

# Chapter 1

## Introduction

Complex system is a collection of a large number of units, they can interact with each other and because of the interaction some collective behaviour can emerge. The properties of the system can not be predicted from behaviour of one individual.

Statistical physics attempt to describe behavior of large number interacting particles, atoms and molecules and macroscopies properties for example magnetisation is explained from interaction between particles. Also as complex system we can consider people in society, population of fishes showing flocking pattern, traffic on the roads.

### 1.1 Complex networks

The first mathematical problem solved using graph theory was *Konigsberg*, Kaliningrad in Russia, the problem of seven bridges. The city *Konigsberg* in that time had seven bridges, that were connecting the parts of the city across the river and the island in the middle. The question was is it possible to find a walk that crosses all seven bridges only once. Representing the problem as a graph, Euler managed to simplify the problem, the parts of the land are represented as nodes while bridges between them are links. Crossing each bridge only once is possible if each part of the land has an even number of connections. Thus, in this case it was not possible, as each piece of land was connected with an odd number of bridges.

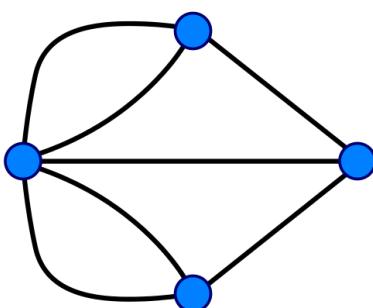
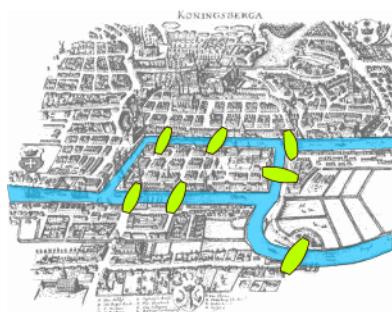


Figure 1.1: The Kronigsber problem of seven bridges.

Despite the differences among complex systems, they can be represented in unique way, using graph theory. The real nature of the components is neglected, and we only represent

the interaction among them. This approximation, allow us to treat equally social (graph of actors), biological (network of proteins) or even technological system (internet, traffic). In recent years complex network theory finds the application in different fields, and it's development is encored by availability of big data.

## 1.2 The structure of complex networks

The complex system can be represented by complex network  $G = (V, E)$ , where the elements of system (atoms, proteins, people) map to set of  $N$  nodes  $V = \{1, 2, \dots, N\}$ . The interactions between elements map to  $L$  links between nodes,  $E = \{e_1, e_2, \dots, e_L\}$ . The **adjacency matrix**  $A = N \times N$  has value 1 if there is connection between two nodes, otherwise it is 0 [1].

### degree distribution

The network properties directly depend on the connectivity between nodes. In the case of regular networks, such as grids, each node has equal number of first neighbors. In general case, the networks have more complicated structure. Thus the important measure is network **degree**  $k$ . The degree of node  $i$  gives the number of nodes attached to node  $i$ ,  $k_i = \sum_j A_{ij}$ . The **degree distribution**  $P(k)$  is probability that randomly chosen node has degree  $k$ . It can be calculated as fraction of  $k$  degree nodes  $N_k$ ,  $p(k) = N_k / N$ . The degree distribution in random network, where all nodes have the same connecting probability follows Poisson distribution  $P(k) = \frac{(Np)^k e^{-Np}}{k!}$ , where  $k$  is mean degree distribution. In real networks degree distribution follows power law. Therefore, real networks have scale-free structure with emergence of the hubs [2].

### assortativity

The **degree-degree correlations** in the network are measured by **assortativity**. If correlations are positive, networks are assortative; there is a tendency that connections exist between similar degree nodes. The negative correlations indicate that large degree nodes have preference to connect nodes with small degree; dissasortative networks. The average first neighbor degree  $k_{nn}$  can be calculated as  $k_{nn} = \sum_{k'} k' P(k'|k)$ . The  $P$  is conditional probability that an edge of degree  $k$  points to node with degree  $k'$ . The norm is  $\sum_{k'} P(k'|k) = 1$ , and detailed balance conditions [1],  $kP(k'|k)P(k) = k'P(k|k')P(k')$  [1]. If the node degrees are uncorrelated,  $k_{nn}$  does not depend on the degree, otherwise increasing/decreasing function indicates on positive/negative correlations in the network.

The Newman defined the assortativity index  $r$  in slightly different way:  $r = \sum_{kl} kl(e_{kl} - q_l q_k) / \sigma_q^2$ , where  $e_{kl}$  is the probability that randomly selected link connect nodes with degrees  $k$  and  $l$ ,  $q_k$  is probability that randomly choosen node is connected to node  $k$  and equals  $q_k = kp_k / \langle k \rangle$ , while  $\sigma_q$  is varience of the distribution  $q_k$ .

## clustering coefficient

The **clustering coefficient** is measure that describe the structure of neighborhood. In networks exist tendency of forming triangles or clusters. This is common in friendship networks where two friends of one person have high probability to be friends. The clustering can be measured by computing the number of links between neighbours of one node,  $c_i = 2e_i / (k_i(k_i - 1))$ . Averaging it over all nodes in the networks we can calculate mean clustering coefficient. It ranges from  $\langle c \rangle = 0$  where connections between neighbouring nodes do not exist, network has structure of three. On the other hand  $\langle c \rangle = 1$  indicates fully connected network.

## real-world networks

Real world networks share similar properties. Mean distance between nodes is small and it is much smaller than number of nodes in the network  $l \ll N$ , this property is called small world phenomena. This cause the fast spread of information or even diseases in the complex systems. In small world networks number of vertices grow exponentially with distance, thus  $l$  increase as  $\log(n)$  or slower. Logarithmic scaling can be proved from various network models, also it is observed in real world complex systems. Clustering coefficient in the real world networks is usually high. Real world networks have one important feature; power-law degree distribution; such networks are called scale-free networks.

## 1.3 The dynamics of complex networks

### The random graph model

The random graph model was introduced by Erdos and Renyi in 1959. This model has  $N$  disconnected nodes. With probability  $p$  each pair of nodes can be connected. The network is characterized only by number of nodes and links,  $G(n, p)$ .

As this process is stochastic, the network with same parameters  $n$  and  $p$ , does not have to be same structure, so it is necessary to consider the ensemble of networks. Then the mean number of links depends on the model parameters:  $\langle m \rangle = n(n - 1)p/2$ . The expected value of node degree can be predicted as  $\langle k \rangle = (n - 1)p$ . The probability  $p$  is defined as density of the network. The probability  $p(k)$  follows the binomial distribution of the form:

$p(k) = p^k(1 - p)^{n-1-k}$ . For large values of  $n$ , this becomes  $p(k) = e^{-k}k^{-k}/k!$ , which is Poisson distribution.

In the case of a large random networks, the average path length is given as  $l = \frac{\ln n - \gamma}{\ln(pn)} + \frac{1}{2}$ . This means that random graph has a very small average path length. This is characteristic of many large networks.

The clustering coefficient  $C = p$ , so for sparse ER graphs the clustering is very small, much smaller than in real world networks.

## 1. Introduction

---

Increasing the probability  $p$ , the giant component may appear. This is sub-graph whose size is proportional to the size of the network. Such change in the network is phase-transition and it is important as small change in the probability  $p$  leads to fundamental change in the system properties. There are two limits of this model, when  $p = 0$ , network is disconnected. If  $p = 1$ , then network is fully connected, and giant component is with size  $O(1)$ . This phenomena is related to percolation phase transition. On the threshold  $p_c$ , the component whose size is proportional to  $n^{2/3}$  emerges. Average path length between two nodes, at critical point is proportional to the  $\ln N$ . The small, logarithmic distance is the origin of the "small-world" phenomena.

The interesting behavior of this model is that, with increasing  $p$  nodes tend to organise in giant component. The subcritical,  $k < 1$  where all components are dimple and small. The size of largest component is  $s = O(\ln n)$ . In critical regime  $k = 1$ , the size of largest component is  $s = O(n^{2/3})$ . Supercritical  $k > 1$ , where the probability of having giant component is 1.

## Small world networks

In the 1999. Watts and Strogatz introduces "small-world" model. This model can generate the networks with small diameter and high clustering coefficient. Their idea is to start from grid like network, where all nodes have same number of neighbors, like ring-lattice or hexagonal lattice, where each node is connected to  $k$  nearest neighbors. Such network has high clustering coefficient, as any pair of consecutive neighbors are connected forming a triangles, while in contrast the network has high average shortest path, as nodes on the opposite sides of the lattice are not connected. The goal of this model is to connect distant nodes and reduce the average path length in the network. This can be simply done by randomly rewiring nodes in the network, with probability  $p$ . Model interpolates between regular network  $p = 0$  and random graph  $p = 1$ , and for some critical probability we can achieve small world networks.

The average shortest-path length from the model is close to that of an equivalent network, and much lower than that of the lattice. The clustering coefficient from the model is still close to that one in the lattice and much larger than in random network.

The degree distribution of this model obviously is not power-law. In regular network, all nodes have equal degree, while in random networks degree distribution becomes Poisson.

## Barabasi-Albert model

The random network model differs from real networks in the two characteristics, growth and preferential attachment. In static models, number of nodes is fixed, while in growing models we try to simulate the continuous change in the system. More important ingredient, are linking rules. In real networks, new nodes tend to link to more connected nodes.

This model is defined as follows, we start from  $m_0$  nodes, randomly connected, and at each timestep we add new node with  $m$  links that will connect to  $m$  nodes already present in the network. The probability that new node connects to node  $i$  depends on node degree  $k_i$  as

$$P(k_i) = \frac{k_i}{\sum_j k_j} \quad (1.1)$$

New node can connect to any node in the network, however nodes with larger degree have higher probability to link new nodes. After time  $t$  the model generates network with  $N = t + m_0$  nodes and  $m_0 + mt$  links. Degree distribution is power-law with exponent  $\gamma = 3$ . As network grows nodes with larger degree becomes bigger, so we end up with few nodes with many links, called hubs. Two simple mechanisms are responsible for emergence of scale-free networks.

#### *degree distribution*

To understand the emergence of scale-free properties we need to analyze the evolution of degree distribution. The rate at which an existing node get new links as result of new nodes connecting to it is

$$\frac{dk_i}{dt} = mP(k_i) = m \frac{k_i}{\sum_j k_j} \quad (1.2)$$

each new node arrives with  $m$  links. The sum is  $2mt - m$  so the equation for large  $t$  becomes:

$$\frac{dk_i}{k_i} = \frac{1}{2} \frac{dt}{t} \quad (1.3)$$

solving this equation we get that degree of node in time step  $t$  is  $k_i(t) = m(\frac{t}{t_i})^\beta$ , where  $\beta = 1/2$ .

We note that degree of each node increase following power-law; the growth in degrees is sub linear, as each new node has more nodes to link than previous. The earlier node  $i$ , the higher is its degree. Hubs are large as they arrived early in the network.

In summary, the analytical calculations predict that the Barabási-Albert model generates a scale-free network with degree exponent 3. The degree exponent is independent of the  $m$  and  $m_0$  parameters. The degree distribution is stationary explaining how different systems have similar structural properties.

In summary, the absence of preferential attachment leads to a growing network with a stationary but exponential degree distribution. In contrast the absence of growth leads to the loss of stationarity, forcing the network to converge to a complete graph. This failure of Models A and B to reproduce the empirically observed scale-free distribution indicates that growth and preferential attachment are simultaneously needed for the emergence of the scale-free property.

In the past decade we witnessed the emergence of two philosophically different answers. The first one views preferential attachment as the interplay between random events and some structural property in the network. The second assumes that each new node or link balances conflicting needs.

The BA model postulates the presence of preferential attachment. Yet, we can build models that generate scale free networks without preferential attachment. The link selection

## 1. Introduction

---

model offers the simplest mechanism that generates a scale-free network. At each time step we add new nodes to the network, we select link at random and connect the new node to one of the two nodes at the end. The higher is degree of the node, the higher is chance that node is located at the end of chosen link. The more k-degree nodes are there, the more likely is that k node is at the end of chosen link. Probability that node at the end of randomly chosen link has degree k is  $q_k = Ckp_k$ . The fact that bias is linear with k indicates that the link selection model builds scale-free networks. Copying model can also generate scale-free networks. In each time step a new node is added to the network. To decide where it connects we randomly select node u. Then with probability p new node links to u, otherwise with probability  $1 - p$  we randomly choose an outgoing link of node u and link the new node to its target. The likelihood that new node connects to degree-k node is  $P(k) = \frac{p}{N} + \frac{1-p}{2L}k$ , the second part is equivalent in selecting a node to randomly selected link. The popularity of the copying model lies in its relevance in real systems. It is common in social networks, citation networks or even protein interactions. in optimization, when new nodes balance conflicting criteria as they decide where to connect

*diameter* The network diameter, represents the maximum distance in the BA model,  $d \sim \frac{\ln N}{\ln \ln N}$ . The diameter grows slower than  $\ln N$ , making the distances in BA model smaller than in random graph. The difference is found for large N. *clustering* The clustering coefficient of the BA model follows  $C \sim \frac{\ln N^2}{N}$ . It is different from clustering found in random networks, and BA networks are in general more clustered.

## Nonlinear BA model

In summary, nonlinear preferential attachment changes the degree distribution, either limiting the size of the hubs ( $\alpha < 1$ ), or leading to super-hubs ( $\alpha > 1$ ). Consequently,  $P(k)$  needs to depend strictly linearly on the degrees for the resulting network to have a pure power law  $p \propto k^{-\gamma}$ . While in many systems we do observe such a linear dependence, in others, like the scientific collaboration network and the actor network, preferential attachment is sublinear. This nonlinear  $P(k)$  is one reason the degree distribution of real networks deviates from a pure power-law. Hence for systems with sublinear attachment the stretched exponential (5.23) should offer a better fit to the degree distribution.

In real systems preferential attachment can be more influenced by the age of the node. If parameter alpha is negative, ageing effect overcomes the role of preferential attachment, and scale-free properties are lost. For large negative alpha, the network turns into the chain, where the youngest nodes are the most attractive. On the other hand for a positive alpha, new nodes will link to older nodes. Positive alpha makes the network more heterogeneous, and scale-free nature still exists but exponent gamma is different from 3. for the high alpha all nodes will tend to connect to oldest node.

In the general ageing model, we have linking rules where rules connecting probability depends on both of node degree and age difference between new and old node. With parameters alpha and beta we can control the structure of generated networks. I already

talked about some limits of the general model. We saw that for specific parameters there are SF networks, BA model, if we move from that point SF behaviour with power-law with exponent 3 is lost. And other classes of networks can appear. In general, model, when alpha and beta are both positive, rich get richer phenomena is more promoted. On the other hand, the region where beta is positive and alpha negative can be interesting, because SF networks can appear only along the critical line.

In growing network models is considered that at each time step one node is added to the network. The remaining question is if there is any change if network growth is not linear anymore and how does it influence the structure of obtained networks. In this work, we use numerical simulations to explore the case when  $M(t)$  is a correlated time-varying function and study how these properties influence the structure of generated networks for different values of parameter  $-\infty < \alpha \leq -1$  and  $\beta \geq 1$  and constant  $L$ .

## 1.4 Network structures

### Bipartite networks

A bipartite network has two partitions,  $U$  and  $V$ . The nodes in the same partition are not connected while links exist only between nodes of a different kind. Bipartite networks represent the membership of people or items in groups. For example, we can define the network of actors as a bipartite graph. In one partition are actors and in other movies. There are no edges between actors or movies, but the actor is connected to the film if it plays in that movie. Another example is a recommender network, such as a network of people and items they like.

The equivalent representation of bipartite network is incidence matrix  $B$ . If  $n$  is number of people and  $g$  number of groups, this matrix is  $gxn$ , having elements  $B_{ij}$  1 if person i belongs to group j.

Even bipartite networks give realistic representation of the system, there is often need to analyze the single type of nodes. From a bipartite network, we can generate two projections. The first one connects nodes partition  $V$  if they point to node  $u$ . Similarly, we can project the network on  $U$  partition, connecting  $u$  nodes. The one mode projection between actors and movies onto actors is undirected network of actors. Actors are connected if they appear in the same movie. We can also create one-mode projection onto movies, where two movies are connected if they share the same actor.

The projections are useful in some manner, but they also lose some important information, for example how many groups nodes share in common. This information can be propagated adding the weight to the edges, equal to the number of common groups.

The product  $B_{ki}$  and  $B_{kj}$  is 1 if  $i$  and  $j$  belong to the same group  $k$ . Thus the total number of groups to which nodes  $i$  and  $j$  belong is

$P_{ij} = \sum_{k=1}^g B_{ki}B_{kj} = \sum_{k=1}^g B_{ik}^T B_{kj}$ . The matrix  $P$  is matrix of one-mode projection. The diagonal elements are non-zero, and represent the number of groups node  $i$  belongs to. To

derive the weighted adjacency matrix, the diagonal elements are set to 0. The adjacency matrix of unweighted projection, each non-zero element needs to be replaced with 1.

### Core-periphery networks

Core-periphery structure describes a network whose nodes are divided into two community, densely connected core and less connected periphery. If we consider the average probabilities of edges within each group as  $p_{11}$  and  $p_{22}$ , and between groups  $p_{12}$ , instead of traditionaly assortative or dissasortative structure we can define core-periphery structure  $p_{11} > p_{12} > p_{22}$ . In the principle core-periphery structure does not have to be limited to only two groups, and we can define layered, onion, structure. The network can have more cores, that are not directly connected to each other.

The simple method for finding core-periphery structure is to assume that nodes in core have higher degree in the core than in the periphery. Another simple method is to construct k-cores. K core is group of nodes that each has connection to at least k other members of the group. K-cores form a nested set, and become denser with higher k. The core-periphery structure can be detected optimizing the measure similar to modularity, as defined by Borgatti and Everett. Their goal is to find the division that minimizes the number of edges in the periphery. So they define the score function that is equal to number of edges in the periphery minus the expected number of such edges placed at random.  $\rho = \frac{1}{2} \sum_{ij} (A_{ij} - p) g_i g_j$ . They used genetic algorithm to minimize this function.

The another way to detect core-periphery structure is to use the inference method based on fits to a stochastic block model. In this method we fit observed network to a block model with two groups, such that edge-probabilities have form  $p_{11} > p_{12} > p_{22}$ . The only downside of this model is that method is going to find the structure that optimize likelihood, and we can not say weather it is core-periphery or community structure.

### Communities

Thus the ability to find groups or clusters in a network can be a useful tool for revealing structure and organization within networks at a scale larger than that of a single node or a few nodes. The occurrence of groups or communities is not limited to social networks. Clusters of nodes in a web network, for instance, might indicate groups of related web pages. Clusters of nodes in a metabolic network might indicate functional units within the network. The ability to find groups also has another practical application: it allows us to take a large network and break it apart into smaller subsets that can be studied separately. The network in Fig. 14.1 is quite small, but others may be much larger, millions of nodes or more, making their analysis and interpretation challenging. Breaking such networks into their component clusters is a useful technique for reducing them to a manageable size. One example of this approach is in network visualization. A network with a million or more nodes can rarely be visualized in its entirety, even with the aid of good visualization software. Such networks are simply

too big to be represented usefully on the screen or on paper. If the nodes of the network divide naturally into groups, however, then we can make a simpler but still useful picture by representing each group as a single node and the connections between groups as edges. An example is shown in Fig. 14.2. This simplified representation allows us to see the large-scale structure of the network without getting bogged down in the details of the individual nodes. If one wanted to see the individual nodes, one could then “zoom in” to a single group and look at its internal makeup. The problem of finding groups of nodes in networks is called community detection. Simple though it is to describe, community detection turns out to be a challenging task, but a number of methods have been developed that return good results in practical situations.

## Stochastic block model

The network or graph is the structure of nodes and edges, where each edge connects two nodes. Nodes can be organized into groups, called communities. Identifying these hidden blocks can lead to interesting insights into the network. However, the community detection problem does not give a precise definition of what a community is. As a consequence, many approaches try to recover such structural patterns in the network [3].

A common definition of a community is that it is densely connected subgraph [4]. We can find these subgraphs by optimizing an objective function, such as modularity function. It measures the difference in the number of edges between the given network and the network with the same number of nodes but randomly connected. In this approach, we try to maximize the density of connections inside a group by focusing more on assortative<sup>1</sup> group structures.

Another type of networks is the bipartite network that has two disjoint sets of nodes. The edges exist only between nodes from different sets. Networks of this class can appear in real-world data, such as users-movies preference, collaboration network for scientists and papers, etc. Application of density-based approach requires to first project bipartite network to one of its partitions and then find communities in that projection. With this, some information is lost. On the other hand, the method that is directly applicable to bipartite networks is Stochastic Block model, from which the models considered in this paper are derived.

Stochastic block model (SBM) is based on connection probabilities between nodes. It is a generative model which includes existence of communities. Parameters that describe SBM for network  $G$  with  $N$  nodes are:

- $k$ : number of groups
- group assignment vector,  $g$ :  $g_i \in \{1, 2..k\}$ , gives the group index of node  $i$ .
- SBM matrix,  $p_{k \times k}$ , whose elements  $p_{ij}$  are the probabilities that edges between groups  $g_i$  and  $g_j$  exist.

---

<sup>1</sup>Networks where nodes tend to connect with other nodes of a similar degree. Edges are more likely inside blocks than out of them.

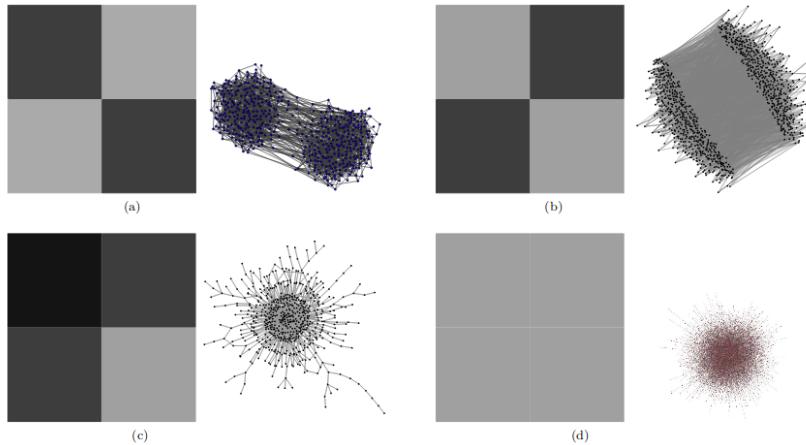


Figure 1.2: Stochastic Block model for different networks structures. (a) assortative. (b) disassortative. (c) core-periphery. (d) Erdos Renyi random graph.

Note that nodes within one group have the same connection probabilities.

SBM can generate and describe different types of network structures. Figure 1.2 [4] shows how the model matrix corresponds to resulting networks with two communities. First, for the assortative network (1.2 a), diagonal elements of the matrix have higher probabilities. This indicates dense connections inside the group, just like in classic community structures. In disassortative structure, (1.2 b), more connections exist between two partitions than inside them, i.e. off-diagonal elements have higher probabilities. Bipartite networks can be represented like this.

Figure (1.2 c) shows how the model represents core-periphery networks. Nodes of one block (core) are well connected with itself and with other partition (periphery). From the last case, we can note that SBM with one group is the Erdos Renyi random graph (1.2 d) because all probabilities inside and between groups are equal.

The benefit of this model is that we can generate many networks with similar group structure. The model can fit real data, which results in finding network communities. For the given network  $G$  and number of groups  $k$ , the best nodes partition  $g$  is found by maximizing the likelihood function. Beside inferring communities, SBM has application in prediction of missing links. This simply formulated model has many variants, motivated by specific properties of real data. For example, for networks which are degree heterogeneous, there is degree corrected SBM. In some social networks, users can belong to more than one group, and this can be modelled with mixed membership SBM. Other extensions include application to bipartite, weighted network, hierarchical model, etc. Also, several algorithms for optimization of likelihood function are proposed. The overview of these versions and methods are given in [5]. In this paper, we will focus on Single and Mixed Membership SBM applied on bipartite networks.

## 1.5 Graph isomorphism

Weisfeiler-Lehman Test

## 1.6 Distributions

### power-law

Power-law distributions characterize many social and biological systems. Power-law distributions are also easy to generate.

the distributions: basic definitions and properties

The nonnegative random variable  $X$  is said to have a power law distribution if

$$Pr[X > x] \sim cx^{-\alpha} \quad (1.4)$$

for constants  $c > 0$  and  $\alpha > 0$ . In power-law distribution asymptotically the tails fall according to power  $\alpha$ . Such distribution leads to much heavier tails than other common models, as exponential distribution.

One specific power-law distribution is Pareto distribution which satisfies

$$Pr[X > x] = \frac{x^{-\alpha}}{k} \quad (1.5)$$

for some  $\alpha > 0$  and  $k > 0$ . The Pareto requires  $X > k$ . The density function of Pareto distribution is  $f(x) = \alpha k^\alpha x^{-\alpha-1}$ . For power law distribution  $\alpha$  is in the range  $0 < \alpha < 2$ , in which case  $X$  has infinite variance, if  $\alpha < 1$   $X$  has infinite mean.

if  $X$  has power law distribution, then  $a$  in log-log plot  $Pr(x)$  will behave as straight line. For the specific case of Pareto distribution, the behaviour is exactly linear as  $\ln(Pr(x)) = -\alpha(\ln x - \ln k)$ . Similarly the density function is also straight line.  $\ln(f(x)) = (-\alpha - 1)\ln(x) + \alpha\ln(k) + \ln(a)$

### log-normal

Many measurements in the nature show a more or less skewed distribution. They are common when mean values are low, variances are large and values can not be negative as example in distribution of mineral resources in the Earth. Such skewed distributions often closely fit to log-normal distribution.

What is the difference between normal and lognormal variability? A major difference is that effect can be additive or multiplicative, leading to normal or lognormal distribution. Basic principles of additive and multiplicative effects can be easily demonstrated with the help of two dices. Adding the two numbers, which is the principle of the most games, leads to values from 2 to 12 with mean of 7, and symmetrical frequency distribution. Multiplying the two numbers, leads to values from 1 to 36 with highly skewed distribution. Although

## 1. Introduction

---

these examples are not normal or lognormal they give us clear difference how different distributions can emerge.

Log-normal distributions are usually characterized in the term of the log-transformed variable, using as parameters the expected value, or the mean, and the standard deviation. This characterization can be advantageous as by definition log-normal distributions are simetrical at the log level.

The basic properties of the lognormal distributions

Random variable  $X$  if  $\log(X)$  is normally distributed., if  $Y = \ln X$  has normal Gaussian distribution.

Only positive values are possible for the variable and distribution is skewed. Two parameters are needed to specify lognormal distribution. Tradionaly the mean  $\mu$  and standard deviation  $\sigma$  or the variance of the  $\sigma^2$  of  $\log(X)$  are used. However there are clear advantages of using transformed data,  $\mu^* = e^\mu$ ,  $\sigma^* = e^\sigma$ . The median of this lognormal distribution is  $\text{med}(X) = \mu^* = e^\mu$ , since  $\mu$  is median of the  $\log(X)$ .

$$f(x) = \frac{1}{x\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2\sigma^2}(\log(x) - \mu)^2\right) \quad (1.6)$$

The mean is  $\exp(\mu + \sigma/2)$  and variance is  $(\exp(\sigma^2) - 1)\exp(2\mu + \sigma^2)$ .

Estimation: The asymptotically most efficient (maximum likelihood) estimators are

$$x^* = \exp\left(\frac{1}{n} \sum_{i=1}^n \log(x_i)\right) = \left(\prod_{i=1}^n x_i\right)^{\frac{1}{n}} \quad (1.7)$$

$$s^* = \exp\left(\left[\frac{1}{n-1} \sum [\log(\frac{x_i}{x^*})]^2\right]^{\frac{1}{2}}\right) \quad (1.8)$$

The lognormal distribution is skewed with mean  $e^{\mu+\frac{1}{2}\sigma^2}$ , median  $e^\mu$  and mode  $e^{\mu-\sigma^2}$ . It has finite mean and variance, in contrast to the power-law distribution.

Despite it has finite moments, the lognormal distribution can be similar to power-law. If  $X$  has a lognormal distribution then loglog plot of density function can appear as straight line for a large portion of a body of distribution. If the variance is large, the distribution may appear linear on log-log plot for several orders of magnitude. The variance of the corresponding normal distribution is large, the distribution may appear linear on a log-log plot. To see this we can check the logarithm of density function. If  $\sigma$  is large then the quadratic term will be small for large range of  $x$  values, so the logarithm of the density function will appear almost linear for large range of values. Recall that normal distribution have property that the sum of two independent normal variables is normal variable. It follows that product of two lognormally distributed random variables also has a lognormal distribution.

## 1.7 Scale-free networks

The study of scale ivariance has a long tradition. Among the fields where this property was analysed were the theory of critical phenomena, percolations and fractal geometry. One of

the first examples considered eas the price fluctuations of cotton in commodities market (Mandelbort, 1963). The future price can not be obtained with arbitary precision from past series., still this series have some form of regularity. The curves for daily, weakly and montly price fluctuations are statistically similar. The fact that some features are found at different time scales is typical sign of fractal behaviour. Similarly in the case of coastline lenght we find fractal behavior. If we try to measure the total lenght, the real shape is so complicated that we always miss some part.

Fractal behaviour might refer to different properties. In some systems scale-free structure is in shape. In this class the fractal shape can be robust, as in the case of branched patterns or electric breakdown. We say robust because these phenomena happen for varaity of external conditions. In the same class we have other systems that are more fragile, in the sense that they arise after precise tuning of some physical quantity. This is in the case of percolations and critical phenomena. Scale-free invariance may be related with dynamics or evolution of the system. The time activity of the system may display self-similar behaviour. The only sign of fractal behaviour is the mathematical form, power-law fluctuations of the time-series.

The self-similarity can be present in the way the different parts of a system interact with each other. This is the case with self-similar graphs and the power-law scaling appers in the distribution of topological quantities like the number of interactions per part of the system. These phenomena are fractals in the topology.

## 1.8 Preferential attachment

One of the most successful applications of multiplicative processes is given by preferential attachment. To date, this is the most successful mechanism adopted in the study of growing networks. Interestingly, the idea that we are going to explain has been independently rediscovered several times in different fields and ages. Precisely for this reason it has also been given several names. For example: Yule Process, Matthew effect, Rich gets richer, Preferential Attachment, Cumulative advantage. In the community there is some agreement (Mitzenmacher, 2004; New- man, 2005) that the first to present this idea has been G. Yule (1925) in order to explain the relative abundance of species and genera in biological taxonomic trees. As shown in Chapter 8 when considering a set of biolog- ical species we have that the classification (taxonomic) tree has scale-free properties. The null hypothesis consists in considering that the set of species arises from a common evolution. Therefore we consider one parent species and after mutation we obtain a new one that very likely can be grouped in the same genus. Every now and then though, speciated species (the new ones) can be so different from the parent one that they can form a new genus on their own (or be grouped in an existing different one). The probability of speciating will be larger for genera that are already large, since mutation rate is constant for any individual. This explanation allow us to focus on the two ingredients of the model. Firstly you have to invoke a certain a priori dynamics (hereafter called growth). Secondly, this dynamics selects successful elements and

makes them even more successful (hereafter called preferential attachment). In detail, take a set of elements each of which is characterized by a certain number  $N_{i,t}$ . As a possible example this could be the number of different genera that have  $i$  species per genera. The set can also be a set of vertices in a graph and the number  $N_{i,t}$  can represent the number of pages whose in-degree is  $i$ . Now let us introduce a rule that introduces new elements in the set; these elements will not be shared equally between the older ones, but rather will be assigned more to those that already have many. Let us consider that  $N_{i,t}$  gives the number of vertices with certain degree  $i$  (the total

## 1.9 In This Thesis

# Chapter 2

## Driving signals

The complex networks grow through the addition of new nodes, and growing networks models consider that growth is constant over time. This approximation is sufficient for explaining how properties of complex networks can emerge; for example, in the Barabasi-Albert model such as in real systems, we find scaling of degree distribution. Models mostly focus on linking rules and their influence on the topology of complex networks.

Still, the growth of real systems changes over time. In online social networks, new users join on daily basis and the users' activity might have bursty nature. We can consider a co-authorship network, where links are created between scientists when they publish a paper. The dynamic of real networks can be complex and highly influenced by non-linear signals. The growth signal; the number of new nodes in each time step; has cycles and trends. Circadian cycles are directly reflected into growth signals and we also find long-range correlations and multifractal properties.

In this chapter, we explain the properties of growth signals, both real and computer-generated. We analyze networks created with a growing network model where the interplay between ageing and preferential attachment shape their structure. We are interested to incorporate non-constant growth signals into the model and measure their impact on the complex networks. Differences between networks with the same number of nodes and links can be observed through connectivity patterns. Figure 2.1 describe used model.

### 2.1 Growing signals

#### Long range correlated signals

The main characteristic of long-range correlated time series is power law decay of autocorrelation function,  $C(s) = \langle x_i x_{i+1} \rangle = s^{-\delta}$ . Instead of using correlation function to directly determine type correlations in the signal, in practice is more common to calculate Hurst exponent.

Hurst exponent is used for estimating self-similarity of the time series described with relation  $x(ct) = cHx(t)$ . Hurst exponent and autocorrelation coefficient  $\delta$  are connected as

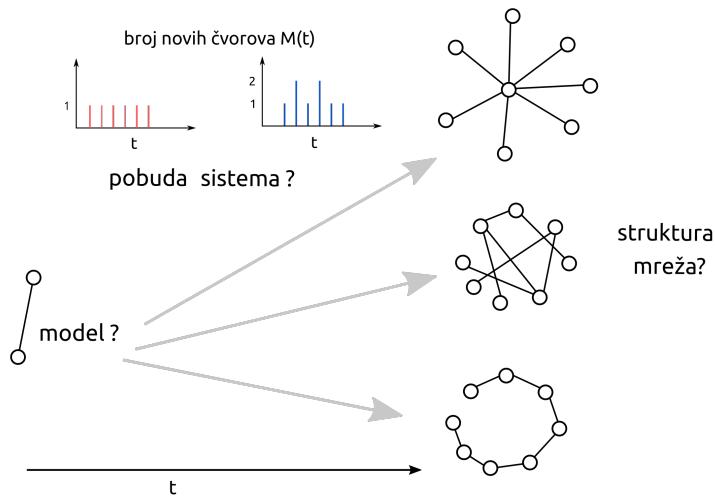


Figure 2.1: Growing network model schema.

$H = 1 - \frac{\delta}{2}$ . When  $H = 0.5$  signal has short range correlations and is considered to be white noise, while for  $H = 1.0$  signal is pink noise. Between this limits  $0.5 < H < 1.0$ , signal has long range correlations.

Monofractal signals can be generated using Fourier transform method [6]:

- first generate one-dimensional sequence of uncorrelated random numbers  $u_i$  from Gaussian distribution with  $\sigma = 1$
- calculate the Fourier transform of the generated sequence,  $u_q$
- filter signal  $x_q = u_q s$ , where  $s$  is Fourier transform of autocorrelation function  $C(s)$
- the inverse Fourier transform  $x_i$  is signal with specific long range correlations

Figure 2.2 shows artificial signals generated using Fourier transform method for different values of Hurst exponents. The obtained signals are round to integers, as in real time series integer values are present. The mean values of signals are close to 4.

For estimation of Hurst exponent from non-stationary signal can be used detrended fluctuation analysis (DFA) [7] [8]. This method removes trends and cycles from the signal, while Hurst exponent is estimated based on residual fluctuations. Signals from real world have usually multifractal structure and can not be described with only one value of Hurst exponent [9]

### Multifractal analysis

Multifractal detrended fluctuation analysis (MF-DFA) [9, 10] to estimate multifractal Hurst exponent  $H(q)$ . For given time series  $\{x_i\}$  with length  $N$ , first we define global profile in the

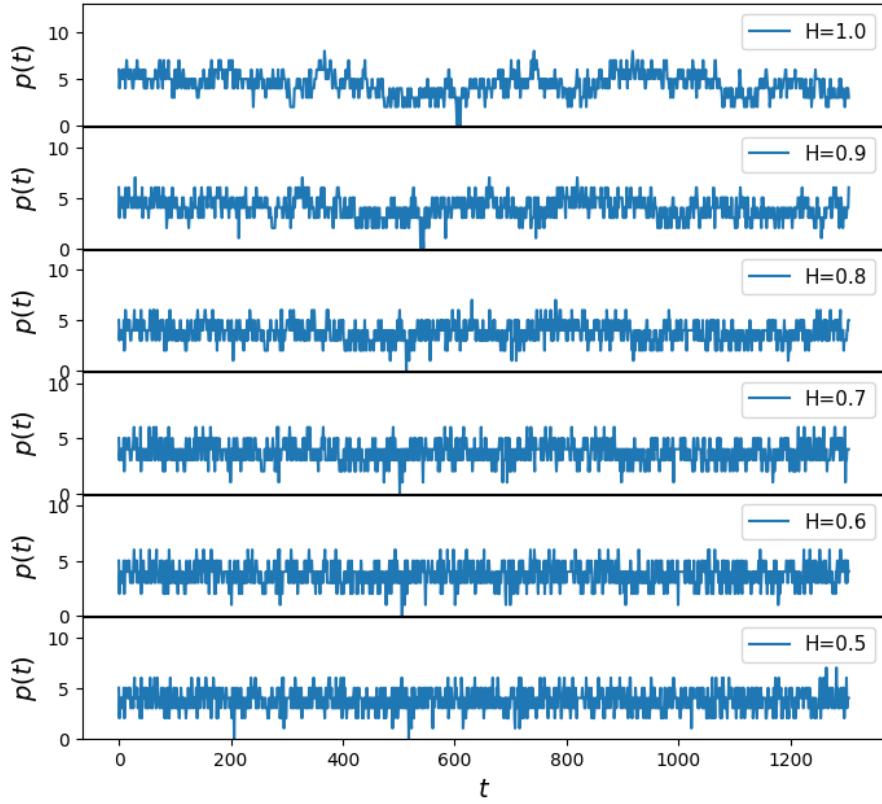


Figure 2.2: Monofractal signals

form of cumulative sum, equation 2.1, where where  $\langle x \rangle$  represents average of the time series:

$$Y(j) = \sum_{i=0}^j (x_i - \langle x \rangle), \quad j = 1, \dots, N \quad (2.1)$$

Subtracting the mean of the time series is supposed to eliminate global trends. The profile of the signal Y is divided into  $N_s = \text{int}(N/s)$  non overlapping segments of length s. If N is not divisible with s the last segment will be shorter. This is handled by doing the same division from the opposite side of time series which gives us  $2N_s$  segments. From each segment  $\nu$ , local trend  $p_{\nu,s}^m$  - polynomial of order m - should be eliminated, and the variance  $F^2(\nu, s)$  of detrended signal is calculated as in equation 2.2:

$$F^2(\nu, s) = \frac{1}{s} \sum_{j=1}^s [Y(j) - p_{\nu,s}^m(j)]^2 \quad (2.2)$$

Then the q-th order fluctuating function is:

$$F_q(s) = \left\{ \frac{1}{2N_s} \sum_{\nu}^{2N_s} \left[ F^2(\nu, s) \right]^{\frac{q}{2}} \right\}^{\frac{1}{q}}, q \neq 0$$

$$F_0(s) = \exp \left\{ \frac{1}{4N_s} \sum_{\nu}^{2N_s} \ln \left[ F^2(\nu, s) \right] \right\}, q = 0 \quad (2.3)$$

The fluctuating function scales as power-law  $F_q(s) \sim s^{H(q)}$  and the analysis of log-log plots  $F_q(s)$  gives us an estimate of multifractal Hurst exponent  $H(q)$ . Multifractal signal has

## 2. Driving signals

---

different scaling properties over scales while monofractal is independent of the scale, i.e.,  $H(q)$  is constant.

### Real signals

In this work, we use two different growth signals from real systems figure 1: (a) the data set from TECH community from Meetup social website [36] and (b) two months dataset of MySpace social network [37]. TECH is an event-based community where members organize offline events through the Meetup site [36]. The time unit for TECH is event since links are created only during offline group meetings. The growth signal is the number of people that attend the group's meetings for the first time. MySpace signal shows the number of new members occurring for the first time in the dataset [37] with a time resolution of one minute. The number of newly added nodes for the TECH signal is  $N = 3217$ , and the length of the signal is  $T = 3162$  steps. We have shortened the MySpace signal to  $T = 20221$  time steps to obtain the network with  $N = 10000$  nodes.

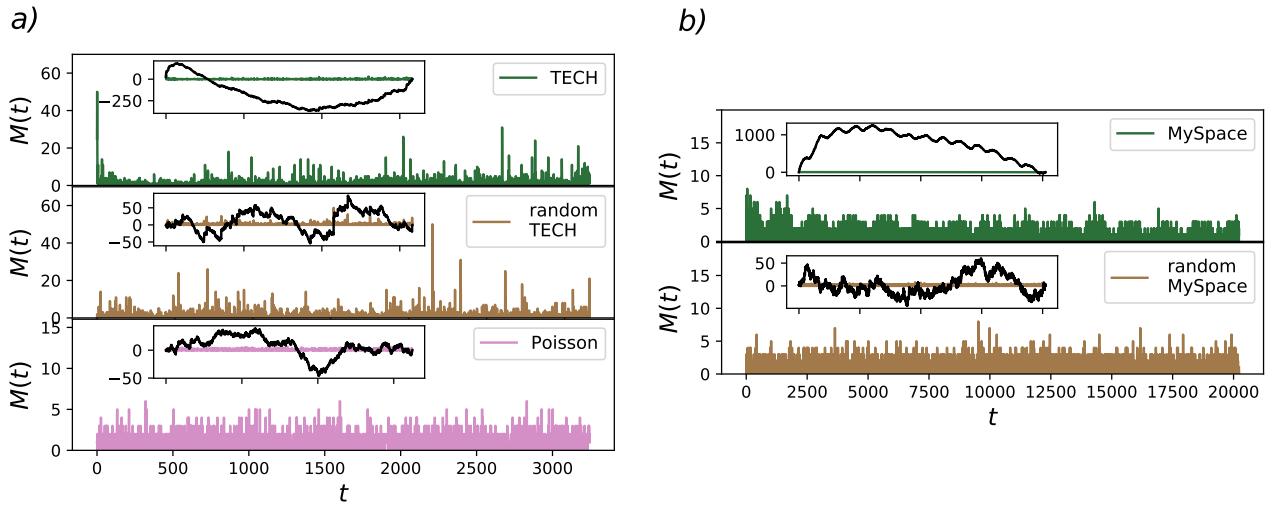


Figure 2.3: Growth signals for TECH (a) and MySpace (b) social groups, their randomized counterparts, and random signal drawn from Poissonian distribution with mean 1. The cumulative signals are shown in insets.

Real growth signals have long-range correlations, trends and cycles [37, 27, 25]. We also generate networks using randomized signals and one computer-generated white-noise signal to explore the influence of these signal's features on the structure of evolving networks. We randomize real signals using reshuffling procedure and keep their length and mean value, the number of added nodes, and probability density function of fluctuations intact, but destroy cycles, trends, and long-range correlations. Besides, we generate a white-noise signal from a Poissonian probability distribution with a mean equal to 1. The length of the signal is  $T = 3246$ , and the number of added nodes in the final network is the same as for the TECH signal.

Figures 2.3 (a) and 2.4 show that the TECH signal has long trends and a broad probability density function of fluctuations. The trends are erased from the randomized TECH signal,

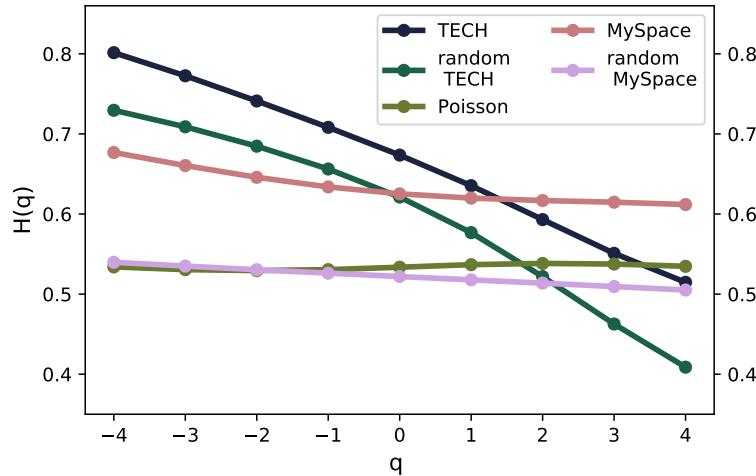


Figure 2.4: Dependence of Hurst exponent on parameter  $q$  for all five signals shown in figure 2.3 obtained with MFDFA.

but the broad distribution of the signal and average value remain intact. MFDFA analysis shows that real signals have long-range correlations with Hurst exponent approximately 0.6 for  $q = 2$ , figure 2.4. The TECH signal is multifractal, the consequence of both broad probability distribution for the values of time series and different long-range correlations of the intervals with small and large fluctuations. Shuffling of the time series does not destroy the broad distribution of values, the reason for the persistent multifractality of the TECH randomized signal, figure 2.4.

MySpace signal has a long trend with additional cycles that are a consequence of human circadian rhythm, figure 2.3(b). It is multifractal for  $q < 0$ , and has constant value of  $H(q)$  for  $q > 0$ , figure 2.4. In MFDFA, with negative values of  $q$ , we put more emphasis on segments with smaller fluctuations, while for positive  $q$  emphasis is more on segments with larger fluctuations [10]. Segments with smaller fluctuations have more persistent long-range correlations in both real signals, see figure 2.4. Randomized MySpace signal and Poissonian signal are monofractal and have short-range with  $H = 0.5$  correlations typically for white noise.

## 2.2 Growing network model with aging nodes

The model starts with small number of nodes randomly connected. Further, at each time step new node arrives in the network and makes connection with one old node, already present in the network. The way in which new nodes are linked is governed by various mechanisms. They can have preference to nodes with high degree (preferential attachment), or preference to nodes with specific age. In the network model with aging nodes the probability that link is created between two nodes depends on the node degree and age [11]:

$$\Pi_i(t) \sim k_i(t)^\beta \tau_i^\alpha \quad (2.4)$$

## 2. Driving signals

---

where  $k_i(t)$  is a degree of a node  $i$  at time  $t$ , and  $\tau_i$  is age difference between node  $i$  and newly added node.

The values of model parameters  $\beta$  and  $\alpha$  control the topology of networks. For example if we fix  $\beta = 1$  and  $\alpha = 0$  generated networks are scale-free; degree distribution is  $P(k) \sim k^{-\gamma}$  with  $\gamma = 3$ , while in the case of nonlinear preferential attachment  $\beta \neq 1$  and  $\alpha = 0$  scale-free properties disappear. Scale-free property can be produced along the critical line  $\beta(\alpha^*)$  in the  $\alpha - \beta$  phase diagram, see Figure 2.5. For  $\alpha > \alpha^*$  networks have gel-like small world behavior, while for  $\alpha < \alpha^*$  but close to line  $\beta(\alpha^*)$  networks have stretched exponential shape of degree distribution [11].

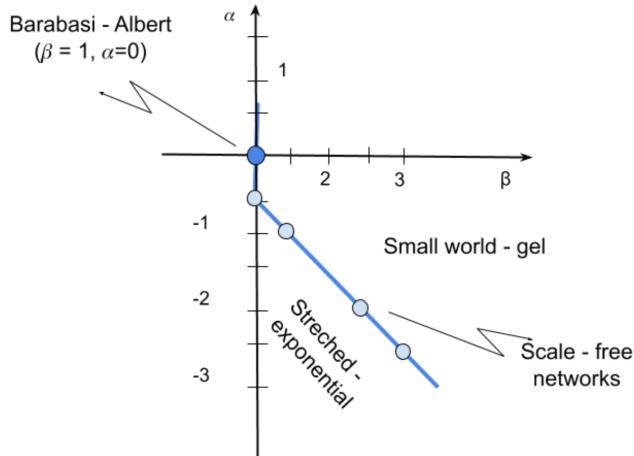


Figure 2.5: Phase diagram of aging network model

The networks generated with constant growth signal are uncorrelated trees. To enable formation of clusters in the network new nodes need to create more than one link. We adapt the original model such that at each time step we add  $M \geq 1$  new nodes that make  $L \geq 1$  links with existing nodes in the network corresponding to probability 2.4. The master equation for  $N_k$ ,  $k$  degree nodes can be written as:

$$\partial_t N_k = \sum_{j=1}^{M(t)} r_{k-j \rightarrow k} N_{k-j} - \sum_{j=1}^{M(t)} r_{k \rightarrow k+j} N_k + M(t) \delta_{k,L}. \quad (2.5)$$

At each time step we add  $M(t)$  nodes with  $L$  links. As multiply links between two nodes are not allowed, we'll get  $M(t)$  new nodes with degree  $L$ , that describes third term in the equation. Old nodes can increase their degree from 1 to  $M(t)$ , as same node can be chosen by different new nodes. The first term in the equation describes nodes with degree  $k \in \{k - M(t), \dots, k - 1\}$  that getting degree  $k$ , while in second term nodes with degree  $k$  entering degree  $k \in \{k + 1, \dots, k + M(t)\}$ . The quantities  $r_{k-j \rightarrow k}$  and  $r_{k \rightarrow k+j}$  are the rates that express the transitions of a node from class with degree  $k - j$  to one with degree  $k$  and from class with degree  $k$  to class with degree  $k + j$  respectively.

The equation 2.5 is not solvable in a general case. It was solved for the case  $M(t) = 1$  and specific values of parameters  $\alpha$  and  $\beta$  using continuous approach [12]. In this work, we use

numerical simulations to explore the case when  $M(t)$  is a correlated time-varying function and study how these properties influence the structure of generated networks for different values of parameter  $-\infty < \alpha \leq -1$  and  $\beta \geq 1$  and constant  $L$ .

## 2.3 Structural differences between networks

### D-measure

Between two nodes in the network, we can define different paths, but the most important one are the shortest paths,  $d_{ij}$ . Diameter defines the largest shortest path found in the network. For each node  $i$  we can define the distribution of the shortest paths between node  $i$  and all others nodes in the network,  $P_i = \{p_i(j)\}$ , where  $p_i(j)$  is percent of nodes at distance  $j$  from node  $i$ . The connectivity patterns can efficiently describe difference between two networks. To specify how much  $G$  and  $G'$  are similar we use D-measure [13]

$$D(G, G') = \omega \left| \sqrt{\frac{J(P_1, \dots, P_N)}{\log(d)}} - \sqrt{\frac{J(P'_1, \dots, P'_N)}{\log(d')}} \right| + (1 - \omega) \sqrt{\frac{J(\mu_G, \mu_{G'})}{\log 2}}.$$

D-measure calculates Jensen-Shannon divergence between  $N$  shortest path distributions,  $J(P_1, \dots, P_N) = \sum_{i,j} p_i(j) \log\left(\frac{p_i(j)}{\mu_j}\right)$ , where  $\mu_j = (\sum_{i=1}^N p_i(j))/N$  is mean shortest path distribution. The first term in equation 2.6 compares local differences between two networks, and Jensen-Shannon divergence between  $N$  shortest path distributions  $J(P_1, \dots, P_N)$  is normed with network diameter  $d$ . The second part determines global differences, computing  $J(\mu_G, \mu_{G'})$  between mean shortest path distributions. We consider equally important local and global properties of the networks, and parameter  $\omega$  is set to 0.5. The D-measure ranges from 0 to 1. The lower D-measure is, networks are more similar and for D-measure  $D = 0$ , structures are isomorphic.

The advantage this measure has is that it can distinguish between networks generated with the same model parameters. To examine how different growth signals influence the network structure, we use D-measure and compare networks generated with the same model parameters  $\alpha, \beta$  and fixed number of links per new node  $L$ , but different growth signals. The growth of first network is driven by fluctuating signal  $M_1 = M(t)$ , while the other one grows by constant rate  $M_2 = \langle M(t) \rangle = const$ .

We focus on the region of model phase diagram with negative  $\alpha$  and positive  $\beta$  as there is found the transition line from stretched-exponential across scale-free to the small world- gel networks. We take range of parameters  $-3 \leq \alpha \leq -0.5$  and  $1 \leq \beta \leq 3$  with steps 0.5 and we also vary the the number of links each new node can create  $L \in 1, 2, 3$ . For each combination of  $(\alpha, \beta, L)$  we generate the sample of 100 networks, and compare the structure of network grown with fluctuating and the constant signal. The results represented by D-measure are obtained averaging the D-measure between all possible pairs of generated networks.

## 2. Driving signals

---

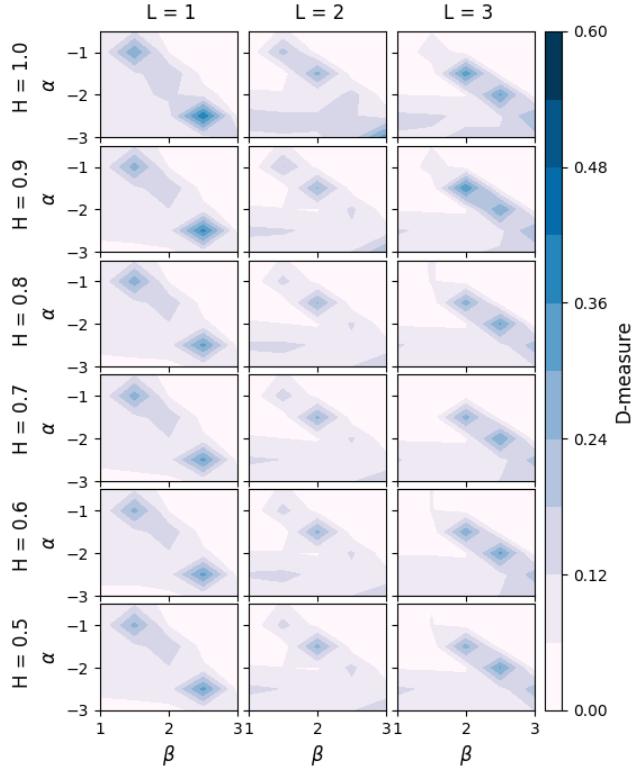


Figure 2.6: D-distance between networks generated with different long-range correlated signals with fixed value of Hurst exponent and networks generated with constant signal  $M=4$ .

First, we explore how monofractal signals, see Figure 2.2 shape the structure of complex networks. The D-measure between networks grown with monofractal signal, with  $H \in \{0.5, 0.6, 0.7, 0.8, 0.9, 1.0\}$  and constant signal  $M = 4$  are shown in figure 2.6. The higher values of D-measure are found in the region of critical line  $\beta(\alpha^*)$ . The most considerable influence is on networks with scale-free distribution. Comparing D-distance in only one point of phase diagram, for example  $L = 1, \alpha = -2.5, \beta = 2.5$ , we find correlations in the signal (Hurst exponent is larger), make bigger impact on the network structure. D-measure between networks grown with signal with Hurst exponent  $H = 1.0$  and constant signal is  $D(H = 1.0, M = 4) = 0.405$ , while between networks grown with signal with  $H = 0.8$  and constant signal is  $D(H = 0.8, M = 4) = 0.316$ . For  $\alpha > \alpha^*$  networks have similar structural properties and D-measure is close to 0. In the region of networks with stretched exponential degree distribution  $\alpha < \alpha^*$  differences are small.

For signals from real communities we find non-zero values of D-measure 2.7. The largest difference between networks is as before along critical line  $\beta(\alpha^*)$ , for scale free network. For values  $\beta < \beta(\alpha^*)$  the structural differences exist, but they become smaller. In the region of gel small world networks  $\alpha > \alpha^*$  structural differences are small and close to zero. In the region around critical line we find that D-measure depends on the properties of the signal. Multifractal signals TECH has the largest impact on network structure; the maximum obtained value of D-measure is  $D_{max} = 0.552$ . Similar behavior we discover for other multifractal signals, random TECH and MySpace. For networks generated with

uncorrelated signals, random mySpace and Poisson, difference exists but it is much smaller and comparable with monofractal signals.

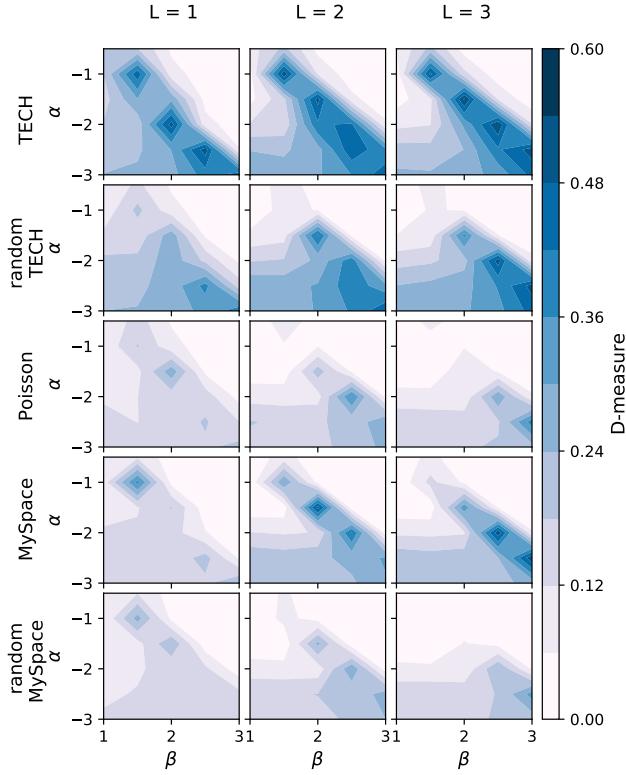


Figure 2.7: The comparison of networks grown with growth signals shown in figure 2.3 versus ones grown with constant signal  $M = 1$ , for value of parameter  $\alpha \in [-3, -1]$  and  $\beta \in [1, 3]$ .  $M(t)$  is the number of new nodes, and  $L$  is the number of links added to the network in each time step. The compared networks are of the same size.

The position of the critical line slightly moves toward larger  $\beta$  with higher link density  $L$ . The addition of more than one node does not influence its position. Although, for fixed network density, we find a critical line independent of the growth signal's properties as can be seen in Figures 2.7, 2.6.

We can note that D-measure rises for lower  $\alpha$ . In the case of constant signal, number of nodes added to the network is equal for each time step, so at time interval  $T$  the network has  $MT$  nodes. In fluctuating signal the number of nodes added during time interval  $T$  vary with time. In signals, such as TECH, where are present peaks in the number of new users, emergence of hubs happens faster. As we decrease the parameter  $\alpha$ , fluctuations present in the signal become more important and emergence of the hubs happens even for uncorrelated signals. The trends present in the real signals further promote the emergence of hubs in the network.

## The assortativity and clustering

We further explore the assortativity index and clustering coefficient of networks generated with monofractal signals with different values of Hurst exponent. We show results for several

## 2. Driving signals

---

ageing model parameters to show the difference between network this model can produce, 2.8. All networks are disassortative, with a negative degree-degree correlation index. For the values of parameters below critical line,  $\alpha = -2.5, \beta = 1.5$   $r$  does not depend on the Hurst exponent. Above the critical line are small-world networks, and they are disassortative with a minimum value of assortativity index  $r = -1$ , for  $L = 1$ , indicating the presence of a hub that connects to many nodes. The assortativity index slightly grows with link density.

In the region of critical parameters, the assortativity index depends on the value of the Hurst exponent. The larger influence on the assortativity index have correlated signals, with Hurst exponent  $H > 0.8$ , so networks become more disassortative, see line for parameters  $L = 1, \alpha = -2.5, \beta = 2.5$  in Figure 2.8. The long-range correlations have a stronger effect on the evolution of networks with lower density.

We calculate the mean clustering coefficient, Figure 2.8. For  $L = 1$  networks are uncorrelated trees, with clustering coefficient 0. For network density  $L > 1$ , nodes are organized into clusters. Under the critical line, for parameter  $L = 3, \alpha = -2.5, \beta = 1.5$ , clustering coefficient is constant and low. Similar values are obtained for clustering coefficient for critical parameters  $L = 3, \alpha = -1.5, \beta = 2.0$ , but for Hurst exponent  $H > 0.8$  clustering coefficient increase. Small world networks,  $L = 3, \alpha = -1.5, \beta = 2.5$  are clustered, the value of  $\langle c \rangle$  is high. The value of clustering for networks created with the constant signal is 0.8. Networks grown with white noise signal and signal with  $H=0.6$  have higher values of the clustering, while networks grown with signals that have Hurst exponent larger than 0.6 have the same value of clustering, which is below 0.8.

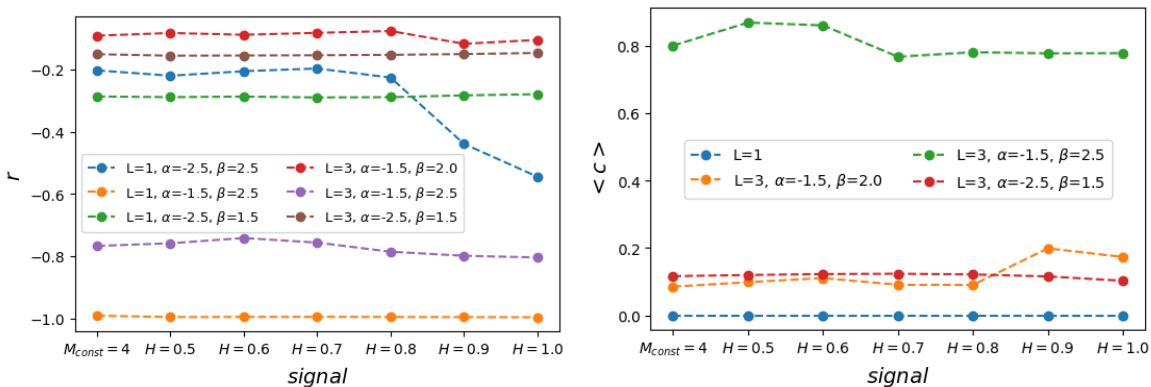


Figure 2.8: Aindex

We examine degree distribution, degree correlations and clustering coefficient of networks generated by real signals, as researchers has shown that these measures provide the sufficient set for describing structure of complex network. D-measure showed that multifractals have larger influence on networks than monofractals, especially on scale-free networks.

Figure 2.9 shows properties of networks generated with model parameters  $L = 2, \alpha = -1.0, \beta = 1.5$ , that lies on critical line. The degree distributions  $P(k)$  of networks generated with real signals TECH and MySpace have emergence of super-hubs. Degree distributions generated with randomized signals and white noise signal do not differ from

degree distribution of networks generated with constant signal. Networks generated with real signals average neighbouring degree  $\langle k \rangle_{nn}(k)$  and clustering coefficient  $c(k)$  depend on node degree, while in networks generated with constant and randomized signals they weakly depend on the degree  $k$ .

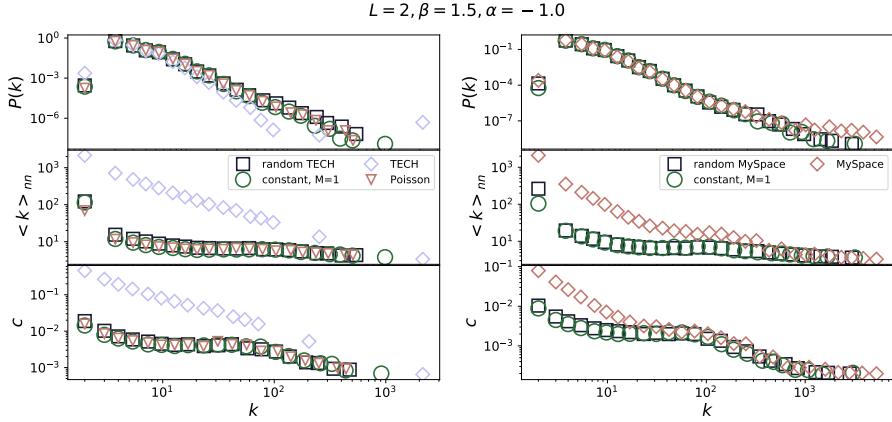


Figure 2.9: Degree distribution, the dependence of average first neighbor degree on node degree, dependence of node clustering on node degree for networks grown with different time-varying and constant signals. Model parameters have value  $\alpha = -1.0$ ,  $\beta = 1.5$  and  $L = 2$  for all networks. The networks are from scale-free class.

We also find structural differences between networks, obtained with model parameters under the critical line  $\alpha < \alpha^*$ , see Figure 2.10. The difference is mostly found for TECH signal. Degree distribution  $P(k)$  shows emergence of hubs in networks grown with TECH signal, while the randomized and Poisson signal are more similar to networks grown with constant signal. MySpace signal; whose generalized Hurst exponent  $H(q)$  weakly depends on scale parameter  $q$  and whose long-range correlations and trends are easily destroyed; do not influence the structure of networks more than constant or randomized signal.

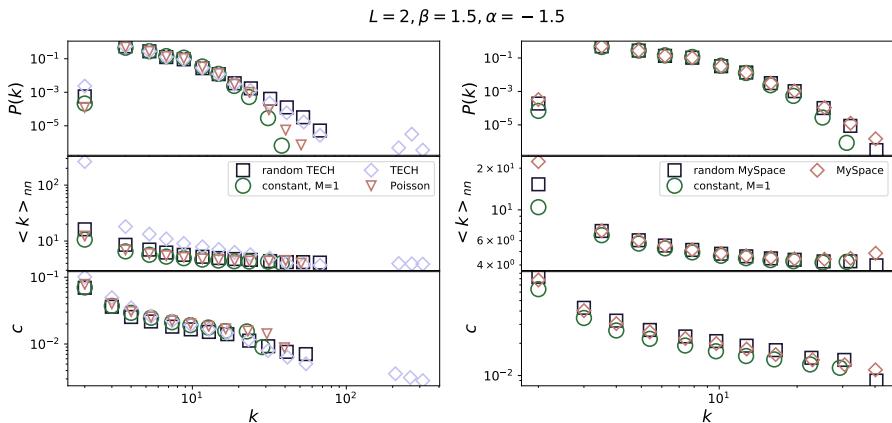


Figure 2.10: Degree distribution, the dependence of average first neighbor degree on node degree, dependence of node clustering on node degree for networks grown with different time-varying and constant signals. Model parameters have value  $L = 2, \alpha = -1.5, \beta = 1.5$ . The networks have stretched exponential degree distribution.

## 2. Driving signals

---

The properties of time-varying signal do not influence the topological properties of small-world gel networks, Figure 2.11. Here model promote existence of hubs. As this is mechanism through which the fluctuations alter the structure of evolving networks, the properties of the signal are not relevant.

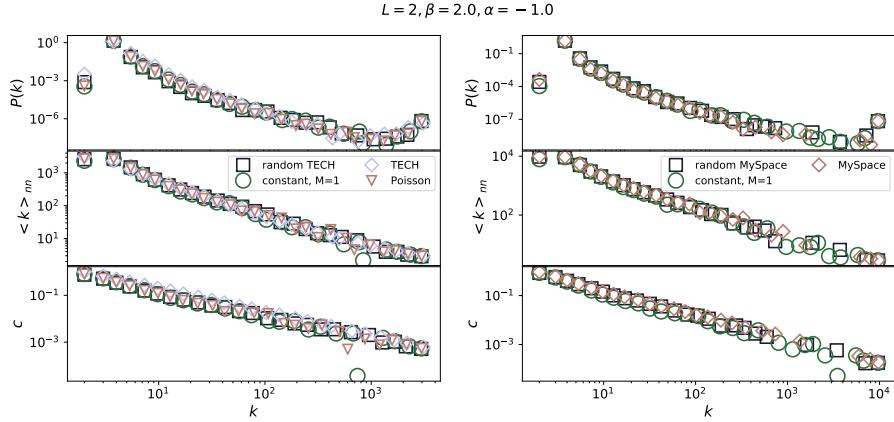


Figure 2.11: Degree distribution, the dependence of average first neighbor degree on node degree, dependence of node clustering on node degree for networks grown with different time-varying and constant signals. Model parameters have value  $L = 2, \alpha = -1.0, \beta = 2.0$ . Generated networks have scale-free properties.

## 2.4 Conclusions

We demonstrate that the resulting networks' structure depends on the features of the time-varying signal that drives their growth. The previous research [14, 15] indicated the possible influence of temporal fluctuations on network properties. Our results show that the temporal properties of growth signals generate networks with power-law degree distribution, non-trivial degree-degree correlations, and clustering coefficient even though the local linking rules, combined with constant growth, produce uncorrelated networks for the same values of model parameters [11].

We observe the most substantial dissimilarity in network structure along the critical line, the values of model parameters for which we generate networks with broad degree distribution. Figure 2.7 shows that dissimilarity between networks grown with time-varying signals and ones grown with constant signals always exists along this line regardless of the features of growth signal. However, the magnitude of this dissimilarity strongly depends on these features. We observe the largest structural difference between networks grown with multifractal TECH signal and networks that evolve by adding one node in each time step. The identified value of D-measure is similar to one calculated in the comparison between sub-critical and super-critical Erdős–Rényi graphs [13] indicating the considerable structural difference between these networks. Our findings are further confirmed in figure ??(b). The networks generated with signals that have trends and long-range temporal correlations differ the most from those grown with the constant signal. Our results show that even white-noise type signals can generate networks significantly different from ones created with constant signal for low values of  $\alpha^*$ .

The value of D-measure declines fast as we move away from the critical line, figure 2.7. The main mechanism through which the fluctuations influence the structure of evolved networks is the emergence of hubs and super hubs. For values of  $\alpha \ll \alpha^*$ , the nodes attach to their immediate predecessors creating regular networks without hubs. For  $\alpha \sim \alpha^*$  graphs have stretched exponential degree distribution with low potential for the emergence of hubs. Still, multifractal signal TECH enables the emergence of hub even for the values of parameters for which we observe networks with stretched-exponential degree distribution in the case of constant growth figure ??(a). By definition, small-world gels generated for  $\alpha > \alpha^*$  have super-hubs [11] regardless of the growth signal, and therefore the effects that fluctuations produce in the growth of networks do not come to the fore for values of model parameters in this region of  $\alpha - \beta$  plane.

Evolving network models are an essential tool for understanding the evolution of social, biological, and technological networks and mechanisms that drive it [1]. The most common assumption is that these networks evolve by adding a fixed number of nodes in each time step [1]. So far, the focus on developing growing network models was on linking rules and how different rules lead to networks of various structural properties [1]. Growth signals of real systems are not constant [15, 14]. They are multifractal, characterised with long-

## 2. Driving signals

range correlations [15], trends and cycles [16]. Research on temporal networks has shown that temporal properties of edge activation in networks and their properties can affect the dynamics of the complex system [17]. Our results imply that modeling of social and technological networks should also include non-constant growth and that its combination with local linking rules can significantly alter the structure of generated networks.

# Chapter 3

## Groups growth model

### 3.1 Introduction

Social groups, informal or formal, are mesoscopic building elements of every socio-economic system that direct its emergence, evolution, and disappearance []. The examples span from countries, economies and science to society in general. Settlements, villages, towns and cities are formal and highly structured social groups of countries. Their organisation and growth determine the functioning and sustainability of every society [18]. Companies are the building blocks of an economic system and their dynamics are important indicators of the level of its development [19]. Scientific conferences, as scientific groups, enable fast dissemination of the latest results, exchange and evaluation of ideas as well as a knowledge extension, and thus are an integral part of science [20]. The membership of individuals in various social groups, online and offline, can be essential when it comes to the quality of their life [21, 22, 23]. Therefore, it is not surprising that the social group emergence and evolution are at the center of the attention of many researchers [24, 25, 26, 27].

Along with massive data sets comes the need to develop methods and tools for their analysis and modeling. Methods and paradigms from statistical physics have proven to be very useful in studying the structure and dynamics of social systems [28]. The main argument for using statistical physics to study social systems is that they consist of many interacting individuals. Due to this, they exhibit different patterns in their structure and dynamics, commonly known as *collective behavior*. While building units of a social systems can be characterized by many different properties, only few of them enforce collective behavior in the systems. The phenomenon is known as *universality* in physics and is commonly observed in social systems such as in voting behavior [29], or scientific citations [30]. It indicates the existence of the universal mechanisms that govern the dynamics of the system [].

The availability of large-scale and long-term data on various online social groups has enabled the detailed empirical study of their dynamics. The focus was mainly on the individual groups and how structural features of social interaction influence whether

### 3. Groups growth model

---

individuals will join the group [31] and remain its active members [20, 32]. The study on LiveJournal [31] groups has shown that decision of an individual to join a social group is greatly influenced by the number of her friends in the group and the structure of their interactions. The conference attendance of scientists is mainly influenced by their connections with other scientists and their sense of belonging [20]. The sense of belonging of an individual in social groups is achieved through two main mechanisms [32]: expanding of the social circle at the beginning of joining the group and strengthening of the existing connections in the later phase. The dynamics of social groups depend on their size [32]. Analysis of the evolution of large-scale social networks has shown that edge locality plays a critical role in the evolution of social networks [33]. Small groups are more cohesive with continued membership, while large groups tend to change their active members constantly [33]. These findings help us understand the growth of a single group, the evolution of its social network, and the influence of the network structure on the group growth. However, how the growth mechanisms influence the distribution of members of one social system among groups is still anecdotal.

Furthermore, it is not clear whether the growth mechanisms of social groups are universal or system-specific. The size distribution of social groups has not been extensively studied. Rare empirical evidence of the size distribution of social groups indicates that it follows power-law behavior [34]. However, the distribution of company sizes follows log-normal behavior and remains stable over decades [35, 36]. Analysis of the sizes of the cities shows that the distribution of all cities also follows a log-normal distribution. In contrast, the distribution of the largest cities resembles Zipf's distribution [37].

A related question that should be addressed is whether we can create a unique yet relatively simple microscopic model that reproduces the distribution of members between groups and explains the differences observed between social systems. French economist Gibrat proposed a simple growth model to reproduce the observed log-normal size distribution of companies and cities. However, the analysis of the growth rate of the companies [35] has shown that growth mechanisms are different from ones assumed by Gibrat. In addition, the analysis of the growth of three online social networks showed that population growth is not determined by the population size and spatial factors, and it deviates from Gibrat's law [38]. Other mechanisms, for instance, growth through diffusion, have been used for modeling and prediction of rapid group growth [39]. However, the growth mechanisms of various social groups and the source of the scaling observed in socio-economic systems remain hidden.

Here we analyze the size distribution of formal social groups in two different systems: Meetup online platform and subreddits on Reddit. We are interested in the scaling behavior of size distributions and the distribution of growth rates. Empirical analysis of the dependence of growth rates, shown in this work, indicates that growth cannot be explained through

Gibrat's model. Here we contribute with a simple microscopic model that incorporates some of the findings of previous research [31, 34]. We show that the model can reproduce size distributions and growth rate distributions for both studied systems. Moreover, the model is flexible and can produce a broad set of size distributions depending on the value of model parameters.

The paper is organized as follows: in Section 3.2 we describe the data, while in Section 3.3 we present our empirical results. In Section 3.4 we introduce model parameter and rules. In section 3.5 we demonstrate that model can reproduce the growth of social groups in both systems and show the results for different values of model parameters. Finally, in Section 3.6, we present concluding remarks and discuss our results.

## 3.2 Data

We analyse the growth of social groups from two widely used online platforms: Reddit and Meetup. Reddit<sup>1</sup> enables sharing diverse web content and members of this platform interact exclusively online through posts and comments. The Meetup<sup>2</sup> allows people to use online tools to organize offline meetings. The building elements of the Meetup system are topic-focused groups, such as food lovers or ICT and data science professionals. Due to their specific activity patterns - events where members meet face-to-face - Meetup groups are geographically localised.

We compiled the Reddit data from <https://pushshift.io/>. This site collects data daily and, for each month, publishes merged comments and submissions in the form of JSON files. Specifically, we focus on subreddits - social groups of Reddit members interested in a specific topic. We select subreddits active in 2017 and follow their growth from their beginning until 2011 – 12. The considered dataset contains 17073 subreddits with 2195677 active members, with the oldest originating from 2006 and the youngest being from 2011. For each post under a subreddit, we extracted the information about the member-id of the post owner, subreddit-id, and timestamp. As we are interested in the subreddits growth in the number of members, for each subreddit and member-id we selected timestamp when member made a post for the first time. Finally, in the dataset we include only subreddits active at least two months.

The Meetup data were downloaded in 2018 using public API. The Meetup platform was launched in 2003, and at the moment we accessed the data, there were more than 240 000 active groups. For each group, we extracted information about the date it had been founded, its location, and the total number of members. We focused on the groups founded from 2003 until 2017 in big cities London and New York, where Meetup platform achieved considerable popularity. We considered groups active at least two months. There were 4673 groups with

---

<sup>1</sup><https://www.reddit.com/>

<sup>2</sup>[www.meetup.com](http://www.meetup.com)

### 3. Groups growth model

---

831685 members in London and 4752 groups with 1059632 members in New York. In addition, we extracted the id of each member in the group and the information about organised events. This allowed us to obtain the date when a member joined a group, which is the first time she attended group event.

In both systems, we approximated the timestamp when the member joined the group. Based on this information, we can calculate the number of new members per month  $N_i(t)$ , the group size  $S_i(t)$  at each time step, and the growth rate for the group in both systems. The time step for both systems is one month. The size of the group  $i$  at time step  $t$  is the number of members that joined that group ending with the month, i.e.,  $S_i(t) = \sum_{k=t_{i0}}^{k=t} N_i(t)$ , where  $t_{i0}$  is the time step in which the group  $i$  was created. We do not consider when a member leaves a group or subreddit since this kind of information is not available to us. For these reasons, the size of considered groups is a non-decreasing function. The growth rate  $R_i(t)$  at step  $i$  is obtained as logarithm of successive sizes  $R = \log(S_i(t)/S_i(t - 1))$ .

While the forms of communication between members and activities that members engage in differ in those two systems, some common properties exist between them. Members can form a new groups and join existing ones in both systems. Furthermore, each member can belong to an unlimited number of groups. For these reasons, we can use the same methods to study and compare the formation of groups in Reddit and Meetup.

## 3.3 Empirical analysis of social group growth

Figure 3.1 summarize properties of the groups in Meetup and Reddit systems. The number of groups grows exponentially over time. Nevertheless, we notice that Reddit has substantially larger number of groups than Meetup. The Reddit groups are prone to engage more members in a shorter period of time. Size of the Meetup groups is in the range from several members up to several tens of thousands of members, while sizes of subreddits are between a few tens of members up to several millions. The distributions of group sizes follows the lognormal distribution

$$P(S) = \frac{1}{S\sigma\sqrt{2\pi}} \exp\left(-\frac{(\ln(S) - \mu)^2}{2\sigma^2}\right), \quad (3.1)$$

where  $S$  is the group size and  $\mu$  and  $\sigma$  are parameters of the distribution. We used package [40] to fit Eq. 3.1 to Reddit and Meetup data and found that distribution of groups sizes for Meetup groups in London and New York follow similar distributions with the values of parameters  $\mu = -0.93$ ,  $\sigma = 1.38$  and  $\mu = -0.99$  and  $\sigma = 1.49$  for London and New York respectively. The distribution of sizes of subreddits also has the log-normal shape with parameters  $\mu = -5.41$  and  $\sigma = 3.07$ . Even though these distributions are from the same class, for subreddits we find broader distribution that may resemble power-law distribution. Our analysis shown in Supportive Information (SI) confirms that the distribution exhibits a log-normal behavior, see SI-Table 1 and SI-Fig. 1.

The log-normal distributions can be generated by multiplicative processes [41]. If there is a quantity with size  $S_i(t)$  at time step  $t$ , it will grow so after time period  $\delta$  the size of

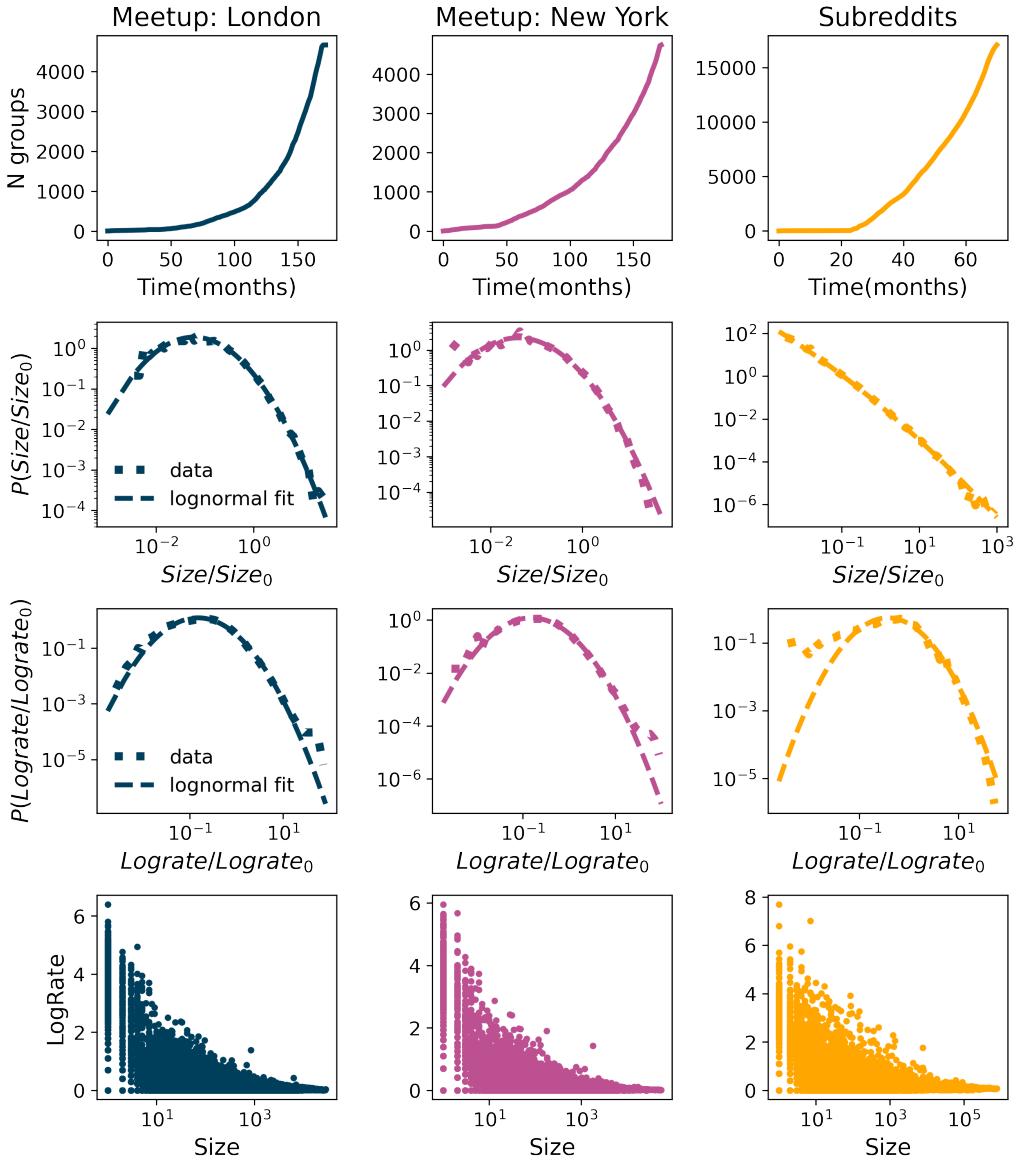


Figure 3.1: The number of groups over time, normalized sizes distribution, normalized log-rates distribution and dependence of log-rates and group sizes for Meetup groups created in London from 08-2002 until 07-2017 that were active in 2017 and subreddits created in the period from 01-2006 to the 12-2011 that were active in 2017.

the quantity is  $S(t + \Delta t) = S(t)r$ , where  $r$  represents a random process. The Gibrat law states that growth rates  $r$  are uncorrelated and do not depend on the current size. In order to describe the growth of social groups, we calculate the logarithmic growth rates defined as  $R = \log \frac{S_t}{S_{t-\Delta t}}$ . According to Gibrat law, the distribution of sizes follow log-normal distribution. For logarithmic growth rates expected distribution is normal, or as it is shown in many studies it is better explained with Laplacian ("tent shaped") distribution [42], [43]. In figure 3.1 we calculate distributions of log-rates. For both systems, log-rates are very well approximated with log-normal distribution. The Fig. 3.1 shows that log-rates depend on the groups size, especially for the smaller and medium size groups. Our empirical analysis implies that the growth of Meetup and Reddit groups violates the basic assumptions of the Gibrat's law [44, 45], and thus, this growth can not be explained as a simple multiplicative

### 3. Groups growth model

---

process.

We are considering a relatively large time period for online groups. The fast expansion of Information Communications Technology (ICT) led to change of how members access online systems. With the use of smartphones the online systems became more available, which led to exponential growth of ICTs systems, figure 3.1 and potential change in the mechanisms that influence growth of social groups in them. For these reasons we aggregate groups according to year they were founded for each of the three data sets and look at the distributions of these sizes in the year 2017 for Meetup groups and 2011 for Reddit. For each year and each of the three data sets we calculate the average size of the groups that were created in a year  $y$   $\langle S^y \rangle$ . We normalize the size of the groups created in year  $y$  with corresponding average size  $s_i^y = S_i^y / \langle S^y \rangle$  and calculate the distribution of the normalized sizes for each year. The distribution of normalized sizes for all years and all data sets is shown in figure 3.2. All distributions exhibit log-normal behavior. Furthermore, the distributions for the same data set and different years follow a universal curve with same value of parameters  $\mu$  and  $\sigma$ . The universal behavior is observed for distribution of normalized log-rates as well, see Fig. 3.2 (bottom panel). These results indicate that growth of the social groups did not change due to increased growth of members in systems. Furthermore, it implies that the growth is independent of the size of the whole data set.

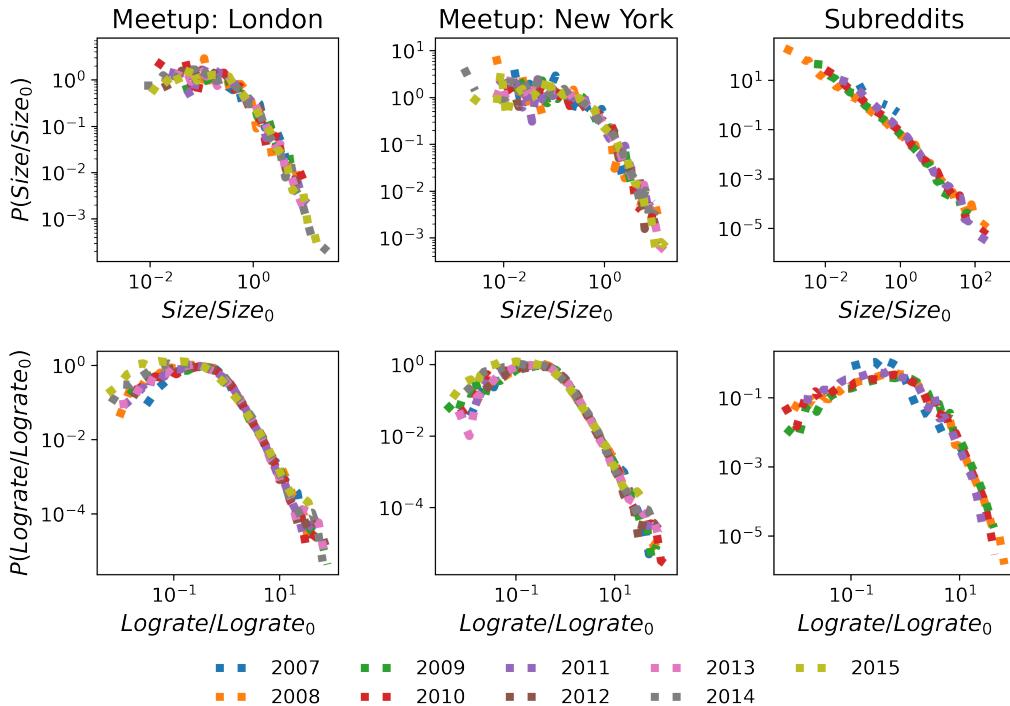


Figure 3.2: The figure shows the groups' sizes distributions and log-rates distributions. Each distribution collects groups founded in the same year and is normalized with its mean value. The group sizes are at the end of 2017 for meetups and 2011 for subreddits.

## 3.4 Model

Growth of social groups can not be explained with the simple rules of Gibrat's law. Previous research on group growth and longevity has shown that social connections with members of a group influence individual's choice to join that group [39, 34]. Moreover, individual's interests and the need to discover new content or activity also influence diffusion of individuals between groups. Furthermore, social systems constantly grow since new members join every minute. The properties of the growth signal that describes the arrival of new members influence both dynamics of the system [46, 47] and the structure of social interactions [48]. Furthermore, number of social groups in the social systems is not constant. They are constantly created and destroyed.

In Ref. [34] authors propose the co-evolution model of the growth of social networks. In this model, authors assume that social system evolves through co-evolution of two networks: network of social contacts between members and network of members' affiliations with groups. This model addresses the problem of growth of social networks that includes both linking between members and social group formation. In this model, a member of a social system selects to join a group either through random selection or according to her social contacts. In the case of random selection, there is a selection preference toward larger groups. If member chooses to select a group according to her social contacts, the group is selected randomly from the list of groups with which her friends are already affiliated.

While the co-evolution model [34] was not created with the intent of studying the growth and size distribution of social groups, authors show that their model is able to reproduce distribution of group sizes for several online social networks that follow power-law distribution. Our empirical analysis, shown in Sec. 3.3 shows that distribution of group sizes is not always power-law, indicating that certain mechanisms proposed in co-evolution model are not universal for all social systems. To fill the gap in understanding how social groups in social system grow, we propose a model of group growth that combines random and social diffusion between groups but following different rules than co-evolution model [34].

Figure 3.3 shows a schematic representation of our model. Similar to co-evolution model [34], we represent social system with two evolving networks, see Fig. 3.3. One network is bipartite network which describes the affiliation of individuals to social groups  $\mathcal{B}(V_U, V_G, E_{UG})$ . This network consists of two partitions, members  $V_U$  and groups  $V_G$ , and set of links  $E_{UG}$ , where a link  $e(u, g)$  between a member  $u$  and a group  $g$  represents the member's affiliation with that group. Bipartite network grows through three activities: arrival of new members, creation of new groups, and through members joining groups. By definition, in bipartite networks links only exist between nodes belonging to different partitions. However, as we explained above, social connections affect whether a member will join a certain group or not. In the simplest case, we could assume that all members belonging to a group are connected with each other. However, previous research on this subject

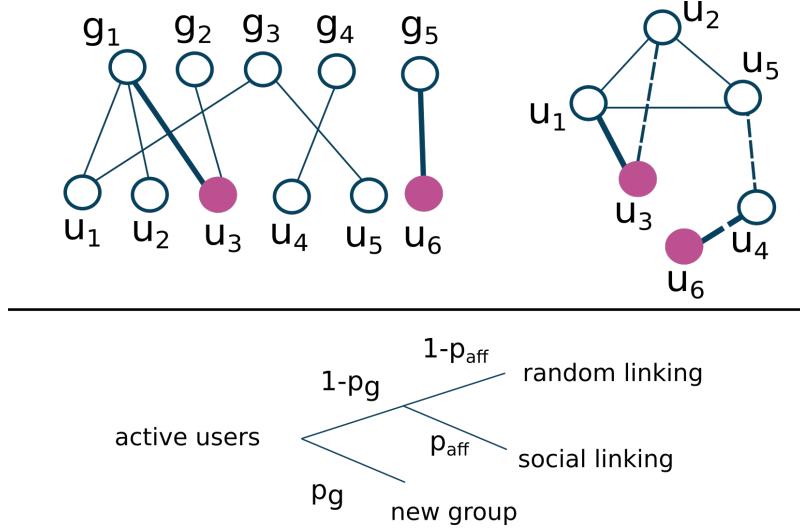


Figure 3.3: The top panel shows bipartite (member-group) and social (member-member) network. Filled nodes are active members, while thick lines are new links in this time step. In the social network dashed lines show that members are friends but still do not share same groups. The lower panel shows model schema. **Example:** member  $u_6$  is a new member. First it will make random link with node  $u_4$ , and then with probability  $p_g$  makes new group  $g_5$ . With probability  $p_a$  member  $u_3$  is active, while others stay inactive for this time step. Member  $u_3$  will with probability  $1 - p_g$  choose to join one of old groups and with probability  $p_{aff}$  linking is chosen to be social. As its friend  $u_2$  is member of group  $g_1$ , member  $u_3$  will also join group  $g_1$ . Joining group  $g_1$ , member  $u_3$  will make more social connections, in this case it is member  $u_1$ .

[32, 31, 34] has shown that the existing social connections of members in a social group are only a subset of all possible connections. For these reasons, we introduce another network  $\mathcal{G}(V_U, E_{UU})$  that describes social connections between members. The social network grows through addition of new members to the set  $V_U$  and creation of new links between them. The member partition in bipartite network  $\mathcal{B}(V_U, V_G, E_{UG})$  and set of nodes in members' network  $\mathcal{G}(V_U, E_{UU})$  are identical.

For convenience, we represent bipartite and member networks with adjacency matrices  $B$  and  $A$ . The element of matrix  $B_{ug}$  equals one if member  $u$  is affiliated with group  $g$ , and zero otherwise. In matrix  $A$ , the element  $A_{u_1u_2}$  equals one if members  $u_1$  and  $u_2$  are connected and zero otherwise. The neighbourhood of member  $u$   $\mathcal{N}_u$  is a set off groups that member is affiliated with. On the other hand, the neighbourhood of group  $g$   $\mathcal{N}_g$  is a set of members affiliated to that group. The size of set  $\mathcal{N}_g$  equals to the size of the group  $g$   $S_g$ .

In our model, the time is discrete and networks evolve through several simple rules. In each time step we add  $N_U(t)$  new members and increase the size of the set  $V_U$ . For each newly added member we create the link to a randomly chosen old member in the social network  $G$ . This condition allows each member to perform social diffusion [39], i.e., to choose a group according to her social contacts. Not all members from set  $V_U$  are active in

each time step. Only a subset of existing members is active in one time step. Activity of old members is a stochastic process and is determined by parameter  $p_a$ ; every old member is activated with probability  $p_a$ . Old members activated in this way and new members make a set of active members  $\mathcal{A}_U$  at time t.

The group partition  $V_G$  grows through creation of new groups. Each active member  $u \in \mathcal{A}_U$  can decide with probability  $p_g$  to create a new group, or to join an already existing one with probability  $1 - p_g$ .

If the active member  $u$  decides that she will join an existing group, she first needs to a choice of this group. A member  $u$  with probability  $p_{aff}$  decides to select a group based on her social connections. For each active member, we look at how many social contacts she has in each group. The number of social contacts  $s_{ug}$  that member  $u$  has in group  $g$  equals to the overlap of members affiliated with a group  $g$  and social contacts of member  $u$ , and is calculated according to

$$s_{ug} = \sum_{u_1 \in \mathcal{N}_g} A_{uu_1}. \quad (3.2)$$

Member  $u$  selects an old group  $g$  to join according to probability  $P_{ug}$  that is proportional to  $s_{ug}$ . Member only considers groups with which it has no affiliation. However, if an active member decides to neglect her social contacts in the choice of the social group, she will, with probability  $1 - p_{aff}$ , select a random group from the set  $V_G$  with which she is not yet affiliated.

After selecting the group  $g$ , a member joins that group and we create a link in bipartite networks between a member  $u$  and a group  $g$ . At the same time, member selects X members of a group  $g$  which do not belong to her social circle and creates social connections with them. As a consequence of this action, we create X new links in network  $\mathcal{G}$  between member  $u$  and X members from group  $g$ .

The evolution of bipartite and social networks, and consequently growth of social groups, is determined by parameters  $p_a$ ,  $p_g$  and  $p_{aff}$ . Parameter  $p_a$  determines the activity level of members and takes values between 0 and 1. Higher values of  $p_a$  result in higher number of active members and thus faster growth of number of links in both networks, as well as the size and number of groups. Parameter  $p_g$  in combination with parameter  $p_a$  determines the growth of the set  $V_G$ .  $p_g = 1$  means that members only create new groups, and the existing network consists of star-like subgraphs with members being a central nodes and groups as leafs. On the other hand  $p_g = 0$  means that there is no creation of new groups and the bipartite network only grows through addition of new members and creation of new links between members and groups.

Parameter  $p_{aff}$  is especially important. It determines the importance of social diffusion.  $p_{aff} = 0$  means that social connections are irrelevant and the choice of group is random. On the other hand,  $p_{aff} = 1$  means that only social contacts become important for group

selection.

Our model is different from co-evolution model Ref. [34]. In our model  $p_{aff}$  is constant and the same for all members. In the co-evolution model this probability depends on members degree. The members are activated in our model with probability  $p_a$ , while in co-evolution model members are constantly active from the moment they are added to a set  $V_U$  until they become inactive after time  $t_a$ . Time  $t_a$  differs for every member and is drawn from exponential distribution with rate  $\lambda$ . In co-evolution model the number of social contacts that member has within the group is irrelevant for the group selection. On the other hand, in our model members tend to choose more often groups in which there is a greater number of their social contacts. While in our model, in the case of random selection of a group, member selects a uniformly at random a group that she is not affiliated with, in the co-evolution model the choice of group is preferential.

## 3.5 Results

The differences between our and co-evolution model, described in previous sections, at first glance may appear small. However, they lead to huge differences in the distribution of the size of social groups. The distribution of group sizes in co-evolution model is a power-law. Our model adds flexibility to produce groups with log-normal size distribution. This expands classes of social systems that can be modeled.

### Model description

First, we explore the properties of size distribution depending on parameters  $p_g$  and  $p_{aff}$ , and fixed value of activity parameter  $p_a$  and constant number of members added in each step  $N(t) = 30$ . The parameter  $X$  is set to value 25 for all simulations presented in this work. Our detailed analysis of the results for different values of parameter  $X$  shows that these results are independent of the value of parameter  $X$ .

Figure 3.4 shows some of the selected results and their comparison with power-law and log-normal fits. We see that values of both  $p_g$  and  $p_{aff}$  parameters, influence the type and properties of size distribution. For low values of parameter  $p_g$ , left column in Fig. 3.4, the obtained distribution is log-normal. The width of the distribution depends on  $p_{aff}$ . Higher values of  $p_{aff}$  lead to a broader distribution.

As we increase  $p_g$ , right column Fig. 3.4, the size distribution begins to deviate from log-normal distribution. The higher the value of parameter  $p_g$ , the faster grows the number of groups available to members. For the value of parameter  $p_g = 0.5$ , every second active member creates a group in each time step, and the number of groups increases fast. How members are distributed in these groups depends on the value of parameter  $p_{aff}$ . When  $p_{aff} = 0$ , social connections are irrelevant for the choice of the group and members choose groups at random. The obtained distribution slightly deviates from log-normal, especially

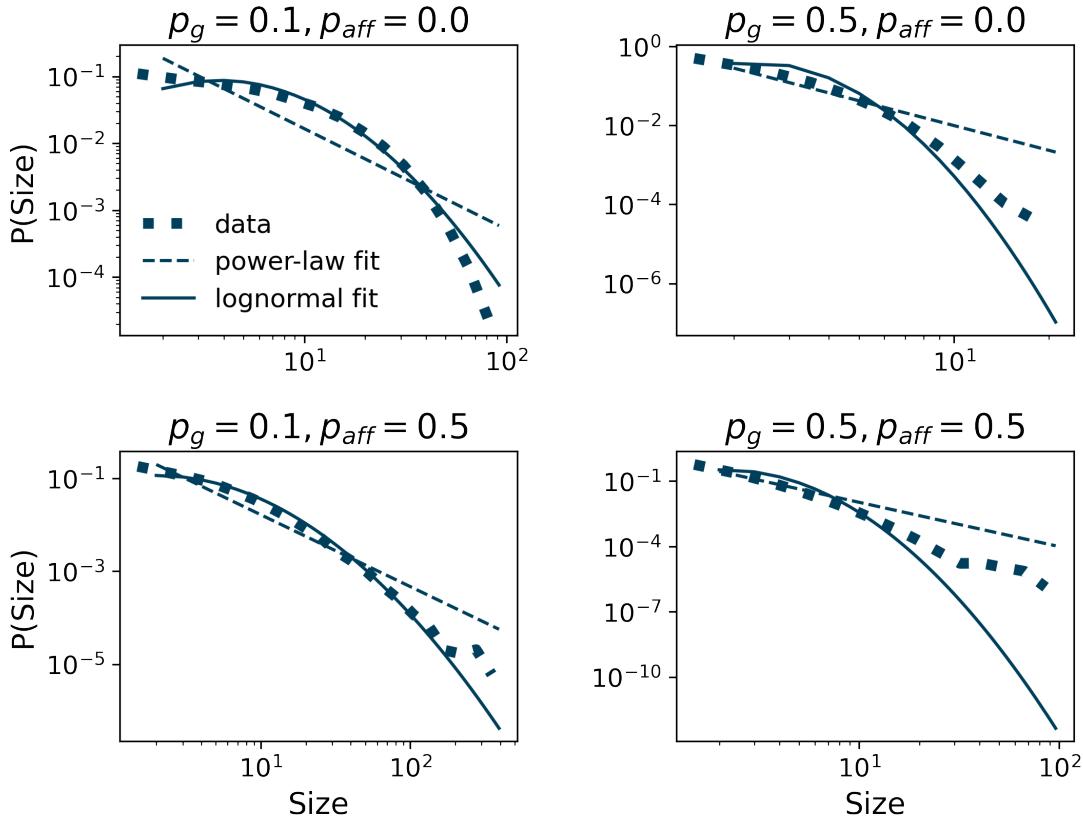


Figure 3.4: The distribution of sizes for constant growth of members,  $N = 30$ . The probability that members are active is fixed to  $p_a = 0.1$ , while we vary the probability for the creation of groups  $p_g$  and affiliation linking  $p_{\text{aff}}$

for large group sizes. In this case large groups sizes become more probable than in the case of log-normal distribution. The non zero value of parameter  $p_{\text{aff}}$  means that the choice of group becomes dependent on social connections. When member chooses a group according to her social connections, larger groups have higher probability to be affiliated with social connections of active members, and thus this choice resembles preferential attachment. For these reasons, the obtained size distribution has more broad tail than log-normal distribution, and begins to resemble power-law distribution.

## Modeling real systems

The social systems do not grow at constant rate. In Ref. [48] authors have shown that features of growth signal influence the structure of social networks. For these reasons we use the real growth signal from Meetup groups located in London and New York, and Reddit community to simulate the growth of the social groups in these systems. Figure 3.5 top panel shows the time series of the number of new members that join each of the three systems each month. All three systems have relatively low growth at the beginning, and than the growth accelerates as the system becomes more popular.

### 3. Groups growth model

---

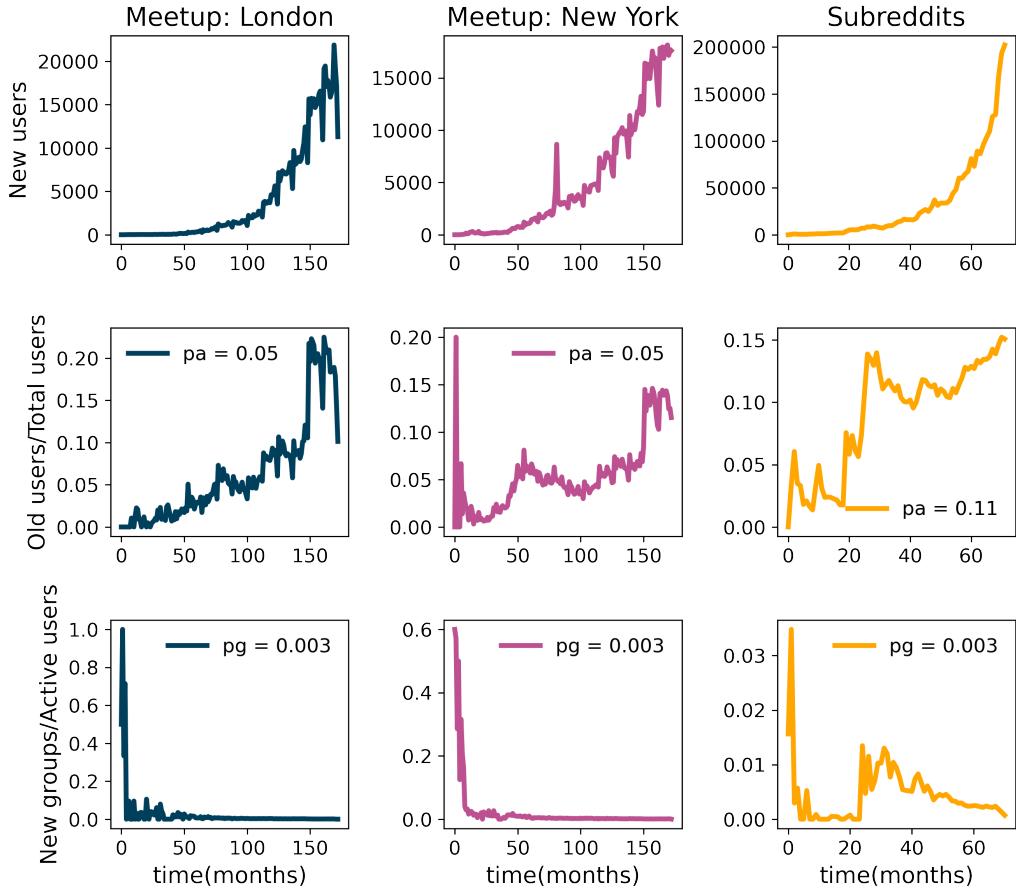


Figure 3.5: The time series of number of new members (top panel), ratio between old members and total members in the system (middle panel), and ratio between new groups and active members (bottom panel) for Meetup groups in London, Meetup groups in New York, and subreddits.

We also use empirical data to estimate  $p_a$ ,  $p_g$  and  $p_{aff}$ . Probabilities that old members are active  $p_a$  and that new groups are created  $p_g$  can be approximated directly from the data. Activity parameter  $p_a$  is the ratio between the number of old members that were active in month  $t$  and the total number of members in the system at time  $t$ . Figure 3.5 middle row shows the variation of parameter  $p_a$  during the considered time interval for each system. The values of this parameter fluctuates between 0 and 0.2 for London and New York based Meetup groups, while its value is between 0 and 0.15 for Reddit. To simplify our simulations we assume that  $p_a$  is constant in time, and estimate its value as its median value during the 170 months for Meetup systems, and 80 months of Reddit system. For Meetup groups based in London and New York  $p_a = 0.05$ , while Reddit members are more active on average and  $p_a = 0.11$  for this system.

Figure 3.5 bottom row shows the evolution of parameter  $p_g$  for the three considered systems. The  $p_g$  in month  $t$  is estimated as the ratio between the groups created in month  $t$   $N_{gnew}(t)$  and the total number of groups that month  $N_{gnew}(t) + N_{gold}(t)$ , i.e.,  $p_g(t) = \frac{N_{gnew}(t)}{N_{new}(t) + N_{old}(t)}$ . We see from Fig. 3.5 that  $p_g(t)$  has relatively high values at the beginning of the system's existence. This is not surprising. At the beginning these systems have relatively small number

$p_{aff}$	JS cityLondon	JS cityNY	JS reddit2012
0.1	0.0161	0.0097	0.00241
0.2	0.0101	0.0053	0.00205
0.3	0.0055	0.0026	0.00159
0.4	0.0027	<b>0.0013</b>	0.00104
0.5	<b>0.0016</b>	0.0015	0.00074
0.6	0.0031	0.0035	0.00048
0.7	0.0085	0.0081	0.00039
0.8	0.0214	0.0167	<b>0.00034</b>
0.9	0.0499	0.0331	0.00047

Table 3.1: Jensen Shannon divergence between group sizes distributions from model (in model we vary affiliation parameter  $p_{aff}$ ) and data.

of groups and often cannot meet the needs for content of all their members. As the time passes, the number of groups grows, as well as content offerings within the system, and members no longer have a high need to create new groups. Figure 3.5 shows that  $p_g$  fluctuates less after the first few months, and thus we again assume that  $p_g$  is constant in time and set its value to median value during 170 months for Meetup and 80 months for Reddit. For all three systems  $p_g$  has the value of 0.003

The affiliation parameter  $p_{aff}$  is not possible to estimate directly from the empirical data. For these reasons, we simulate the growth of social groups each of the three systems with the time series of new members obtained from the real data and estimated values of parameters  $p_a$  and  $p_g$ , while we vary the value of  $p_{aff}$ . For each of the three systems, we compare the distribution of group sizes obtained from simulations for different values of  $p_{aff}$  with ones obtained from empirical analysis using Jensen Shannon (JS) divergence. The JS divergence [?] between two distributions  $P$  and  $Q$  is defined as

$$JS(P, Q) = H\left(\frac{P + Q}{2}\right) - \frac{1}{2}(H(P) + H(Q)) \quad (3.3)$$

where  $H(p)$  is Shannon entropy  $H(p) = \sum_x p(x)\log(p(x))$ . The JS divergence is symmetric and if  $P$  is identical to  $Q$ ,  $JS = 0$ . The smaller the value of JS divergence, the better is the match between empirical and simulated group size distributions. The Table 3.1 shows the value of JS divergence for all three systems. We see that for London based Meetup groups the affiliation parameter is  $p_{aff} = 0.5$ , for New York groups  $p_{aff} = 0.4$ , while the affiliation parameter for Reddit  $p_{aff} = 0.8$ . Our results show that social diffusion is important in all three systems. However, Meetup members are more likely to join groups at random, while for the Reddit members their social connections are more important when it comes to choice of the subreddit.

Figure 3.6 shows the comparison between the empirical and simulation distribution of group sizes for three considered systems. We see that empirical distributions for Meetup groups based in London and New York are perfectly reproduced by the model and chosen values of parameters. In the case of Reddit, the distribution is very broad, and the tail of distribution is well reproduced by the model.

### 3. Groups growth model

The bottom row of Fig. 3.6 shows the distribution of logarithmic values of growth rates of groups obtained from empirical and simulated data. We see that the tails of empirical distributions for all three systems are well emulated by the ones obtained from the model. However, there are deviations which are the most likely consequence of using median values of parameters  $p_a$ ,  $p_g$ , and  $p_{aff}$ .

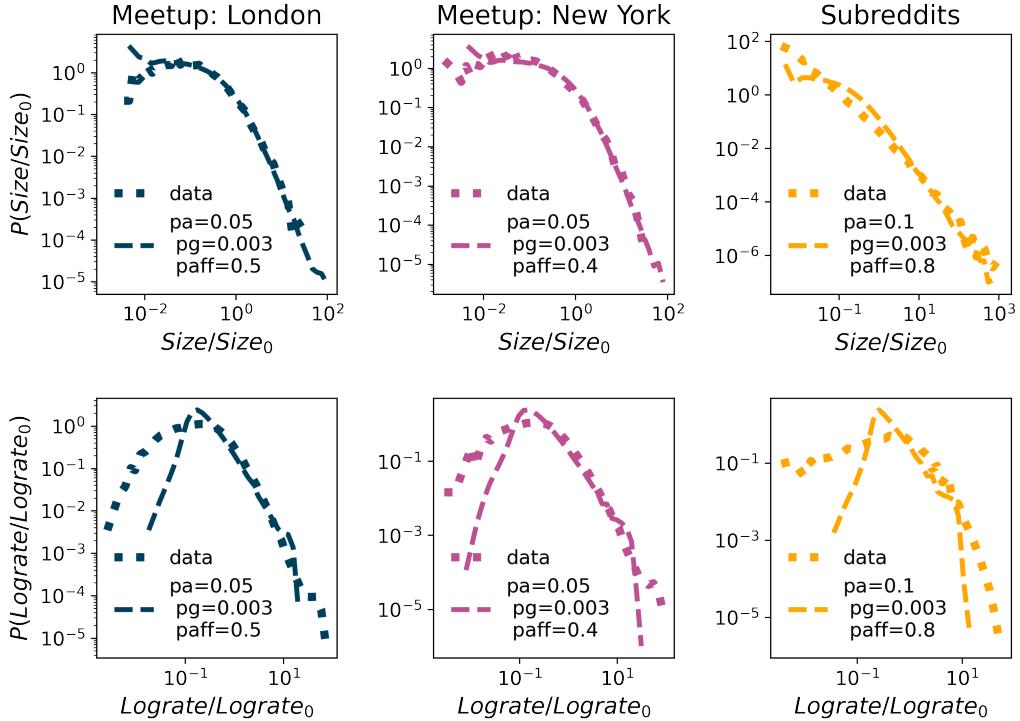


Figure 3.6: The comparison between empirical and simulation distribution for group sizes (top panel) and logrates (bottom panel).

## 3.6 Discussion and conclusions

The growth of the cities and companies has attracted the attention of researchers in the previous few decades [37, 35, ?]. It is not surprising if we keep in mind that their growth determines other processes essential for the functioning of cities and economies. In the cities, economic and innovation growth scale with their size [?]. The understanding of growth and segmentation of economic systems are essential for the long-term prediction of their evolution, and risk assessment []. The growth of social groups and social system segmentation has been slightly overlooked since the focus was mainly on the structure of social networks and their evolution.

The results of our empirical analysis and theoretical modeling show that there are universal growth rules that govern the growth of social groups in these systems. Through rigorous empirical analysis of the growth of social groups in three systems, Meetup groups located in London and New York, and Reddit, we show that the distribution of group sizes in these

systems has log-normal normal behavior. The empirical distributions of normalized sizes of the groups created in different years fall on top of each other and have the same values of parameters for the same system. Furthermore, the distributions for Meetup groups located in London and New York have similar model parameter values, suggesting that groups' growth in these two systems are similar. Numerical simulations further confirm these findings. By tuning our model's parameters, we can reproduce the distribution of group sizes in all three systems.

Our results show that while the processes that govern the growth of social groups in studied social systems are the same, their importance varies among systems. The analysed groups grow through two mechanisms [34]: members join a group that is chosen according to their interests or social relations with the group's members. The number of members in the system is growing, as well as the number of groups. The empirical distribution of growth rates differs for Meetup and Reddit. The observed differences can be explained by different modalities of interactions between their members. Meetup members need to invest more time and resources to interact with their peers. The events are localised in time and space, and thus the influence of peers in selecting another social group may be limited. On the other hand, Reddit members do not have these limitations. The interactions are online, asynchronous, and thus not limited in time. The influence of peers in choosing new subreddits and topics thus becomes more important. The inspection of numerical simulations confirms these observations. The values of  $p_{aff}$  parameters for Meetup and Reddit imply that social connections in diffusion between groups are more critical in Reddit than in Meetup.

Gibrat's law is the first empirical law used by researchers to describe and explain the growth and segmentation of various socio-economical systems, including cities and firms. The possibility of application of common law to the growth of social groups in different systems indicates the existence of universal growth patterns and mechanisms that govern that growth []. Detailed and rigorous empirical analysis of the growth of the cities and firms showed that it goes beyond Gibrat's law []. Our and the work of other researchers [34] confirm that these findings also hold for the growth of social groups. The analysis of monthly growth rates shows that these rates are log-normally distributed and depend on the size of a group. Furthermore, we cannot reduce the model proposed in this work to the law of proportional growth. Although our analysis shows that Gibrat's law does not apply to the growth of social groups, our findings confirm that universal patterns characterise this growth.

The results presented in this paper contribute to our knowledge of the growth and segmentation of socio-economical systems. Our rigorous analysis shows that the distribution of sizes of groups for studied systems follows a log-normal distribution. The findings of the previous research suggested the power-law behavior of this distribution. A detailed and comprehensive analysis of distributions of group sizes in social systems is needed. These and future results will help us better understand the growth and segmentation of social systems and predict their evolution and sustainability.



# **Chapter 4**

## **The role of trust in knowledge based communities**

Information and communications technologies (ICTs) have enabled faster and easier creation and sharing of knowledge. Furthermore, they have provided access to a large amount of data which enabled a detailed study of their emergence and evolution [47], as well as user's roles [49], patterns of their activity [50, 51, 52]. However, relatively small attention was given to sustainability of SE communities. Most of the research was focused on the activity and factors that influence the increase of the users' activity in these communities. Factors such as need for experts and the quality of their contributions have been thoroughly investigated [53]. It was shown that growth of communities and mechanisms that drive it may depend on the topic around which the community was created [54].

### **4.1 The Stack Exchange**

The Stack Exchange is a network of question-answer websites on diverse topics. In the beginning, the focus was on computer programming questions with StackOverflow<sup>1</sup> community. Its popularity led to the creation of the Stack Exchange network that these days counts more than 100 communities on different topics. The SE communities are self-moderating, and the questions and answers can be voted, allowing users to earn Stack Exchange reputation and privileges on the site.

The new site topics are proposed through site Area51<sup>2</sup>, and if the community finds them relevant, they are created. Every proposed StackExchange site needs interested users to commit to the community and contribute by posting questions, answers and comments. After a successful private beta phase site reaches the public beta phase, other members are allowed to join the community. The site can be in the public beta phase for a long time until it meets

---

<sup>1</sup>More information about StackOverflow is available at: <https://stackoverflow.co/> and broad introduction to StackExchange network is available at: <https://stackexchange.com/tour>.

<sup>2</sup>Visit <https://area51.stackexchange.com/faq> for more details about closed and beta SE communities and the review process.

#### 4. The role of trust in knowledge based communities

---

specific SE evaluation criteria for graduation. Otherwise, it may be closed with a decline in users' activity.

We focused analysis on four pairs of SE communities with the same topic. Astronomy, Literature and Economics are active communities<sup>3</sup> The first time, these communities were unsuccessful and thus closed. We also compare closed Theoretical Physics with the Physics site, considering that those two topics engage similar type of users.

### Data

Stack Exchange data are public and regularly released. As closed communities were active between 180 and 210 days, we extracted only first 180 days of data. Given that first few months can be crucial for further development of the community [55], we are interested in early evolution of Stack Exchange sites.

Detailed information about questions, answers, and comments are available for each SE community. Each post is labelled with a unique ID, the user's ID who made the post, and creation time. On Stack Exchange, users interact on several layers: Those interactions are considered positive.

- posting an answer on the question; for every question, we extract IDs of its answers
- posting a comment on the question or answer; for every question and answer, we selected IDs of its comments
- accepting answer; for each question, we selected the accepted answer ID

Even though posts can be voted and downvoted, information about a user who voted is absent, so we do not consider these interactions between users. Comments can not be downvoted, while we find only around 3% negatively voted answers and questions, Table 4.1.

Table 4.1: Percentage of negatively voted interactions

Site	Status	Questions	Answers
Physics	Beta	5%	4%
	Closed	1%	2%
Astronomy	Beta	3%	3%
	Closed	2%	1%
Economics	Beta	4%	4%
	Closed	7%	4%
Literature	Beta	2%	5%
	Closed	2%	1%
<b>Average</b>		3.2%	3%

---

<sup>3</sup>Astronomy, Literature and Economics graduated on December 2021 and during our research, they were still in the public beta phase.

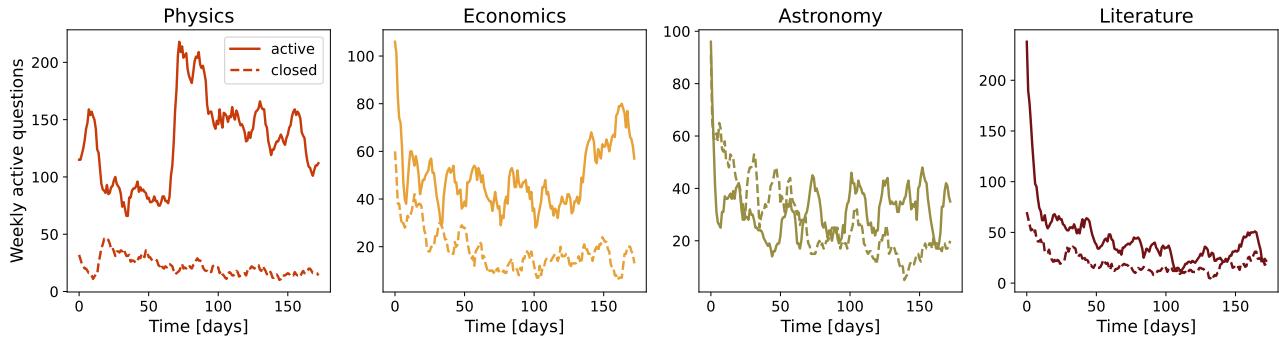


Figure 4.1: Number of active questions within 7 days sliding windows. Solid line - active sites; dashed lines - closed sites.

## Comparison between active and closed SE communities

Table 4.2 compares the first 180 days between closed and active communities. When it comes to basic statistics, active communities had larger number of users, questions, answers and comments. Another simple indicator if community is going to graduate or decline can be time series of active questions for period of 7 days in Figure 4.1. The question is active if had at least one activity, posted answer or comment during previous seven days. We find that live communities have larger number of active questions after first three months. Still, this difference is smaller for literature and astronomy. For astronomy we observe that closed community had higher number of active questions in the early period of community life.

Table 4.2: Community overview for first 180 days, Number of users  $n_u$ , number of questions  $n_q$ , number of answers  $n_a$ , number of comments  $n_c$

Site	Status	First Date	$n_u$	$n_q$	$n_a$	$n_c$
Astronomy	Closed	09/22/10	336	474	953	1444
	Beta	09/24/13	405	644	959	2170
Economics	Closed	10/11/10	275	368	458	1253
	Beta	11/18/14	648	1024	1410	3553
Literature	Closed	02/10/10	284	318	523	1097
	Beta	01/18/17	478	910	907	3301
Physics	Closed	09/14/11	281	349	564	2213
	Launched	08/24/10	1176	2124	4802	15403

Similarly, the official Stack Exchange community evaluation process considers simple metrics <sup>4</sup>. To determine the success of sites they measure how many questions are answered, how many questions are posted per day, and how many answers are posted per question. There are two measures: the number of avid users and the number of visits that

<sup>4</sup><https://stackoverflow.blog/2011/07/27/does-this-site-have-a-chance-of-succeeding/>

#### 4. The role of trust in knowledge based communities

---

are not easily interpreted from the data. The site is *healthy* if it has 10 questions per day, 2.5 answers per question and more than 90% of answered questions. For less than 80% of answered questions, 5 questions per day and 1 question per answer site *needs some work*.

We calculated Stack Exchange statistics for astronomy, economics, literature and physics and results are presented in the Table 4.3. After observed period of 180 days only live physics is healthy site while other live communities are at least in two criteria labeled as *okay*. Closed sites mostly *need some work*, the exception is closed astronomy. For example it has *excellent* percent of answered questions and *okay* answer ratio.

Table 4.3: Community overview for first 180 days according to SE criteria

Site	Status	Answered	Questions per day	Answer ratio
Astronomy	Closed	<u>95</u> %	2.62	<u>2.02</u>
	Beta	<u>96</u> %	3.57	<u>1.49</u>
Economics	Closed	68 %	2.04	<u>1.25</u>
	Beta	<u>84</u> %	<u>5.66</u>	1.37
Literature	Closed	79 %	1.77	<u>1.65</u>
	Beta	74 %	<u>5.04</u>	<u>1.10</u>
Physics	Closed	83 %	1.93	<u>1.64</u>
	Beta	<u>93</u> %	<b>11.76</b>	<b>2.74</b>
Stack Exchange criteria	excellent	> 90 %	> 10	> 2.5
	needs some work	< 80 %	< 5	< 1

This simple measurements presented in tables 4.2 and 4.3 and in figure 4.1 do not provide us clear indications about community sustainability. Only for physics topic the difference between active and closed community is evident, while for other communities it is not so clear. Thus, we need deeper insights into structure and dynamics of these communities to understand. The structure of social interactions within communities and dynamics of collective trust may provide better explanation why some communities succeed and other died.

## 4.2 Network properties of Stack Exchange data

On Stack Exchange sites, the interaction between users happens through posts. As we are interested in examining the characteristics of the users, we map interaction data to the networks. Using complex network theory, we can quantify the properties of obtained networks and compare different SE communities, e.g. alive and closed SE sites. In the user interaction network, the link between two nodes, user  $i$  and  $j$ , exists if user  $i$  answers or comments on the question posted by user  $j$ , or user  $i$  comments on the answer posted by user  $j$ . The created network is undirected and unweighted, meaning that we do not consider multiply interactions between users or the direction of the interaction.

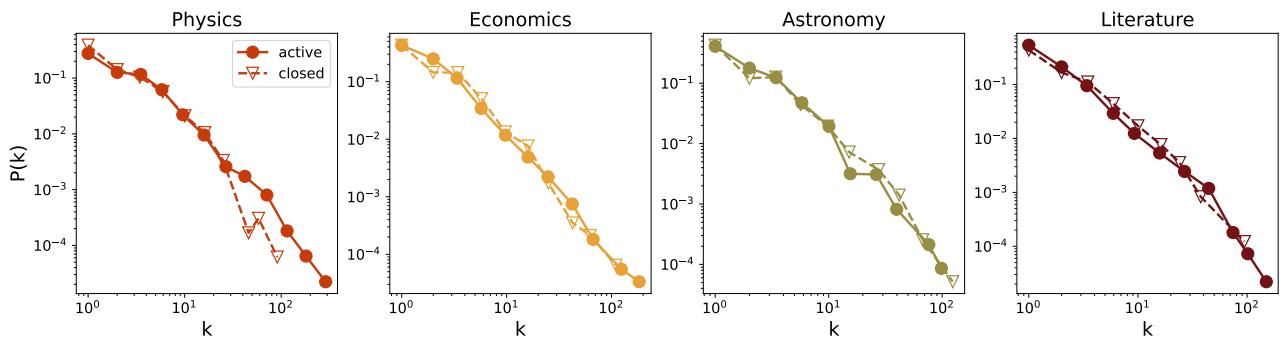


Figure 4.2: Degree distribution.

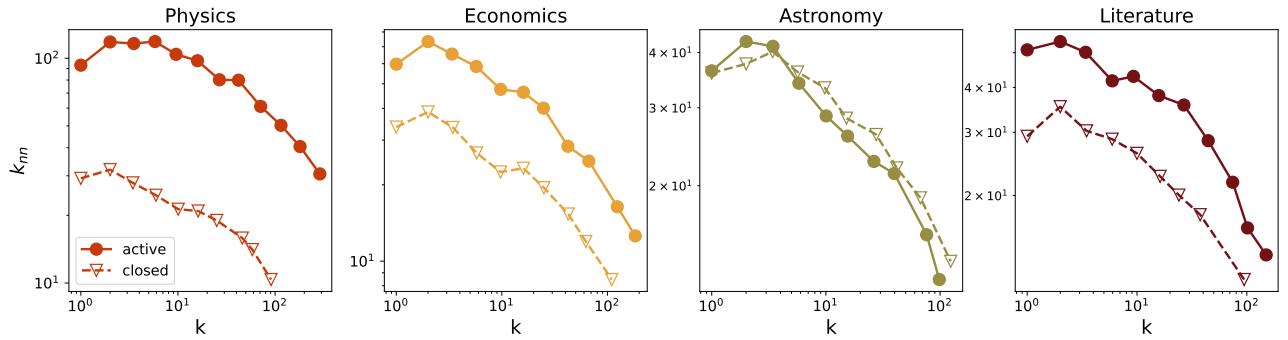


Figure 4.3: Neighbour degree.

Instead of creating a static network from the data in the first 180 days of community life, we study how network snapshots evolve. At each time step  $t$ , we create network snapshot  $G(t, t + \tau)$ , for time window of the length  $\tau$ . We fix the time window to  $\tau = 30$  days and slide it by  $t = 1$  day through time. Discussion of how the length of the sliding window influences the results is given in appendix A. Sliding the time window by one day, we can capture changes in the network structure daily, as two 30 days consecutive networks overlap significantly.

We calculate different structural properties of observed networks and study their evolution. There are many local and global measures of the network properties [56]. They are

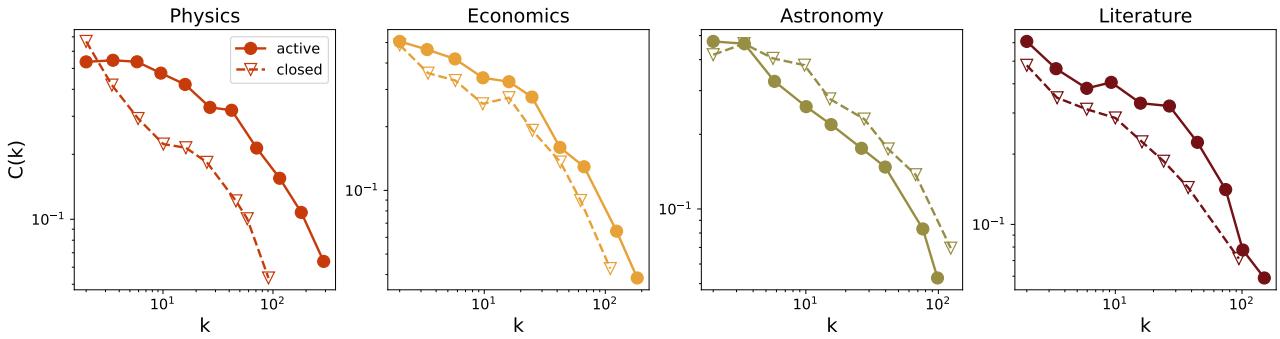


Figure 4.4: Neighbour degree.

also dependent, still it was shown that degree distribution, degree-degree correlations and clustering coefficient can describe the properties of the most complex networks [57].

## Clustering

The clustering coefficient of a node quantifies the average connectivity of between its neighbours and cohesion of its neighbourhood [56]. It is a probability that two neighbours of a node are also neighbours, and is calculated using the following formula:

$$c_i = \frac{e_i}{\frac{1}{2}k_i(k_i - 1)} . \quad (4.1)$$

Here  $e_i$  is the number of links between neighbours of the node  $i$  in a network, while  $\frac{1}{2}k_i(k_i - 1)$  is the maximal possible number of links determined by the node degree  $k_i$ . The clustering coefficient of network  $C$  is the value of clustering averaged over all nodes. Here we investigate how clustering coefficient in a SE community is changing with time by calculating its value for all network snapshots. We compare the behavior of clustering for active and closed communities on the same topic in order to better understand how cohesion of these communities is changing over time. Members' clustering coefficient measures the probability that other members connected to them are also connected. Study on dynamics of social group growth shows that that links between one's friends that are members of a social group increase the probability that that individual will join the social group [31]. Furthermore, successful social diffusion typically occur in networks with high value of clustering coefficient [58]. These results suggest that high local cohesion should be a characteristic of sustainable communities.

We first analyse structural properties of Stack Exchange communities and examine the difference between successful and unsuccessful ones. We calculate the mean clustering coefficient for 30-days window networks and examine how it changes with time. Figure 4.2 shows the evolution of mean clustering coefficient for all eight communities. All communities that are still alive are clustered, with the value of mean clustering coefficient higher than 0.1. Physics, the only launched community, has the value of clustering coefficient above 0.2 for the first 180 days. During larger part of the observed period, the clustering coefficient

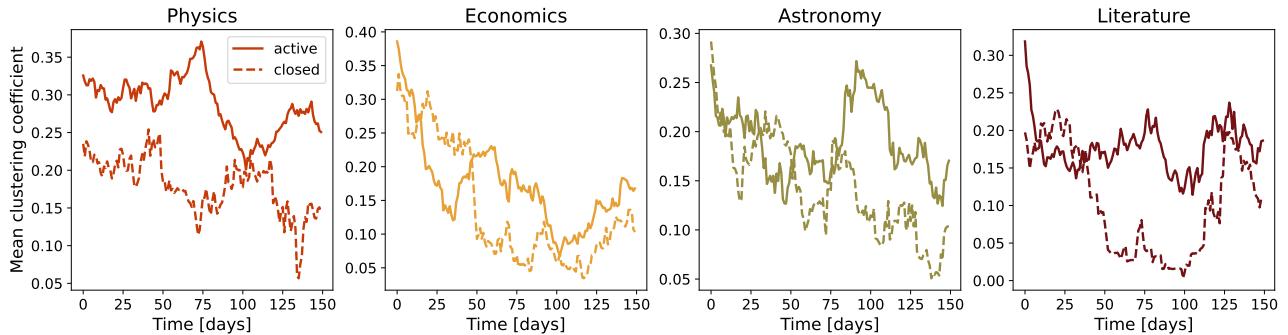


Figure 4.5: Mean clustering coefficient.

of an active community is higher compared to the clustering coefficient of its closed pair. If we compare active communities with their closed counterpart, the closed communities have higher value of the mean clustering coefficient in the early phase while later communities that are still active have higher values of clustering coefficient. These results suggest that all communities have relatively high local cohesiveness, and that lower values of clustering coefficient in the later phase of community life may be an indicator of its decline.

## Core-periphery structure

Real networks, including social networks, have a distinct mesoscopic structure [59, 60]. Mesoscopic structure is manifested either through community structure or core-periphery structure. Networks that have community structure consist of a certain number of group of nodes that are densely connected with each other, with sparse connections between groups. Networks with core-periphery structure consist of two groups of nodes, with higher edge density within one group and between groups but low edge density in the second group [60]. Research on dynamics of user interaction in SE communities shows that there is a small group of highly active members, similar to core, that have frequent interactions with casual or low active members of community [50, 54]. These results indicate that we should expect a core-periphery structure in SE communities.

Core-periphery pattern means that network consists of two components: a core, densely connected group of nodes, and periphery, a second group of nodes that are loosely connected with the core and with each other. Classification of nodes into one of these two groups provide us with information about their functional and dynamical roles in the network. Active users typically belong to core, while periphery consists of less active users.

To investigate core-periphery structure of SE communities and how it evolves trough time, we analyse the core-periphery structure of every 30 days network snapshot. We use Stochastic Block Model (SBM) adapted for core-periphery inference of network structure [60] to determine the core-periphery structure.

**SBM** is model where each node, in given network  $G$ , belongs to one of  $B$  blocks. Vector  $\theta_i = r$  indicates that node  $i$  is in block  $r$ , while SBM matrix  $\{p\}_{B \times B}$ , specify the probability  $p_{rs}$  that nodes from group  $r$  are connected to nodes in group  $s$ . The SBM model is looking for the most probable model that can reproduce a given network  $G$ . Probability of having model parameters  $\theta, p$  given network  $G$  is proportional to likelihood of generating network  $G$ , prior of SBM matrix and prior on block assignments:

$$P(\theta, p|G) = P(G|\theta, p)P(p)P(\theta) \quad (4.2)$$

$$P(G|\theta, p) = \prod_{i < j} p_{r_i s_j}^{A_{ij}} (1 - p_{r_i s_j})^{1 - A_{ij}} \quad (4.3)$$

where  $A_{ij}$  is number of edges between nodes  $i$  and  $j$ .

Prior on  $p$  is modified for core-periphery model such that  $P(p) = 3! I_{0 < p_{22} < p_{12} < p_{11} < 1}$ , while prior on  $\theta$  consists of three parts: probability of having  $l$  blocks; given the number of layers probability  $P(n|l)$  of having groups of sizes  $n_1..n_l$  and the probability  $P(\theta|n)$  of having particular assignments of nodes to blocks.

For fitting model in the work [60] authors use Metropolis-within-Gibbs algorithm. The likelihood of SBM model increase with number of blocks and model itself does not define optimal number of communities. Inferring minimum description length of the model is one approach to decide which model is more likely.

For each 30 days snapshot network we run 50 iterations and choose the model parameters  $\theta$  and  $p$  according to minimum description length. MDL does not change much among inferred core-periphery structure, see Fig. A8, while looking into adjusted rand index we can notice that difference exists. Still, ARI between pair-wised compared partitions is large ( $ari > 0.9$ ) indicating stability of inferred structures.

## core-periphery structure of knowledge-sharing networks

Furthermore, we examine core-periphery structure of these communities and its evolution. Specifically, we are interested in the evolution of connectivity in the core. Figure 4.3 shows the number of links between nodes in the core per node  $\frac{L}{N}(t)$ .  $\frac{2L}{N}$  is the average degree of the node in the core, and thus,  $\frac{L}{N}$  is the half of the average degree. Again, Physics community has the much higher value of this quantity than Theoretical physics during the whole observed period, indicating higher connectivity between core members. Higher connectivity between core members in the active community is also characteristic for Literature, although this quantity has the same value for active and closed communities at the end of the observation period. The differences between active and closed communities are not that evident for Economics and Astronomy, see Fig. 4.3. Active and closed Economics communities have similar connectivity in the core during the first 50 days. After this period, the connectivity in the core of the active community the twice as large as in the closed community and the difference grows at the end of observation period. The connectivity in the core of closed Astronomy community is higher than the connectivity in the core of the active community during the first 50 days. But as the time progresses, this difference changes in the favor of live community, while at the end of the observation period the difference disappears.

The difference between active and closed communities is more prominent if we consider average number of core-periphery edges per core node. The connectivity between core and periphery is higher for the still active communities than for the closed ones, see Fig. 4.3. This is very obvious if we compare Physics and Theoretical physics community. Moreover, Physics community has the highest connectivity compared to all other communities. When it comes to active communities that are still in the beta phase, they either have the same core-periphery connectivity as their closed counter part, or as in the case of Astronomy, their periphery is weaker connected to the core during the first 50 days of their life, see Fig. 4.3.

On average, the cores of the active communities have higher number of nodes in the core than the closed communities, Fig. A11. However, the relative size of the core compared to the size of the whole network is similar when we compare closed and active communities on the same topic. This is even true for communities on physics topic. The size of the core fluctuates with time for active and closed communities. The core membership also changes with time. This core membership is changing more for the closed communities. We quantify this by calculating the Jaccard index between the cores of the subnetworks in the moment  $t_i$  and  $t_j$ . Figure A9 in Supplementary Information shows the value of Jaccard index between any two of the 150 subnetworks. The highest value of the Jaccard index is around the diagonal and

#### 4. The role of trust in knowledge based communities

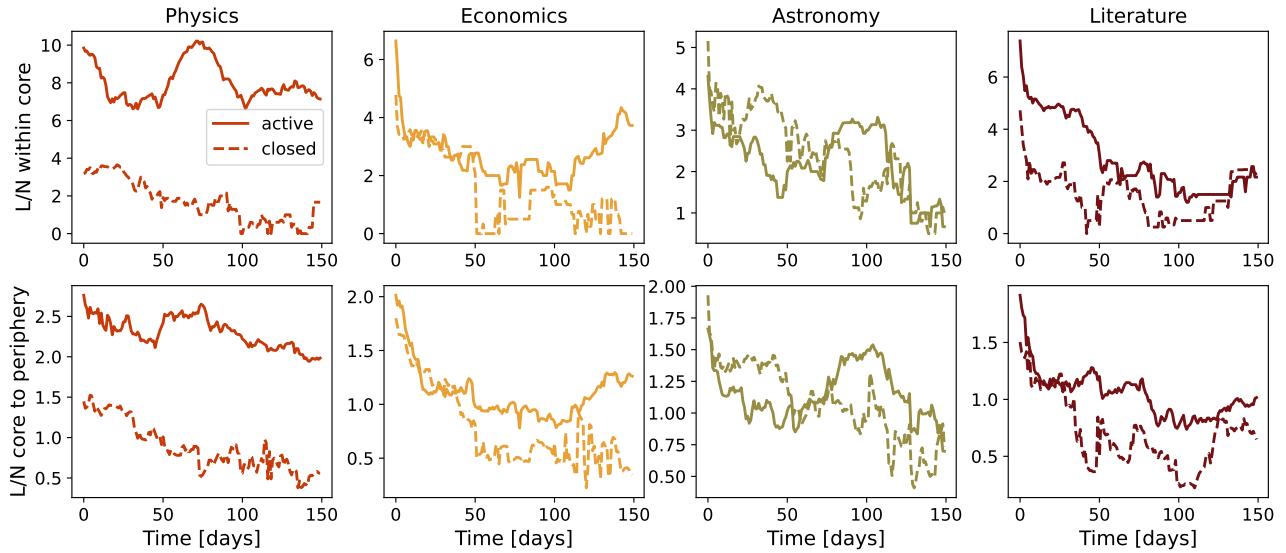


Figure 4.6: Links per node in core and links per node between core and periphery.

has value close to 1. This is expected, since these subnetworks are for consecutive days and the difference between them is smaller. The value of Jaccard index decreases with number of days between two subnetworks  $|t_i - t_j|$  faster in closed communities Fig. A10. This difference is the most prominent for the literature communities, while this difference is practically non-existent for Astronomy. The relatively high overlap between cores of even more distant subnetworks for active communities, further confirms that the core is more stable in these communities than in their closed counterparts.

### Core-periphery structure of the interaction networks

In Q-A communities are common two types of users: popular and casual users. Popular users tend to generate the majority of interactions - they are likely to post more questions, also take part in answering questions and tend to engage discussions through comments. For popular users we consider 10 of most active users. We analyse interactions between popular and casual users and among popular users in the sub-networks of 30 days [ $t+30$ ]. In both cases the number of links per nodes in active sites are larger than in closed communities (figure 4.8).

Although this separation of users puts an emphasis on differences between closed and active sites, it does not guarantee that all popular users are in the top 10. To solve this dilemma we use the SBM (Stochastic Block Model) algorithm to detect the core and the periphery of each 30 days sub-network. Such a split of users leads us to similar conclusions as before. (see figure A.3 - 2nd column)

Stochastic models start from random configuration and the algorithm can converge to different local stable states. For each 30 days sub-network we run 50 iterations of SBM and choose the model parameters  $\theta, p$  according to minimum description length. As example we show analysis of inferred sample of core-periphery structures for 30 days area51 astronomy

networks, Figure B.1. We represent mean minimum description length (MDL) and mean number of nodes in the core with standard deviation. MDL does not change much among inferred core-periphery structures, still difference between obtained configurations is notable in the number of nodes in the core. To investigate in more details similarity between obtained core-periphery configurations in the sample we calculate several measures between pair-wised partitions such as normalized mutual information, adjusted rand index, F1 measure and Jaccard index. Those measures are larger than 0.5, and in most cases higher than 0.9 indicating stability of the inferred core-periphery structures.

To study the stability of the core across the time we compute Jaccard's coefficient between core users in  $[t+30)$  networks selected at times  $t_1$  and  $t_2$ , (figure 4.4). Higher values of the Jaccard index indicate that core users tend to stay in the core. The detected cores experience a lot of change over time and the highest overlap between core users is in the network closer in the time. The average Jaccard index between core users in all sub-networks separated by time interval  $|t_1 - t_2|$  with standard deviation confidence interval is presented in figure 4.5. Compared to closed sites, active sites show more stability in the core. Even the number of core users obtained in the launched and closed communities are comparable 4.6 (a high difference is found only for physics ), the ratio between total core and periphery reputation is evidently higher in the active than in closed sites, figure 4.7.

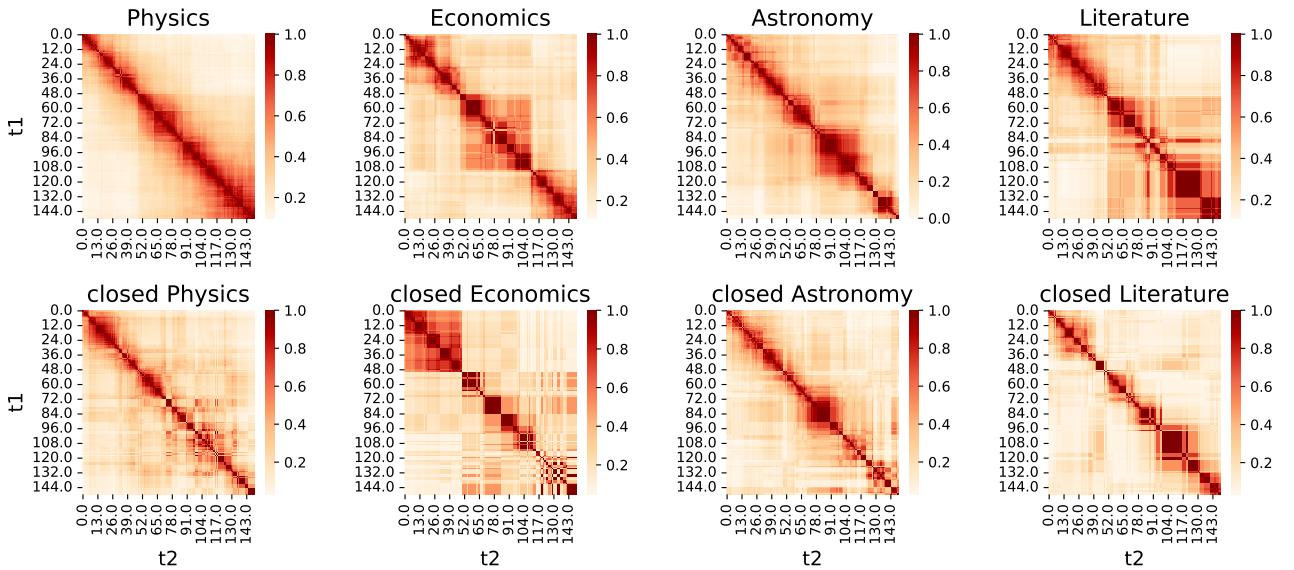


Figure 4.7: Jaccard index between core users in sub-networks at time points  $t_1$  and  $t_2$

#### 4. The role of trust in knowledge based communities

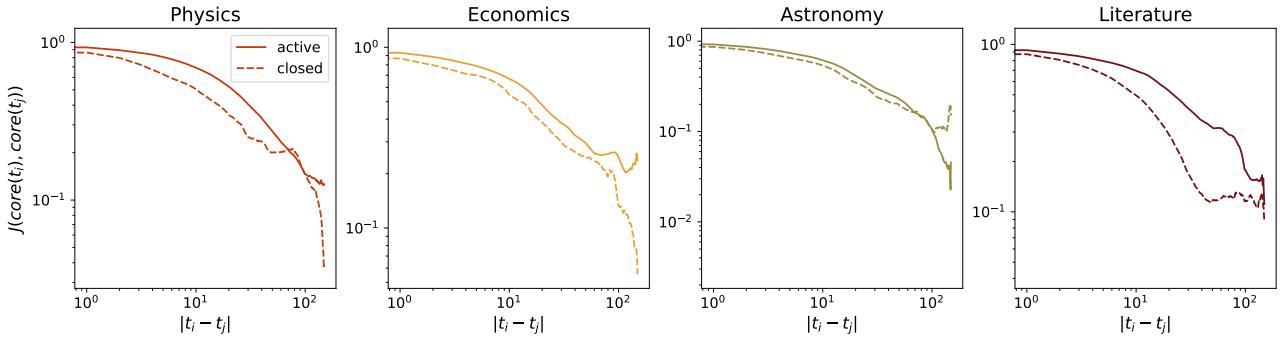


Figure 4.8: Jaccard index between core users in 30days sub-networks for all possible pairs of 30 days sub-networks separated by time interval  $|t_i - t_j|$

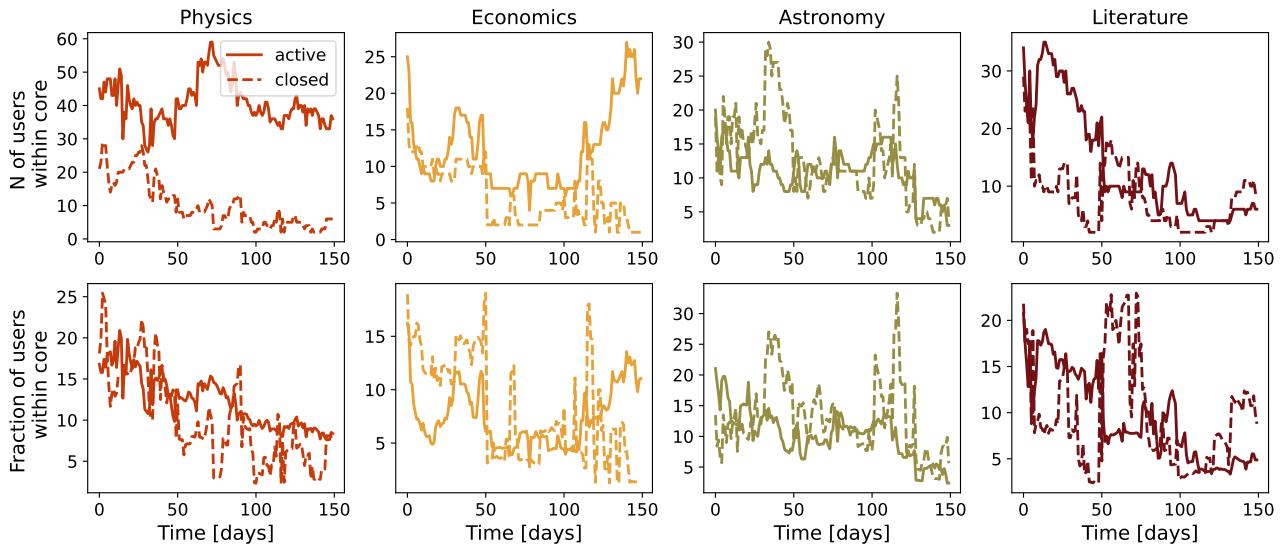


Figure 4.9: Just for reference size of the core (top) and fraction of users in core (bottom). Solid lines - active sites; dashed lines - closed sites.

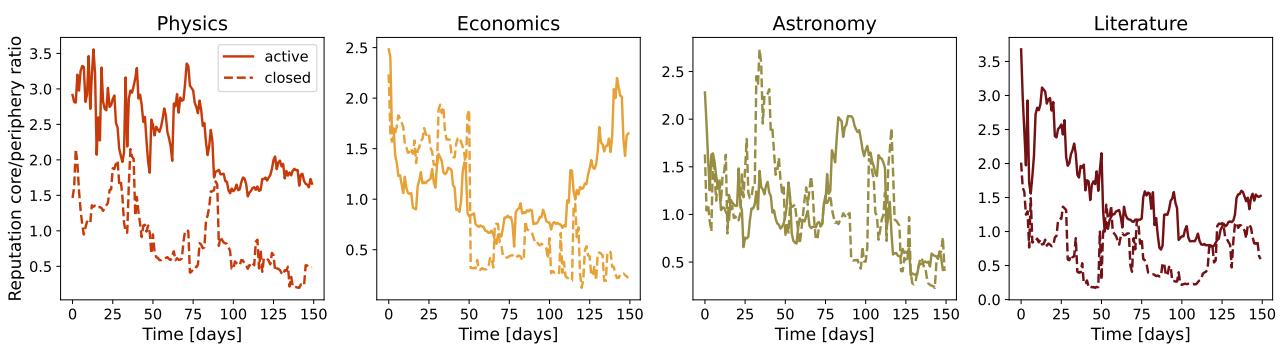


Figure 4.10: Ratio between the total reputation within network core and periphery. Solid lines beta communities, dashed lines area 51 communities.

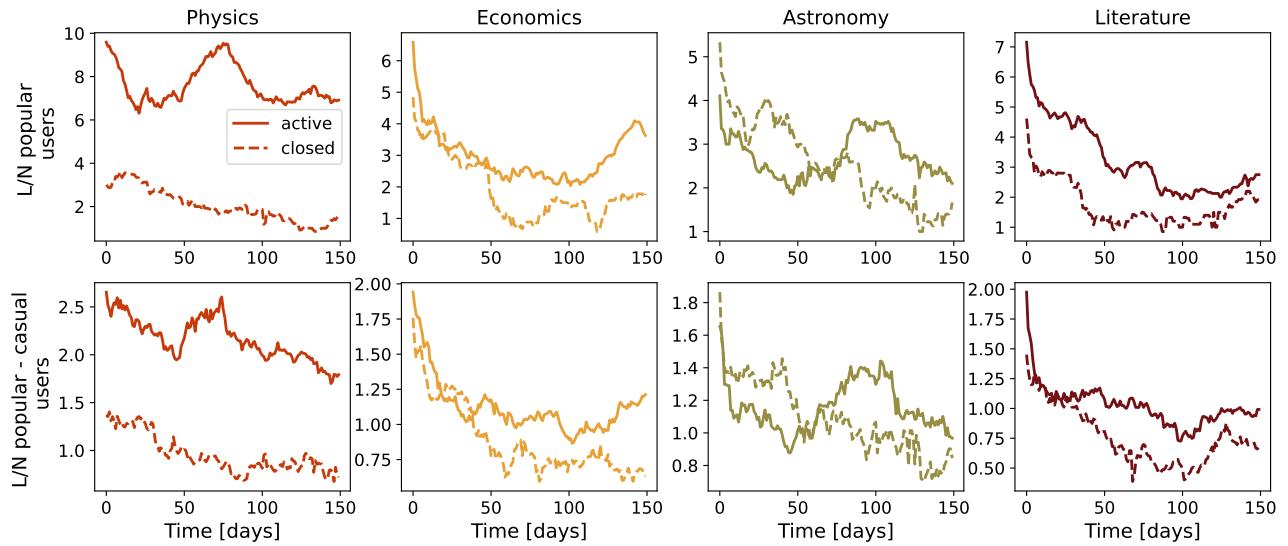


Figure 4.11: Links per node among popular users (top 10% of users) and between popular and casual users (everyone but popular users). Reminder: only 3rd and 5th columns should stay and only for reference to previous research, while our point is to this selection via core/periphery decomposition without thresholding.

### 4.3 Dynamic reputation model

Any dynamical trust or reputation model has to take into account distinct social and psychological attributes of these phenomena in order to estimate the value of any given trust metric [61]. First of all, the dynamics of trust is asymmetric, meaning that trust is easier to lose than to gain. As part of asymmetric dynamics, in order to make trust easier to loose the trust metric has to be sensitive to new experiences (recent activity or the absence of the activity of the agent), while still maintaining nontrivial influence of old behavior. The impact of new experiences has to be independent of the total number of recorded or accumulated past interactions, making high levels of trust easy to lose. Finally, the trust metric has to detect and penalize both the sudden misbehavior and the possibly long term oscillatory behavior which deviates from community norms.

We estimate dynamic reputation of the Stack Exchange users using Dynamic Interaction Based Reputation Model (DIBRM) [62]. This model is based on the idea of dynamic reputation, which means that the reputation of users within the community changes continuously through time: it should rapidly decrease when there is no registered activity from the specific user in the community (reputation decay), and it should grow when frequent, constant interactions and contributions to the community are detected. The highest growth of user's reputation is found through bursts of activity followed by short period of inactivity.

In our implementation of the model, we do not distinguish between positive and negative interactions in the Stack Exchange communities. Therefore, we treat any interaction in the community (question, answer or comment) as potentially valuable contribution. In fact, evaluation criteria for Stack Exchange websites going through beta testing, described in SI, do not distinguish between positive and negative interactions. The percentage of negative interactions in the communities we investigated was below 5%, see Table 1 in SI. Filtering positive interactions would also require filtering out comments because they are not rated by the community, and that would eliminate a large portion of direct interactions between the users of a community, which is essential for estimating their reputation.

In DIBRM, reputation value for each user of the community is estimated combining three different factors: 1) *reputation growth* - the cumulative factor which represents the importance of users' activities; 2) *reputation decay* - the forgetting factor which represents the continuous decrease of reputation due to inactivity; *the activity period factor* - measuring the length of the period of time in which the change of reputation happened. In case of Stack Exchange communities, the forgetting factor has a literal meaning, as we can assume that past contributions provided by a user are being forgotten by active users as their attention is captured by more recent content.

In line with the basic dichotomy of reputation dynamics, which revolves around the varying influence of past and recent behavior, DIBRM has two components: *cumulative factor* - estimating the contribution of the most recent activities to the overall reputation of the user; *forgetting factor* - estimating the weight of past behavior. Estimating the value of

recent behavior starts with the definition of the parameter storing the basic value of a single interaction  $I_{b_n}$ . Cumulative factor  $I_{c_n}$  then captures the additive effect of recent successive interactions. The reputational contribution  $I_n$  of most recent interaction  $n$  of any given user is estimated in the following way:

$$I_n = I_{b_n} + I_{c_n} = I_{b_n} \left(1 + \alpha \left(1 - \frac{1}{A_n + 1}\right)\right) \quad (4.4)$$

Here,  $\alpha$  is the weight of the cumulative part and  $A_n$  is the number of sequential activities. If there is no interaction at  $t_n$ , this part of interactions has a value of 0. Important property of this component of dynamic reputation is the notion of sequential activities. Two successive interactions made by a user are considered sequential if the time between those two activities is less or equal to the time parameter  $t_a$  which represents the time window of interaction. This time window represents maximum time spent by the user to make a meaningful contribution (post a question or answer or leave a comment).

$$\Delta_n = \frac{t_n - t_{n-1}}{t_a} \quad (4.5)$$

If  $\Delta_n < 1$  is less than one the number of sequential activities  $A_n$  will increase by one, which means that the user is continuing to communicate frequently. On the other hand, large values  $\Delta_n$  greatly increase the effect of the forgetting factor. This factor plays a major role in updating the total dynamic reputation of a user in each time step (after every recorded interaction):

$$T_n = T_{n-1} \beta^{\Delta_n} + I_n \quad (4.6)$$

Here,  $\beta$  is the forgetting factor. In our implementation of the model, the trust is updated each day for every user irrespective of their activity status. Therefore, the decay itself is a combination of  $\beta$  and  $\Delta_n$ : the more days pass without recorded interaction from a specific user, the more their reputation decays. Lower values of beta lead to faster decay of trust as shown on figure ??.

## Selection of dynamical reputation model parameters

One of the largest drawbacks of DIBRM is the parameter tuning problem. In previous applications of the model [62, 63] there was no single best set of parameter values for modeling dynamic reputation in Stack Exchange communities. For example, in [63] the best approximation of the official Stack Exchange reputation is obtained with  $t_a = 2, \beta = 1, \alpha = 1.4$  which means there is no active forgetting factor. In our application of DIBRM to SE communities we opted for a different set of parameter values. Details of parameter search and tuning are presented in SI.

For basic reputation contribution of a single interaction we selected  $I_{b_n} = 1$  and at the same time this is the threshold value of an active user. This value is intuitive as every interaction has initial contribution of +1 to user's reputation, although the previous works have used

## 4. The role of trust in knowledge based communities

---

values of +2 and +4. Following the previous work and after examining the median/average time between subsequent interactions of the same user, we selected  $t_a = 1$ , which also means that reputation in our model will be updated every day during the time-window of the analysis, regardless of whether the user is active or not. To emphasize the bursts of activity and frequent recent interactions, cumulative factor has a larger value  $\alpha = 2$ . Finally, the most delicate parameter is the forgetting factor, which at the same time determines the weight of past interactions and the reputational punishment due to user inactivity. Here we need to select the value of parameter  $\beta$  so we include the forgetting due to inactivity but not to penalize it too much. In Fig. A1 we show how different values of parameter  $\beta$  influence the time needed for user's reputation to fall on value  $I_n = 1$  due to user's inactivity and value of dynamical reputation in the moment of the last activity. The higher the value of parameter  $\beta$  and initial dynamical reputation of users, the longer time it takes for user's reputation to fall on baseline value. For parameter  $\beta = 0.9$  and  $I_n = 5$ , user's reputation falls on value  $I_n = 1$  after less than 20 days, while this time is doubled for  $\beta = 0.96$ . We see, that for higher values of parameter  $\beta$  the time needed for  $I_n$  to fall on value 1 becomes longer, and that the initial value of reputation becomes less important.

Figure A2 in SI shows the difference between the number of users that had at least one activity in the window of 30 days and number of users with reputation higher than 1 during the same period for different values of parameter  $\beta$ . The minimal difference between these two variables is observed for the values of  $\beta$  between 0.94 and 0.96 for both live and closed communities. Since we want to compare communities, we select  $\beta = 0.96$  after verifying that this level of reputational decay does not reduce the number of active users (based on their dynamic reputation) below the actual number of users who have been active (interacted with the community) in the time window of 30 days.

To summarize, our model of dynamical reputation has three parameters: 1) basic reputation contribution  $I_{bn} = 1$ ; 2) cumulative factor  $\alpha = 2$ ; 3) forgetting factor  $\beta = 0.96$ . The selected values of parameters are used for measuring dynamical reputation of user in all four pair SE communities. Given these values of parameters, the minimal reputation achieved by the user immediately after they have made an interaction in the SE community is 1. This reputation will decay below 1 if the user does not perform another interaction within the one-day time window. For any user in a community, when their reputation drops below 1, we consider this user inactive which means that the user at that time is not "visible" in the community and their past contributions at that time are unlikely to impact other users. The number of active users and mean user reputation for different Stack Exchange communities are shown in Fig. 4.14.

### Dynamic reputation - $\beta$ parameter

Our implementation of dynamic reputation model was based on  $\beta = 0.96$ . There are several reasons for selecting this value.

In Dynamic reputation model, the  $\beta$  parameter controls the strength of the forgetting factor

of the model. The value of this parameter should reflect the core feature of the reputational systems and make reputation easier to loose. Due to user's inactivity, any level of reputation will eventually decay to below 1. Dependence of time needed for reputation to drop below this level and the  $\beta$  parameter, as well as reputation before inactivity is shown on Figure 4.9. Here  $I_n$  is equal to the raw number of interactions in the community without forgetting or cumulative factor at work.

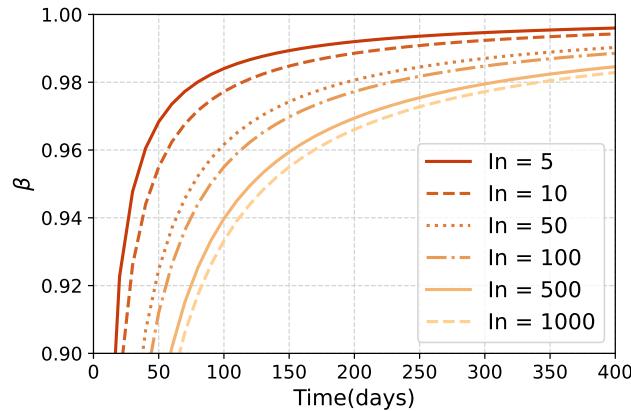


Figure 4.12: Dependence of parameter  $\beta$  and number of days  $\Delta$  needed for reputation  $I_n$  to drop to  $I_{n_0} = 1$ . Dependence of parameter  $\beta$  and number of days when reputation due inactivity decreases from  $I_n$  to  $I_0$  is given as  $\beta = \left(\frac{I_{n_0}}{I_n}\right)^{(1/\Delta)}$

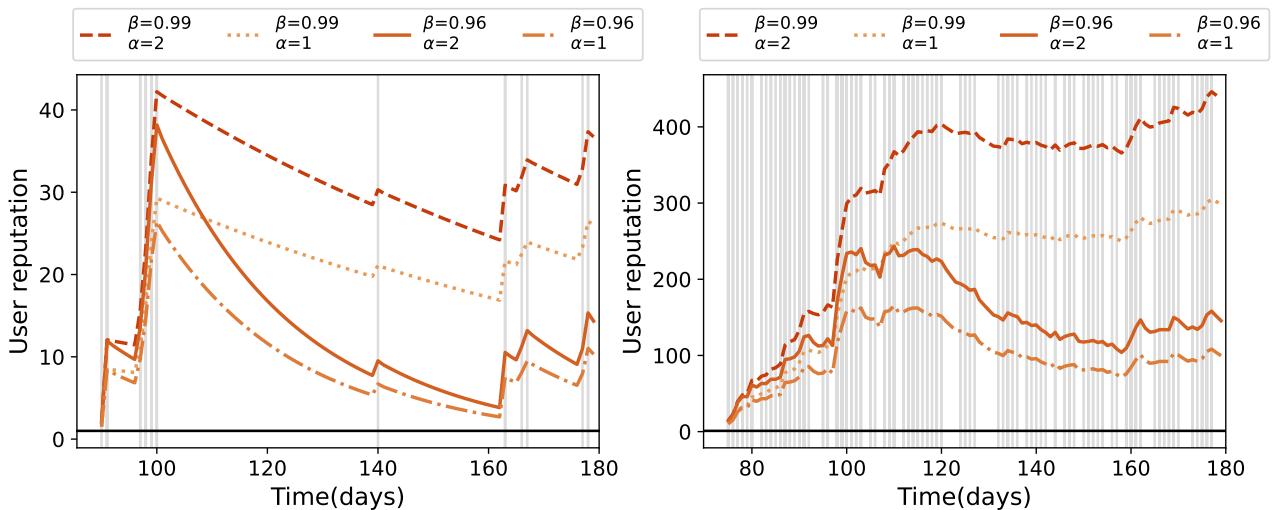


Figure 4.13: Single users reputations

For  $\beta$  values below 0.96, the decay is fast and within two to four months of inactivity even high values of reputation are reduced below the threshold. On the other hand, with  $\beta$  values the decay process is more differentiated and high reputation becomes harder to loose, surviving up to a year of inactivity. For  $\beta$  equal to 0.96, it takes a month for reputation based on 5 interactions to decay and around five months for high reputation based on 500 or 1000 interactions to decay below the threshold.

**30 days sliding window** We compared the number of users with estimated reputation higher than 1 for different parameters  $\beta$  and concluded that  $\beta$  close to 0.96 approximates the number of users with recorded interactions in a given 30 days sliding window. For each pair of communities we calculated number of users with at least one interactions in every 30 days sliding window and then we estimated several times series expressing the number of users with reputation higher than 1 for fixed  $\beta$ . Then we calculated the root mean square error (RMSE) between those time series for the first 200 days. Values of RMSE are shown on Figure 4.11. For each community, we can find parameter  $\beta$  that minimizes RMSE. Although  $\beta$  does not have a unique value across communities, it varies between 0.95 and 0.96.

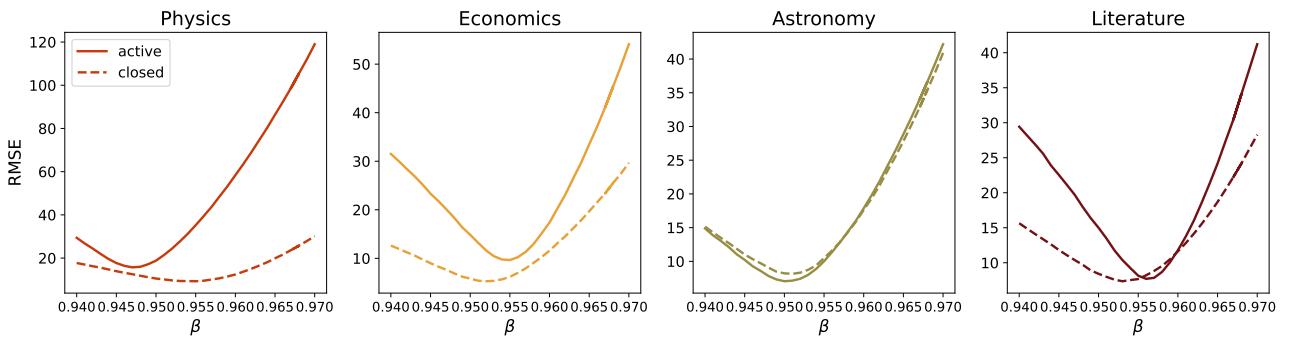


Figure 4.14: RMSE between number of active users in sliding window of 30 days and number of users with reputation  $> 1$  for  $0.94 < \beta < 0.97$  with step 0.001.

Figure 4.12 shows comparison between number of users in 30 days sliding window, number of users for these optimal values  $\beta = 0.954$  and  $\beta = 0.96$ . For  $\beta = 0.96$  we observe that in most communities estimated number of active users consistently slightly higher than the actual number of users which have made at least one interaction in that sliding window. This means that dynamic reputation model in some cases overestimates the reputation of the user, but far more important is that it never underestimates the real number of active users. Since we base our calculations of total and average reputation within the community only on users whose reputation is higher than the threshold this is important as no active users are disregarded by the model due to the value of the decay parameter.

Finally, it's important that our dynamic reputation captures the trend of long-term user activity. In Figure 4.13 solid lines show the time series of estimated dynamic reputation for  $\beta = 0.96$  while dashed lines show the number of users who were active in a given sliding window and continued to be active in the next one. Although the total estimated number of active users is expectedly higher, two time series follow similar trends in different communities.

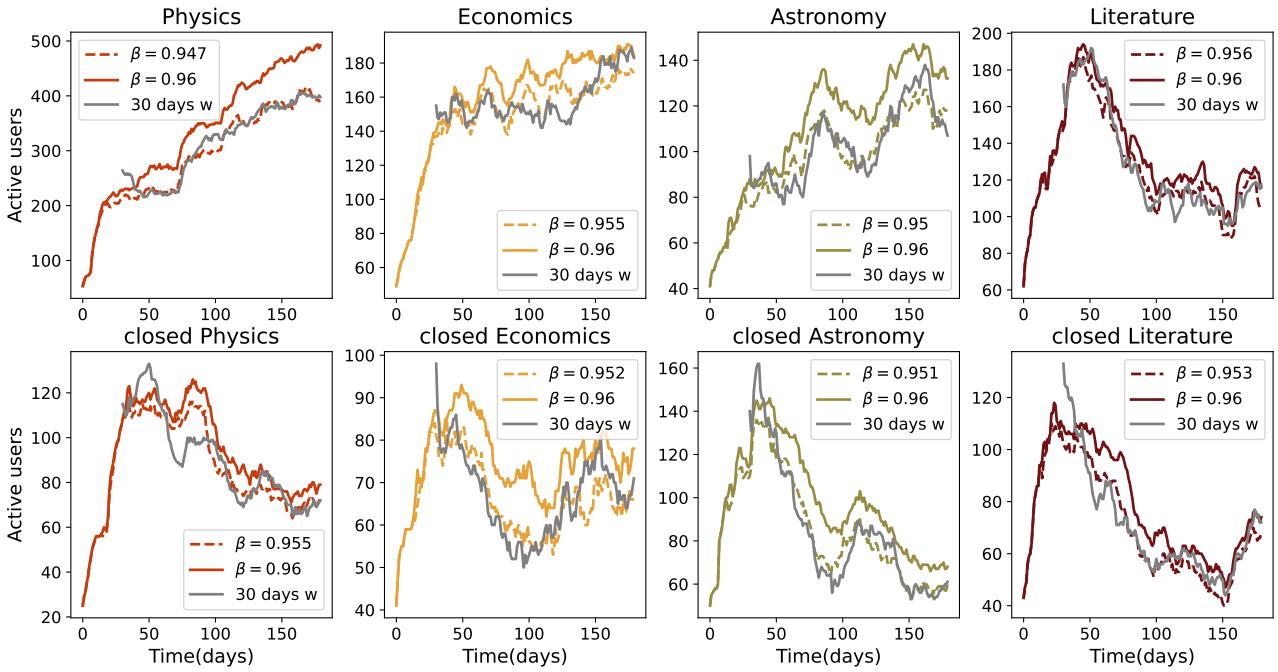


Figure 4.15: Number of active users in a sliding window of 30 days and number of users with dynamic reputation higher than 1 for  $\beta = 0.954$  and  $\beta = 0.96$  which provide the best fit to the number of users in 30 days sub-networks for each community

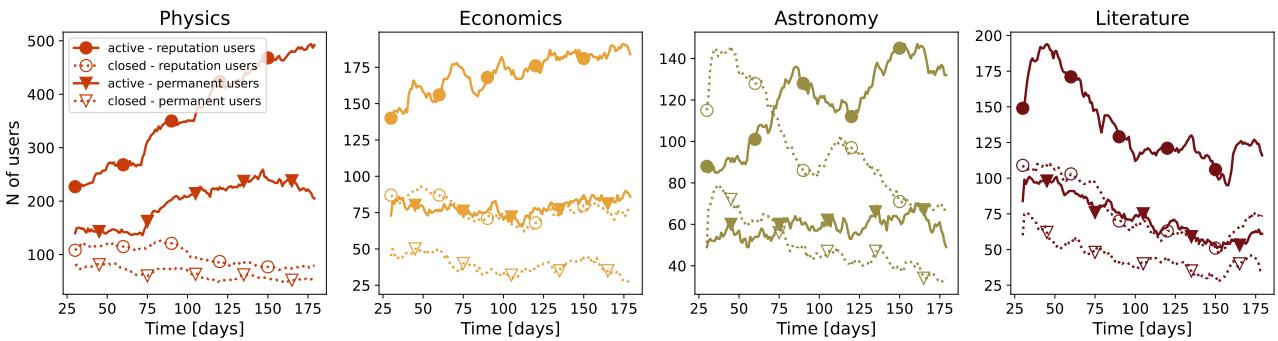


Figure 4.16: Solid lines represent number of users with dynamic reputation higher than 1 for  $\beta = 0.96$  while dashed lines are number of users within 30 days sliding window who were active and remained to be active. Blue lines are beta, while red lines are area51 communities.

## 4. The role of trust in knowledge based communities

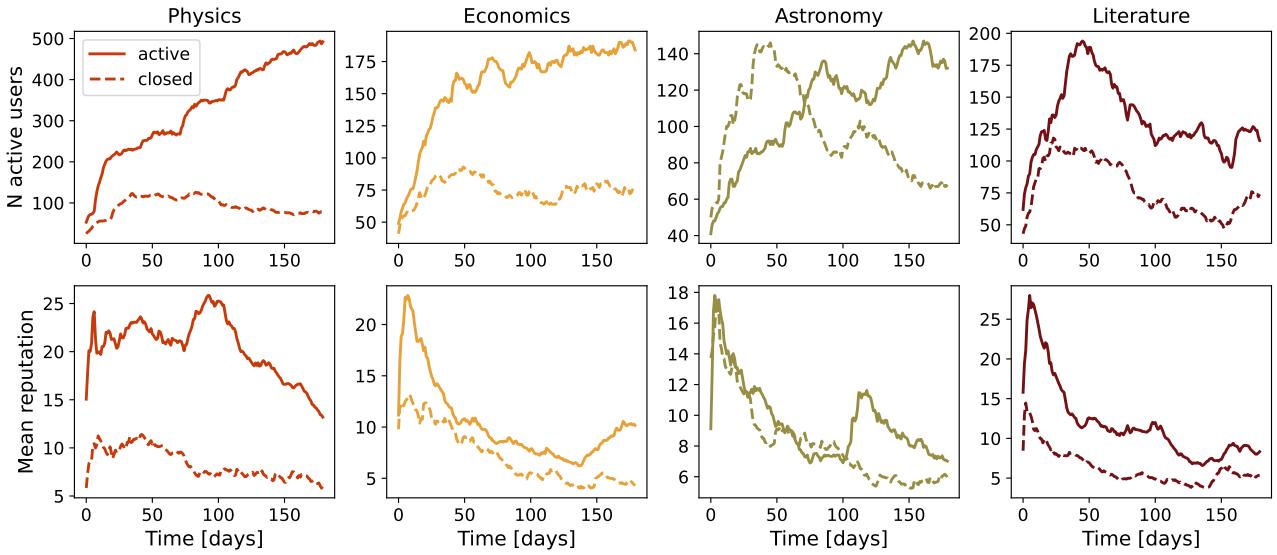


Figure 4.17: Dynamic Reputation on the four pairs of Stack Exchange websites: Astronomy, Literature, Economics, Physics and Theoretical Physics.

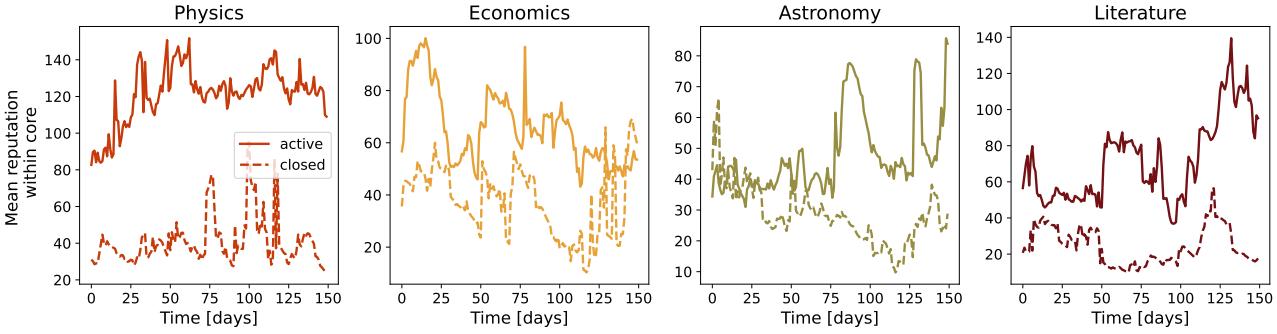


Figure 4.18: Dynamical reputation within core.

### Dynamic reputation of users within the network of interactions

Examined network properties suggest that there are structural differences between active and closed communities. Active communities have higher and more stable local cohesiveness compared to their closed counterparts. The overlap of the set of nodes in the core for active communities shows a significant overlap even for distant subnetworks, meaning that the membership of the core in active communities is more stable.

To further explore the differences between active and closed communities, we focus on dynamical reputation which is our proxy for collective trust in these communities. We investigate whether and how core-periphery structure is related to collective trust in the network. Figure 4.15 shows the mean dynamical reputation in the core of active and closed communities and its evolution during the observation period. There are clear differences between active and closed communities when it comes to dynamical reputation. The mean dynamical reputation of core users is always higher in active communities than in closed. As expected, the largest difference is observed between Physics and Theoretical Physics community. The difference between active communities which are still in the beta phase

and their closed counterparts is not as prominent, however, the active communities have higher mean dynamical reputation especially in the later phase of community life. The only difference in the pattern is observed for astronomy communities at the early phase of their life, when closed community has a higher value of dynamical reputation than active community. This is in line with similar patterns in the evolution of mean clustering and core-periphery structure.

By definition, the core consists of very active individuals and thus we expect higher total dynamical reputation of users in the core in comparison to the the total reputation of users belonging to subnetworks periphery. Figure A12 shows the ratio between the total reputation of core and periphery for closed and active communities and its evolution. The ratio between total reputation of core and periphery in Physics is always higher than in the Theoretical physics community. Similar pattern can be observed for literature communities, although the difference is not as clear as in the case of physics. Ratio of total dynamical reputation between core and periphery is higher for closed community than active one on the economics topic in the early days of community life. However, in the later stage of their lives this ratio becomes higher for active communities. Communities around astronomy topic deviate from this pattern, which once again shows the specificity of these communities.

To complete the description of the evolution of dynamic reputation active and closed communities, we examine the evolution of Gini index of dynamical reputation in the whole network which is shown in Fig. A5 in Supplementary Information. The Gini index is always higher for active communities than for closed ones, especially for later times in observation period. Only pattern of Astronomy communities deviates from the pattern observed for other three pairs during the early days. These results indicate that the dynamical reputation is distributed in the population more unequally in the active than in closed communities. The evolution of assortativity coefficient that measures correlations between dynamical reputation of connected users in the subnetworks, shown in Fig. A6, shows that networks are disassortative for the largest part of the observation period. These results suggest that users with high dynamical reputation have tendency to connect with users with low value of dynamical reputation.

In Figure ?? we show mean user reputation in core and in periphery over time (30 day sliding windows as before). We see that the mean user reputation in core is greater in the currently active sites (solid lines, top panels) than in their closed pairs (dashed lines). In the bottom panels, we see that the mean reputation on the network periphery has substantially lower values, and the difference between active and closed sites is less pronounced.

For reference in Fig 4.6 we show core sizes in all sites. We show these in absolute numbers (total number of nodes) and as a fraction of network size through time.

**Gini coefficient** Besides the number of active users (who at given moment of observation have reputation higher than the threshold) and the population mean value of dynamical reputation, we have investigated in more details the distribution of dynamical reputation within discussed communities. We have observed that the distributions are often skewed

## 4. The role of trust in knowledge based communities

---

which prompted us to compare the communities in terms of their Gini coefficient. The gini coefficient is a simple measure that shows us the degree of reputation inequality within the community. We calculate the value based on the dynamic reputation values of users at every time step (day) and report the values in Fig. 4.16. We see that all communities (both still active and closed ones) have gini coefficient values higher than 0.5 throughout first six month period. Interestingly, except in the case of Astronomy, currently active communities had higher reputation inequality every day during first six month period. As in many other measures, in the case of astronomy, closed community started as more unequal one (signalled by higher gini coef values), but after around two months the situation changed.

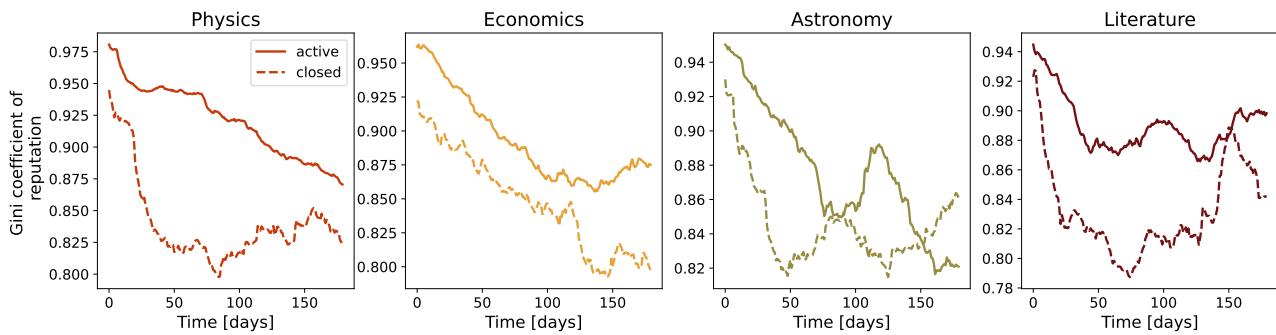


Figure 4.19: Gini index of dynamic reputation within population

## Dynamic reputation in the network of interactions

In the few figures below, we investigate whether users' dynamic reputation is related with users' position within the network.

### Dynamic Reputation assortativity

We first look at user interaction patterns, e.g. we investigate whether users connect with others of similar or different reputation (positive/negative assortativity). We operationalize this by measuring assortativity of dynamic reputation on interaction network. Practically this is a measure of correlation between dynamic reputation of users who are linked in the interaction network. These results are shown in Fig. 4.17. We look at 30 day unweighted undirected networks of interactions (questions, answers and comments) and calculate assortativity by using users' reputation on the last day of observed time window. We see small values of assortativity that are mostly negative, signaling weak correlations between reputation levels of interacting users. The fact that the values are mostly negative are expected, users of different dynamic reputation interact, e.g. active, high reputation users respond to the questions of new, less reputable users. Exceptions are closed astronomy and literature sites that occasionally had positive assortativity values, signaling existence of links between users of similar reputation levels.

### DynRep & Degree DynRep & BC

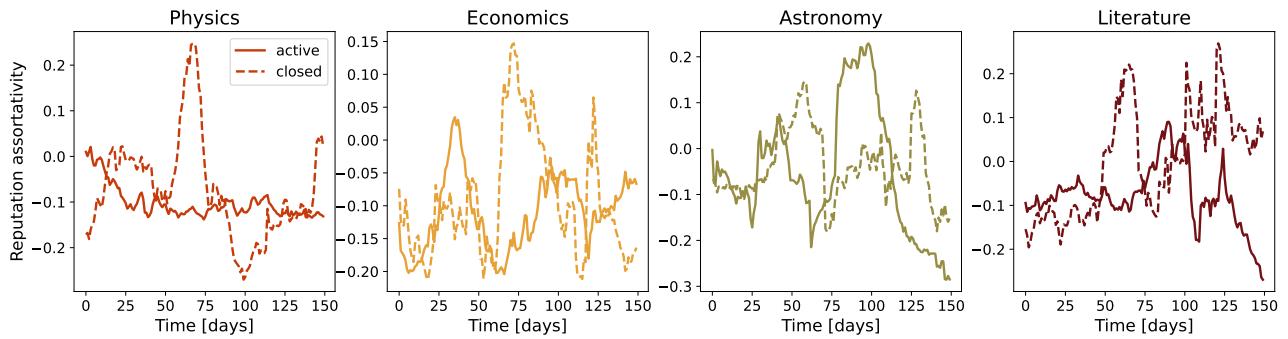


Figure 4.20: Dynamic Reputation assortativity in the network of interactions (questions, answers, comments, unweighted, undirected network). Solid lines - active sites; dashed lines - closed sites.

We continue to investigate whether the user's reputation correlates with typical network centrality measures calculated at user's node in the interaction network. As previously, we compare node's centrality in the 30 day network with the node's dynamic reputation on the last day of the period, repeat the process every day for the first six months. Correlation coefficient between dynamic reputation and degree in the network is very high, as expected, as most of the interactions that contributed to user's reputation are also present as links in the network. We show these results in Fig. 4.18(top). However, we again see the distinction between active and closed communities where this correlation is higher in active communities, except in the first month of sliding windows. Astronomy is an exception here as well as we see that the correlations were similar in both closed and still active sites throughout observed period. In the bottom panels of Fig. 4.18 we present correlation coefficients of dynamic reputation and user's betweenness centrality in the interaction network. These correlations are also high and most of the time higher in the later networks of active than closed communities. This is particularly interesting due to global nature of betweenness centrality measure and less obvious relation of it to user's dynamic reputation.

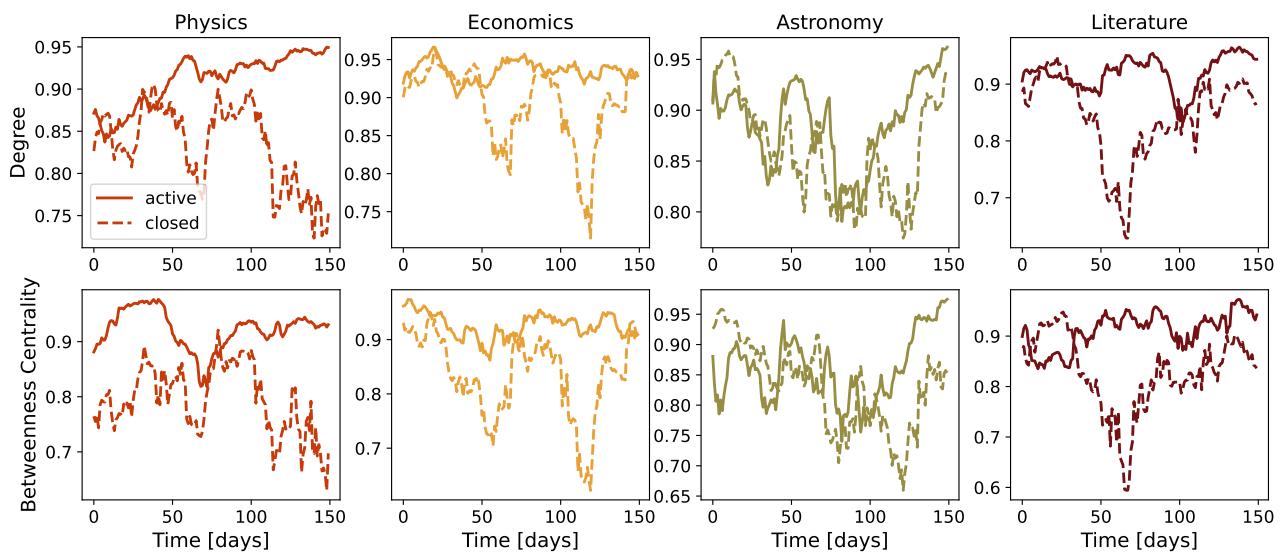


Figure 4.21: Coefficient of correlation between users' Dynamic Reputation and users' network degree (top) and users's betweenness centrality (bottom). Solid lines - active sites; dashed lines - closed sites.

# Appendix A

## The choice of the sliding window

There is no well-specified procedure for the choice of sliding window  $\tau$ . Previous studies showed that if  $\tau$  is small sub networks become sparse, while for too large sliding windows some important structural changes can not be observed [64, 65]. We analyse how networks properties and properties of dynamical reputation change with the window size, see SI for more details. Figure A13 in SI shows how considered network properties and dynamical reputation depend on the time window size for active and closed communities on the astronomy. We observe that fluctuations of all measures are more pronounced for time window of 10 days than for 30 and 60 days. However, we find that while the structural properties of networks evolve at different paces over varied time windows the trends remain very similar. The observed qualitative difference between closed and live communities is independent of  $\tau$ , especially if we compare time window size of 30 and 60 days. The time window size of 30 days ensures enough amount of interaction, even for closed communities, while the number of observation points remains relatively high. For these reasons, we choose a sliding window of 30 days.

In this section, we investigate how the size of sliding windows affect network properties over time. Figure A.3 summarize results for one pair of communities, area51 and beta astronomy, but similar conclusions can be observed for other pairs of sites. We show the network properties for sub-networks of 10, 30, and 60 days sliding windows. For a sliding window of 10 days, results may be too noisy and we may not observe some important trends in the community. The number of users for beta astronomy seems to fluctuate around some mean value. On the larger scale, 30 days window, it is more clear that the number of users slightly increase over time. Contrary, for too large an aggregation window (60 days), important information about the time series can be lost, such as the local minimum of the number of users around time step 80 that is observed for the 30-day sliding window. Looking into other network characteristics such as L/N and clustering we conclude that differences between closed and active sites are more transparent with a larger aggregation window, still, on each scale, beta sites show a higher number of nodes, number of links per node and clustering coefficient.

As before we study the structure of created sub-networks through the lens of core-

## A. The choice of the sliding window

---

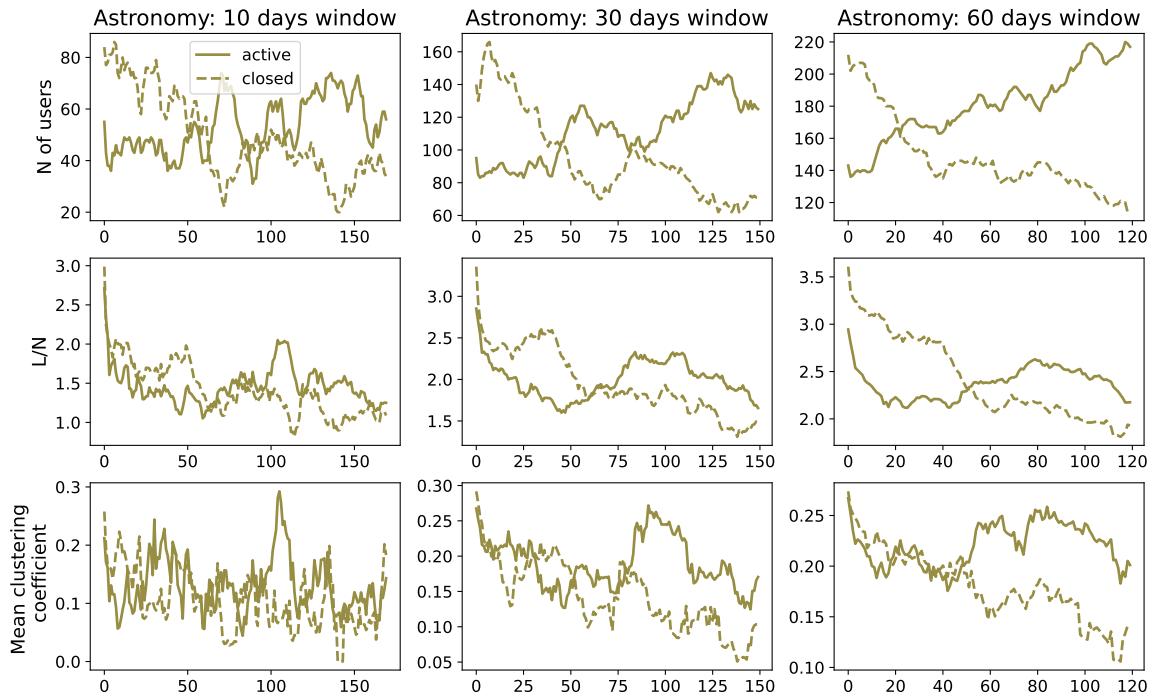


Figure A.1: Results for different sliding windows. Example is for astronomy, blue solid lines - active, orange dashed lines - closed site.

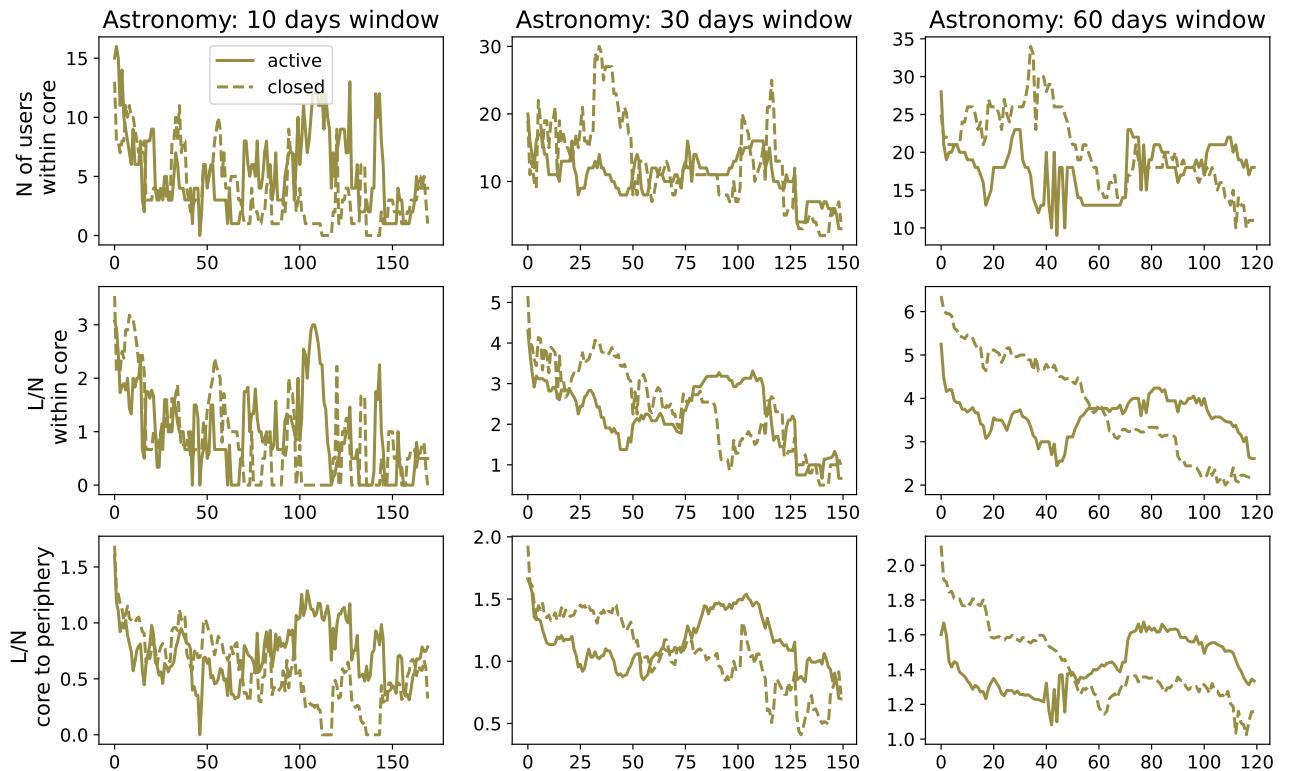


Figure A.2: Results for different sliding windows. Example is for astronomy, blue solid lines - active, orange dashed lines - closed site.

periphery structure. On small scales, the window of 10 days, there are often few, or even no nodes in the core and it can affect the calculation of other measures of interest. Such behaviour is more typical for closed communities. With the size of the sliding window number of nodes in the core increases and results of core-periphery measures become smoother. Finally, the choice of the sliding window does not change conclusions that core users in the beta communities produce more activity and make the strong core. However, our main results are shown for a sliding window of 30 days, as it makes a good compromise between large and small time scales.

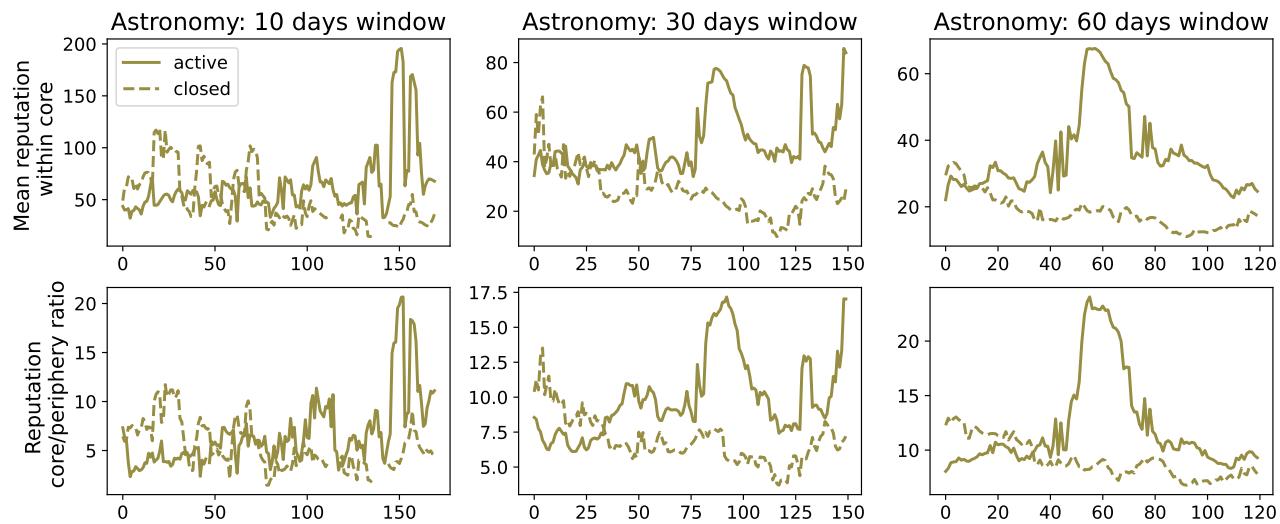


Figure A.3: Results for different sliding windows. Example is for astronomy, blue solid lines - active, orange dashed lines - closed site.



# Appendix B

## Robustness of core-periphery algorithm

### Precision and recall

Consider the network  $G(V, L)$ , with a set of nodes  $V$  and a set of links between them  $L$ . The stochastic community detection algorithms may converge to different configurations. To quantify the similarity between the obtained structures and robustness of the algorithm, we run 50 iterations and calculate several similarity measures between pairwise partitions  $C$  and  $C'$ .

The core-periphery structure has two groups so confusion matrix [66] can be defined as:

		partition C	
		core	periphery
partition $C'$	core	$n_{TP}$	$n_{FN}$
	periphery	$n_{FP}$	$n_{TN}$

The diagonal elements correspond to the number of nodes found in the same class in both node configurations. The number of nodes in the core found in  $C$  and  $C'$  is denoted as true positive  $n_{TP}$ , while the number of nodes in the periphery in  $C$  and  $C'$  is denoted as true negative  $n_{TN}$ . The off-diagonal elements of the confusion matrix indicate the number of nodes differently classified. We can define the number of nodes found in the first configuration  $C$  in the core but in  $C'$  in the periphery as a false positive,  $n_{FP}$ , similarly the number of nodes found in the periphery in the partition  $C$ , and in the core in partition  $C'$  as a false positive,  $n_{FN}$ .

From the confusion matrix, we can write the precision  $P = n_{TP} / (n_{TP} + n_{FP})$  and recall  $R = n_{TN} / (n_{TN} + n_{FN})$ . These measures range from 0 to 1. The precision (recall) corresponds to the proportion of instances predicted to belong (not belong) to the considered class and which indeed do (do not) [66].

## F1 measure

The **F1 measure** is the harmonic mean of precision and recall [66]:

$$F_1 = 2 \frac{P \cdot R}{P + R} = \frac{2n_{TP}}{2n_{TP} + n_{FN} + n_{FP}} \quad (\text{B.1})$$

It can be interpreted as a measure of overlap between true and estimated classes; it is 0 for no overlap to 1 if overlap is complete.

## Jaccard coefficient

The **Jaccard's coefficient** is the ratio of two classes' intersection to their union [66]. It can also be expressed in terms of confusion matrix:

$$J = \frac{C_{core} \cap C'_{core}}{C_{core} \cup C'_{core}} = \frac{n_{TP}}{n_{TP} + n_{FP} + n_{FN}} \quad (\text{B.2})$$

## Normalized mutual information

**Normalized mutual information (NMI)** is similarity measure between two partitions  $C$  and  $C'$  based on information theory [67]:

$$NMI(C, C') = \frac{MI(C, C')}{(H(C) + H(C'))/2} \quad (\text{B.3})$$

where  $MI$  is mutual information between sets  $C$  and  $C'$ , while  $H(C)$  is entropy of given partition. The entropy is defined as  $H(C) = -\sum_{i=1}^{|C|} P(i) \log(P(i))$ , where  $P(i) = |U_i|/N$  is the probability that an object is randomly classified as  $i$  (in this special case  $i = 0$ , the node belongs to the core, or  $i = 1$ , the node belongs to the periphery). The mutual information between sets  $C$  and  $C'$  measures the probability that the randomly chosen node is a member of the same group in both partitions:

$$MI(C, C') = \sum_i \sum_j P(i, j) \log\left(\frac{P(i, j)}{P(i)P'(j)}\right) \quad (\text{B.4})$$

where  $P(i, j) = |U_i \cap U_j|/N$

$NMI$  ranges from 0 when the partitions are independent to 1 if they are identical.

## Adjusted rand index

**Adjusted rand index.** For the set of nodes  $V$ , with  $n$  nodes, consider all possible combination of pairs  $(v_i, v_j)$ . We can select the number of the pairs where nodes belong to the same group in both partitions,  $C$  and  $C'$ , denoted as  $a$ . Similarly, as  $b$ , we can define the number of pairs whose nodes belong to different groups in partitions. Then, unadjusted rand index [68] is given as  $RI = \frac{a+b}{\binom{n}{2}}$ , where  $\binom{n}{2}$  is number of all possible pairs. The RI between two randomly

assigned partitions is not close to zero; for that reason, it is common to use the adjusted rand index [69], defined as:

$$ARI = \frac{RI - E[RI]}{\max(RI) - E[RI]} \quad (\text{B.5})$$

where  $E[RI]$  is expected value of  $RI$ , and  $\max(RI)$  is maximum value of  $RI$ .

As example we show analysis of inferred sample of core-periphery structures for 30 days closed Astronomy, Stack Exchange networks, Figure B.1. We represent the mean minimum description length (MDL) and the mean number of nodes in the core with standard deviation. MDL does not change much between inferred core-periphery structures; the difference between obtained configurations is still notable in the number of nodes in the core. To investigate in more details similarity between obtained core-periphery configurations in the sample we calculate several measures between pair-wise partitions such as normalized mutual information, adjusted rand index, F1 measure and Jaccard index. These measures are greater than 0.5 and, in most cases, greater than 0.9, indicating stability of the inferred core-periphery structures.

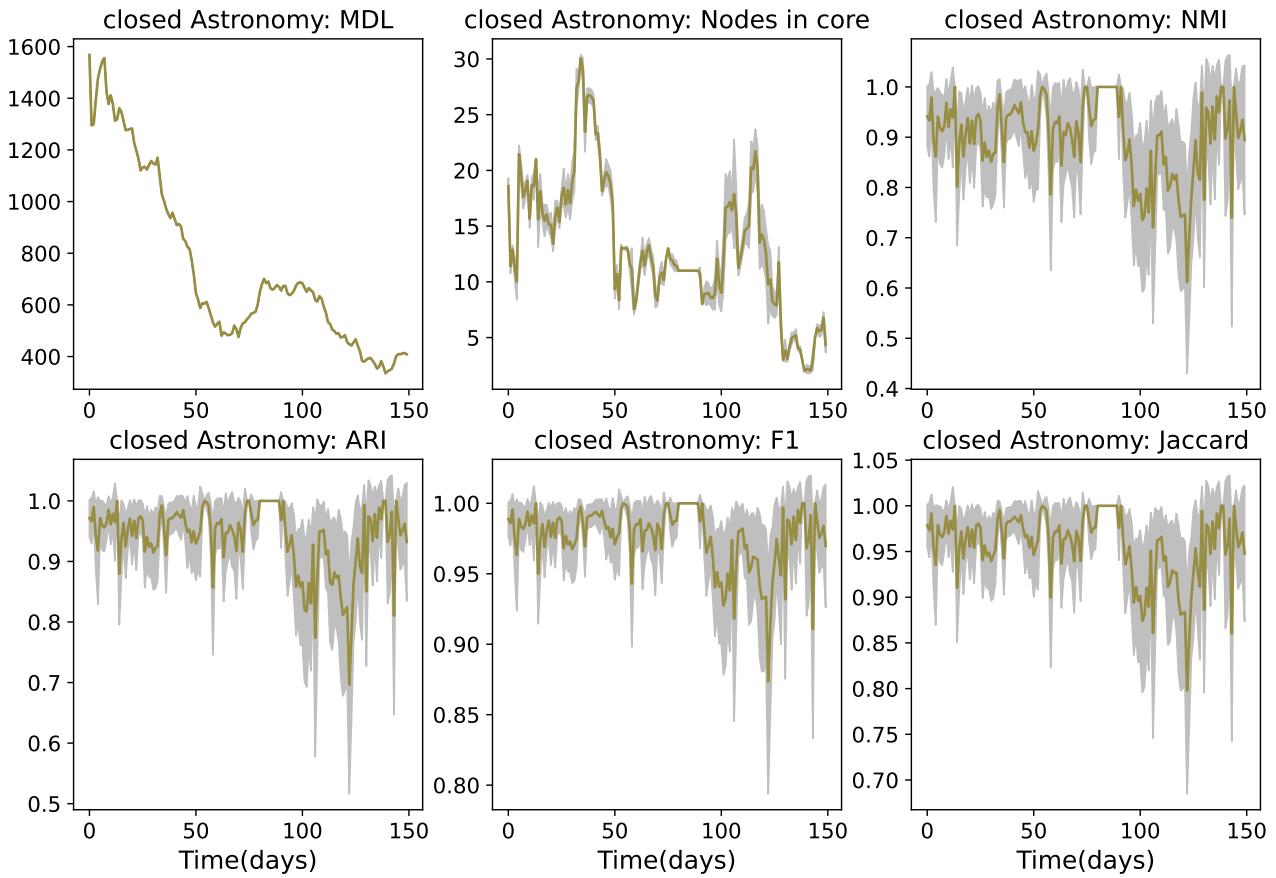


Figure B.1: Minimum description length, number of nodes in core, normalized mutual information, adjusted rand index, F1 measure and Jaccard index, among 50 samples for 30-days sub-networks. Results are given for closed astronomy.



# Bibliography

- [1] Stefano Boccaletti, Vito Latora, Yamir Moreno, Martin Chavez, and D-U Hwang. Complex networks: Structure and dynamics. *Physics reports*, 424(4-5):175–308, 2006.
- [2] M. E. J. Newman. *Networks: An Introductio*. Oxford University Press, 2010.
- [3] Martin Rosvall, Jean-Charles Delvenne, Michael T. Schaub, and Renaud Lambiotte. Different approaches to community detection. *CoRR*, abs/1712.06468, 2017.
- [4] Santo Fortunato and Darko Hric. Community detection in networks: A user guide. *Physics Reports*, 659:1 – 44, 2016. Community detection in networks: A user guide.
- [5] Thorben Funke and Till Becker. Stochastic block models: A comparison of variants and inference methods. *PLOS ONE*, 14(4):1–40, 04 2019.
- [6] Hernán A Makse, Shlomo Havlin, Moshe Schwartz, and H Eugene Stanley. Method for generating long-range correlations for large systems. *Physical Review E*, 53(5):5445, 1996.
- [7] Jan W Kantelhardt, Eva Koscielny-Bunde, Henio HA Rego, Shlomo Havlin, and Armin Bunde. Detecting long-range correlations with detrended fluctuation analysis. *Physica A: Statistical Mechanics and its Applications*, 295(3-4):441–454, 2001.
- [8] C-K Peng, Sergey V Buldyrev, Shlomo Havlin, Michael Simons, H Eugene Stanley, and Ary L Goldberger. Mosaic organization of dna nucleotides. *Physical review e*, 49(2):1685, 1994.
- [9] Jan W Kantelhardt, Stephan A Zschiegner, Eva Koscielny-Bunde, Shlomo Havlin, Armin Bunde, and H Eugene Stanley. Multifractal detrended fluctuation analysis of nonstationary time series. *Physica A: Statistical Mechanics and its Applications*, 316(1-4):87–114, 2002.
- [10] Espen Alexander Fürst EAFI Ihlen. Introduction to multifractal detrended fluctuation analysis in matlab. *Frontiers in physiology*, 3:141, 2012.
- [11] Kamalika Basu Hajra and Parongama Sen. Phase transitions in an aging network. *Physical Review E*, 70(5):056103, 2004.
- [12] Sergey N Dorogovtsev and José FF Mendes. Scaling properties of scale-free evolving networks: Continuous approach. *Physical Review E*, 63(5):056125, 2001.

## Bibliography

---

- [13] Tiago A Schieber, Laura Carpi, Albert Díaz-Guilera, Panos M Pardalos, Cristina Masoller, and Martín G Ravetti. Quantification of network structural dissimilarities. *Nature communications*, 8(1):1–10, 2017.
- [14] Marija Mitrović and Bosiljka Tadić. Emergence and structure of cybercommunities. In *Springer Optimization and Its Applications*, volume 57, pages 209–227. Springer International Publishing, 2012.
- [15] Marija Mitrović Dankulov, Roderick Melnik, and Bosiljka Tadić. The dynamics of meaningful social interactions and the emergence of collective knowledge. *Scientific reports*, 5(1):1–10, 2015.
- [16] Milovan Suvakov, Marija Mitrovic, Vladimir Gligorijevic, and Bosiljka Tadic. How the online social networks are used: dialogues-based structure of myspace. *Journal of The Royal Society Interface*, 10(79):20120819, 2013.
- [17] Petter Holme and Jari Saramäki. Temporal networks. *Physics reports*, 519(3):97–125, 2012.
- [18] Marc Barthelemy. *The structure and dynamics of cities*. Cambridge University Press, 2016.
- [19] César A Hidalgo and Ricardo Hausmann. The building blocks of economic complexity. *Proceedings of the national academy of sciences*, 106(26):10570–10575, 2009.
- [20] Jelena Smiljanić, Arnab Chatterjee, Tomi Kauppinen, and Marija Mitrović Dankulov. A theoretical model for the associative nature of conference participation. *PloS one*, 11(2):e0148528, 2016.
- [21] Ali Montazeri, Soghra Jarvandi, Shahpar Haghigat, Mariam Vahdani, Akram Sajadian, Mandana Ebrahimi, and Mehregan Haji-Mahmoodi. Anxiety and depression in breast cancer patients before and after participation in a cancer support group. *Patient education and counseling*, 45(3):195–198, 2001.
- [22] Kathryn P Davison, James W Pennebaker, and Sally S Dickerson. Who talks? the social psychology of illness support groups. *American Psychologist*, 55(2):205, 2000.
- [23] Wendy K Tam Cho, James G Gimpel, Daron R Shaw, et al. The tea party movement and the geography of collective action. *Quarterly Journal of Political Science*, 7(2):105–133, 2012.
- [24] Sinan Aral and Dylan Walker. Identifying influential and susceptible members of social networks. *Science*, 337(6092):337–341, 2012.
- [25] Sandra González-Bailón, Javier Borge-Holthoefer, and Yamir Moreno. Broadcasters and hidden influentials in online protest diffusion. *American behavioral scientist*, 57(7):943–965, 2013.

- [26] János Török, Gerardo Iñiguez, Taha Yasseri, Maxi San Miguel, Kimmo Kaski, and János Kertész. Opinions, conflicts, and consensus: modeling social dynamics in a collaborative environment. *Physical review letters*, 110(8):088701, 2013.
- [27] Taha Yasseri, Robert Sumi, András Rung, András Kornai, and János Kertész. Dynamics of conflicts in wikipedia. *PloS one*, 7(6):e38869, 2012.
- [28] Claudio Castellano, Santo Fortunato, and Vittorio Loreto. Statistical physics of social dynamics. *Reviews of modern physics*, 81(2):591, 2009.
- [29] Arnab Chatterjee, Marija Mitrović, and Santo Fortunato. Universality in voting behavior: an empirical analysis. *Scientific reports*, 3(1):1–9, 2013.
- [30] Filippo Radicchi, Santo Fortunato, and Claudio Castellano. Universality of citation distributions: Toward an objective measure of scientific impact. *Proceedings of the National Academy of Sciences*, 105(45):17268–17272, 2008.
- [31] Lars Backstrom, Dan Huttenlocher, Jon Kleinberg, and Xiangyang Lan. Group formation in large social networks: membership, growth, and evolution. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 44–54, 2006.
- [32] Jelena Smiljanić and Marija Mitrović Dankulov. Associative nature of event participation dynamics: A network theory approach. *PloS one*, 12(2):e0171565, 2017.
- [33] Jure Leskovec, Lars Backstrom, Ravi Kumar, and Andrew Tomkins. Microscopic evolution of social networks. In *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 462–470, 2008.
- [34] Elena Zheleva, Hossam Sharara, and Lise Getoor. Co-evolution of social and affiliation networks. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1007–1016, 2009.
- [35] Luís A Nunes Amaral, Sergey V Buldyrev, Shlomo Havlin, Heiko Leschhorn, Philipp Maass, Michael A Salinger, H Eugene Stanley, and Michael HR Stanley. Scaling behavior in economics: I. empirical results for company growth. *Journal de Physique I*, 7(4):621–633, 1997.
- [36] Michael HR Stanley, Luis AN Amaral, Sergey V Buldyrev, Shlomo Havlin, Heiko Leschhorn, Philipp Maass, Michael A Salinger, and H Eugene Stanley. Scaling behaviour in the growth of companies. *Nature*, 379(6568):804–806, 1996.
- [37] Giorgio Fazio and Marco Modica. Pareto or log-normal? best fit and truncation in the distribution of all cities. *Journal of Regional Science*, 55(5):736–756, 2015.
- [38] Konglin Zhu, Wenzhong Li, Xiaoming Fu, and Jan Nagler. How do online social networks grow? *Plos one*, 9(6):e100023, 2014.

## Bibliography

---

- [39] Sanjay Ram Kairam, Dan J Wang, and Jure Leskovec. The life and death of online groups: Predicting group growth and longevity. In *Proceedings of the fifth ACM international conference on Web search and data mining*, pages 673–682, 2012.
- [40] Jeff Alstott, Ed Bullmore, and Dietmar Plenz. powerlaw: A python package for analysis of heavy-tailed distributions. *PLOS ONE*, 9(1):1–11, 01 2014.
- [41] Michael Mitzenmacher. A brief history of generative models for power law and lognormal distributions. *Internet mathematics*, 1(2):226–251, 2004.
- [42] Hernan Mondani, Petter Holme, and Fredrik Liljeros. Fat-tailed fluctuations in the size of organizations: the role of social influence. *PLoS One*, 9(7):e100527, 2014.
- [43] Dongfeng Fu, Fabio Pammolli, Sergey V Buldyrev, Massimo Riccaboni, Kaushik Matia, Kazuko Yamasaki, and H Eugene Stanley. The growth of business firms: Theoretical framework and empirical evidence. *Proceedings of the National Academy of Sciences*, 102(52):18801–18806, 2005.
- [44] Gerald F Frasco, Jie Sun, Hernán D Rozenfeld, and Daniel Ben-Avraham. Spatially distributed social complex networks. *Physical Review X*, 4(1):011008, 2014.
- [45] Jiang-Hai Qian, Qu Chen, Ding-Ding Han, Yu-Gang Ma, and Wen-Qing Shen. Origin of gibrat law in internet: Asymmetric distribution of the correlation. *Physical Review E*, 89(6):062808, 2014.
- [46] Marija Mitrović, Georgios Paltoglou, and Bosiljka Tadić. Quantitative analysis of bloggers’ collective behavior powered by emotions. *Journal of Statistical Mechanics: Theory and Experiment*, 2011(02):P02005, 2011.
- [47] Marija Mitrović Dankulov, Roderick Melnik, and Bosiljka Tadić. The dynamics of meaningful social interactions and the emergence of collective knowledge. *Scientific reports*, 5(1):1–10, 2015.
- [48] Ana Vranić and Marija Mitrović Dankulov. Growth signals determine the topology of evolving networks. *Journal of Statistical Mechanics: Theory and Experiment*, 2021(1):013405, 2021.
- [49] Akrati Saxena and Harita Reddy. Users roles identification on online crowdsourced q&a platforms and encyclopedias: a survey. *Journal of Computational Social Science*, pages 1–33, 2021.
- [50] Tiago Santos, Simon Walk, Roman Kern, Markus Strohmaier, and Denis Helic. Activity archetypes in question-and-answer (q&a) websites—a study of 50 stack exchange instances. *ACM Transactions on Social Computing*, 2(1):1–23, 2019.

- [51] Rogier Slag, Mike de Waard, and Alberto Bacchelli. One-day flies on stackoverflow-why the vast majority of stackoverflow users only posts once. In *2015 IEEE/ACM 12th Working Conference on Mining Software Repositories*, pages 458–461. IEEE, 2015.
- [52] Anamika Chhabra and S RS Iyengar. Activity-selection behavior of users in stackexchange websites. In *Companion Proceedings of the Web Conference 2020*, pages 105–106, 2020.
- [53] Himmel Dev, Chase Geigle, Qingtao Hu, Jiahui Zheng, and Hari Sundaram. The size conundrum: Why online knowledge markets can fail at scale. In *Proceedings of the 2018 World Wide Web Conference*, pages 65–75, 2018.
- [54] Tiago Santos, Simon Walk, Roman Kern, Markus Strohmaier, and Denis Helic. Self-and cross-excitation in stack exchange question & answer communities. In *The World Wide Web Conference*, pages 1634–1645, 2019.
- [55] Yaniv Dover, Jacob Goldenberg, and Daniel Shapira. Sustainable online communities exhibit distinct hierarchical structures across scales of size. *Proceedings of the Royal Society A*, 476(2239):20190730, 2020.
- [56] Stefano Boccaletti, Vito Latora, Yamir Moreno, Martin Chavez, and D-U Hwang. Complex networks: Structure and dynamics. *Physics reports*, 424(4-5):175–308, 2006.
- [57] Chiara Orsini, Marija M Dankulov, Pol Colomer-de Simón, Almerima Jamakovic, Priya Mahadevan, Amin Vahdat, Kevin E Bassler, Zoltán Toroczkai, Marián Boguná, Guido Caldarelli, et al. Quantifying randomness in real networks. *Nature communications*, 6(1):1–10, 2015.
- [58] Damon Centola, Víctor M Eguíluz, and Michael W Macy. Cascade dynamics of complex propagation. *Physica A: Statistical Mechanics and its Applications*, 374(1):449–456, 2007.
- [59] Santo Fortunato. Community detection in graphs. *Physics reports*, 486(3-5):75–174, 2010.
- [60] Ryan J Gallagher, Jean-Gabriel Young, and Brooke Foucault Welles. A clarified typology of core-periphery structure in networks. *arXiv preprint arXiv:2005.10191*, 2020.
- [61] Claudiu Duma, Nahid Shahmehri, and Germano Caronni. Dynamic trust metrics for peer-to-peer systems. In *16th International Workshop on Database and Expert Systems Applications (DEXA'05)*, pages 776–781. IEEE, 2005.
- [62] A. Melnikov, J. Lee, V. Rivera, M. Mazzara, and L. Longo. Towards dynamic interaction-based reputation models. In *2018 IEEE 32nd International Conference on Advanced Information Networking and Applications (AINA)*, pages 422–428, 2018.
- [63] Ekaterina Yashkina, Arseny Pinigin, JooYoung Lee, Manuel Mazzara, Akinlolu Solomon Adekoju, Adam Zubair, and Luca Longo. Expressing trust with temporal frequency of

## Bibliography

---

- user interaction in online communities. *Advances in Intelligent Systems and Computing*, pages 1133–1146, Cham, 2020. Springer International Publishing.
- [64] Gautier Krings, Márton Karsai, Sebastian Bernhardsson, Vincent D Blondel, and Jari Saramäki. Effects of time window size and placement on the structure of an aggregated communication network. *EPJ Data Science*, 1(1):1–16, 2012.
  - [65] Naomi A Arnold, Benjamin Steer, Imane Hafnaoui, Hugo A Parada G, Raul J Mondragon, Félix Cuadrado, and Richard G Clegg. Moving with the times: Investigating the alt-right network gab with temporal interaction graphs. *Proceedings of the ACM on Human-Computer Interaction*, 5(CSCW2):1–17, 2021.
  - [66] Vincent Labatut and Hocine Cherifi. Accuracy measures for the comparison of classifiers. *arXiv preprint arXiv:1207.3790*, 2012.
  - [67] Leon Danon, Albert Diaz-Guilera, Jordi Duch, and Alex Arenas. Comparing community structure identification. *Journal of statistical mechanics: Theory and experiment*, 2005(09):P09008, 2005.
  - [68] Jorge M Santos and Mark Embrechts. On the use of the adjusted rand index as a metric for evaluating supervised classification. In *International conference on artificial neural networks*, pages 175–184. Springer, 2009.
  - [69] Lawrence Hubert and Phipps Arabie. Comparing partitions. *Journal of classification*, 2(1):193–218, 1985.