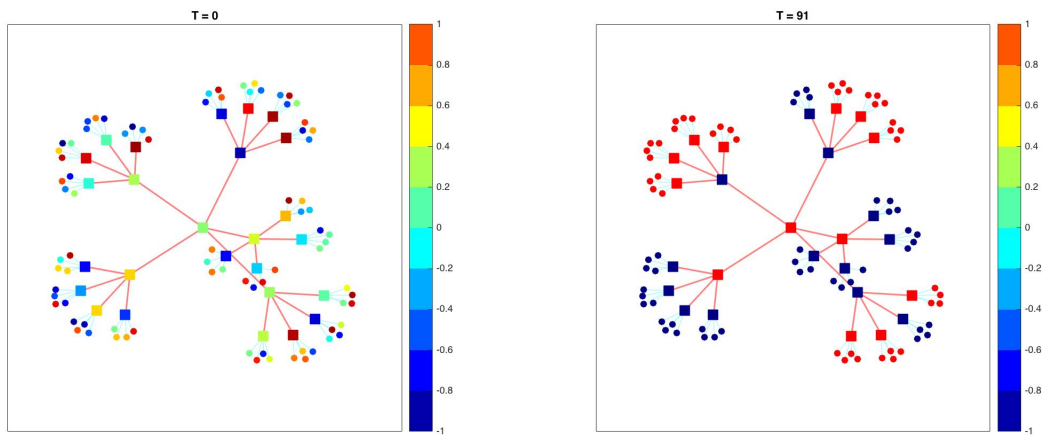


Figure 23: Initial and final opinions of one run out of 1000 in a 5-regular Cayley tree with 106 individuals showing bipolarization