



UNIVERSITÀ  
DEGLI STUDI  
FIRENZE

**Dottorato in Informatica, Sistemi e Telecomunicazioni**

Indirizzo: Dinamica Non Lineare e Sistemi Complessi

Ciclo XXVI

Coordinatore: Prof. Luigi Chisci

*Local Dynamics in Complex Networks:*

*from community detection to epidemic spreading*

Settore Scientifico Disciplinare FIS/03

**Dottorando:**

Dott. Emanuele Massaro

**Tutore:**

Dott. Franco Bagnoli

**Co-Tutore:**

Dott. Andrea Guazzini

**Referente**

Prof. Stefano Ruffo

Anni 2011/2013



*“I do not feel obliged to believe that the same God who has endowed us with sense, reason, and intellect has intended us to forget their use..”*

Galileo Galilei

# *Abstract*

## **Local Dynamics in complex networks: from community detection to epidemic spreading**

by Emanuele MASSARO

The great attention around complex networks in the last ten years is due to two main factors: on one hand interactions in many real systems can be generalized in the form of graphs or networks; on the other hand the significant amount of data coming from the web and analysed with the methods of the network science may allow to predict future behavior in many disciplines, especially in the social context. We faced the problem of community detection and epidemic spreading in complex networks, considering also the role of a local dynamics, a factor often neglected in literature. Regarding community detection, we addressed the issue of modifying existing well performing algorithms by incorporating elements from the domain application fields, i.e., domain-inspired. We focused on an approach inspired by psychology, which may be useful for further strengthening the link between social network studies and the mathematics of community detection. We introduced a community-detection algorithm derived from the van Dongen's Markov Cluster algorithm (MCL) method by considering networks' nodes as agents capable to take decisions. In this framework we introduced a memory factor to mimic typical human behaviors such as the oblivion effect. The method is based on information diffusion and includes a non-linear processing phase. Our approach has three important features: the capacity of detecting overlapping communities, the capability of identifying communities from an individual point of view and the fine tuning of the community detectability with respect to a prior knowledge of data. We tested our method in both synthetic and real-world (dynamic and static) networks and we showed that the adaptation of more complex heuristics allows to reach very efficient results comparing with other well-known community detection algorithms. Turning towards epidemic spreading in complex networks, we studied the influence of global, local and community-level risk perception on the extinction probability of a disease in several models of social networks. In particular, we studied the infection progression in a susceptible-infected-susceptible (SIS) model on several modular networks, formed by a certain number of random and scale-free communities. We confirmed the expectation that in scale-free networks the progression is faster than in random ones with the same average connectivity degree. In this thesis we found that the knowledge of the infection level in one's own neighborhood is the most effective property in stopping the spreading of a disease, but at the same time the more expensive one in terms of the quantity of required information. Therefore,

the cost/effectiveness optimum resulted in a tradeoff among several parameters. Finally we developed a self-organized percolation method for modeling the risk perception in epidemic spreading in both single-layer and multiplex networks.

## *Acknowledgements*

One of the joys of completion is to look over the journey past and remember all the friends and family who have helped and supported me along this long but fulfilling road.

I would like to express my heartfelt gratitude to Dr.Franco Bagnoli who is not only a mentor but a dear friend. I could not have asked for better role models, each inspirational, supportive, and patient. I could not be prouder of my academic roots and hope that I can in turn pass on the research values and the dreams that he has given to me.

I would also like to thank my co-advisor and friend, Dr.Andrea Guazzini who passed on to me important notions about cognitive science and psychology and important tips for overcoming this PhD.

My PhD was funded by European research project RECOGNITION, and I would like to thank the all the members of the project for their support.

As a member of both the Department of Information Engineering (DINFO) and the Department of Physics and Astronomy at the University of Florence, I have been surrounded by wonderful colleagues; both communities have provided a rich and fertile environment to study and explore new ideas. At DINFO, I would like to thank Prof. Luigi Chisci. Life at the Department of Physics and Astronomy is ever fantastic, and, I would like to thank all the people who supported this thesis: Prof. Stefano Ruffo and Prof. Duccio Fanelli for their scientific advices, Claudia and Alessandro who started their PhD with me, Giovanna Pacini for her patience and solace and many others, in particular Francesca, Gianluca and Aurelio.

Many thanks to the staff at ABC group of the Max Planck Institute of Human Development (2012), with kind regards to Henirk Olsson. I am grateful for the chance to visit and be a part of the lab. Thank you for welcoming me as a friend and helping to develop the ideas in this thesis.

I would not have contemplated this road if not for my parents, Carla and Gabriele, and my sisters, Margherita, Ilaria and Cinzia.

# Contents

<b>Acknowledgements</b>	<b>vi</b>
<b>List of Figures</b>	<b>ix</b>
<b>List of Tables</b>	<b>xi</b>
<b>1 Introduction</b>	<b>1</b>
<b>2 Complex Networks</b>	<b>3</b>
2.1 Introduction . . . . .	3
2.2 Some basic concepts . . . . .	4
2.2.1 Average path length . . . . .	4
2.2.2 Clustering coefficient . . . . .	5
2.2.3 Degree distribution . . . . .	5
2.3 Complex Network models . . . . .	6
2.3.1 Random Networks . . . . .	7
2.3.2 Small-World Effect . . . . .	7
2.3.3 Scale-free Networks . . . . .	9
2.3.4 Hierarchical Networks . . . . .	10
<b>3 Community Detection</b>	<b>11</b>
3.1 Introduction . . . . .	11
3.2 The Method . . . . .	12
3.3 Information Dynamics Algorithm . . . . .	14
3.4 Static Networks . . . . .	15
3.4.1 Zachary "Karate Club" . . . . .	17
3.4.2 Bottlenose dolphin Network . . . . .	19
3.4.3 Entropy of information . . . . .	20
3.4.4 Local and long range interaction . . . . .	21
3.5 Dynamic Networks . . . . .	25
3.5.1 Performance evaluation . . . . .	29
3.6 Evaluating cerebral cortex connectivity with local information algorithm . . . . .	33
3.6.1 Net Explorer . . . . .	34
3.6.2 Other algorithms . . . . .	35
3.6.3 Numerical Results . . . . .	36

3.6.4	Model's approximation of the clustering structure and of the adjacency matrix representing the empirical cortical map of the cat . . . . .	37
3.6.5	Qualitative evaluation of the algorithms' performance . . . . .	39
3.6.6	Empirical weighted cortical map approximation by the <i>NetExplorer</i> algorithm . . . . .	42
3.6.7	Discussion and conclusions . . . . .	43
3.7	A Cognitive-inspired Model for Self-organizing Networks . . . . .	44
3.7.1	Evaluation . . . . .	48
3.7.2	Final remarks . . . . .	49
3.8	Considerations on more complex heuristics . . . . .	50
3.8.1	IDA + LTE . . . . .	50
3.8.2	Double pruning . . . . .	57
3.9	Impact of local information in growing networks . . . . .	59
3.9.1	Related work . . . . .	61
3.9.2	The model . . . . .	62
3.9.3	Results . . . . .	66
3.9.4	Conclusions . . . . .	68
3.10	Final remarks . . . . .	68
<b>4</b>	<b>Epidemic Spreading</b>	<b>71</b>
4.1	Single-layer networks . . . . .	71
4.1.1	Introduction . . . . .	71
4.1.2	The networks model . . . . .	73
4.1.3	The risk perception model . . . . .	77
4.1.4	Results and Discussion . . . . .	79
4.2	Multiplex Networks . . . . .	83
4.2.1	Introduction . . . . .	83
4.2.2	The network model . . . . .	87
4.2.3	Multiplex network model . . . . .	89
4.2.4	The self-organized percolation method . . . . .	90
4.2.4.1	Risk perception . . . . .	92
4.2.4.2	Virtual risk perception . . . . .	94
4.2.5	Final remarks . . . . .	97
<b>5</b>	<b>Conclusions</b>	<b>99</b>
<b>A</b>	<b>Algorithms pseudocode</b>	<b>103</b>
	<b>Bibliography</b>	<b>105</b>



# List of Figures

2.1	Example of an undirected graph (network) and its mapping on the adjacency matrix. . . . .	3
2.2	Regular lattice, random network and connectivity distribution of a random network. . . . .	7
2.3	Small-world network. . . . .	8
3.1	An example of an adjacency matrix. . . . .	13
3.2	<i>Simple artificial network</i> . . . . .	15
3.3	Detection of community in a simple generated network. . . . .	16
3.4	Results on the Zachary's karate club network 1. . . . .	17
3.5	Results on the Zachary's karate club network 2. . . . .	17
3.6	Results on the Zachary's karate club network 3. . . . .	18
3.7	Results on the Zachary's karate club network 4. . . . .	18
3.8	Communities detected by our algorithm. . . . .	18
3.9	Results on the Bottlenose dolphin network 1. . . . .	19
3.10	Results on the Bottlenose dolphin network 2. . . . .	19
3.11	Results on the Bottlenose dolphin network 3. . . . .	20
3.12	Temporal evolution of the Shannon Entropy of Information. . . . .	21
3.13	The asymptotic configuration of the entropy as function of the parameters $m$ and $\alpha$ . . . . .	22
3.14	Hierarchical community detection. . . . .	22
3.15	Undirected network. . . . .	23
3.16	Undirected network 2. . . . .	23
3.17	Different results on the Zachary's karate club and Bottlenose dolphin networks. . . . .	24
3.18	Normalized mutual information on the GN benchmarks. . . . .	25
3.19	Dynamic networks. . . . .	30
3.20	Community detection in mobility networks. . . . .	31
3.21	Shannon Entropy of Information during time for different scenarios. . . . .	32
3.22	Comparison between the local entropy of a traveller (blue line) and a normal agent (black line). . . . .	32
3.23	Comparison of different community detection algorithms over the cortical map of the cat. . . . .	38
3.24	Normalized error of prediction for different algorithms. . . . .	39
3.25	Reconstruction of the cortical map of the cat. . . . .	42
3.26	Fitness function and mean energy. . . . .	47
3.27	Energy minimization approach. . . . .	48
3.28	Schematic representation of IDA + LTE algorithm. . . . .	52

---

3.29	IDA+LTE algorithm on the Zachary's Karate Club network. . . . .	53
3.30	Trigger in the <i>Lusseau's network of bottlenose dolphins</i> . . . . .	55
3.31	Myriel in the novel <i>Les Miserables</i> by Virctor Hugo. . . . .	56
3.32	(a) Washington in NCAA footbal college network. . . . .	56
3.33	Adjacency matrix of a network composed by 200 nodes and 3 levels where $p_1 = 0.9$ , $p_2 = 0.2$ and $p_3 = 0$ . . . . .	57
3.34	Evolution of the double pruning heuristic on the Zachary Karate Club network [1]. . . . .	58
3.35	Normalized Mutual Information on LFR benchmarks. . . . .	59
3.36	Schematic description of the method. . . . .	63
3.37	Politcal Blogs and hamsterster.com networks. . . . .	66
3.38	Simulation of the temporal density evolution of the <i>Political blogs network</i> . . . . .	67
4.1	Schematic representation of our network-generation model. . . . .	74
4.2	Connectivity degree distribution of generated networks. . . . .	75
4.3	Modularity values of generated networks. . . . .	76
4.4	Effect of the precaution parameters in different networks. . . . .	80
4.5	Infected individuals and the effect of precaution parameter in community- structured networks. . . . .	81
4.6	Percentage of infected agents for a random network. . . . .	82
4.7	Distribution of connectivity degree of scale-free and random networks. . . . .	86
4.8	Clustering coefficient (y-axis) as function of <i>transitivity parameter</i> $p_t$ . . . . .	88
4.9	Example of multiplexes generated with our method. . . . .	89
4.10	Cumulative frequency distribution of the infected individuals for the SIS dynamics achieved with our self-organized percolation method for differ- ent networks. . . . .	91
4.11	Evolutionary process of the direct percolation. . . . .	92
4.12	Schematic phase diagram and results for the risk perception. . . . .	93
4.13	Schematic phase diagram for the virtual risk perception. . . . .	95
4.14	Risk Perception in multiplex networks. . . . .	96

# List of Tables

2.1	Property of several real networks. Each network has the number of nodes $N$ , the clustering coefficient $C$ , the average path length $L$ and the degree exponent $\gamma$ of the power-law degree distribution. . . . .	6
3.1	Detailed scenario configuration . . . . .	29
3.2	Results from the information dynamics algorithm for the different levels in the network shown in Figure 3.36 (a). $L$ is the number of the level, $n$ is the id of the node and $p$ is the probability to join with the most connected node (bold node). . . . .	65
3.3	Statistics of the social networks. $C$ (mean clustering coefficient), $l$ (average path length) and $d$ (diameter of the network). . . . .	67
4.1	Critical values for the extinction of the epidemic on case of scale-free networks of 500 nodes and 5 communities considering a maximum threshold time $T_{max} = 1000$ , necessary for the extinction of the epidemics. . . . .	82



*For myself. . .*



# Chapter 1

## Introduction

We live in a world of networks. The networks around us and we ourselves, as people, are part of the network of social relations between individuals. Examples of networks in the world are the Internet, the rail network, the subway, neural networks, the telephone network, or less concrete entities, such as the relations of knowledge and collaboration between people. The study of networks is therefore very important, given the wide variety of structures and systems of the real world that can be incorporated into the category of “complex networks”. In general, a graph or network is a very general approximation of a system constituted by many entities, called nodes (which may represent, in various cases, persons, computers, proteins, chemicals, etc.) linked to each other and interacting through connections (which may be, therefore, a cable between computers, hyperlinks between web pages, a collaboration between people, a reaction between chemical substances, etc.). A complex network is a network with non trivial topological features that would not be recognizable as a simple network. The majority of existing networks, being them social, biological or technological (as well as some of the phenomena that depend on or derived from the networks) may be considered complex: this depends on some features such as, e.g., the degree distribution of links, the high coefficient of grouping, the assortativity between the vertices, and often evidence of a hierarchical structure. Another important observed feature of complex networks is the presence of a community structure. Many real networks are not homogeneous and do not consist of a single block of indistinct nodes, rather they exhibit some community structures, i.e. group of vertexes more connected where there are high densities of links within each group, and a relatively low number of connections between the various groups. In recent years was explored in detail the possibility of automatically identify communities within networks giving rise to a new field of research called “community detection”: in this thesis we summarize the previous works and we present a new cognitive-inspired algorithm for detecting communities in many networks from synthetic to real and dynamic networks.

---

Given the impact of community detection in many areas, such as psychology and social sciences, we have addressed the issue of modifying existing well performing algorithms by incorporating elements of the domain application fields, i.e. domain-inspired. We have focused on an approach inspired by psychology and social science which may be useful for further strengthening the link between social network studies and mathematics of community detection. The rest of this thesis is organized as follows: we start by describing some basic features of complex networks in Section 2. In Section 3 we show the applications of our method to different scenarios ranging from static to dynamic networks and from synthetic to real-world networks. Considering psychological notions as mentioned above, we adopted a local algorithm where an individual is simply modeled as a memory and a set of connections to other individuals. In our approach, the “learning” (nonlinear) phase is modeled after competition in chemical/ecological world. Here we want to emphasize not only the good efficiency of the algorithm in detecting community but also its capability to discover overlapping nodes and to reveal the subjective vision of hierarchical levels of the network. Moreover we show that it is very difficult to obtain an objective definition of community and that our method is very efficient for discovering different subjective hierarchical levels using more complicated heuristics as reported in Section 3.8: here we show also that our method is comparable in terms of efficiency with other well-known community detection algorithms. In the last part of this thesis we analyse a process widely studied in complex networks: the epidemic spreading. In particular we have a look to the risk perception in modelling the spread of epidemics in clustered (Section 4.1) and multiplex networks (Section 4.2). This section has a dual purpose: on the one hand we want to analyse networks in order to reproduce real scenarios; on the other side we want to take into account the precautions that people take when they become aware of a certain disease. Moreover we present a self-organized percolation method for both single layer and multiplex networks: this method allow us to detect the critical values for avoiding the epidemic spreading, simulated through the so-called Susceptible-Infected-Susceptible (SIS) dynamics, in a very general way. Finally we discuss the main results achieved during my PhD in the Conclusions.



## Chapter 2

# Complex Networks

### 2.1 Introduction

Many different systems can be effectively described in terms of graphs or networks. The best-known representation of a network is that made by drawing points (nodes) and lines connecting them (arcs). This is also the form of a family of objects known in mathematics as graphs. Formally, a graph  $G = (V, E)$  is defined as a pair of sets: the set  $V$  of nodes, or vertices, and the set  $E$  of the connections between them (also known as edges, links), and each arc is drawn as a line connects two nodes. Graph theory was officially founded in the late eighteenth century with the solution of the famous problem of the bridges of Königsberg by the Swiss mathematician Leonhard Euler [2]. The technique used by him proved to be of much greater utility than simple puzzle solving. The German physicist Gustav Kirchoff analyzed electric circuits in terms of graphs and chemists found a natural correspondence between graphs and structures of atoms and molecules. A graph can describe a transportation network, a neural network of the brain

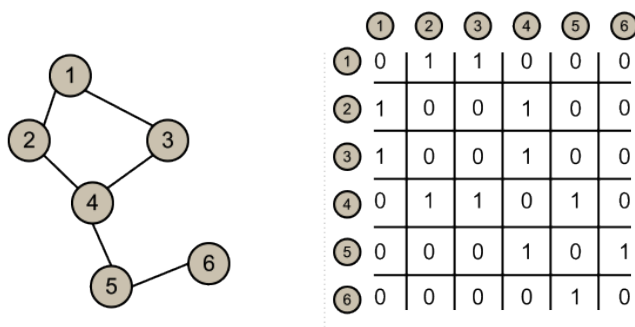


FIGURE 2.1: Example of an undirected graph (network) and its mapping on the adjacency matrix.

or an economy in which companies are the nodes and the transactions between them are edges. In the twentieth century the theory has assumed a predominantly statistical and algorithmic predisposition [3, 4]. Given the representation of the network as a graph, all deductions, measures and indicators of the theory can be applied to this case. The terms graph and network are synonymous in all respects. From the mathematical point of view a graph  $G = (V, E)$  can be represented by a matrix (called adjacency matrix  $A$  Figure 2.1) whose elements different from zero represent a link between two nodes (if different from 1 represents a "weight" of the connection: cost, speed, energy, etc.). intersection if the two nodes are connected. The complex networks observed in real case often share characteristics that include global statistics as short "distance" between nodes or the presence of community structure. In many of these we observe that the relative commonness of vertices with a certain degree greatly exceeds the average. The understanding of the elements and relationships of each class of networks is essential for proper identification and description, and to try to discover the principles of structural construction and groped to predict the behaviour over time (evolution).

## 2.2 Some basic concepts

Although many features and measures of complex networks have been proposed and investigated in the last decades, three spectacular concepts as the average path length, the clustering coefficient and the degree distribution of links play a key role in the recent study and development of complex networks theory. In fact, the original attempt of Watts and Strogatz in their work on small-world networks [5] was to construct a network model with small average path length as a random graph and relatively large clustering coefficient as a regular lattice, which evolved to become a new network model as it stands today. On the other hand, the discovery of scale-free networks was based on the observation that the degree distributions of many real networks have a power-law form, albeit power-law distributions have been investigated for a long time in physics for many other systems and processes [6].

### 2.2.1 Average path length

The shortest paths play an important role for the spreading of information within a network. Suppose that a user needs to send a packet from one computer to another in the Internet network: if your package arrives at its destination by the shortest path between two nodes, ensures a faster transfer and minimizes the use of network resources. The average path length is defined to represent the length of all the shortest paths of a graph  $G$  in a single matrix  $M$ , in which each element  $j$  is the length of the shortest

path between the node  $i$  and  $j$ . An important property of the network is the average length of the shortest paths, known as *average path length* defined as the average of the shortest paths, evaluated on all pairs of nodes in the network:

$$L = \frac{1}{N(N-1)} \sum_{i=1}^N \sum_{j=1, i \neq j}^N m_{ij}. \quad (2.1)$$

The diameter  $D$  of a network is defined as the maximal distance among all the shortest path between any pair of nodes in the network. It was an interesting discovery that the average path length of most real complex networks is relatively small, even in those cases where these kinds of networks have many fewer edges than a typical globally coupled network with a equal number of nodes. This smallness inferred the small-world effect, hence the name of small-world networks.

### 2.2.2 Clustering coefficient

The clustering, or transitivity, of a graph measures the tendency of two nodes that are adjacent to a common node to be connected with one another. We remember that the clustering coefficient is a function of the number of local triples. Following the definition of Barrat et al. [7], the clustering coefficient  $C$  is defined as the average of the local clustering coefficients  $c_i$ :

$$C = \frac{1}{N} \sum_i c_i = \frac{1}{N} \sum_i \left( \frac{1}{k_i(k_i-1)} \sum_{j,h} a_{ij} a_{ih} a_{jh} \right), \quad (2.2)$$

where  $k_i$  is the connectivity degree of node  $i$ . Then  $0 < c_i < 1$  and  $0 < C < 1$ , and therefore in a fully connected network  $C = 1$ . A high value of the coefficient  $C$  indicates that there are many connections between neighbouring nodes. In a completely random network consisting of  $N$  nodes,  $C \sim 1/N$ , which is very small as compared to most real networks. It has been found that most large-scale real networks have a tendency toward clustering, in the sense that their clustering coefficients are much greater than  $O(1/N)$ , although they are still significantly less than one (namely, far away from being globally connected).

### 2.2.3 Degree distribution

The simplest and perhaps the most important feature of a single node is its degree. The degree  $k_i$  of a node  $i$  is the number of links that end on the node, and is defined

considering the adjacency matrix  $A$  as:

$$k_i = \sum_{j=1}^N a_{ij}, \quad (2.3)$$

where  $N$  is the number of nodes in the network. If a network is directed, nodes have two different degrees, the *in-degree*, which is the number of incoming edges, and the *out-degree*, which is the number of outgoing edges. The average of  $k_i$  over all  $i$  is called the average degree of the network, and is denoted by  $\langle k \rangle$ . The degree distribution  $P(k)$  of a network is then defined as the fraction of nodes in the network with degree  $k$ . Thus if there are  $N$  nodes in total in a network and  $n_k$  of them have degree  $k$ , we have  $P(k) = n_k/N$ . As we will see in the next section most of the observed networks in real world seem to follow a power-law degree distribution where  $P(k) \sim k^{-\gamma}$ .

TABLE 2.1: Property of several real networks. Each network has the number of nodes  $N$ , the clustering coefficient  $C$ , the average path length  $L$  and the degree exponent  $\gamma$  of the power-law degree distribution.

Network	N	C	L	$\gamma$
WWW [8]	153127	0.11	3.1	$\gamma_{in}=2.1$ $\gamma_{out}=2.45$
E-mail [9]	56969	0.03	4.95	1.81
Software [10]	1376	0.06	6.39	2.5
Electronic circuits [11]	329	0.34	3.17	2.5
Language [12]	460902	0.437	2.67	2.7
Movie actors [5, 6]	225226	0.79	3.65	2.3
Math. co-authorship [13]	70975	0.59	9.50	2.5
Food web [14]	154	0.15	3.40	1.13
Metabolic system [15]	778	—	3.2	$\gamma_{in} = 2.2$ $\gamma_{out}=2.2$

## 2.3 Complex Network models

Measuring some basic properties of a complex network, such as the average path length  $L$ , the clustering coefficient  $C$ , and the degree distribution  $P(k)$ , is the first step toward understanding its structure. The next step is to develop a mathematical model to reproduce the main features found in real-world networks, thereby obtaining a platform for which mathematical analysis is possible.

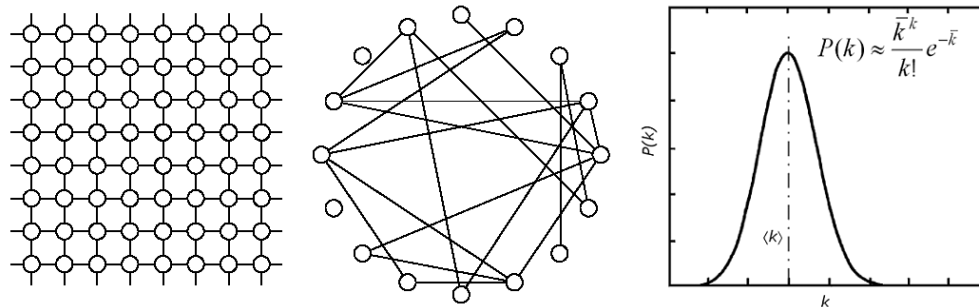


FIGURE 2.2: Regular lattice, random network and connectivity distribution of a random network.

### 2.3.1 Random Networks

The classical theory of networks, first introduced by the studies of Euler (1736), was defined by hungarian mathematicians P. Erdős and A. Rényi [16]. The basic assumption is that each pair of nodes of the network is connected randomly with a probability  $p$ . This leads to a statistically homogeneous network in which, despite the randomness of the model, the majority of nodes has the same number of connections  $\bar{k}$ . In particular, the connectivity follows a Poisson distribution with a peak at the value  $\bar{k}$ , then the probability of finding a strongly connected node decays exponentially, so  $P(k) \sim e^{-k}$  for  $k \gg \bar{k}$  as shown in the Figure 2.2. In this scenario the average path between two nodes is relatively small as well as the frequency of nodes that tend to cluster together (clustering coefficient).

### 2.3.2 Small-World Effect

On the other hand, empirical studies carried out on biological, social or technological networks shown that there are many deviations from this random structure [17]. One of the most important kind of networks is the social network of relations among people where the nodes are individual and there is a link between them if they know each other. A network of this kind is hardly representable as a regular lattice, and, unlike a totally random network, it is inhomogeneous. There are different grouping, i.e. there is a high probability that two nodes connected to a other are connected to each other: friends of my friends are friends (with some probability). The “small world” effect is fairly well known: two friends find they have mutual acquaintances. The idea that the chain of acquaintances between any two people to world was limited had been formulated by the hungarian novelist Frigyes Karinthy in 1929 [18] and scientifically enunciated by the sociologist Stanley Milgram in 1967 [19]. The Milgram experiment was to study the route followed by letters sent from Nebraska to direct acquaintance with a final destination

address of Pittsburgh. The average number of steps turned out to be five to six people involved. The experiment was not perfectly rigorous, but since then the term "six degrees of separation" has become the term to identify the "small world" phenomenon. After studying the networks of partnerships such as that of actors in films or the electricity distribution network or the neural connections of the small worm *C. elegans*, Steven Strogatz and Duncan Watts in 1998 discovered an interesting phenomenon [5]. Almost halfway between a regular lattice and a totally random network, these real networks are characterized by small average distance  $L$  between the nodes and high clustering coefficient  $C$  (measure the local density of connections in groupings). The proposed model describes the construction of a small-world network starting from a regular lattice and replacing the links with other present between different nodes with probability  $p$ . These shortcuts are responsible for the "small world" effect (Figure 2.3).

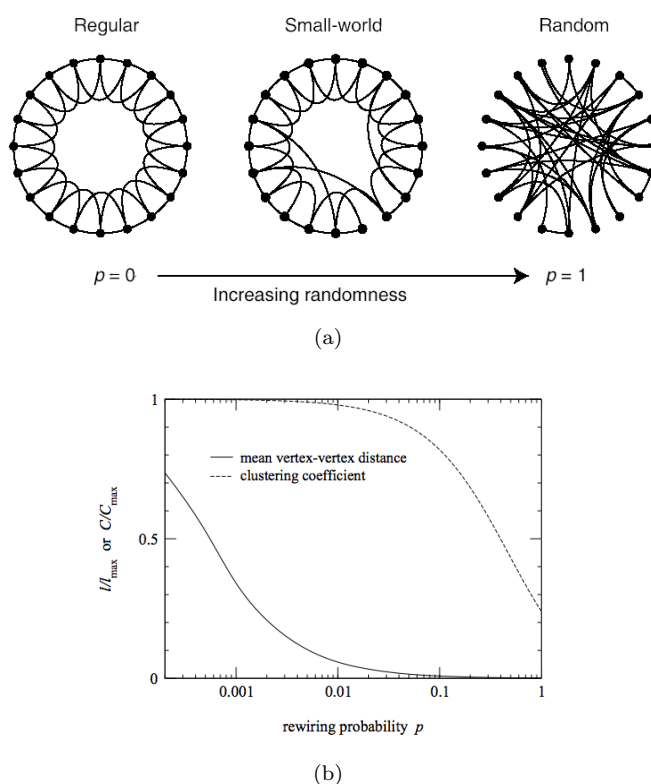


FIGURE 2.3: Small-world network.

(a) The small network includes a handful of extra links thrown in between pair of nodes chosen entirely at random. (b) The clustering coefficient  $C$  and mean vertex-vertex distance  $L$  in the small-world model of Watts and Strogatz [5] as a function of the rewiring probability  $p$ . For convenience, both  $C$  and  $L$  are divided by their maximum values, which they assume when  $p = 0$ . Between the extremes  $p = 0$  and  $p = 1$ , there is a region in which clustering is high and mean vertex-vertex distance is simultaneously low [17].

### 2.3.3 Scale-free Networks

In a small-world network the degree distribution of nodes follow a Poisson law, but many real networks show an entirely different distribution of degrees. This is the case of networks such as the WWW, film actors, metabolic networks or electric ones. Scale-free networks are characterized by a power-law degree distribution:

$$P(k) \sim k^{-\gamma}, \quad (2.4)$$

where  $2 < \gamma < 3$  (observed). This observation indicates that large networks self-organize into a state without scale (hence the term scale-free), a feature not provided by existing models of random networks. A common feature of all the random models (including the ER) is that they provide that the degree distribution ( $P(k)$ ) has an exponential trend with a peak around the mean value  $\hat{k}$ . In order to explain this mismatch, Barabási and Albert have found that there are two aspects of real networks that are not incorporated into these models [6]. First, they assume that previous models start with a fixed number of vertices ( $N$ ) that are connected or rewired randomly, without changing the size. On the contrary, many real world networks are open and are formed by the continuous addition of new nodes to the system, i.e., the number of nodes  $N$  grows during the period of life of the network. For example, the network of actors grows with the addition of new players to the system, the WWW grows exponentially over time with the addition of new web pages, and the scientific literature increases with the publication of new papers. Consequently, a common feature to these systems is that the network continuously expands with the addition of new nodes, which are connected to those already present in the system.

On the other hand random networks models assume that the probability that two nodes are connected to each other is random and uniform. However, many real networks show preferential connectivity: for example, it is more likely that a new player is assigned a support role to the best-known and established actors. The probability that a new actor recites with a more established actor is much greater than that of acting together with other actors less known. In a similar way, it is more likely that a newly created web page contains links to documents already well known and with high connectivity, and more likely that a new manuscript is mentioned in a very well known paper and already frequently mentioned, with respect to a less known article and, therefore, less cited. These examples indicate that the probability with which a new node connects to the existing ones is not uniform; there is a higher probability that it is connected to a node that has already a substantial amount of connections.

In light of this consideration Barabási and Albert they have developed a model for the generation of scale-free networks [20]. The “evolutionary” feature of real self-organized networks is explained through two mechanisms by the so-called Barabási-Albert (BA) model [6]:

1. **Growth:** we start by considering a network of few nodes  $m_0$  and at each time step a new node with connectivity  $m < m_0$  is added to the network.
2. **Preferential attachment:** the main hypothesis of the model is that the more connected a node is, the more likely it is to receive new links. The new node links to an existing node in the network  $i$  with a probability  $\Pi$  which is function of the connectivity degree  $k$  of node  $i$ , then  $\Pi(k_i) = \frac{k_i}{\sum_j k_j}$ , where  $j$  represents all the nodes in the network.

### 2.3.4 Hierarchical Networks

The main problem of the BA model is that the clustering coefficient is vanishing when the size of the network increases, which does not correspond to what is observed in real networks. The idea of the hierarchical network, introduced by Ravasz and Barabási try to eliminate this problem [21]. In fact many real networks are expected to be fundamentally modular, meaning that the network can be seamlessly partitioned into a collection of modules with an underlying hierarchical structure and an high degree of clustering. This fact is also very important in the context of social relationships where each people, following the studies of Dunbar [22], have a sort of hierarchical structure of relationship from family to friend until work colleagues reaching the knowledge of around 150 people. This evolutionary model allow to generate scale-free networks with community and hierarchical structure. A model for generating hierarchical networks will be shown in Section 3 while a model for scale-free networks with community structure will be shown in Section 4.



## Chapter 3

# Community Detection

### 3.1 Introduction

Detecting communities is a task of great importance in many disciplines, namely sociology, biology and computer science [23–27], where systems are often represented as graphs. Community detection is also linked to clustering of data: many clustering methods establish links among representative points that are nearer than a given threshold, and then proceed in identifying communities on the resulting graphs [28, 29]. Given a graph, in a broad sense, a community is a group of vertices “more linked” than between the group and the rest of the graph. This is clearly a poor definition, and indeed, on a connected graph, there is not a clear distinction between a community and the rest of the graph. In general, there is a continuum of nested communities whose boundaries are somewhat arbitrary: the structure of communities can be seen as a hierarchical dendrogram [30].

In general, community detection algorithms rely on global quantities like betweenness, centrality, etc. [30, 31] and most algorithms require the graph to be completely known. This constraint is problematic for networks like the World Wide Web, which for all practical purposes is too large and too dynamic to ever be fully known.

Moreover in complex networks, and in particular in social networks, it is very difficult to give a clear definition of community: it is caused by the fact that nodes often results in overlapping communities because they belong to more than one cluster or module or community. The problem of overlapping communities was discussed in [32] and recently a solution to it were presented in [33]. For instance people usually belong to different communities at the same time, depending on their families, friends, colleagues, etc.: so each people, making a *subjective* community detection algorithm, has its own vision of communities in his social environment.

In social networks, the definition of a community could be linked to the human capability of information processing, particularly the poor evaluation of probabilities. When faced with insufficient data or insufficient time for a rational processing, we humans have developed algorithms, denoted heuristics, that allows to take decisions in these situations. The modern approach to the study of cognitive heuristics defines them as those *strategies that prevent one from finding out or discovering correct answers to problems that are assumed to be in the domain of probability theory*. The ratio of a cognitive algorithm for community detection is based on the fact that humans' networks are the results of the individual strategies of single subjects; on the other hand they are presumably shaped and evolved by the social structures in which they live [34, 35].

The thesis is organized as follows: we start by describing a new algorithm for detecting communities in complex networks in Section 3.2. Considering psychological notions as mentioned above, we adopted local algorithm where an individual is simply modeled as a memory and a set of connections to other individuals. The “learning” (nonlinear) phase is modeled after competition in chemical/ecological world, where resources fighting each other in order not to fall into oblivion. In Section 3.3 we describe the first algorithm in which information about neighbouring nodes is propagated and elaborated locally, but the connections do not change. Here we want to emphasize not only the good efficiency of the algorithm in detecting community but also its capability to discover overlapping nodes and a sort of subjective vision of hierarchical levels of the network. Next, in Section 3.4.3 we give an interpretation of Shannon entropy of information as quality function for estimating models parameters. Finally we discuss our results and we propose future steps in Conclusions.

## 3.2 The Method

We consider  $N$  individuals, labeled from 1 to  $N$ . Let us denote by  $A$  the adjacency matrix,  $A_{ij} = 1$  (0) indicates the presence (absence) of a link from site  $j$  to site  $i$ ; all links have the same weight (Figure 3.1). Each individual  $i$  is characterized by a state vector  $S_i$ , representing his knowledge of the outer world. We interpret  $S$  as a probability distribution, assuming that  $S_i^{(k)}$  is the probability that individual  $i$  belongs to the community  $k$ . Thus,  $S_i^{(k)}$  is normalized on the index  $k$ . We shall denote by  $S = S(t)$  the state of the all network at time  $t$ , with  $S_{ik} = S_i^{(k)}$ . We shall initialize the system by setting  $S_{ij}(0) = \delta_{ij}$ , where  $\delta$  is the Kroneker delta,  $\delta_{ij} = 1$  if  $i = j$  and zero otherwise. In other words, at time 0 each node only knows about itself.

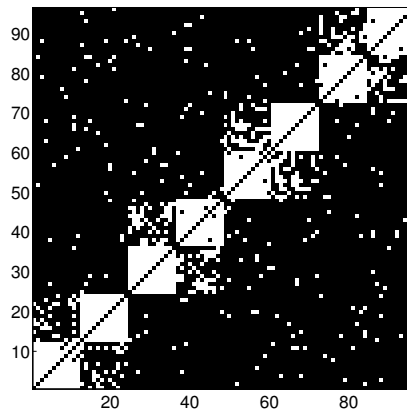


FIGURE 3.1: An example of an adjacency matrix.

It is a three-level matrix composed by 4 blocks of 2 sub-communities of 8 nodes each. The link probability inside a sub-community is 0.98, in first level blocks is 0.3 and among blocks is 0.03. White points indicate the presence of a link between the node  $i$  and the node  $j$ ,  $A_{ij} = 1$ .

As mentioned the competition phase is modeled thinking to a chemical/ecological analogy. Our algorithms are based on the concept of *diffusion and competitive interaction* in network structure introduced by Nicosia et al. [36].

If two populations  $x$  and  $y$  are in competition for a given resource, their total abundance is limited [37]. After normalization, we can assume  $x + y = 1$ , i.e.,  $x$  and  $y$  are the frequency of the two species, and  $y = 1 - x$ . The reproductive step is given by  $x' = f(x)$ , which we assume to be represented by a power  $x' = x^\alpha$ . For instance,  $\alpha = 2$  models the birth of individuals of a new generation after binary encounters of individuals belonging to the old generation, with noneverlapping generations (eggs laying) [38].

After normalization we obtain:

$$x' = \frac{x^\alpha}{x^\alpha + y^\alpha} = \frac{x^\alpha}{x^\alpha + (1 - x)^\alpha}. \quad (3.1)$$

Introducing  $z = (1/x) - 1$  ( $0 \leq z < \infty$ ), we get the map

$$z(t + 1) = z^\alpha(t), \quad (3.2)$$

whose fixed points (for  $\alpha > 1$ ) are 0 and  $\infty$  (stable attractors) and 1 (unstable), which separates the basins of the two attractors. Thus, the initial value of  $x$ ,  $x_0$ , determines the asymptotic value, for  $0 \leq x < 1/2$ ,  $x(t \rightarrow \infty) = 0$ , and for  $1/2 < x < 1$ ,  $x(t \rightarrow \infty) = 1$ .

By extending to a larger number of components for a probability distribution  $P_i$ , the competition dynamics becomes

$$P'_i = \frac{P_i^\alpha}{\sum_j P_j^\alpha}, \quad (3.3)$$

and the iteration of this mapping, for  $\alpha > 1$ , leads to a Kroneker delta, corresponding to the larger component. However, the alternation between information and competition can generate interesting behaviors.

### 3.3 Information Dynamics Algorithm

The dynamics of the network is given by an alternation of communication and elaboration phases. Communication is implemented as a simple diffusion process, with memory  $m$ . The memory parameter  $m$  allows us to introduce some limitations in human cognition such as the mechanism of oblivion and the timing effects: the most recent information has more relevance than the previous one [39, 40].

We assume that each individual spends the same amount of time in communication, so that people with more connections dedicate less time to each of them. Since the amount of available time is limited, we normalize the adjacency matrix on the columns (i.e., we assign at each link the inverse of the output degree of the incoming node), forming a Markov matrix  $M$

$$M_{ij} = \frac{A_{ij}}{\sum_k A_{kj}}. \quad (3.4)$$

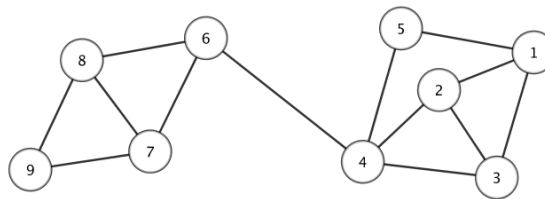
Note that in many mathematical texts the indices are inverted, so that the Markov matrices are normalized on the rows. We prefer the “physics” notation so that matrix multiplication with a probability distribution  $P$  takes the usual form  $P' = MP$ . Then in the communication phase, the state of the system evolves as

$$S(t + \frac{1}{2}) = mS(t) + (1 - m)MS(t). \quad (3.5)$$

As described in the Eq. 3.3, the competition phase is modeled thinking to a competitive interaction between the nodes in the network [36].

In this way the dynamic of the model is given by a sequence  $S(t) \rightarrow S(t + \frac{1}{2}) \rightarrow S(t + 1)$ :

$$\begin{aligned} S_{ik} \left( t + \frac{1}{2} \right) &= mS_{ik}(t) + (1 - m) \sum_j M_{ij} S_{jk}(t), \\ S_{ik}(t + 1) &= \frac{S_{ik}^\alpha(t + \frac{1}{2})}{\sum_j S_{ij}^\alpha(t + \frac{1}{2})}. \end{aligned} \quad (3.6)$$

FIGURE 3.2: *Simple artificial network*

The graph is composed by of 9 nodes and 13 links divided in 2 communities. It is possible to identify two different communities: the first one composed by nodes 1-2-3-4-5 and the second one by 6-7-8-9.

We assume that individuals' memory is large enough so that they can keep track of all information about all other individuals. In a real case, one should limit this memory and apply an input/output filtering. Individuals do not change their connectivity. For testing purposes we use three networks and analyzing and discussing our model peculiarities. The three case studies, of growing or different complexity, are: a simple artificial network used to show the typical output of our algorithm, the Zachary *karate club* network [1] and the bottlenose dolphins network [41].

### 3.4 Static Networks

In this first case study, as shown in Figure 3.2 represented by its adjacency matrix  $A$ :

$$A = \begin{pmatrix} 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 & 1 & 1 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 0 \end{pmatrix}$$

the algorithm face with a very simple task and converges to an optimal solution in few iterations and for a wide range of model's parameters  $m$  and  $\alpha$ . Analyzing state matrix  $S(t)$ , it is possible to identify two different communities marked by nodes 5 and 9.

$$S(T) = \begin{pmatrix} 0 & 0 & 0 & 0 & 0.9999 & 0 & 0 & 0 & 0.0001 \\ 0 & 0 & 0 & 0 & 0.9982 & 0 & 0 & 0 & 0.0018 \\ 0 & 0 & 0 & 0 & 0.9982 & 0 & 0 & 0 & 0.0018 \\ 0 & 0 & 0 & 0 & 0.9383 & 0 & 0 & 0 & 0.0617 \\ 0 & 0 & 0 & 0 & 0.9975 & 0 & 0 & 0 & 0.0025 \\ 0 & 0 & 0 & 0 & 0.1309 & 0 & 0 & 0 & 0.8691 \\ 0 & 0 & 0 & 0 & 0.0054 & 0 & 0 & 0 & 0.9946 \\ 0 & 0 & 0 & 0 & 0.0054 & 0 & 0 & 0 & 0.9946 \\ 0 & 0 & 0 & 0 & 0.0061 & 0 & 0 & 0 & 0.9939 \end{pmatrix}$$

In Figure 3.3(b) it is possible to identify two different communities highlighted by upper values in the graph. The first community is composed by node 1-2-3-4-5 and the second one by 6-7-8-9. Our algorithm is capable also to detect overlapping nodes (4 and 6) as "middle" values between blue lines. In this way each node knows exactly its role in the network.

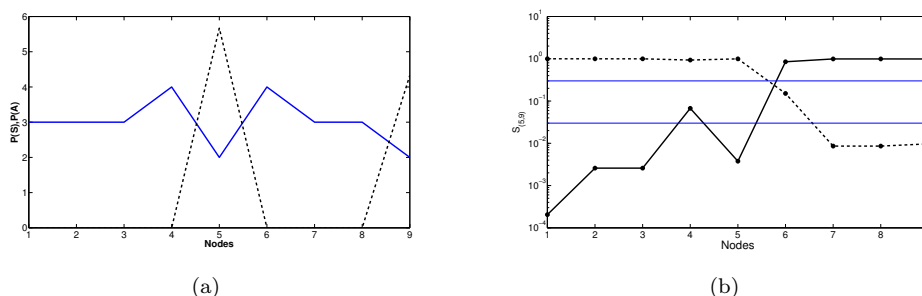


FIGURE 3.3: Detection of community in a simple generated network.

(a) On the x-axis of both figures there are is number of nodes. On the y-axis: the cumulative distribution  $P_j^{(S)}$  (dashed black line,  $P_j^{(S)} = \sum_i S_{ij}$ , multiplied by five) and  $P_j^{(A)}$  (blue line,  $P_j^{(A)} = \sum_i A_{ij}$ , connectivity). The information propagation algorithm identifies communities by leaves (nodes 5 and 9 with lower connectivity) with  $m = 0.3$  and  $\alpha = 1.4$ . (b) The value of state vectors, at the final asymptotic time, of node 5 (dashed black line) and node 9 (black line). We can observe upper values indentifying communities: the first one composed by nodes 1-2-3-4-5 and the second one by nodes 6-7-8-9. The algorithm is capable also to detect the *communication nodes* 4 and 6 between the blue lines. In this way we can indentify the overlap between the communities and also define a sort of *objective vision* of nodes. It is clear that the upper nodes know very well which is their community as well as nodes 4 and 6 that know that they are in a middle state between two communities.

### 3.4.1 Zachary "Karate Club"

The second test case is a typical network literature example: the network proposed by Zachary in the 1977 , and known as "karate club" [1]. Although this network (Figure 3.4(a)) is rather small, our algorithm shows interesting results. With  $m = 0.25$  and  $\alpha = 1.4$  the algorithm has detected three communities in few steps as described in the Figure 3.4(b). In figures Figure 3.5(a-b), Figure 3.6(a-b) and Figure 3.7(a-b) we show that each community has the own subjective vision of network structures. Analyzing these results it is possible it is possible to detect three communities in the network and two overlapping nodes (node 1 and 3) as shown in Figure 3.8.

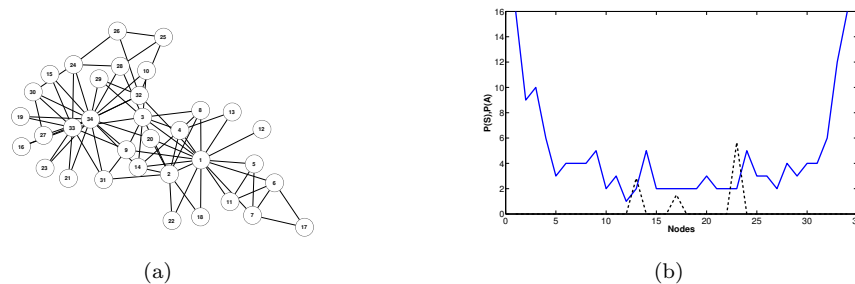


FIGURE 3.4: Results on the Zachary's karate club network 1.

(a) Zachary's karate club network. (b) On the x-axis of the figure there is the number of nodes. On the y-axis: the connectivity (blu line)  $P(A)$  and the cumulative distribution (dashed black line)  $P(S)$  are reported at final asymptotic time with  $m = 0.2$  e  $\alpha = 1.4$ . The  $P(S)$  reveals three underlying substructures labelled by nodes 13,17 and 23.

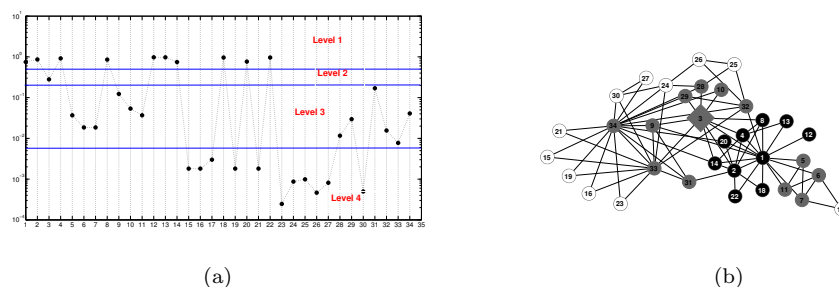


FIGURE 3.5: Results on the Zachary's karate club network 2.

(a) Values of state vector 13: the community is defined by black nodes in Figure 3.5(b) corresponding to upper values in the Level 1. As well as in the simple artificial network network the algorithm has detected not only the righth community but also the overlapping node corresponding to big gray diamond in Figure 3.5(b), in the Level 2. The Level 3 and the Level 4 (respectively gray and white nodes in Figure 3.5(b)) correspond to the different level of knowledge the others. In the Level 3 it is possible to find "direct" contacts while in the fourth level it is possible to detect "friends of my friends".

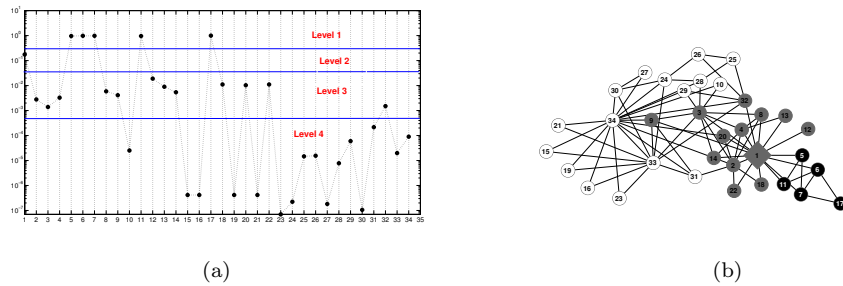


FIGURE 3.6: Results on the Zachary's karate club network 3.

(a) Values of state vector 17: the community is defined by black nodes in Figure 3.6(b) corresponding to upper values in the Level 1. For the description see the caption of the Figure 3.5.

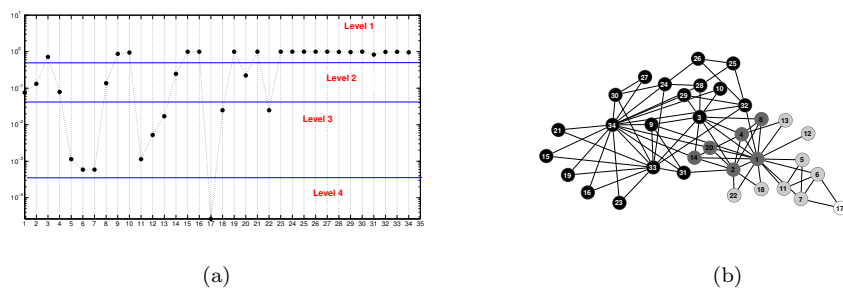


FIGURE 3.7: Results on the Zachary's karate club network 4.

(Values of state vector 17: the community is defined by black nodes in Figure 3.7(b) corresponding to upper values in the Level 1. In this case we haven't found overlapping but is also possible to detect a sort of scale of friendship inside the network labeled by the 4 levels in Figure 3.7(a) and by black-to-white color scale in Figure 3.7(b) .

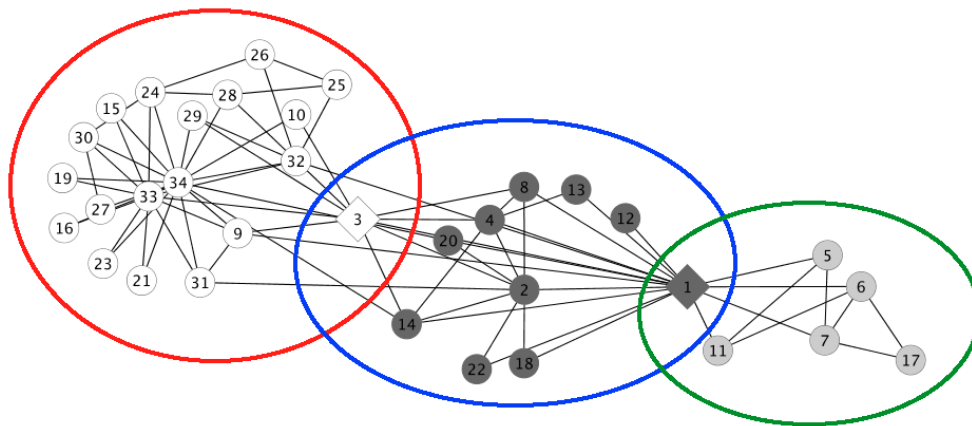


FIGURE 3.8: Communities detected by our algorithm.

Communities are marked by dark gray (shown in Figure 3.5), light gray (shown in Figure 3.6) and white nodes (shown in Figure 3.7). On the other hand three circles represent the overlap between communities because of the role of node 3 and 1 as explained in Figure 3.5 and Figure 3.6.



### 3.4.2 Bottleneck dolphin Network

The third case study concerns a community network of dolphins. The network we study was constructed from observations of a community of 62 bottlenose dolphins over a period of seven years from 1994 to 2001 [41] as shown in Figure 3.9(a). With  $m = 0.5$  and  $\alpha = 1.03$  our algorithm detects two communities (Figure 3.9(b)). Results are reported in Figure 3.10(a-b) and Figure 3.11(a-b) where we show the capability of our algorithm to detect the right communities but also 7 overlapping nodes between them.

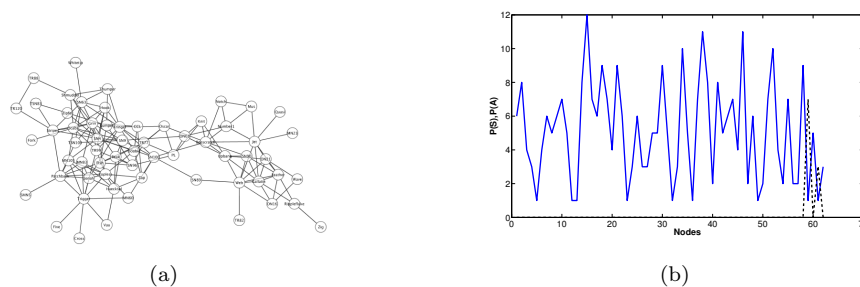


FIGURE 3.9: Results on the Bottleneck dolphin network 1.

(a) Bottleneck dolphin network. This network has a size of 62 nodes and from direct observation it is known that it has two communities. (b) On the x-axis of the figure there is the number of nodes. On the y-axis: the connectivity (blue line)  $P(A)$  and the cumulative distribution (dashed black line)  $P(S)$  are reported at final asymptotic time with  $m = 0.5$  and  $\alpha = 1.03$ . The  $P(S)$  reveals two underlying substructures labeled by nodes 59 and 61

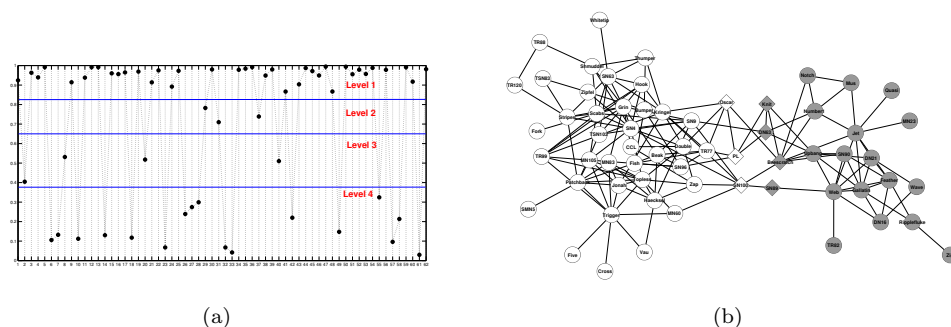


FIGURE 3.10: Results on the Bottleneck dolphin network 2.

(a) Values of state vector number 59: the four levels indicate respectively nodes of community (white nodes), overlapping nodes (white diamond nodes), closer nodes of the other community (gray diamonds) and finally nodes of the other community (gray nodes). (b) Communities detected by our algorithm.

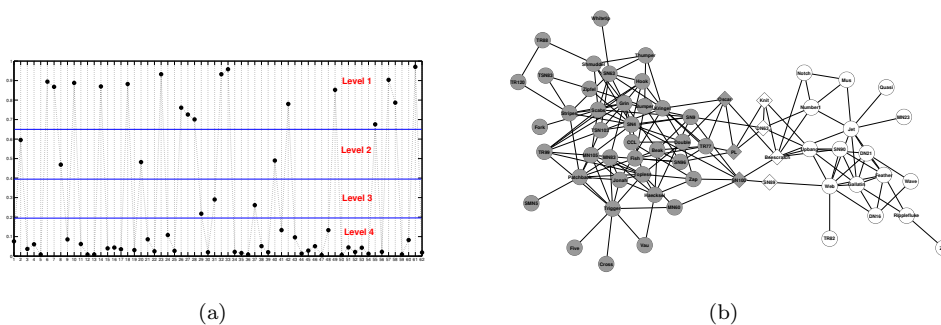


FIGURE 3.11: Results on the Bottlenose dolphin network 3.

(a) Values of state vector number 61: the four levels indicate respectively nodes of community (white nodes), overlapping nodes (white diamond nodes), closer nodes of the other community (gray diamonds) and finally nodes of the other community (gray nodes). (b) Communities detected by our algorithm.

### 3.4.3 Entropy of information

In order to present the temporal results in a compact way, we computed the entropy  $E$  of a configuration  $S$ , using the cumulative distribution over the non-normalized index,

$$\begin{aligned}
 P_i^{(S)} &= \sum_j S_{ij}, \\
 E^{(S)} &= - \sum_i P_i^{(S)} \log(P_i^{(S)}) - \sum_i (1 - P_i^{(S)}) \log(1 - P_i^{(S)}).
 \end{aligned} \tag{3.7}$$

The entropy  $E$  is maximal for the flat distribution, when each node knows only itself, and minimal (zero) what all the network has only one label (or has become just one star for the rewiring algorithm). If the population is evenly distributed among  $n$  clusters, the entropy is  $E = \log(n)$ . Let us to study the artificial complex network illustrated in Figure 3.1. This network is composed by three levels with different probability to have a link in a region.

As we observed our algorithm is able to observe all levels of a hierarchical network. In Figure 3.12(a) it is possible to identify the final level of the artificial network. The value of Entropy  $E(t)$  can help us understand the structure of the network at *priori*. In fact, different levels of a hierarchical structure are identified by the plateau as we can observe in Figure 3.12(b). This result is emphasized by the entropy's first derivative where we can observe three distinct peaks (Figure 3.12(c)). The final monocluster, using the adjacency matrix  $A$  in the Eq. 3.6, is identified by the major hub in the network (Figure 3.12(d)). Role of parameters  $m$  and  $\alpha$  has shown in Figure 3.13 where it is possible to observe the value of Entropy  $E$  at the state of convergence for each combination of parameters. Then if we change the parameters we can discover different

communities structures: in Figure 3.14 we show the final asymptotic state with  $m = 0.7$  and  $\alpha = 1.4$  on the same network used before.

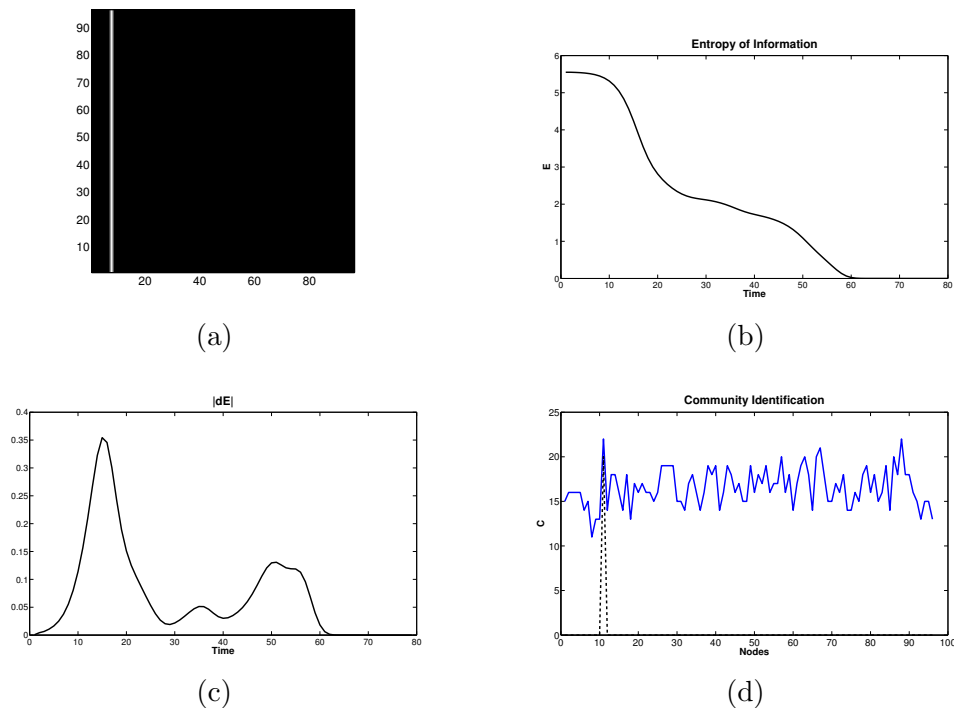


FIGURE 3.12: Temporal evolution of the Shannon Entropy of Information.

Referring to the artificial hierarchical network illustrated in Figure 3.1: (a) the asymptotic configuration observed with our algorithm using  $m = 0.27$  and  $\alpha = 1.25$ . (b) The Entropy of Information ( $E$ ) vs Time: we can observe three different *plateaux*. The final configuration,  $E = 0$ , corresponds to the monocluster shown in (a); (c) plot of the first derivative of the Entropy which show the three *plateau* with three different peaks. (d) We observe that the final community is identified by the higher connectivity.

### 3.4.4 Local and long range interaction

We have generated matrices as defined in Section 3.2, with  $K = 5$ ,  $N = 120$ ,  $G = 3$ ,  $C = 2$  (3 groups of 40 nodes, and communities of 20 nodes). After having generated the networks with uniform local connectivity  $K$ , links are rewired with probability  $p_r$ . If rewired, the site is connected to another one (possibly already connected) in the same community with probability  $p_c$ , in the same group with probability  $(1 - p_c)p_g$  and to a random node with probability  $(1 - p_c)(1 - p_g)$ .

An example of such matrix is reported in Figure 3.15-a. The rewiring probabilities ( $p_r = 0.2$ ,  $p_c = 0.9$ ,  $p_g = 0.7$ ) are such that the local structure is extremely evident, followed by the community structure. The group structure is almost invisible.

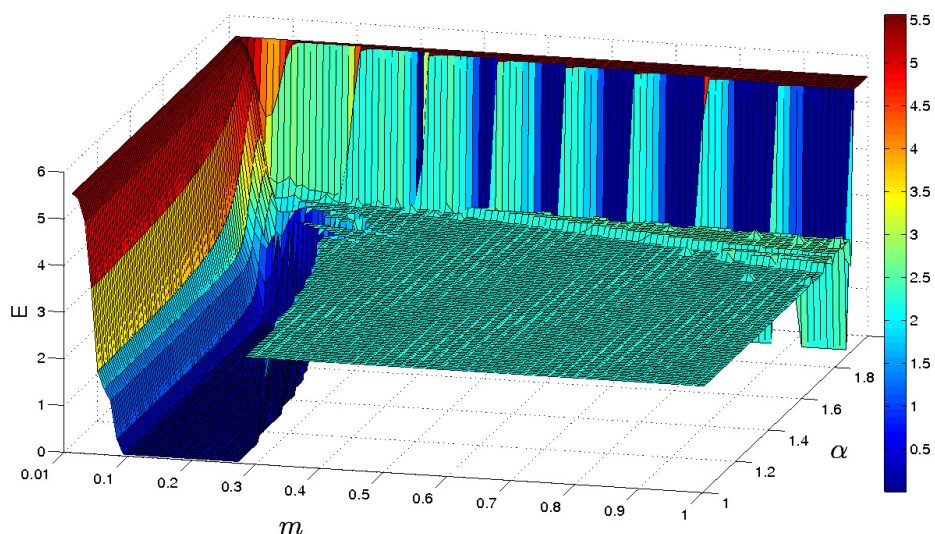


FIGURE 3.13: The asymptotic configuration of the entropy as function of the parameters  $m$  and  $\alpha$ .

We can clearly observe many different surfaces in this 3D graph: the surfaces correspond to different asymptotic configurations.

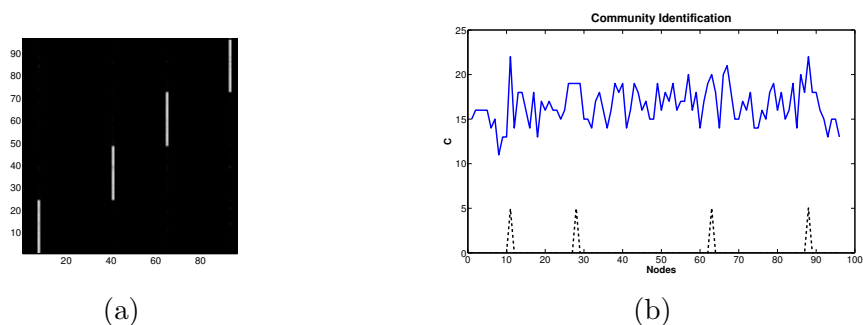


FIGURE 3.14: Hierarchical community detection.

(a) Asymptotic configuration of the matrix  $S$  using  $m = 0.7$  and  $\alpha = 1.4$  (this configuration corresponds to the large green surface in Figure 3.13). We observe that the asymptotic result corresponds to the middle layer (four communities) of the hierarchical structure of the network; (b) We have identified all the different levels in a hierarchical complex network changing the parameters  $m$  and  $\alpha$ .

In order to reveal all structures of communities, we have slowly varied  $m$ , for a given value of  $\alpha$ . The community structure for  $\alpha = 1.04$  is reported in Figure 3.15-b. The levels of  $\exp(H)$  corresponding to the group and community structures are marked. By changing the value of the memory  $m$ , nodes tend to accumulate more knowledge about the external world, and, due to the competition phase, their memory becomes dominated in general by just one label, that marking the community the node belongs to. There are occasional transitions when two communities fuse together, in the sense that a label from one community invades the other. It is possible to see these transitions. In particular,

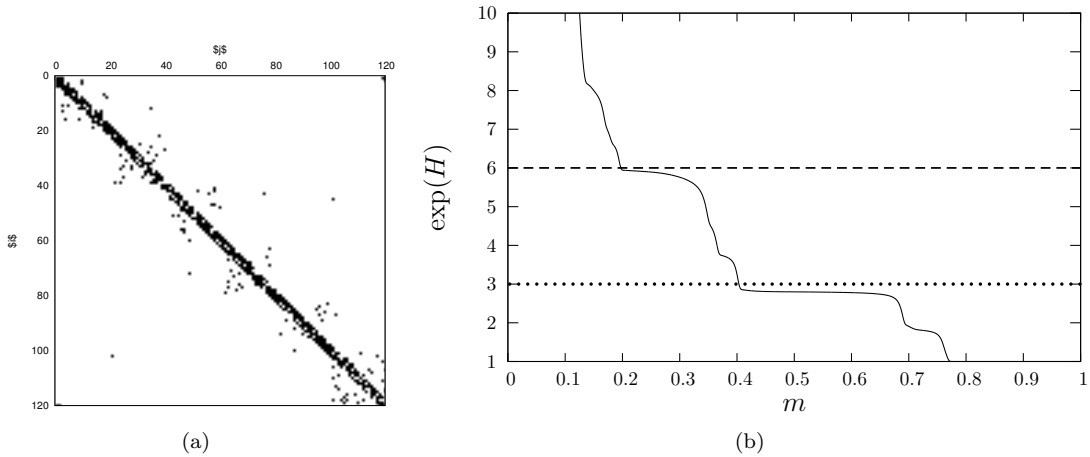


FIGURE 3.15: Undirected network.

(a) One adjacency matrix for  $K = 5$ ,  $N = 120$ ,  $G = 3$ ,  $C = 2$ ,  $p_r = 0.2$ ,  $p_c = 0.9$ ,  $p_g = 0.7$ . Black dots corresponds to  $A_{ij} > 0$ . (b) The plot of  $\exp(H)$  vs  $m$  for the network of Figure 3.15(a). The dotted line marks the value of  $\exp(H)$  corresponding to the three groups level and the dashed line marks the value corresponding to the six communities level.

plateaus corresponding to a structure in six and three communities are evident for large intervals of  $m$ . The final state corresponds to just one community (we have not reported the trivial initial phase, with  $H \simeq \log(120)$ ).

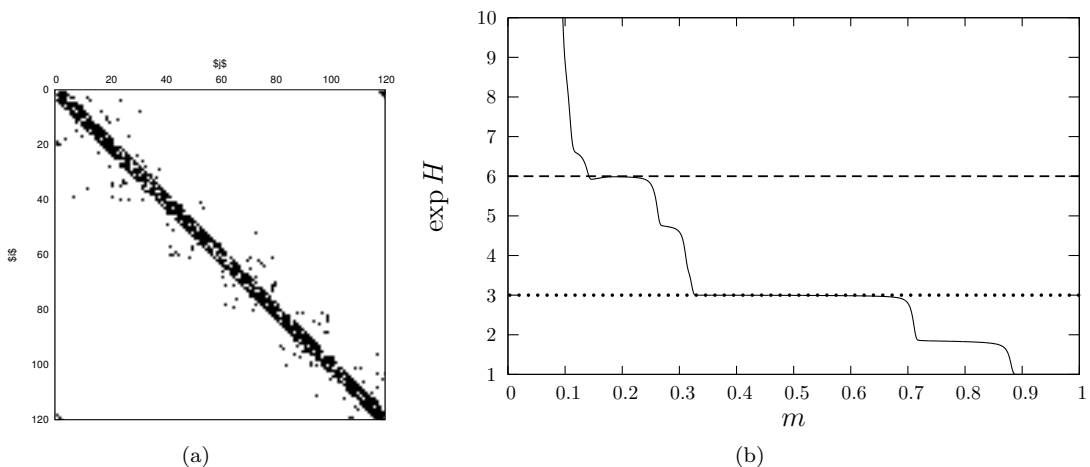


FIGURE 3.16: Undirected network 2.

(a) One adjacency matrix for  $K = 7$ ,  $N = 120$ ,  $G = 3$ ,  $C = 2$ ,  $p_r = 0.2$ ,  $p_c = 0.9$ ,  $p_g = 1.0$ . Black dots corresponds to  $A_{ij} > 0$ . (b) The plot of  $\exp(H)$  vs  $m$  for the network of Figure 3.16(a). The dotted line marks the value of  $\exp(H)$  corresponding to the three groups level and the dashed line marks the value corresponding to the six communities level.

Since the matrices are generated stochastically, it may happen that two communities are

more connected in one realization, and therefore the plateaus may happen for slightly different values of  $H$ .

Actually, the long-range connections at the network levels are not strictly needed: due to the local connectivity all nodes are connected, and we can set  $p_g = 1$  and still have the transition to a single community, but this is favored by a larger local connectivity  $K$ . See for instance the Figure 3.16 .

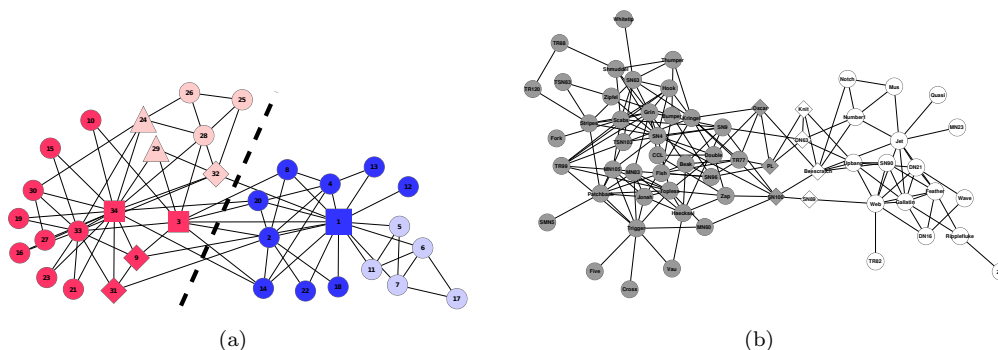


FIGURE 3.17: Different results on the Zachary's karate club and Bottlenose dolphin networks.

(a) Zachary's karate club network ( $m = 0.2$  and  $\alpha = 1.4$ ). (b) Bottlenose dolphin network ( $m = 0.5$  and  $\alpha = 1.03$ ).

We have also applied our method to two real networks, the well-known Zachary *karate club* network, Figure 3.17-a [1], and the social interaction of bottlenose dolphins observed by Leusseau, Figure 3.17-b [41].

For the Zachary club, our algorithm identifies four communities with different overlapping nodes between them. Considering the hierarchical structure of the network it is possible to merge together two sub-communities. Diamonds denote the overlapping nodes between the two principal communities. Triangles mark the overlapping nodes between the two sub-communities while square are the overlapping nodes between both subcommunities and communities.

The bottlenose dolphin network has a size of 62 nodes and was obtained by direct observation. Our algorithm detects 2 principal communities but also 7 overlapping nodes (diamonds) between them.

Finally, in order to evaluate our algorithm's performance we computed the normalized mutual information (NMI) on a Girvan-Newman (GN) benchmark graph [45] varying the mixing parameter  $\mu$ , see Figure 3.18. The graph consists of 128 nodes, each with expected degree 16, which are divided into four groups of 32. The GN benchmark is regularly used for testing algorithms for community detection.

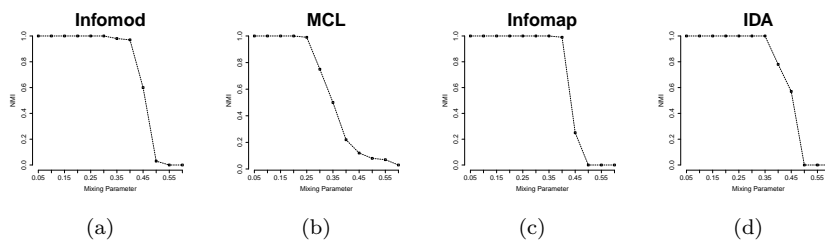


FIGURE 3.18: Normalized mutual information on the GN benchmarks.

Test of the algorithm on the GN benchmark based on normalized mutual information (NMI) on the y axes, and the mixing parameter  $\mu$  on the x axes. (a) Infomod [42], (b) MCL [43] (c) Infomap [44], (d) our model.

We created 11 networks varying the mixing parameters  $\mu$ , and compared the results with other well known community-detection algorithms. We performed simulations with different values of parameters  $m$  and  $\alpha$ . Results (Figure 3.18-(d)) show that our algorithm achieves very good performance: in fact, up to  $\mu = 0.35$  it always finds the predefined partition in four clusters. In the Figure 3.18-(a)-(b)-(c) we reported the results achieved by Lancichinetti and Fortunato [46] on three well-known community detection algorithms.

A final remark concerns the memory requirement of our cellular automata. We have chosen here the simplest implementation by furnishing to all nodes enough memory to contain the whole network (*i.e.*,  $S$  is a  $N \times N$  matrix), but in practice the number of entries different from zero are always quite few. It is therefore possible to assume that the nodes have bounded memory, as required by the “prescriptions” of human heuristics.

### 3.5 Dynamic Networks

Nowadays, the closer and closer interaction between devices and their users is a clear expression of the increasing tightness among the cyber world and the physical one. Let us consider, for example, mobile devices that are in charge of autonomously accomplish tasks like discern, collect and redistribute the important information (for their users) that circulate in the environment. On the one hand, devices exploit through their applications the information coming from the physical world to adapt and optimise their behaviour in the cyber world and, on the other hand, the status of the applications in the cyber world can affect the behaviour of their human users in the physical world (such as the users behaviour in social gaming or other social-oriented applications). This strong interaction has not only the quite obvious effect of generating a huge amount of information that flow from one world to the other, but it also triggers between the two a deeper connection, leading to the so called Cyber-Physical World (CPW) convergence

scenario [47]. In this context, mobile devices play an important role because they are the actual representation of their users in the cyber world or in other terms, mobile devices act as proxies of their human counterparts in the cyber world. The challenge here is to devise methodologies that make devices able to properly mine the acquired knowledge in order to let them being aware about their environment so that they can make proper decisions for specific tasks. Opportunistic Networks (OppNets) and the problems connected to them, represent a very well fitting example of the this general concept.

OppNets [48] are dynamic, delay tolerant wireless networks made by mobile nodes (e.g. human users equipped with smartphones) where the connectivity between them is not granted at any time instant. In OppNets the communication between nodes can occur only upon contacts, (i.e. when nodes are in reciprocal transmission range) and the information spreading is mainly done according to the *store carry and forward* paradigm: nodes exploit any contact with other peers to exchange messages under the condition that the other peer is deemed a good candidate to bring the message closer to the destination. The efficient delivery or spread of information to the interested users in this kind of networks is currently an open research problem. To this end, researchers, not only have to consider the typical physical problems of wireless networks but also the aspects connected with the humans' behaviour like their mobility patterns, their natural tendency to aggregate in social communities, etc. The ability of catching and understanding such social information, in order to predict and exploit human behaviour, has a great relevance for the development of effective solution for the above mentioned problems in OppNets. Let us consider for example the message forwarding problem in OppNets: due to the high mobility of devices, the challenge for a forwarding method is to quickly forward the message from the source to the destination, without introducing too many duplicate messages or overhead information. Here, the nodes' awareness about information like the social relationships, the aggregation habits and the community structure of their human users (all information coming from the physical world and exploited in the cyber world), can help to select suitable forwarders while containing the delivery costs.

In this work we focus on the community detection problem, or in other terms, we want to identify the underlying social structure emerging from the habits and the social relationships that nodes (human users) exhibit in a dynamic network by repeatedly entering in contact and exchange information with each other. The existing approaches of community detection can be divided into *centralised* and *decentralised*. The former consider the network as a static weighted graph that does not evolve over time and they use the whole global information about physical contacts and interactions between users over a



given amount of time to accomplish the task. In the literature are present many centralized community detection approaches as reported in [49]. Moreover, both in *physics* and *computer science* fields, many well-performed and improved algorithms for detecting communities in complex networks have been presented in the last decade, among the others we refer to the so-called OSLOM [50], INFOMAP and HIERARCHICAL INFOMAP [44, 51], MODULARITY OPTIMIZATION [52], LOUVAIN METHOD [53] and the LABEL PROPAGATION METHOD [54]. Though they are very useful for offline data analysis on mobility traces to define a priori strategies of forwarding, data dissemination, energy saving, etc, they are completely unsuitable to be used in online contexts, i.e. to define online, distributed algorithms that nodes can use to understand the structure of the communities of the network where they move. We recall that in our scenario nodes must be able to make proper decision without relying on centralised information so it is very important that nodes autonomously build a local representation of their surrounding environment. Several decentralised approaches have been proposed for community detection. Differently from the centralized ones, they consider the network as a time-evolving entity and, more important, they do not rely on a global vision of the network but only on a local one, i.e. every node in the network builds and updates its own representation of the existing social communities over time. For example, in [55, 56] the authors presented three community detection algorithms (i.e., SIMPLE, k-CLIQUE, and MODULARITY) while another improved approach can be found in the work of Borgia et al. [57]. All these methods, are based only on the contact duration to build the representation of the nodes' social structure, however, we tackle the problem from another point of view. In particular in our method the representation of community structure emerges spontaneously due the elaboration of information carried on by the agents locally and over the time without any intervention from the outside.

Here we focus on the development of a completely local approach of community detection considering also some social and psychological aspects: humans' communities are large and varied, we recognize several levels of grouping, sometimes dependent on the context, and we have probably developed our language as a tool for faster communication and discovering of social relationships. Therefore in social networks it is very difficult to have a precise definition of community because of people often belong to different communities at the same time and there is not a clear distinction between a community and a rest of the graph. In general, there is a continuum of nested communities whose boundaries are somewhat arbitrary. A community-detection algorithm should therefore return different "views", according to the value of some control parameters. At a superficial level, most of our information processing concerns the evaluation of probabilities. When faced with insufficient data or insufficient time for a rational processing, humans have developed algorithms, called heuristic in the cognitive psychology area, that allow

us to take decisions in these situations. The modern approach to the study of Cognitive Heuristics defines them as those *strategies that prevent one from finding out or discovering incorrect answers to problems that are assumed to be in the domain of probability theory*. Basically the Cognitive Heuristics Program proposed by Goldstein and Gigerenzer suggests to start from fundamental psychological mechanisms in order to design the models of heuristics [34]. These models have to satisfy the following constraints: (a) *Ecologically rational* (i.e., they exploit structures of information in the environment), (b) *Founded in evolved psychological capacities* such as memory and the perceptual system, (c) *Fast and frugal*, and simple enough to operate effectively when time, knowledge, and computational power are limited.

Our main goal is to explore the behaviour of exploratory methods inspired by human heuristics, in the hope of exploiting the “social knowledge” of human mind and also for developing more “natural” human-computer interfaces. Clearly, we do not pretend to simulate the real human behaviour, but only to study the behaviour of simplified models inspired by it. In particular, we deal with the task of identifying communities in an existing graphs, using a local and dynamic algorithm where an individual is simply modelled as a memory and a set of connections to other individuals. In this approach the information about neighbouring nodes is propagated and elaborated locally over the time as function of the previous meetings. In this way we are able to simulate a process in which the agents, through an alternation of communication and elaboration phases, have their local subjective representation of network in which the communities exploration is given by the probability to belong to one or more clusters at the same time. This method, already tested for detecting communities in static networks [58–60], is now applied in dynamic environments.

We now show how the algorithm performs when used in nodes moving according to one of the reference mobility models in the opportunistic networking literature, i.e. HCMM [61], already used in several works like [62, 63] to evaluate the performance of both forwarding and data dissemination approaches for OppNets. This allows us to show that our algorithm can be used to dynamically detect the structure of communities of users in mobile social networking environments. Mobility traces generated by HCMM incorporate temporal, social and spacial notions in order to obtain a proper representation of the real user movements. More precisely, nodes move in an area of  $1000m^2$  divided in a  $6 \times 6$  grid where a single grid’s cell represent a physical location that gathers nodes belonging to a single social community. In this synthetic scenario, communities are placed far from each other so as to avoid any border effect e.g. involuntary communication between groups. In each community we find two kind of moving nodes: travellers and non-travellers. Non-travellers roam only inside their community, while travellers, from time to time, use to visit other social communities different from the one they belong to.

With this configuration we want to simulate different social communities where usually people stay, apart for few of them that due to their social relationships can meet people from different social communities. In this context, the only way to exchange information between nodes is through nodes mobility, and travellers play an important role because they are the unique bridge between communities. In our experimental setup, we consider a network of  $N = 90$  nodes, divided in 3 separated communities and we study the performance of the algorithm through incrementally increasing the number of travellers for each community. We want to evaluate the average discovery time of the underlying community structure together with the goodness of the detection itself. Indeed, increasing the number of travellers contributes to ease the information flow from one community to another but also to blur the actual community boundaries making the community detection problem more and more challenging. The detailed scenario configuration can be found in Table 3.1.

TABLE 3.1: Detailed scenario configuration

Parameter	Value
Node speed	Uniform in $[1, 1.86m/s]$
Transmission range	$20m$
Simulation Area	$1000 \times 1000m$
Number of cells	$6 \times 6$
Number of nodes	90
Number of communities	3
Number of travellers per community	3, 4, 7, 13
Simulation time	50000s

### 3.5.1 Performance evaluation

In the Figure 3.19 we show the results of the algorithm in the case with 4 travellers for each community. In Figure 3.19(a) we show the snapshot regarding the community structure revealed by our algorithm in which we can observe the 3 principal clusters but also the overlapping nodes between the community that are the corresponding travellers. As we told above our model defines through the state matrix  $S$  the probability for a node to belong to a certain community: this result is reported in Figure 3.19(b)-(c)-(d) where each bar of the histograms correspond to the probability for the nodes to belong to the different communities. For instance, looking at Figure 3.19(b) we can observe that nodes 1, 4, 7, 10 and 13 have an high probability to stay into the community 1 but the first four nodes have also a little probability to belong to other communities. In fact node 1 (blue bar in Figure 3.19(b)) is a member of community 1 with  $p \sim 0.78$  and of the community 2 with  $p \sim 0.22$  because it is a traveller between the two communities. While the node 13 has a probability  $p \sim 1$  to belong to the community 1: in this way

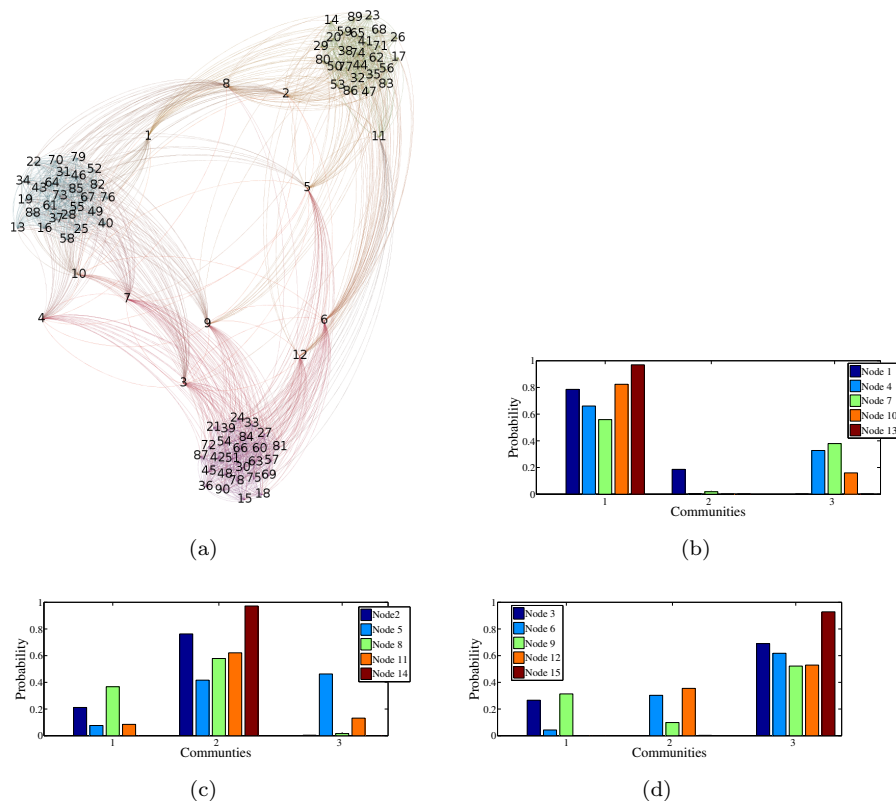
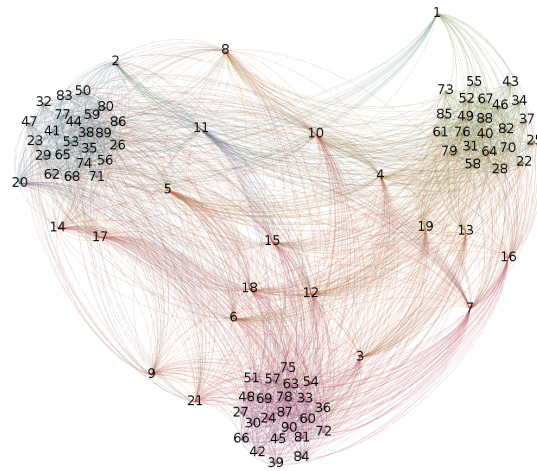


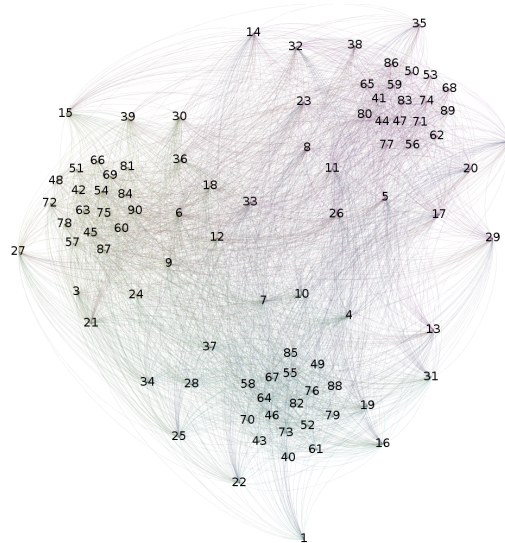
FIGURE 3.19: Dynamic networks.

Case with  $N = 90$  nodes, 3 communities and  $N_{tr} = 4$  travellers for each community. (a) Community structure of the network revealed by our algorithm where the link represent the encounters between the agents during time while the travellers are the overlapping nodes between the three principal communities. (b) Probability to belong to the principal communities in the case with 4 travellers for each community. Local vision of nodes 1, 4, 7, 10 and 13. These nodes are in the same community labelled as 4. As we can observe the nodes 1, 4, 7, 10 are the travellers of the community 1. In fact they belong with a certain probability also to other communities. On the contrary the node 13 has a very high probability to belong only to its community.

each node is aware of its role inside its community. In Figure 3.20(a)-(b) we report the snapshots of the final community structure detected by our algorithm considering 7 and 13 travellers, respectively: also here the algorithm is able to detect not only the three principal clusters but also the travellers as the overlapping nodes between the communities. In Figure 3.21 we show the different plots of the Shannon Entropy of Information for different cases considering different number of travellers. Here we can observe, not only the three plateaus corresponding to three principal clusters, but also different times for reaching the final state: in fact increasing the number of travellers the time for reaching the asymptotic state decrease as we can observe in Figure 3.21. In Figure 3.22 we report the local entropy for a traveller (black line) and for a normal agent (blue line) during time. The local entropy  $E^i$  is simply define as  $E^{(i)} = -\sum_j^N S_{ji} \log S_{ji}$



(a)



(b)

FIGURE 3.20: Community detection in mobility networks.

Final community structures detected by the algorithm considering (a)  $N_{tr} = 7$  and (b)  $N_{tr} = 13$  travellers respectively.

that represents the knowledge of the single node about the surrounding world. As we can observe the traveller tend to update its knowledge because of sometimes it jumps between other communities.

In this section we want to explore the performance of our method taking into account two different heuristics in order to allow the method to be free by the parameters  $m$  and  $\alpha$ . The main feature of these two heuristics called respectively *IDA + LTE* and *Double pruning* is to reveal the hierarchical community structure of complex networks from a

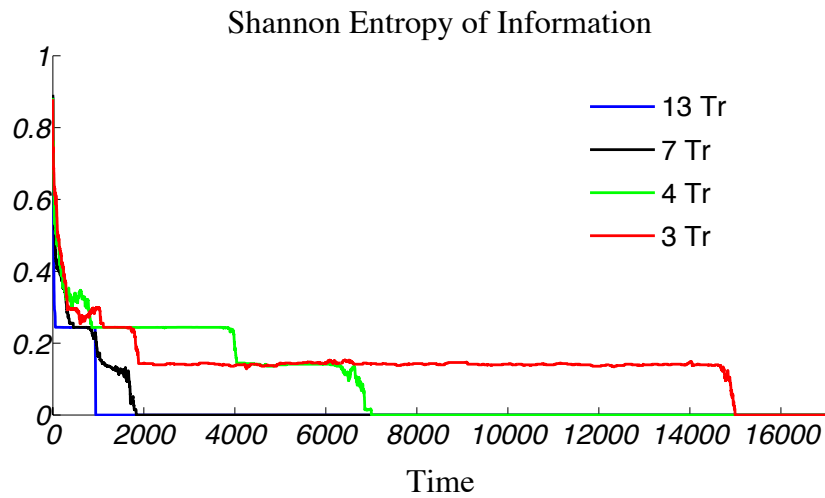


FIGURE 3.21: Shannon Entropy of Information during time for different scenarios. Here  $N_{tr} = 13$ ,  $N_{tr} = 7$ ,  $N_{tr} = 4$  and  $N_{tr} = 3$  travellers.

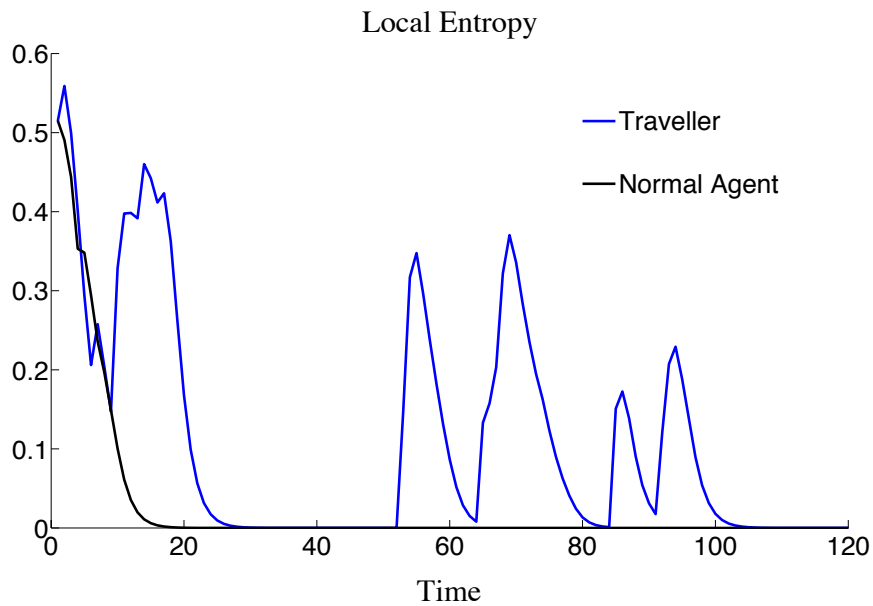


FIGURE 3.22: Comparison between the local entropy of a traveller (blue line) and a normal agent (black line). The peaks correspond to the jumps in other communities.

local viewpoint. We demonstrate that a node can discover the network's multi-levels giving its subjective vision of the world until it wants or can explore the graph without any global information about the entire network. The experimental results show that the proposed algorithms are efficient in both real-world and synthetic networks.

### 3.6 Evaluating cerebral cortex connectivity with local information algorithm

Brain's intricate network of neurons is one of the most prototypical examples of a complex system. Since the end of the 80s, with the development of techniques that allow to study anatomical connectivity, a considerable amount of data has become available to researchers [64]. Classical studies have focused on the rat, cat and macaque brains obtaining the adjacency matrix of brain structures or the weighted connection matrix (in which the weights reflect the strength of the connection) [64, 65].

In order to find similarities between anatomical and functional/cognitive properties of the brain, clustering organization of the cerebral cortex has been investigated in several works. Limiting to the cat cerebral cortex, the first studies conducted on this subject [66, 67] have employed both traditional multidimensional scaling methods and evolutionary optimization algorithm. The latter method identified clusters maximizing the intra-cluster connections and minimizing the inter-cluster ones, i.e. areas that belong to the same cluster should be more densely connected among them than areas in other clusters. This is a very general principle commonly used in community detection algorithm [45, 68]. The cluster organization found with this algorithm broadly agrees with functional distinctions of the cerebral cortex. These results have been confirmed in subsequent works [69, 70]. Summarizing, these studies found that this kind of data analysis is able to detect local clusters of densely interconnected areas that correspond to main functional properties of the brain.

So, the classical approach to the cortical network clustering has been almost always based on evolutionary algorithms that globally process the entire pattern of connectivity. The principal aim of this section is to complement the results of such clustering approach, using a local algorithm based on a bio-inspired model of diffusion of the information, equipped with a bounded memory and processing capacity. We name our algorithm *NetExplorer*. It is based on a modified version of the Van Dongen algorithm [43, 58, 59]. Three principal aspects distinguish our use of the Van Dongen algorithm from the ordinary community detection algorithms. The former is represented by an intrinsic local approach of the method, which allows the method to run even without a perfect knowledge of the network. The second aspect relies on the stochastic (i.e instead of deterministic) nature of the algorithm, which randomizes, using a precise recipe, the parameters of the model in order to optimize the exploration of the network [59]. The latter is that our method represents a microscopical model describing the communication dynamics of a complex network, representing in this way a model potentially useful to study more deeply the actual structure/organization of real neural networks in the brain.

On these basis in this work we employ graph theory to investigate how different brain areas communicate. In particular, we will apply the *NetExplorer* algorithm to detect the cat cerebral cortex communities. The accuracy and effectiveness of our model have been carried out comparing the numerical estimation of *NetExplorer* with the results provided by six other classical/relevant clustering algorithms.

In the next sections, we describe our model and the other six algorithms and then we present the results of the numerical simulations. In the last section finally the conclusions and the discussion are reported.

### 3.6.1 Net Explorer

Here we introduce *NetExplorer*, a community-detection algorithm derived from the van Dongen's Markov Cluster algorithm (MCL) method [43, 60]. The MCL algorithm simulates a sort of diffusion process over the graph, followed by a prouning phase in which the competition among the links is performed in order to eliminate the weakest. In the model the graph is expressed by the correspondent adjacency matrix  $A$ : specifically, the adjacency matrix of a finite graph  $\mathbf{G}$  on  $n$  vertices is the  $n \times n$  matrix where the non-diagonal entry  $A_{ij} = 1(0)$  indicates the presence (absence) of a link from the node  $i$  to the node  $j$ .

In our model the number of plateaus of the function  $E_i^t$  correspond to the different  $L$  levels: If we evaluate the first derivative of the entropy we can identify  $L$  peaks, while in the second derivative, we observe  $L$  changes of sign. For this reason, we evaluate the first and the second derivative of the local entropy for each node.

Analogously for the entropy defined above, it is possible to introduce the concept of local entropy for each node in order to study the local view of agents. Similarly it is possible to detect different plateaus corresponding to the different network sub-clusters that the single node discovers during time. We observed that we can use a fixed value of the parameter  $m$ , while we have to change the value of  $\alpha$  in order to find the community and in particular the hubs that labels each community. For this reason, we repeat several times an exploration phase of the network in which the nodes save their state vector  $S_i^t$ , in a *temporary memory box*, until they observe a change in the sign of the second derivative. If the following condition is satisfied

$$\text{sign} \left( \frac{d^2 E_i^{t-1}}{dt^2} \right) \neq \text{sign} \left( \frac{d^2 E_i^t}{dt^2} \right), \quad (3.8)$$

the state vector  $S_i^t$  is stored into the temporary long term memory together with the value of the first derivative of the local entropy and the entropy. When a node meets an



impasse (e.g. its state vector entropy do not evolve any more) its  $\alpha$  is changed by the following mechanism, if

$$\left| \frac{E_i^{t-1}}{K_i} \right| - \left| \frac{E_i^t}{K_i} \right| < \epsilon, \quad (3.9)$$

where  $K_i$  is the connectivity degree of the node  $i$ , and the parameter  $\epsilon$  is equal to  $10^{-5}$ . Then, a counter  $\tau$  is increased by 1, and if  $\tau_i$  becomes greater than a given threshold (say  $\tau^*$ ), the parameter  $\alpha$  is updated in the following way:

$$\alpha_i = 1.5|\eta\sigma| + 1, \quad (3.10)$$

where  $\eta$  is a random Gaussian variable with mean 1 and standard deviation  $\sigma$ , and  $\alpha$  ranges in the interval  $(0, \infty)$ .

### 3.6.2 Other algorithms

In order to explore different results and different scenarios we have tested other six well-known algorithms for detecting communities in complex networks, in particular: *Oslom*, *Infomap*, *Hierarchical Infomap*, *Label Propagation Method*, *Modularity Optimization* via simulated annealing and the *Louvain Method*.

The *Order Statistics Local Optimization Method (Oslom)* [50] is a local optimization method applied to a fitness to measure the statistical significance of individual communities. The statistical significance is defined as the probability to find a a very similar community in a random model without community structure.

*Infomap* [44] is a dynamic algorithm that compresses the information of a random walk on the graph. It then defines a function which is optimized when it reaches the minimum description length of the random walk. Such optimization is carried out with a combination of greedy search and simulated annealing. There is also a hierarchical version of this algorithm (*Hierarchical Infomap* [51]) that is also able to find the multilevel organization of several complex networks.

*Modularity Optimization* [52], via simulated annealing, is essentially a fast implementation of the Girvan and Newman (*GN*) algorithm [45]. The *GN* algorithm is a divisive algorithm in which links are iteratively removed based on the value of their betweenness. The procedure of link removal ends when the modularity of the resulting partition reaches a maximum. The *Modularity Optimization* [52] via simulated annealing starts from a set of isolated nodes and then some links of the original graph are added to these nodes in order to obtain the maximum value of the modularity.

*Louvain Method* [53] defines a fitness function in order to optimize the modularity of Girvan and Newman in neighbourhood through a local optimization method. Then, after finding the first partition, the clusters are agglomerate together giving rise to super nodes, and so on until the modularity doesn't increase any more.

*Label Propagation* [54] method uses the concept of the spreading of the information over networks considering the node neighbourhood in order to identify the communities. Similar to our model, initially, each node is labelled by an unique value then trough the diffusion of information nodes become aware of the surrounding world. The process ends when the knowledge of each node does reach a stationary state. The resulting communities are labelled by the last label values.

### 3.6.3 Numerical Results

We applied the *NetExplorer* and the other community detection algorithms to a connection matrix of the cat cerebral cortex comprising 52 areas (clustered into 4 different systems, respectively the visual, auditory, somatosensory-motor, and frontolimbic one) linked by 818 pathways [66, 67]. This database about cortico-cortical connections was produced employing information of a large set of anatomical studies. The authors reviewed these studies in order to assign connection weights between cortical areas. They employed strict criteria such as using only data from cat's brain or giving priority to results obtained from higher resolution techniques. Connections weights were measured from 0 (i.e. absent or unreported) to 3 (i.e. strong or dense connections). Since our goal was to evaluate and compare both the quantitative and qualitative predictive ability of the community detection algorithms considered, we organize this section accordingly. First of all we define a standardized error function to estimate the degree of agreement between each model's prediction and the empirical dataset considered. As a first approximation of such an error, we compute the normalized difference (3.11) between the vectors and the adjacency matrices representing, respectively, the models' attribution of the cortical areas to the principal cortical clusters/communities, and the boolean representation of the entire map.

$$N_E^V = \frac{\sum_{i=1}^N 1 - \delta(E_i^V - P_i^k)}{N} \quad (3.11)$$

where  $E^V$  is the vector containing for each area  $i$  the  $id$  of the clusters empirically containing it.  $P_i^k$  is the  $id$  of the cluster coupled by the community detection models  $k$  to the area  $i$ . Finally the  $\delta(x)$  function assumes the value of 1 whenever  $x = 0$  (i.e.  $E_i^V$  is equal to  $P_i^k$ ), and 0 otherwise.

$$N_E^M = \frac{\sum_{i=1}^N \sum_{j=1}^N 1 - \delta(E_{ij}^M - P_{ij}^k)}{N(N-1)} \quad (3.12)$$

where  $E^M$  is the real empirical cortical map, in which the element  $E_{ij}^M$  assumes the value 1 if the two areas share the same hub, and is 0 otherwise.  $P_{ij}^k$  has the same shape of  $E^M$ , and represent the resulting matrix of the  $k$  models' prediction about the entire cortical map.

Both the normalized errors clearly range within the interval  $(0, 1)$ , and represents the effectiveness of the approximation of the empirical map made a method in terms of percentage of correct attribution.

### 3.6.4 Model's approximation of the clustering structure and of the adjacency matrix representing the empirical cortical map of the cat

In order to facilitate a visual inspection of the approximation accuracy of the different methods, the *empirical structure* (Figure 3.23(a)) has been reported bigger than the others in the center of the first row, and the principal clusters have been coloured just to make easier the distinction and comparison of the clusters, respectively within and between the sub figures (i.e. the same principal clusters are coloured the same in the different subfigures). The resulting images are reported in Figure 3.23(b-c-d-e-f-g-h), and the diagonal of the represented matrices is the *cluster attribution vector* defined above as  $(P_i^k)$ .

Starting from the vectors  $P^k$  the normalized errors for each method have been computed, and reported in Figure 3.24. The Figure 3.24(a) reports the normalized error of the *principal clusters composition* approximated by the community detection models under scrutiny, with respect to the empirical structure coming from the literature. Each method has been applied directly on the weighted cortical network used as reference dataset ( $E^M$ ), and each algorithms have provided as output an approximated clustering structure, and an approximated reproduction of the entire cortical map (Figure 3.23). A first inspection of the coloured representation of the clustering structure, suggests that despite a slightly difference in the total number of clusters detected (i.e. sometimes the error regards very few nodes/areas), all the methods approximate the clusters with an error smaller than the 30% (Figure 3.24(a)). Nevertheless the goodness of such an approximation range between the 5% to the 30%, with the *NetExplorer* algorithm that made an error of misattribution of only 3 areas on 52 with respect to the empirical reference. Noteworthy even the 3 misattributions of the *NetExplorer* method represent actually a meaningful information, suggesting the possible existence of two important

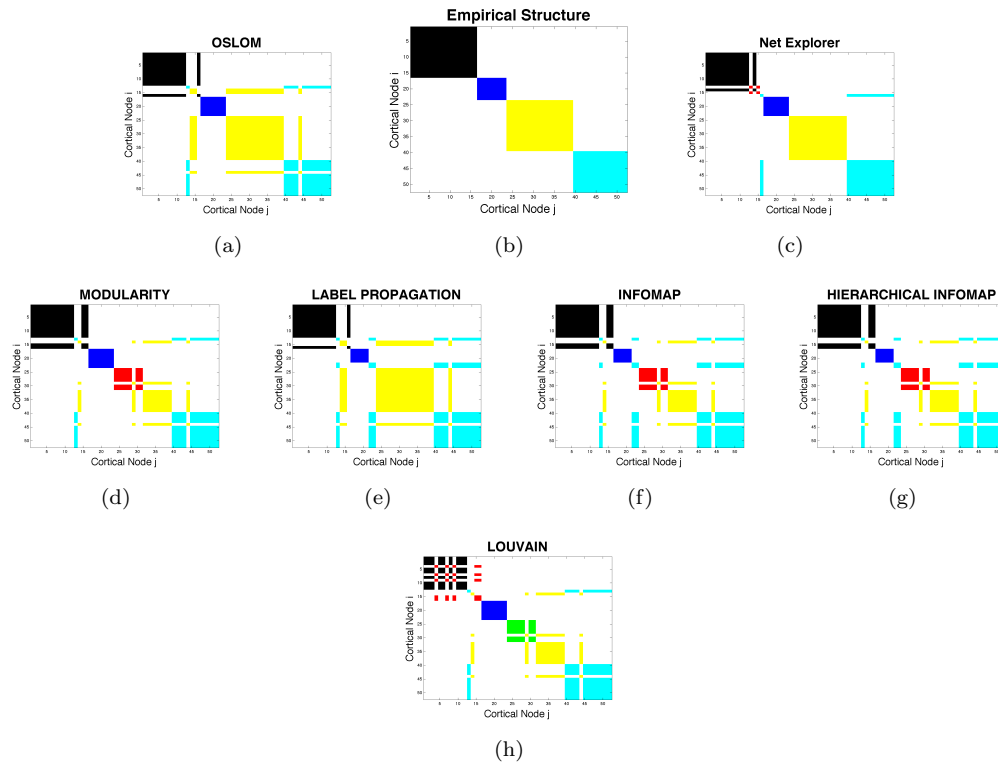


FIGURE 3.23: Comparison of different community detection algorithms over the cortical map of the cat.

In the figure are reported the different clustering structures resulting from the application to the different methods of community detection to the weighted cortical map of the cat. On the two axes are indicated the ID of the cerebral nodes considered. Nodes 1-16 are visual areas: 1. area 17, 2. area 18, 3. area 19, 4. PMLS (posteromedial lateral suprasylvian area), 5. PLS (posterolateral lateral suprasylvian area), 6. AMLS (anteromedial lateral suprasylvian area), 7. ALLS (anterolateral lateral suprasylvian area), 8. VLS (ventrolateral suprasylvian area), 9. DLS (dorsolateral suprasylvian area), 10. area 21a, 11. area 21b, 12. area 20a, 13. area 20b, 14. area 7, 15. AES (anterior ectosylvian sulcus), 16. PS (posterior suprasylvian). Nodes 17-23 are auditory areas: 17. AI area, 18. AII area, 19. AAF (anterior auditory field), 20. P (posterior auditory field), 21 VP (ventroposterior auditory field), 22. Epp (posterior part of the posterior ectosylvian gyrus), 23. Tem (temporal auditory field). Nodes 24-39 are somatosensory-motor areas: 24. area 3a, 25. area 3b, 26. area 1, 27. area 2, 28. area SII, 29. area SIV, 30. area 4g, 31. area 4, 32. area 6l, 33. area 6m, 34. area 5Am, 35. area 5Al, 36. area 5Bm, 37. 5Bl area, 38. SSAi area, 39. SSAo area. Nodes 40-52 are frontolimbic areas: 40. PFCMil (infralimbic medial prefrontal cortex), 41. (dorsal medial part of prefrontal cortex), 42. PFCL (lateral part of prefrontal cortex), 43. Ia (agranular insula), 44. Ig (granular insula), 45. CGA (anterior cingulate cortex), 46. (posterior cingulate cortex), 47. RS (retrosplenial cortex), 48. area 35, 49. area 36, 50. pSb (presubiculum), 51. Sb (subiculum), 52. Enr (entorhinal cortex). (a) Osлом, (b) Empirical data, (c) NetExplorer, (d) Modularity Opt., (e) Label Propagation, (f) Infomap, (g) Hierarchical Infomap and (h) Louvain.

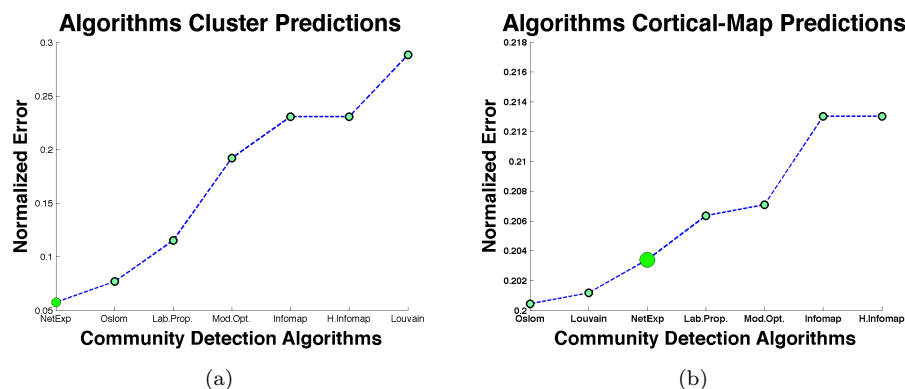


FIGURE 3.24: Normalized error of prediction for different algorithms.

In the figure are reported the normalized error of prediction ( $N_E^{V,M} \in (0, 1)$ ) for all the computational methods of community detection taken into account. The subfigures report respectively (a) the goodness of the principal clusters composition and (b) the goodness of the approximation of the entire cortical map. In the subfigure (a) the NetExplorer algorithm provides the best approximation of the clustering structure of the map ( $N_E^V \simeq 0.05$ ), and Oslom gets the closest performance to the NetExplorer with an error of  $N_E^V \simeq 0.075$ . In the subfigure (b) the best prediction is carried out by the Oslom algorithm, with an error of ( $N_E^M \simeq 0.2$ ), whereas NetExplorer ranks in the third position with an error of ( $N_E^M \simeq 0.206$ ).

hubs (i.e. the nodes 16 as *bridge* between the cluster 1 and the cluster 4, and the nodes 15 and 13 that together could be interpreted as a sensorial integration sub-cluster again bridging the cluster 1 with the others, see next section).

The Figure 3.24(b) reports the goodness of the methods' approximations of the entire cortical map under investigation taking into account both the vector  $P^K$  and the entire matrices  $P^M$ . To estimate such a function the matrices  $P^M$  have been compared with the unweighted adjacency matrix representing the known cortical network of the cat. The average normalized error for such a challenge increases, ranging within the interval (20%, 22%) probably indicating a lack of information to estimate the crossconnections between different cluster. Beyond the possible interpretation of such a result, it is noteworthy that our algorithm performs better than 4 of the 6 competitors, getting a good result with respect to the best method which here is *Oslom*.

### 3.6.5 Qualitative evaluation of the algorithms' performance

With regard to the functional and cognitive evaluation of the detected clusters we observe that misattributions of cortical areas of each algorithm are not equally explainable.

*NetExplorer* detects five clusters three of whom perfectly correspond to auditory, somatosensor - motor and frontolimbic systems. Visual cortex is divided into a big cluster

that comprehends the majority of visual regions (areas 17, 18, 19, *PMLS*, *PLLS*, *AMLS*, *ALLS*, *VLS*, *DLS*, 21*a*, 21*b*, 20*a*, 7) and a smaller cluster of only two visual regions (area 20*b* and *AES*). *AES* is a multimodal (visual and auditory) area that receives information from area 20*b* (a retinotopically organized region). The detection of this cluster can be explained with the multimodal nature of *AES* area (and its connection with area 20*b*) that makes them sufficiently different from other visual cortical areas to create another cluster (although there isn't any physiological evidence of such a role of area 20*b*). The only misattribution of the *NetExplorer* solution regards *PS* area, one of the few cortical areas (the others are *AES* and area 7) that is well connected to some frontolimbic structures such as Ig (granular insula) and area 36 (perirhinal cortex). Since its role of connection between the "output" of visual system and the more cognitive-oriented frontal regions (in particular, the brain structures responsible of mnemonic functions), it can be reasonable that *NetExplorer* assigns *PS* to the frontolimbic system.

As far as other algorithms are concerned, the worst solution is produced by *Louvain*. This algorithm detects 6 clusters in which only the auditory cortex is correctly individuated. Visual areas are divided into two clusters: the first one is composed by areas 17 – 19, *PLLS*, *AMLS*, *VLS*, area 21*a*, 21*b*, and 20*a*, whereas the second cluster is composed by *PMLS*, *ALLS*, *DLS*, *AES* and *PS*. A functional criterion that allows to find something in common within the areas of this two clusters, is very hard to develop. *PMLS*, *PLLS*, *AMLS*, *ALLS*, *VLS*, and *DLS* are visual areas that are (to a certain extent) overlapping (their functional distinction is not so clear). Moreover, there are visual areas that are misattributed to other clusters: 20*b* area (a visual retinotopic area as we noted before) and area 7 (a visual area that responds also to other sensorial modalities such as somatosensorial and auditory). Area 20*b* is attributed to the frontolimbic cluster: this area might possibly have a role in the integration of visual information with the cognitive and emotional role of the frontolimbic system. The misattribution of area 7 to the somatosensory-motor cluster is much easily explainable because it responds also to somatosensorial and auditory stimulus beyond visual. The *Louvain* algorithm also individuates two clusters in the somatosensory-motor group of cortical areas. The first one is composed by five somatosensory areas (areas 3*a*, 3*b*, 1, 2, and *SII*) and two motor areas (areas 4 and 4*g*) whereas the second one is formed by five somatosensory areas (*SIV*, 5*Am*, 5*Al*, *SSAi* and *SSAo*) and four motor areas (6*l*, 6*m*, 5*Bm*, and 5*Bl*). Again, interpreting this result is not easy because each cluster contains cortical areas with very different function (somatosensorial, motor, visualmotor, somatosensorial integrated with other sensorial modalities) without an apparent order. Moreover, the most natural distinction of the somatosensory-motor cluster (that is between somatosensorial areas and motor areas) is not respected. With regard to the frontolimbic cluster, there

is only one misattribution: *Ig* (granular insula) area is attributed to the second cluster of the somatosensory-motor system (together with areas *SIV*, *5Am*, *5Al*, *SSAi* and *SSAo*, *6l*, *6m*, *5Bm*, and *5Bl*). A possible explanation of this misattribution is that *Ig* area responds to stimuli of different sensorial modalities. However, there are other multimodal areas (such as *CGA* and *CGP*) that are assigned to the frontolimbic cluster by the algorithm (and not to a specific sensorial modality).

*Infomap*, *Hierarchical Infomap* and *Modularity Optimization* find five clusters. *Infomap* and *Hierarchical Infomap* give the same solution. In particular, the visual cluster is correctly individuated with the exception of areas *20b* and *7*. As we said before, the misattribution of area *7* is easier to explain compared to area *20b*. With regard to the auditory function, we have that two regions that are incorrectly attributed to the frontolimbic cluster (*Epp* and *Tem* areas). These two areas belong to the so-called “auditory belt” that has a role in the multimodal visual/auditory processing of the stimuli. It would be more reasonable that the algorithm assigned these areas to the visual cluster. Nevertheless, due to their multimodal nature it can be reasonable to have close relationship with the frontolimbic structures. The algorithms find again two clusters in the sensorimotor part of the brain with the same structures as *Louvain*’s (and so not so easy to explain from a functional point of view). Lastly, *Ig* (granular insula) area is attributed to the second cluster of the somatosensory-motor areas (as in *Louvain*’s solution). *Modularity Optimization* gives a similar result compared to *Infomap* and *Hierarchical Infomap* but it is able to correctly detect all the areas that belong to the auditory cluster.

*Label Categorization* and *Oslom* find correctly only 4 clusters. Three areas that belong to the visual cluster, are both attributed to the somatosensory-motor system (area *7* and *AES*) and the frontolimbic cortex (area *20b*). We have already discussed the misattribution of area *7* and area *20b*. *AES* is very similar to area *7* (a visual area that responds also to other sensorial modalities such as somatosensorial and auditory); an explanation of this misattribution can rely on its multimodal nature. As in *Infomap* (and *Hierarchical Infomap*), the two areas of the auditory belt (*Epp* and *Tem* areas) are attributed to the frontolimbic cluster. Somatosensory-motor and frontolimbic systems are correctly clustered with the exception of *Ig*, as we have seen before. Lastly, *Oslom* algorithm gives the best solution from a qualitative point of view. Auditory, somatosensory-motor and frontolimbic clusters are correctly identified by the algorithm with the exception of *Ig*. With regard to the visual system, as in *Label Categorization* area *7* and *AES* are attributed to somatosensory-motor system whereas area *20b* is assigned to the frontolimbic cortex.

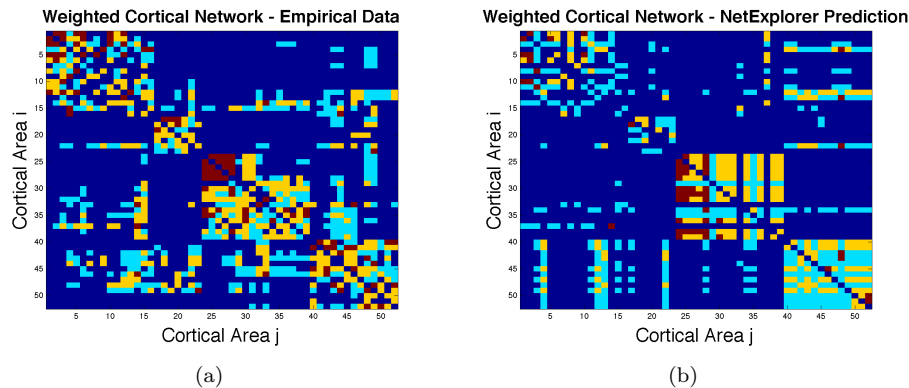


FIGURE 3.25: Reconstruction of the cortical map of the cat.

In the figure are represented the weighted cortical networks ( $W^E$  and  $W^{NetExp}$ ) deriving respectively from the empirical data (a), and the NetExplorer approximation (b). Noteworthy the error characterizing such a prediction ( $N_E^W \simeq 0.13$ ) is smaller than the error of the best methods' performance (i.e. *Oslo*) in the prediction of the simple adjacency matrix which was of  $N_E^M \simeq 0.2$ .

To sum up, we can conclude that *Infomap*, *Hierarchical Infomap*, *Modularity Optimization* and *Louvain* detect more than four clusters but it is no easy to explain the distinction within somatosensory-motor system (for all these five algorithms) and within visual cortex areas (only *Louvain*). *Oslo* and *Label Categorization* correctly detect four clusters and misattribution of areas to other systems can have a cognitive or functional explanation. *NetExplorer* is the only algorithm that is able to correctly detect the cortical areas that belong to auditory, sensorimotor and frontolimbic structures. Although it finds a second smaller cluster in the visual system areas, these misattributions are explainable on the basis of the cognitive role of these areas.

### 3.6.6 Empirical weighted cortical map approximation by the *NetExplorer* algorithm

In this section we have taken into account the weighted representation of the cortical map provided by our model to estimate its predictive efficiency with respect to the empirical reference weighted representation of the map. It is important to notice that, despite the empirical weighted map is contemporary the seeding structure of the algorithms and the final reference for the results, is absolutely not trivial to be able to explain all the complex features of a cortical/neural network just considering its *visible* connectivity pattern. The goodness of approximation of the weighted cortical map provided by our algorithm is of 13%, and it is even better than that provided for the unweighted network, which was around the 20% (Figure 3.25). The degree of approximation of our method



(Figure 3.25a vs Figure 3.25b) is quite good also because the not symmetric structure of the final matrix produced.

Finally the *NetExplorer* algorithm individuate even the hubs (i.e. the most important nodes within a community) for each principal cortical cluster revealed. Such nodes are respectively the number 3 (black), the number 15 (Red and new with respect to the literature), the number 17 (Blue), the number 28 (Yellow), and the number 48 (Cyan).

### 3.6.7 Discussion and conclusions

The general results of this work confirm the goodness of the complex network analysis to the study of the cerebral cortex connectivity and to give a valuable support to the neurocomputational modeling [64]. The particular complexity implied by the study of the cortical/neural dynamics, frequently requires a differential approach to the interpretation of the results. In particular for the cognitive neuroscience the goodness of an approximation of a cortical/functional structure, have to be always considered as composed by a qualitative as well as a quantitative aspects. The former relies on the coherence of the classification produced by the community detection method, while the latter is related to the quantitative consistency/robustness of the solution provided. Both the previous indicators have been considered in the results section and will be merged together in the present section in order to get a most general overview of the system under scrutiny.

The overall forecasting efficiency of the analyzed methods suggests that the *NetExplorer* algorithm is the best to predict the actual positioning of the nodes within the cluster structure proposed in literature. The normalized error ranges within the interval (5%, 30%), nevertheless all the methods detect correctly the main backbone of the analyzed map. For what concern the entire cortical map approximation, the normalized error increases (20%, 22%), indicating how the complexity of the cross connections between the principal clusters increases with respect to the backbone. In the approximation of the entire map the best method appears to be *Oslom*, but even in this case the performance of our algorithm is quite close to the best one, ranking as third in general.

The clusters (and the hubs for each clusters) detected by *NetExplorer* can give interesting clues on how information flows through the connections of cat cerebral cortex. We can interpret the cluster individuated by the algorithm (with the exception of the small *AES-20b* sub-cluster within the visual system) as a further confirmation of the importance of small range connections within distinct functional systems [71]. From a neuroscientific point of view, neurons that are responsible of the communication among distant areas are especially important because they break the rigid cluster organization of cortical areas

[71]. The hubs individuated with our algorithm could represent the best candidate for individuating the regions in which such neurons are present. In particular, we found that for the visual system the hub is area 19 whereas for the auditory system the hub is *AI* region. This means that the single areas that belong to the visual system speak mainly with area 19 that can be an ideal candidate to send long range projections in order to share information. The same holds for auditory system, in which *AI* area could have this role too. With regard to the other two clusters, the main hub of the somatosensory-motor system is area *SII* whereas the main hub of frontolimbic structures is area 35 (perirhinal cortex). The fifth cluster detected by the algorithms (composed only by *AES* and area 20b) can be considered as another interesting candidate for neurons with peculiar communicative role within visual cortex.

Concluding, a very interesting perspective, opened by our bio-inspired and local approach to the community detection, is represented by the possibility of studying the nested network structure characterizing small range connections, long projections and hub in cerebral networks.

### 3.7 A Cognitive-inspired Model for Self-organizing Networks

Among the capabilities of the human cognitive system that is attracting most computer scientists, there is the ability of humans to develop local algorithms able to exploit what might be called “collective human computation”. As collective human computation, we refer to the natural synchronization between the cognitive elaborations made by a person which is immersed into group dynamics. In such a condition, human beings analyze only some relevant information coming from the group, giving to the group only some relevant contributions for the general problem which is faced. In this way, the group can be described as more than just the sum of its single components [72].

When being faced with insufficient data or insufficient time for rational processing, humans have developed strategies that allow to take decisions in these situations. In general, such an effect has been well described in the cognitive heuristics program proposed by Goldstein and Gigerenzer, which suggest starting from fundamental psychological mechanisms in order to design models of heuristics [34].

Only some relevant information are extracted from the environment while the rest is interpolated for building our knowledge. Among others, Milgram et al. have shown experimentally how humans are able to adopt effective strategies to solve very complex problems, exploiting optimally their partial knowledge of their environment [19, 22].

This kind of human distributed computing has been studied deeply only from the perspective of disciplines such as social cognition and social psychology, while it is not yet well known in other domains.

The social cognition domain studies human cognition as characterized by the use of “fast and frugal” solutions, that are specialized for a social context in which we live using a bounded rationality and limited computational resources [73]. Therefore, the aim of our work is to assemble a working computational scheme inspired by the operating principles of human cognition, based on general assumptions about cognitive high-level functions.

This approach promises to embed some fundamental aspects of the human cognitive system in a computational model in order to obtain a minimization of computational resources needed for the task and the evolution of a dynamic knowledge network capable of generating strategies suitable for networks like the Internet, which are too large and too dynamic to ever be fully/perfectly known [58].

The fundamental aspects on which we focused our modeling, involves the spread of information through a human network, and the knowledge representation arising from the dynamics of short-term memory (STM) and long-term memory (LTM). The passage of information between STM and LTM occurs through a simple cognitive heuristic approach, which compatibly with their computational capacity reduces the dimensionality of the information required to represent the environment in a dynamic manner.

In our previous work [58], we applied such an approach to the community detection problem, which can be considered as a task of great importance in many disciplines [23–27], where systems can be represented as graphs. The first version of the algorithm was characterized by a two step procedure (e.g. discovering and elaboration phases), in which the effect of the nodes’ connectivity on the information spreading was exploited by nodes to assess a first approximation of the topology of the network. In this work, we present a second version of the algorithm in which the third phase was added in order to refine the topology detection by a cognitive inspired strategy which embeds the cognitive dissonance theory [74].

In general, as the Internet nowadays, human social networks have to be considered as a continuum of nested communities whose boundaries are somewhat arbitrary [33].

Here, we propose such a tool for detecting communities in complex networks using a local algorithm, applied as a cellular automaton. In this approximation, a node is just modeled as a memory and a set of links to other nodes. The information about neighbouring nodes is propagated using a standard diffusion process, and elaborated locally using a non-linear competition process among the information. This process can be considered an implementation of the “take the best” heuristic [35], which relies on

the assumption that the most relevant or easily detectable information gives an accurate estimate of the frequency of the related event/contents in the population. The result of the algorithm equips each node with information about possible hubs or super-nodes present in its environment, and such information can be used by the node to rewire its connections whenever its fitness does not satisfy some given requirements.

In real-world applications, such a process can be engineered within the ICT domain. Consider for instance resilience and scalability effects in service ecosystems. There, one important factor is to decentralize services. This can be done with the help of creating overlay networks on top of large-scale ones such as the Internet. An adaptive, intelligent or even resource-optimizing algorithm plays a crucial role for the (self-)maintenance of such systems.

In that way, we could tackle the first steps to create an intelligent, semi-structured peer-to-peer overlay network from an unstructured one, e.g. like a *self-optimizing* FastTrack [75] network. FastTrack itself uses a semi-structured overlay network with a mix of *designated* super-nodes and normal nodes. The latter have to connect to one of the super-nodes in order to minimize redundant communication overhead. There, participating nodes could retrieve content at a certain time with given resource constraints (e.g. bandwidth, energy, latency), detect super-nodes automatically during an operation, and thus change their connections (and therefore the overlay topology) for better conditions.

Analogously for the entropy defined above, it is possible to introduce the concept of local entropy for each node in order to study the local view of agents. Similar as we can observe in Figure 3.12 (d), it is possible to detect different plateaus corresponding to the different network sub-clusters that the single node discovers during time. We observed that we can use a fixed value of the parameter  $m$ , while we have to change the value of  $\alpha$  in order to find the community and in particular the hubs that labels each community. For this reason, we simulate an exploration phase of the network several times in which the nodes save their state vector  $S_i^t$ , in a *temporary memory box*, when they observe a change in sign of the second derivative. If the following condition is satisfied

$$\text{sign} \left( \frac{\delta^2 E_i^{t-1}}{\delta t^2} \right) \neq \text{sign} \left( \frac{\delta^2 E_i^t}{\delta t^2} \right), \quad (3.13)$$

the state vector  $S_i^t$  is stored into the temporary long term memory together with the value of the first derivative of the local entropy and the entropy. When a node meets an impasse (e.g. its state vector entropy and its cognitive dissonance do not evolve anymore) its  $\alpha$  is changed by the following mechanism, if

$$\left| \frac{E_i^{t-1} + D_i^{t-1}}{K_i} \right| - \left| \frac{E_i^t + D_i^t}{K_i} \right| < \epsilon, \quad (3.14)$$

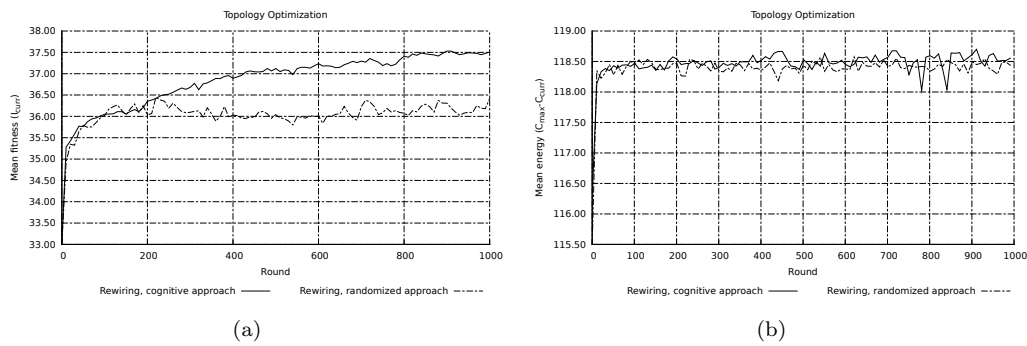


FIGURE 3.26: Fitness function and mean energy.

(a) The fitness function shows the mean number of retrieved items per node. The fitness of the randomized algorithm is represented as the dashed line, our approach as the solid line. (b) Mean energy usage per node of the randomized algorithm (dashed) and our approach (solid).

where  $K_i$  is the connectivity degree of the node  $i$ . Then, a counter  $\tau$  is increased by 1 ( $\tau_i \leftarrow \tau_i + 1$ ), and if  $\tau_i$  becomes greater than a given threshold (say  $\tau^*$ ), the parameter  $\alpha$  is updated in the following way:

$$\alpha_i = 1.5|\eta\sigma| + 1, \quad (3.15)$$

where  $\eta$  is a random Gaussian variable with mean 1 and standard deviation  $\sigma$ . After a typical period of a fixed length ( $\Delta T$ ), the process is stopped for all nodes and a node's long term memory is updated with a new sample respectively experience. The long term memory is characterized by a bound threshold  $B^1$  (here  $B^1 = 5$ ) in order to mimic the ecological limits of such cognitive functions (i.e. bounded rationality). After the node has saved its state vector when the sign of its second entropy derivative changed (Eq. 3.8), it proceeds in structuring its long term memory. First, its first derivatives are decreasingly sorted, and then the first  $B^1$  time positions are recorded. Later, such  $B^1$  element vectors are descendingly sorted with respect to the entropy. Finally, using the time positions, the correspondent state vectors is analyzed and larger elements for each state vector are assumed as potential hubs and therefore stored into the long term memory. At this stage, the long term memory of each node is composed by a list of  $B^1$  sets of potential hubs, ordered following the procedure from the more local to the more global one. Moreover, the long term memory is bounded by another threshold ( $B^2$ ), which represents the long term memory buffer, i.e. the maximum number of the  $B^1$  sets it can consider/contain, so that the long term memory is represented by a  $(B^1, B^2)$  matrix. Finally, each node summarizes its knowledge of the network building a *hub list* obtained by analyzing the frequency in which each hub appears within the long term memory, which is subsequently ordered from the most represented (i.e. the hub with a

larger frequency) to the least represented one. The knowledge of the network (i.e. the hub list) is used by weak nodes in order to increase the fitness.

The nodes' fitness is computed in a general and conservative way following the ratio presented in section II. In the first scenario, the nodes are sorted with respect to the number of objects they collected through their neighbors, while in the second scenario, they are sorted with respect to the amount of energy they spent to collect the maximum number of items. After this phase, the last 9% of the nodes (e.g. the weakest nodes) are chosen for the cognitive rewiring, and in addition 3% of the nodes are chosen for a random rewiring.

Whenever a node does not have a “sufficient” fitness, it eliminates a portion of unnecessary links (i.e. those links which point to nodes detected as non-hubs, in this work just 1 link) and proceeds to try to establish new connections using the hub list it has. Starting from the most relevant hub (i.e. the first from the list) and continuing towards the last one, the rewiring node tries to build new links. Finally, if no hubs have available links, because they have reached the maximum number of connections, the rewiring node adopts a random strategy and establishes a link towards the first available node it finds.

### 3.7.1 Evaluation

We compared our model from section 3.2 with a randomized algorithm. For comparability reasons between the algorithms, they are kept similar, apart that the randomized algorithm is memoryless and therefore nodes have no knowledge about its surrounding and potential hubs they might connect to. Consequently, the randomized algorithm selects the nodes that have to rewire using the same method as the cognitive algorithm does; but where the cognitive algorithm prefers to connect to a hub, the randomized one chooses a random node.

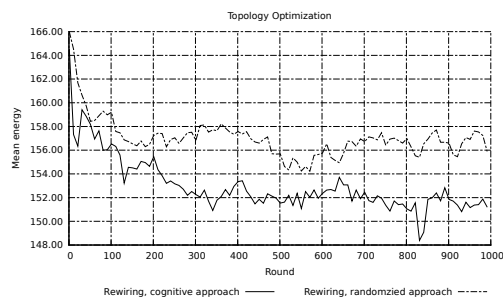


FIGURE 3.27: Energy minimization approach.

Mean energy usage per node of the randomized algorithm (dashed) and our approach (solid).

For the evaluation, we used the two scenarios described before in order to test our algorithms. The initial network topology consists of  $N = 200$  with a mean connectivity per node of 4. A total of  $I = 50$  unique items is distributed among the nodes, where each node needs to retrieve  $I_{max} = 45$  objects from its neighbors. We used this setting in order to analyze a network on a larger scale. Further, we also tested the algorithm for smaller networks, and the results imply a similar behaviour as presented here. We run the simulation 50 times on our Matlab cluster with different random seeds. Figure 3.26 shows an initial result for the first scenario and Figure 3.27 for the second. Both figures show values of the median run regarding final results of fitness and energy.

In the first scenario, the number of retrievable items shall be maximized. Therefore, the “weakest” nodes are determined by the sum of collected items. In Figure 3.26, it is shown that both approaches improve the initial topology significantly at the beginning. After having reached a plateau of 36 items, the randomized approach begins to oscillate, whereas the cognitive approach can exploit its knowledge of potential hubs and steadily micro-optimizes the topology up to more than a *mean* of 1.1 items by not having significant differences in their energy usage. We can also observe that the cognitive approach is less prone to oscillations.

The second scenario shown in Figure 3.27 shows the energy dynamics of both approaches. Each node has unlimited energy available, so that it is able to retrieve all necessary 45 items. The weakest nodes are now defined as nodes who consume the most energy of all. Hence, those are candidates for rewiring in order to minimize the system’s energy. The behaviours of both approaches are quite similar as in the fitness optimization from Figure 3.26. The initial topology improvement significantly reduces the energy consumption of the system. However, oscillation effects occur more often than in the first scenario. Our cognitive approach reduces the *mean* energy consumption of the nodes of more than 4.1 hops per node compared to the randomized algorithm.

### 3.7.2 Final remarks

In this section, we described how we optimize a topology by the means of a cognitive-inspired algorithm. The resulting online optimization problem was tackled with a cognitive model that enables a node to be self-aware about its surrounding community and eventually to detect and distinguish between hubs and non-hubs. This knowledge was exploited by a node to gain a more effective rewiring to other nodes than by random selection. We showed the effectiveness of our approach in two scenarios, in each comparing the achieved results to a randomized algorithm using the same network conditions. In the first scenario, the goal was to find a topology in which a maximum number of unique

items can be retrieved for the system under a given energy constraint that was spent for “hopping”. In the second one, we removed the energy constraint, so that nodes had enough energy for retrieving all items in each round, with the focus on decreasing the system’s overall energy. In both scenarios, the cognitive-inspired algorithm performed significantly better than the random one.

Despite the fact that the algorithm uses global information for the selection of rewiring nodes, the approach shows first steps towards a pure self-organizing network since only local information is used for the hub detection. Overall, we showed first steps that information generated by a cognitive-inspired algorithm can be exploited in order to optimize network topologies. As future work, we plan to (i) deploy the algorithm on a wide range of *large scale* network topologies, (ii) localize the decision making of a node when to rewire or not, and (iii) further elaborate the used scenario by introducing more dynamics into items and nodes. We think that our algorithm is generic enough that it could also be used as a foundation in a wide area of applications beyond the scenario proposed here.

## 3.8 Considerations on more complex heuristics

### 3.8.1 IDA + LTE

Recently, Huang et al. [76] proposed a local algorithm for detecting communities in networks. We report briefly the algorithm procedure. Usually, a network can be represented by a graph  $G = (V, E)$ , where  $V$  is the set of vertices and  $E$  is the set of edges.

**Definition 1** (*Neighborhood*) Let  $G = (V, E, w)$  be a weighted undirected network and  $w(e)$  be the weight of the edge  $e$ . For a vertex  $u \in V$ , the structure neighbourhood of vertex  $u$  is the set  $\Gamma(u)$  containing  $u$  and its adjacent vertices which are incident with a common edge with  $u$  :  $\Gamma(u) = \{v \in V \mid \{u, v\} \in E\} \cup \{u\}$ .

**Definition 2** (*Structural Similarity*) Given a weighted undirected network  $G = (V, E, w)$ , the structure similarity  $s(u, v)$  between two adjacent vertices  $u$  and  $v$  is:

$$s(u, v) = \frac{\sum_{x \in \Gamma(u) \cap \Gamma(v)} w(u, x) \cdot w(v, x)}{\sqrt{\sum_{x \in \Gamma(u)} w^2(u, x)} \cdot \sqrt{\sum_{x \in \Gamma(v)} w^2(v, x)}}. \quad (3.16)$$



When we consider an unweighted graph, the weight  $w(u, v)$  of any edge  $\{u, v\} \in E$  can be set to 1 and the equation above can be transformed to:

$$s(u, v) = \frac{|\Gamma(u) \cap \Gamma(v)|}{\sqrt{|\Gamma(u)| \cdot |\Gamma(v)|}}. \quad (3.17)$$

It corresponds to the so-called edge-clustering coefficient introduced by Radicchi et al. [77].

**Definition 3 (Tightness)** *By employing the structural similarity, we introduce tightness, a new quality function of a local community  $C$ , which is given as follows:*

$$T(C) = \frac{S_{in}^C}{S_{in}^C + S_{out}^C} \quad (3.18)$$

where  $S_{in}^C = \sum_{u \in C, v \in C, \{u, v\} \in E} s(u, v)$  is the internal similarity of the community  $C$  which is equal to two times of the sum of similarities between any two adjacent vertices both inside the community  $C$ ;

$S_{out}^C = \sum_{u \in C, v \in N, \{u, v\} \in E} s(u, v)$  is the external similarity of the community  $C$  which is equal to the sum of similarities between vertices inside the community  $C$  and vertices out of it.

The tightness measure is extended from the weak community definition proposed by F. Radicchi [77]. Similar to other community definitions, the tightness value of a community  $C$ , denoted by  $T(C)$ , will increase when sub-graph  $C$  has high internal similarity and low external similarity. The whole network without outward edges will achieve the maximal value 1, but the problem here is to find the local optimization of the measurement for each community. Suppose a community  $C$  is detected from a certain vertex  $s$ . We explore the adjacent vertices in the neighborhood set  $N$  of  $S$ . So the variant tightness of the community  $C \cup \{A\}$  becomes

$$\begin{aligned} T(C \cup \{A\}) &= \frac{S_{in}^C + 2S_{in}^a}{(S_{in}^C + S_{in}^a) + (S_{out}^C - S_{in}^a + S_{out}^a)} = \\ &= \frac{S_{in}^C + 2S_{in}^a}{(S_{in}^C + S_{in}^a + S_{out}^C + S_{out}^a)}, \end{aligned} \quad (3.19)$$

where  $S_{in}^a = \sum_{\{v,a\} \in E \wedge v \in C} s(v,a)$ ;  $S_{out}^a = \sum_{\{a,u\} \in E \wedge u \notin C} s(a,u)$ . Then the tightness increment of a vertex  $a$  joining in  $C$  is:

$$\begin{aligned} \Delta T_C(A) &= T(C \cup \{A\}) = \\ &= \frac{S_{in}^C + 2S_{in}^a}{(S_{in}^C + S_{in}^a + S_{out}^C + S_{out}^a)} - \frac{S_{in}^C}{S_{in}^C + S_{out}^C} = \\ &= \frac{2S_{in}^a \cdot S_{out}^C - S_{in}^C \cdot S_{out}^a + S_{in}^C \cdot S_{in}^a}{(S_{in}^C + S_{in}^a + S_{out}^C + S_{out}^a)(S_{in}^C + S_{out}^C)} \end{aligned} \quad (3.20)$$

If  $\Delta T_C(A) > 0$ , then  $2S_{in}^a \cdot S_{out}^C - S_{in}^C \cdot S_{out}^a + S_{in}^C \cdot S_{in}^a > 0$  which is equivalent to  $\frac{S_{out}^C}{S_{out}^a} - \frac{S_{out}^a - S_{in}^a}{2S_{in}^a}$ . Then they define the tightness gain in the following [76].

**Definition 4 (Tightness Gain)** The tightness gain for the community  $C$  adopting a neighbor vertex  $a$  can be denoted as

$$\tau_C(A) = \frac{S_{out}^C}{S_{in}^C} - \frac{S_{out}^a - S_{in}^a}{2S_{in}^a}. \quad (3.21)$$

It means that the ratio of external similarity to internal similarity of community  $C$  is greater than the ratio of external similarity increment to internal similarity increment caused by adopting vertex  $a$ . Obviously, this case will result in the increase of the tightness value of community  $C$ . Therefore,  $\tau_C(a)$  can be utilized as a criterion to determine whether the candidate vertex  $a$  should be included in the community  $C$  or not. In the following, they introduce an optional resolution parameter  $\alpha$  to control the scale at which we want to observe the communities in a network.

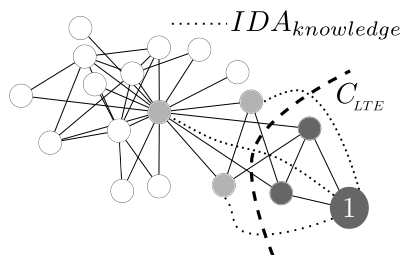


FIGURE 3.28: Schematic representation of IDA + LTE algorithm.

We suppose that  $LTE$  detect for the node 1 the dark-gray nodes as community nodes ( $C_{LTE}$ ): after that we compute the  $IDA$  on the selected nodes. The node takes the information from those accepted as its community creating a virtual link with these (dotted line in figure). The  $LTE$  decides whether to accept the new nodes members of own community or not.

**Definition 5** (*Tunable Tightness Gain*) The tunable tightness gain for the community  $C$  merging a neighbor vertex  $a$  can be denoted as

$$\tau_C^\alpha(A) = \frac{S_{out}^C}{S_{in}^C} - \frac{\alpha S_{out}^a - S_{in}^a}{2S_{in}^a}. \quad (3.22)$$

A parameter  $\alpha \in (0, \infty)$  is introduced as the coefficient of  $S_{out}^a$  which can increase or decrease the proportion of the external similarity of the candidate vertex  $a$ . Here, the criterion for accepting a vertex  $a$  is changed to  $\tau_C^\alpha(A) > 0$ . For  $\alpha = 1$ , the criteria is moderate and can be used in most normal cases. In [76] the authors shows different scenarios for different values of the free parameter  $\alpha$ : setting  $\alpha \in (0, 1)$ , the value of  $S_{out}^a$  is reduced by this coefficient which increases the chance of a candidate vertex  $a$  joining  $C$  and bigger communities will be formed compared to the normal case with  $\alpha = 1$ . On the contrary, it will result in the formation of smaller communities in a network when we set  $\alpha > 1$ . We report the local algorithm proposed by Huang et al. [76] in the following:

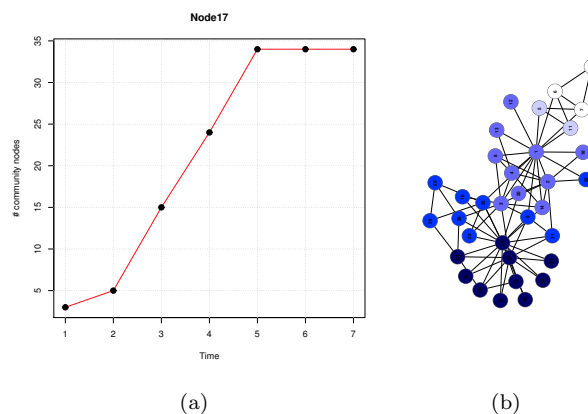


FIGURE 3.29: IDA+LTE algorithm on the Zachary's Karate Club network.

(a) Node 17 in the *Zachary's Karate Club*. In the first step the LTE algorithm accept 3 nodes for the local community. After that we compute the IDA + LTE algorithm: the knowledge of the other nodes permits to discover the structure of the whole network. (b) The color of nodes mark communities of the node 17 during the time. At time  $T=1$  the community nodes are 17-6-7 (white nodes). After that the node is capable to discover other nodes in the network (blue-scale color). At  $T = 2$  the community is detected by nodes 17, 6, 7 plus the nodes 11 and 5. We report the knowledge of the network as communities during the time: starting from the white community to arrive to dark-blue nodes.

1. Pick a vertex  $s \in V$  as the starting vertex.

$$\text{Let } C = \{S\} \text{ and } N = \Gamma(S) - \{S\}$$

2. Select the vertex  $a \in N$  that possess the largest similarity with vertices in  $C$ .

3. If  $\tau_{C^\alpha}(A) > 0$ ,  
 set  $C = C \cup \{a\}$  and  $N = N \cup \Gamma(a) - C$
4. Repeat 2 and 3 until  $N = \emptyset$ .

This algorithm, called local tightness expansion algorithm (LTE), is a very efficient local algorithm for detecting communities in networks from an individual point of view. Anyway when  $N = \emptyset$  we arrive at a static convergence. Here we want to introduce a new algorithm merging together our algorithm that we call Information Dynamics Algorithm (IDA) and LTE for analysing how the information dynamics can improve the performance of LTE and viceversa. We describe the new algorithm in the following (its representation is illustrated in Figure 3.28):

1. Pick a vertex  $s \in V$ .
2. Run LTE and discover  $C_{LTE}$  for the vertex  $s$ .
3. Run IDA and take information of other nodes by community nodes
4. Run LTE and decide to accept new nodes.
5. Repeat 3 and 4 until a fixed memory or when the vertex  $s$  knows the whole network.

For testing purposes we use four real networks analysing and discussing our model peculiarities. The four case studies, of growing or different complexity, are the Zachary *karate club* network [1], the bottlenose dolphins network [41], the network of social interaction in the novel *Les Misérables* by Victor Hugo [78] and the NCAA college football network [41]. For all the simulations we used always the same parameters. For the IDA we assumed  $\alpha = 1.4$  and  $m = 0.2$ . For the LTE we used  $\alpha = 0.3$  for the first step and  $\alpha = 2$  during LTE+IDA dynamics.

### Zachary's karate club

The first test case is the typical network benchmark used in literature: the network proposed by Zachary in the 1977, and known as "karate club" [1]. We studied the dynamics of node 17 during the time. Results are reported in the Figure 3.29 (a) and Figure 3.29 (b). The increasing of information is reported in the Figure 3.29 (a): here we can observe that the  $C_{LTE}$  is a community composed by 3 nodes, the white nodes in the Figure 3.29 (a). After that the node asks to nodes 6 and 7 their knowledge about the networks giving their informations by IDA; the node 17 can decide to accept the new nodes as community nodes through LTE algorithm. At  $T = 2$  it accepts only the node

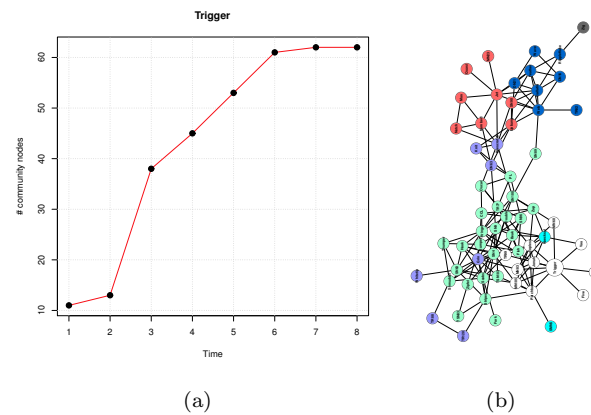


FIGURE 3.30: Trigger in the *Lusseau's network of bottlenose dolphins*.

(a) LTE algorithm in the first step discover 11 nodes for the community. After that we show the increase of the community size trough information dynamics.(b) Nodes' color indicates the community of Trigger during the time. At time  $T = 1$  the community nodes are the white nodes. After that the node is capable to discover other nodes in the network. At  $T = 2$  cyan nodes. At  $T = 3$  green nodes. At  $T = 4$  violet nodes. At  $T = 5$  red nodes. At  $T = 6$  it discovers the whole network with the dark blue nodes except of Zip (gray node) that it knows at  $T = 7$ .

11 and 5 so the community is formed by 5 nodes, and so on. At  $T = 3$  we can observe that the network is splitted in 2 communities that is the same observed by Zachary [1] except only for the node number 22.

### Bottlenose dolphin network

The second real-world network represents the social interactions of bottlenose dolphins living in Doubtful Sound, New Zealand. The network was studied by the biologist David Lusseau, over a period of seven years from 1994 to 2001, who divided the dolphins into two groups according to their age [41]. He we observed the information and community detection dynamics of *Trigger*: results are shown in the Figure 3.30 (a) and in the Figure 3.30 (b). The  $C_{LTE}$  at  $T = 1$  is a community composed by 11 nodes (Figure 3.30 (a)) labeled by white nodes in the Figure 3.30 (b). If we consider *DN63*, *Beescratch* and *Knit* as overlapping nodes between the two principal communities [79], we can say that at  $T = 4$ , *Trigger* discovers the same community structure observed by Lusseau [41].

### Social interaction in the novel *Les Miserables* by Victor Hugo

Results are shown in the Figure 3.31 (a) and in the Figure 3.31 (b).

## NCAA football college network

Results are shown in the Figure 3.32 (a) and in the Figure 3.32 (b).

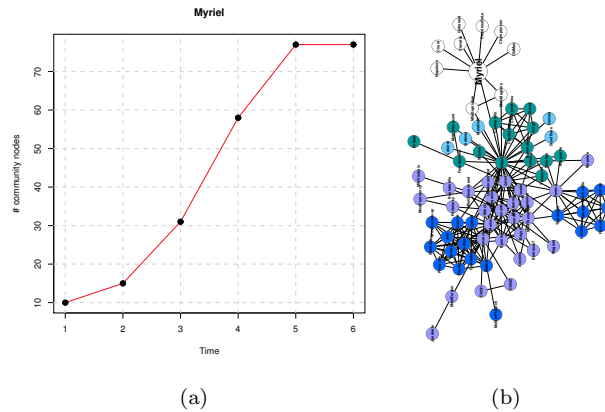


FIGURE 3.31: Myriel in the novel *Les Misérables* by Victor Hugo.

(a) LTE algorithm in the first step discover 10 nodes for the community. After that we show the increase of the community size through information dynamics. (b) Nodes' color indicate the community of *Myriel* during the time. At time  $T = 1$  the community nodes are the white nodes. At  $T = 2$  cyan nodes. At  $T = 3$  green nodes. At  $T = 4$  violet nodes. Finally the blue nodes.

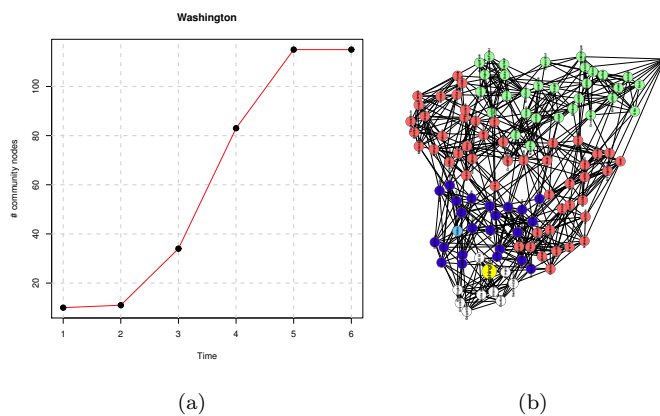


FIGURE 3.32: (a) Washington in NCAA football college network.

(b) Nodes' color indicate the community of *Washington* during the time. At time  $T = 1$  the community nodes are the white nodes. At  $T = 2$  the only light blue node. At  $T = 3$  blue nodes. At  $T = 4$  red nodes. Finally the green nodes.

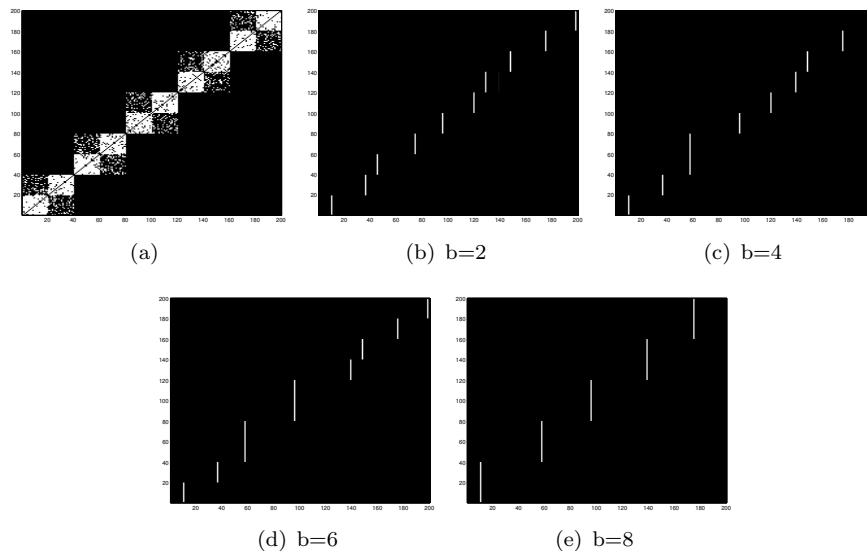


FIGURE 3.33: Adjacency matrix of a network composed by 200 nodes and 3 levels where  $p_1 = 0.9$ ,  $p_2 = 0.2$  and  $p_3 = 0$ .

Trough the double pruning heuristic it is possible to detect the different community levels reaching the discover of the four principal communities for  $b = 8$ .

### 3.8.2 Double pruning

The model has two free parameters  $m$  and  $\alpha$  and in previous works we have shown that it is very difficult to have good values of these parameters for general and different cases. This procedure is very simple: in particular at some moments for each node we evaluate the histogram of the state vector with 4 bins and we set to zero the values that are lower then the second bin values as reported in the Algorithm 2. In this way we are able to generate different view of the clustering levels on hierarchical community structure networks as shown in Figure 3.33 and in Figure 3.34. We use the Lancichinetti-Fortunato-Radicchi (LFR) benchmark graph [46, 80] to evaluate the accuracy of this method. We adopt the normalized mutual information (NMI) to evaluate the quality of detected communities which is currently widely used in measuring the performance of graph clustering algorithms [46]. The accuracy of our method is compared with three others well-known community detection algorithms. The observed results on our hierarchical networks and on the Zachary karate club network as shown in Figure 3.34 allow us to state that for  $b = 4$  we are able to detect the principal communities. For this reason for comparing our method with others we set  $b = 4$ . In order to evaluate and partially validate our approach, we have applied our algorithm comparing its performance with 3 community detection algorithms namely *Infomap* [44], *Infomod* [42] and *MCL* [43]. The input parameters of the benchmark graphs used here are: average degree is 20, the maximum degree 50 for all the networks. While we changed the range

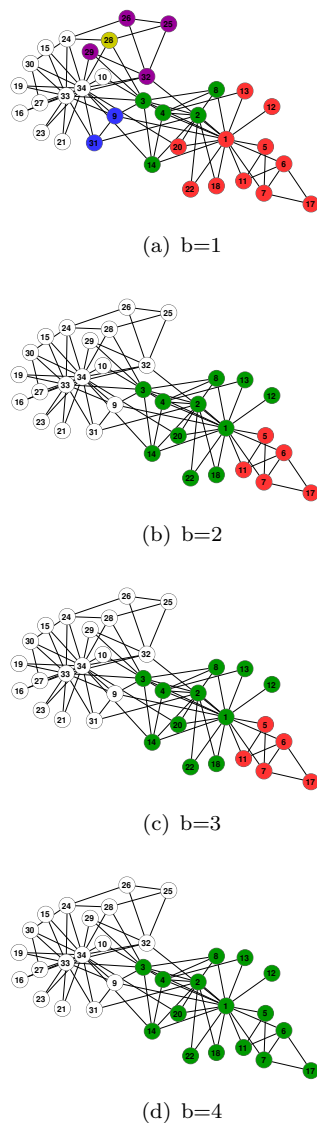
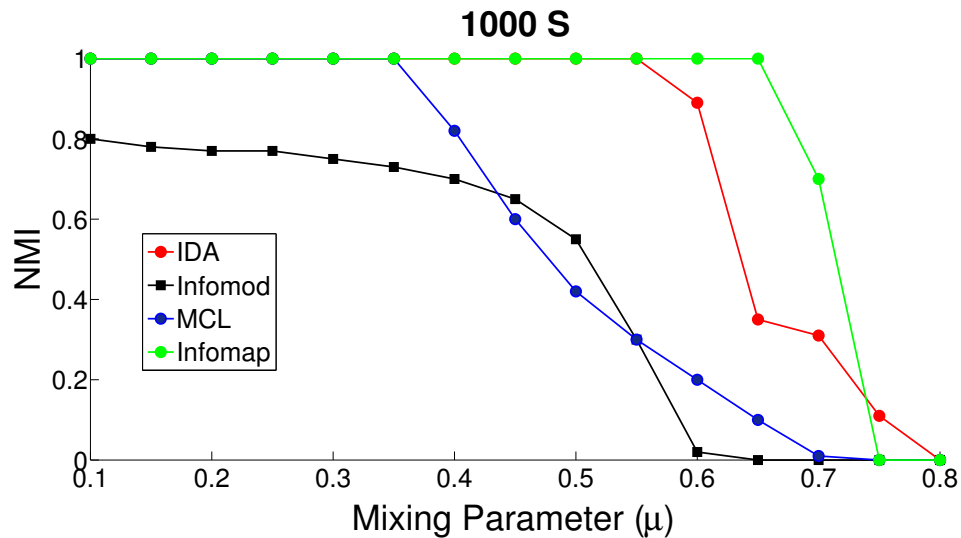


FIGURE 3.34: Evolution of the double pruning heuristic on the Zachary Karate Club network [1].

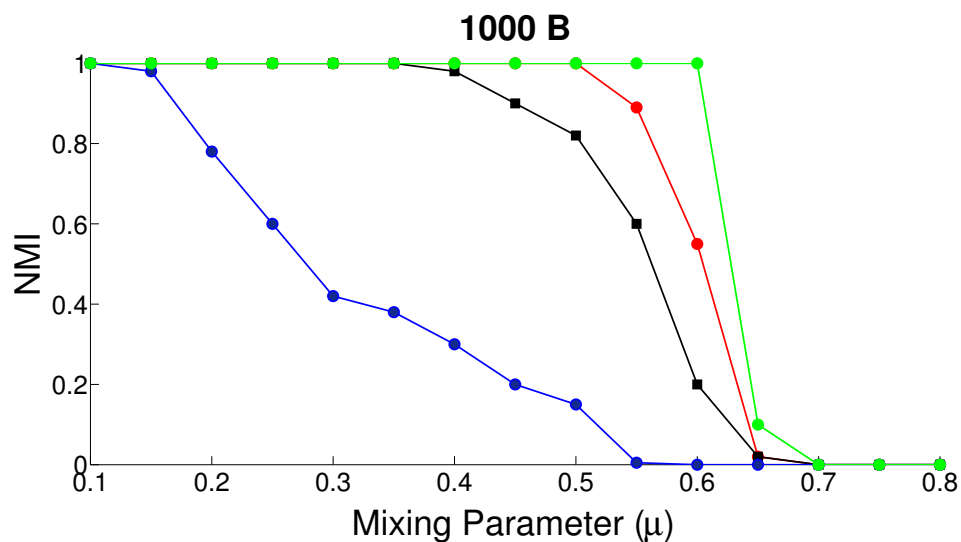
We show that the method is capable to discover different community structure but we reach the observed two communities for  $b = 4$ . Different colors correspond to different communities.

of community size generating two kind of networks of 1000 nodes: 1000*S* (*S* stay for small) means that communities have between 10 and 50 nodes and 1000*B* (*B* stay for big) means that communities have between 20 and 100 nodes. Results of the performance's comparison between our algorithm and the others are reported in Figure 3.35. As we can observe our algorithm with the double pruning heuristic is very competitive with the other algorithms except for the Infomap method which is nowadays the best algorithm for detecting communities in static network even if it can not be easily applied in dynamic environments. Moreover we can observe that our method improves the performance of the MCL algorithm which was our starting point.





(a)



(b)

FIGURE 3.35: Normalized Mutual Information on LFR benchmarks.

(a)-(b) Comparison between our method and other algorithms on LFR Benchmark graphs with different size of communities respectively Small (a) and Big (b) increasing the mixing parameter  $\mu$ .

### 3.9 Impact of local information in growing networks

The emergence and the global adaptation of social networks has influenced human interaction on individual, community, and larger social levels. Most notably, perhaps, is the rise of Facebook, which in October 2011 reached more than half (55%) of the world's global audiences, catching 835.6 millions of users in 2012 [81]. A large number

of models have been proposed that aim at exploring and explaining how local mechanisms of network formation produce global network structure. In the context of social networks it is important to understand why and how people decide to make connections and how they change or modify their own local structure. For this reason it is essential to understand some aspects of how humans behave in social networks: How do people acquire information in on-line social networks and what are the mechanisms that lead people to join together or to visit a specific website?

One of the the most well known mechanism that is used in growing networks is preferential attachment, where new connections are established preferentially to more popular nodes in a network, giving rise to a scale-free network [82]. Moreover, users in on-line social networks tend to form groups, called communities: given a graph, a community is a group of vertices “more linked” among them than between the group and the rest of the graph [45]. This is clearly a poor definition, and indeed, on a connected graph, there is no clear distinction between a community and the rest of the graph. In general, there is a continuum of nested communities whose boundaries are somewhat arbitrary: the structure of communities can be seen as a hierarchical dendrogram [30]. Our communities are large and varied, and we recognize several levels of grouping, sometimes dependent on the context. In recent work we have shown that using information dynamics algorithms where nodes elaborate information locally, we are able to detect such communities in complex networks [58, 59].

Recently, Papadopoulos et al. [83] explored the trade-off between popularity and similarity in growing networks. Nodes in growing networks tend to link not only to the most popular nodes (as in preferential attachment [6]) but also to the closest nodes in terms of affinity. Comparing their results with real-world complex networks, the authors showed that they were able to predict the probability of forming new links with remarkable precision.

In this section we develop a model that emulates the growing of a social network, starting from psychological assumptions that allow us to simulate how people acquire and elaborate information in social networks. We demonstrate the concept of similarity and popularity in growing networks, not by a geometric approach as in Papadopoulos et al. [83], but by using a simple mechanism that explain users’ behaviour in on-line social networks.

The rest of this section is organized as follows: we start by summarizing previous work in section 3.9.1. In section 3.9.2 we describe our model, which uses a local algorithm where an agent is modeled with a memory and a set of connections to other individuals. In the first step the new agent explore the local structure of the network, where it receives information about the neighborhood. The learning (nonlinear) phase is modeled after

competition in the chemical/ecological world, where agents compete with each other. Section 3.9.3 shows the principal results of our simulations. Finally, we discuss our results and propose future steps in the Conclusions.

### 3.9.1 Related work

The goal of much of the research that model the growth of real networks is to reproduce networks with certain properties as well as properties of real-world networks. For instance, we know that many observed networks fall into the class of scale-free networks, meaning that they have power-law (or scale-free) degree distributions. An influential model is the so called Barabási–Albert (BA) model [82]. The main hypothesis of the BA model is that the more connected a node is, the more likely it is to receive new links. From the BA model, however, it is not trivial to generate networks with a community structure or with a high clustering coefficient, something we observe in real social networks. Regarding social networks, Jin et al. [84] presented a model where the friendship between individuals depends on the number of mutual friends and the number of meetings between them. The resulting networks from their simulations show high levels of clustering and a strong community structure in which individuals have more links to others within their community than to individuals from other communities.

In these two models we can already single out two important and psychologically plausible features: (1) the predisposition of people to link with hubs (nodes with higher connectivity degree, where the connectivity degree,  $C(k)$ , of a node in a unweighted network is defined as the sum of its links) and (2) the tendency of people to connect with friends of friends (social community).

A recent paper highlighted the relevance of existing communities in a network [85]. Here, a new node connects to a random node in the network and with a probability  $p$  it links to some nodes in the community of the selected node and with probability  $1 - p$  with some random nodes. Simulation results showed that this model can reproduce features of real world networks, but the authors used a global community detection algorithm, the well-known Louvain Method [75], without considering the local and subjective view of agents in the network.

Popularity is not the only thing that determines if a node will link to another node, social closeness is also important. It has been shown that more similar nodes have higher chances to connect to each other even if they are not popular [86]: this effect is known as *homophily* in social science. As mentioned above, in 2012 Papadopoulos et al. [83] published a paper giving an important contribution to the understanding of the evolutionary properties of complex networks. The authors showed that the popularity

of a node in a network is just one dimension of attractiveness. In their framework, the probability that a new node that arrives in a structured network links to another node is a function of two variables: the connectivity degree of the target node (popularity) and the similarity (affinity) between the new node and the target node. In order to evaluate the trade-off between popularity and similarity they exploited a geometric representation. They place nodes in circles whose distance from the origin depends on the birth time while the angular position define the *social identity* of nodes. Initially the network is empty. At each time  $t > 1$ , a new node labelled  $t$  appears at a random angular position  $\theta$  on the circle, with polar coordinates  $(r_t, \theta_t)$ . In order to implement the trade-off between popularity and similarity, it is assumed that the node evaluates its hyperbolic distance from other nodes with label  $s$  ( $s < t$ ), with coordinates  $(r_s, \theta_s)$ , by means of the function  $x_{st} = r_s + r_t + \ln(\theta_{st}/2)$ , where  $\theta_{st}$  is the relative angle of the two polar coordinates. The node finally connects to the  $m$  nodes with the smallest hyperbolic distance. This approach is interesting, but the mechanism appear rather artificial. In our model, we try to take into account all of these features by considering a simple information dynamics algorithm that allow to us to model not only the preferential attachment mechanism but also the social closeness between individuals.

### 3.9.2 The model

The model is based on a mechanism that emulates people's strategies for acquiring information in social networks, emphasising the local subjective view of an individual and what kind of information the individual can receive when arriving in a new social context.

We start from a simple assumption: we suppose that a new individual, or node, arrives in an already structured network and picks information about it. According to this idea we begin, at time  $t = 0$ , with a fully-connected network of  $N_0$  vertexes and then, at each time step, a new node is added to the network. For explaining our idea we describe a simple case, illustrated in Figure 3.36 (a), representing a new node  $i$  that joins the *Zachary's Karate Club* network [1] (assuming that the node  $i$  is invited by node 17). It discovers  $n$  ordered levels: the nodes in level 1 are those adjacent to the connected node (node 17 in Figure 3.36 (a)), nodes in level 2 are the friends of my friends, and so on. With different probabilities  $p_1 > p_2 > \dots > p_n$  it links to nodes in the different levels. This mechanism implements the assumption that the probability of getting new friends in a social context is strictly correlated with the local structure of the network; it is easier that two persons become friends if they a have a friend in common. Here we

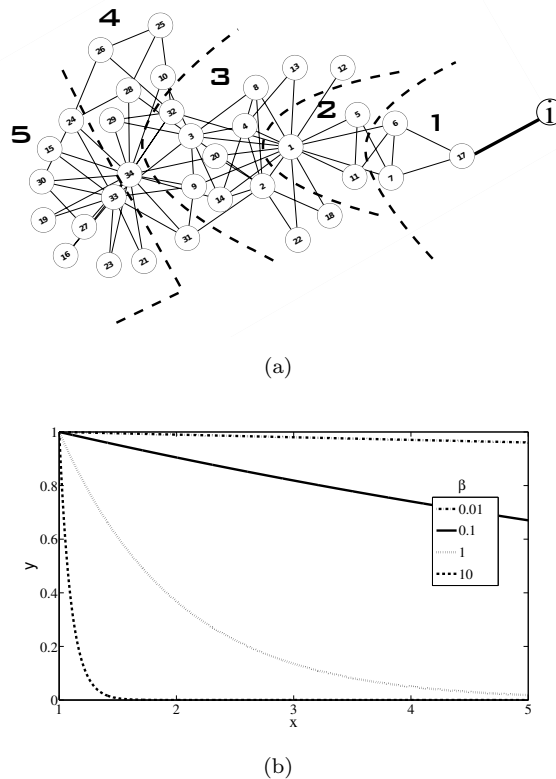


FIGURE 3.36: Schematic description of the method.

(a) Starting configuration of the *Zachary's Karate Club* network [1]: the new node  $i$  links to node 17. The new node has the local hierarchical representation of the network, labelled by levels 1, 2, 3, 4 and 5. With probability  $p_1$  it could link to nodes of level 1, with probability  $p_2$  to nodes of level 2, and so on. (b) Different probability functions derived from Equation 3.23 for different values of parameter  $\beta$  and  $a = 1$ . The  $y$  - axis represents the probability to join with some nodes in the corresponding level shown on the  $x$  - axis.

use a simple exponential function

$$y = ae^{-\beta(x-1)}, \quad (3.23)$$

where  $x$  is the considered level. The  $\beta$  represent the “temperature”: the probability of joining farther nodes, and  $a$  is a normalization constant. In the Figure 3.36 (b) we show the function Eq. ref:exp for some values of  $\beta$  with  $a = 1$ . Assuming  $-\beta(x-1) = z$  and  $e^z = b^z$ , we can express the probability distribution of Eq. 3.23 as  $\sum_{z=0}^{\infty} ab^z = a \frac{1}{b-1}$ . Setting the previous equation equal to 1 we obtain  $a = 1 - b$ .

Then,  $y = (1 - b)b^z = (1 - e^{-\beta})e^{-\beta(x-1)} = P(x)$ . For high values of parameter  $\beta$  the probability to join other levels is very low (e.g., the continuous line in Figure 3.36 (b)).

The second part of the algorithm allows new nodes to locally elaborate the information about the nodes belonging to a given level. The network is represented by the adjacency

matrix  $A_{ij} = 1$  (0), 1 if there is a link between  $i$  and  $j$ , 0 otherwise.

Each individual  $i$  is characterized by a knowledge vector  $S^{(i)}$ , representing his knowledge of the world. The knowledge vector  $S^{(i)}$  is a probability distribution, assuming that  $S_j^{(i)}$  is the probability that individual  $i$  knows about the community  $j$ . It can also be seen as the probability that  $i$  belongs to the community "leaded" by  $j$ , and therefore,  $S_j^{(i)}$  is normalized over the index  $j$ . In order to use a compact notation, we arrange the knowledge vectors for all individuals column by column as  $S_{ij} = S_j^{(i)}$ , forming a knowledge matrix  $S = S(t)$  of the whole network at time  $t$ . We initialize the system by setting  $S_{ij}(0) = \delta_{ij}$ , where  $\delta$  is the Kronecker delta,  $\delta_{ij} = 1$  if  $i = j$  and zero otherwise. In other words, at time 0 each node knows only about itself.

The dynamics of the network is given by an alternation of communication and elaboration phases. The communication is implemented as a simple diffusion process, with memory  $m$ . The memory parameter  $m$  allows us to introduce some important features of the human cognitive system, for example that recently acquired information have more relevance than information gained in the past [39, 40].

In the communication phase, the state of the system evolves as

$$S_{ij}(t + 1/2) = mS_{ij}(t) + (1 - m) \sum_k A_{ik} S_{kj}(t), \quad (3.24)$$

where  $A$  is the adjacency matrix. We assume that nodes talk with each other and we suppose that nodes with high connectivity degree have greater influence in the process of information's diffusion. This is due to the fact that during a conversation it is more likely to know a vertex with high degree instead of one that has few links. For this reason, the information dynamics is a function of the adjacency matrix  $A$ .

The elaboration phase implements elements of fast and frugal heuristics [35]. When people are asked to take a decision, very rarely do they weight all available pieces of information. If there is some aspect that has a higher importance than others, and one item exhibits it, than the decision is taken, otherwise, the second most important factor is considered, etc. In order to implement an adaptive scheme, we exploit a similarity with competition dynamics among species.

If two populations  $x$  and  $y$  are in competition for a given resource, their total abundance is limited. After normalization, we can assume  $x + y = 1$ , i.e.,  $x$  and  $y$  are the frequency of the two species, and  $y = 1 - x$ . The reproduction phase is given by  $x' = f(x)$ , which we assume to be represented by a power  $x' = x^\alpha$ . For instance,  $\alpha = 2$  models birth of individuals of a new generation after binary encounters of individuals belonging to the old generation, with non-overlapping generations [36].

After normalization  $x' = \frac{x^\alpha}{x^\alpha + y^\alpha} = \frac{x^\alpha}{x^\alpha + (1-x)^\alpha}$ . and introducing  $z = (1/x) - 1$  ( $0 \leq z < \infty$ ), we get the map  $z(t+1) = z^\alpha(t)$ , whose fixed points (for  $\alpha > 1$ ) are 0 and  $\infty$  (stable attractors) and 1 (unstable), which separates the basins of the two attractors. Thus, the initial value of  $x$ ,  $x_0$ , determines the asymptotic value, for  $0 \leq x < 1/2$   $x(t \rightarrow \infty) = 0$ , and for  $1/2 < x < 1$   $x(t \rightarrow \infty) = 1$ . By extending to a larger number of components for a probability distribution  $S^{(i)}$ , the competition dynamics becomes

$$S_{ij}(t+1) = \frac{S_{ij}(t+1/2)^\alpha}{\sum_k S_{ik}^\alpha(t+1/2)}, \quad (3.25)$$

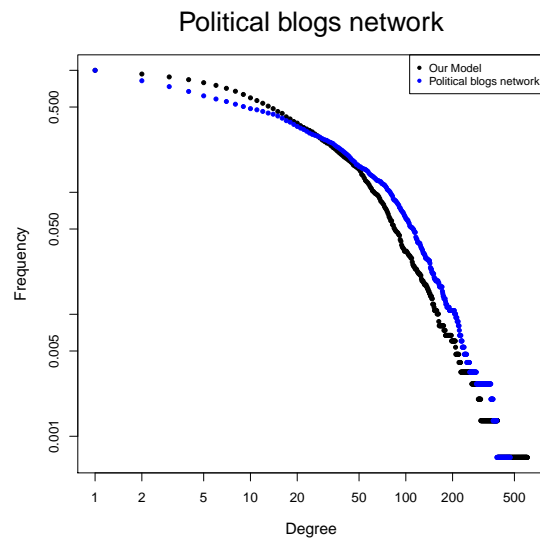
and the iteration of this mapping, for  $\alpha > 1$ , leads to a Kronecker delta, corresponding to the largest component. The parameter  $\alpha$  allows us to model a “pruning effect” of the information, which eliminates unnecessary clutter and a clears the way for more information to enter the field of view of the individuals [87]. The convergence time depends on the relative differences among the components and therefore, when coupled with the information propagation phase, it can produce interesting behaviours. The model has two free parameters, the memory  $m$  and the exponent  $\alpha$ .

Finally, the probability of making a new link ( $P_n$ ) depends on the joint probabilities of two functions  $f(y)$  and  $g(S)$ :  $P_n = f(y) \cdot g(S)$ .

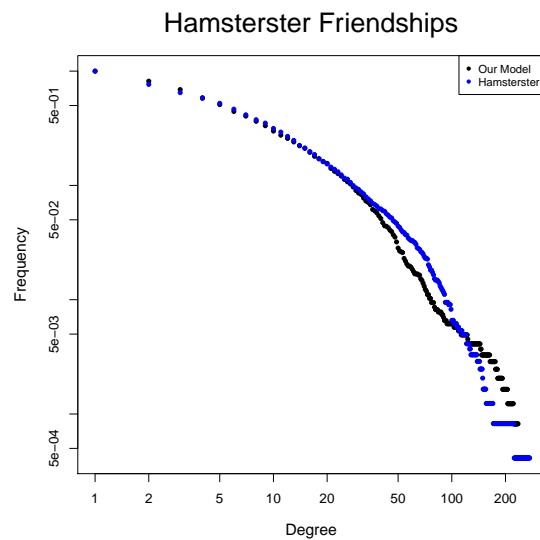
In summary: (1) we start with  $m_0$  nodes ( $m_0 \geq 1$ ); (2) at time  $t$  a new node, labelled by  $t$ , appears in the network; (3) the new node connects with a random node in the network, discovering  $n$  levels; with probability  $p_1$  given by Eq. 3.23 it joins the selected level; (4) the new node links with probability  $p_2$  given by the Eq. 3.25 to the level's nodes. In this way we take into account the social closeness because of the probability to link to nodes in the network depends on the social distance from my *closest* friend *and* the popularity of nodes given by the information dynamics procedure.

TABLE 3.2: Results from the information dynamics algorithm for the different levels in the network shown in Figure 3.36 (a).  $L$  is the number of the level,  $n$  is the id of the node and  $p$  is the probability to join with the most connected node (bold node).

L	n	p
1	6, 7	0.5-0.5
2	<b>1</b> , 5, 11	0.97
3	2, <b>3</b> , 4, 8, 9, 12, 13, 14, 18, 20, 22, 32	0.35
4	10, 25, 26, 28, 29, 31, 33 <b>34</b>	0.67
5	15, 16, 19, 21, 23, <b>24</b> , 27, 30	0.47



(a)



(b)

FIGURE 3.37: Political Blogs and hamsterster.com networks.

(a) Cumulative frequency of node degree distributions ( $\log - \log$  scale) of Political blogs network (blue points) and model generated predictions (black points). (b) Cumulative frequency of node degree distributions ( $\log - \log$  scale) of the network of friendships between users of the website hamsterster.com and model generated predictions (black points). Model predictions are averaged over 10 simulation runs.

### 3.9.3 Results

Results obtained with the information dynamics algorithm, applied to the network represented in Figure 3.36 (a), are shown in Table 3.2. It is more likely that a new node



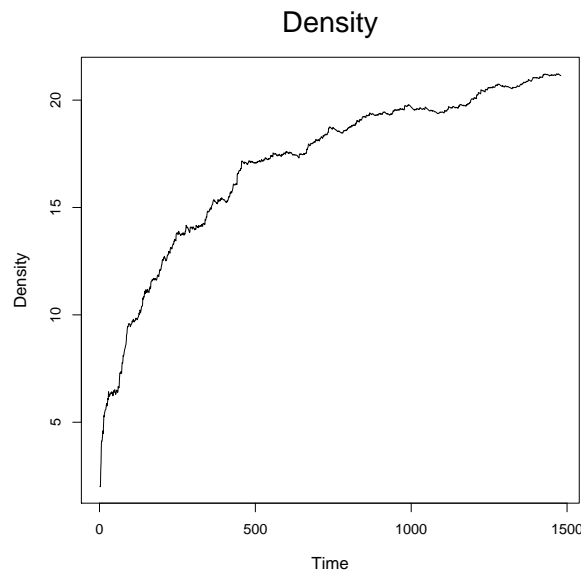


FIGURE 3.38: Simulation of the temporal density evolution of the *Political blogs network*.

TABLE 3.3: Statistics of the social networks.  $C$  (mean clustering coefficient),  $l$  (average path length) and  $d$  (diameter of the network).

	No.vertices	No.edges	$C$	$l$	$d$
Political blogs	1490	19090	0.24	3.39	9
Simulation	1490		0.24	3.23	9
Hamsterster.com	2426	12534	0.09	3.54	10
Simulation	2426		0.09	3.84	11

will be connected to a node with a high degree, as is also predicted by preferential attachment.

In order to validate our model we compare predictions from the model with two real networks. The model predictions are averages over ten simulation runs. The first social network is the *Political blogs network*. It is a directed network of hyperlinks, with  $N = 1490$  nodes, among weblogs on US politics, recorded in 2005 by Adamic and Glance [88]. The degree distribution of the resulting network from a simulation of the model with the same number of nodes is comparable with the real network as shown in Figure 3.37(a). The mean clustering coefficient of the network generated with our simulations is  $C = 0.24$ , the average path length is  $l = 3.23$ , and we obtained a network with the same diameter,  $d = 9$ . The second social network contains friendships between users of the website hamsterster.com [89]. The degree distribution of the network is shown in Figure 3.37(b). The cumulative frequency of the degree distribution is very similar to the network generated with the model (black points in Figure 3.37(b)). The mean clustering coefficient  $C = 0.09$ , the average path length is  $l = 3.84$ , and the

diameter  $d = 11$ . Another important feature of the model is that the final structure of the network spontaneously arises without any constraint on the degree of the new node. In fact in our simulations we don't assume any constraints on the number of links that the new node can do: theoretically the new node can link with all the nodes in the network. In Figure 3.38 we show a simulation of the temporal evolution of network density, defined as the average number of links connected to each node, for the *Political blogs network*. From this result we show that the number of links of incoming node increase over time: in fact the number of link that the new node can do depends on the number of nodes in the network (for instance in a network of 10 nodes the probability to make a certain number of link is less than in a network of 1000 nodes). It is a well known result that the density increases in networks that grow over time.

### 3.9.4 Conclusions

In this section, we introduced a new model of growing complex networks. The model is based on the idea that local structure plays a fundamental role in social networks and may be involved also in the growing process of the network itself. Our model reproduces the main features observed in real networks, such as high clustering coefficient, low characteristic path length, strong division in communities and variability of degree distributions.

Following these encouraging results, future work will compare this model with the model proposed by Papadopoulos et al. [83] and validate new results with larger real world networks. Moreover, we plan to derive analytical predictions from the model in order to fix a priori the model's parameters for forecasting some graph's properties (eg. degree distribution, clustering coefficient, diameter, etc.).

### 3.10 Final remarks

Here we have described an algorithm to identify the communities structures in a network from a local point of view. The method is based on pure information propagation where the nonlinear part of these method, responsible for the actual elaboration of information, is inspired by a chemical/ecological competition model [36]. There is not a unique definition of a community, so an exploratory algorithm, like the one that humans have presumably developed during their evolution, should present different clustering for different values of the parameters, or for different iterations. In this implementation we adopted a frequency-based approach and an unbounded memory at the level of nodes. Unbounded memory means that the node's state vector  $S_i$  has not been limited and it

could potentially reach a size equal to the network size  $N$ . Despite this because of the explanation and normalization phases are sufficient to avoid this problem. Nevertheless it will be very important to limit the computational resources of the node explicitly, as suggested by Simon in 1955 [73], so increasing both the ecological plausibility of the model and the insights which drive the algorithm design. The results that we have obtained are promising. The method under investigation is not competitive with respect to others (see the review [90]), but it provides a natural “scanning” of the various clustering levels. Moreover, our method can be naturally applied to weighted graphs. We have demonstrated, through the definition of *Entropy of Information*, our algorithm is efficient to discover all cluster levels for general networks. We believe that the local algorithm procedure will not only allow to us to study much larger networks but also to mimic single human behavior in social network through specific and simple heuristics decision rules. The model parameters  $m$  and  $\alpha$  play a crucial role for the detection of communities. These results suggest how cognitive heuristics could be designed as those mechanism which allow humans to optimize those parameters in order to maximize the gathered information from the environment. Following this assumption the future works will investigate what kind of computational procedures could be used to mimic this human behavior. We plan to investigate the consequences of bounded memory and computational resources of nodes, in particular in a dynamic environment.



## Chapter 4

# Epidemic Spreading

### 4.1 Single-layer networks

#### 4.1.1 Introduction

Epidemic spreading is one of the most successful and most studied applications in the field of complex networks. The comprehension of the spreading behavior of many diseases, like sexually transmitted diseases (i.e. HIV) or the H1N1 virus, can be studied through computational models in complex networks [7, 91]. In addition to “real” viruses, spreading of information or computer malware in technological networks is of interest as well.

The susceptible-infected-susceptible (SIS) model is often used to study the spreading of an infectious agent on a network. In this model an individual is represented as a node, which can be either be “healthy” or “infected”. Connections between individuals along which the infection can spread are represented by links. In each time step a healthy node is infected with a certain probability if it is connected to at least one infected node, otherwise it reverts to a healthy node (parallel evolution).

The study of epidemic spreading is a well-known topic in the field of physics and computer science. The dynamics of infectious diseases has been extensively studied in scale-free networks [92–95], in small-world networks [96] and in several kind of regular and random graphs.

A general finding is that it is hard to stop an epidemic in scale-free networks with slow tails, at least in the absence of correlations in the network among the infections process and the node characteristics [92]. This effect is essential due to the presence of hubs,

which act like strong spreaders. However, by using an appropriate policy for hubs, it is possible to stop epidemics also in scale-free networks [94, 97].

This network-aware policy is inspired by the behavior of real human societies, in which selection had lead to the development of strategies used to avoid or reduce infections. However, human societies are not structureless, thus a particular focus must be devoted to the community structures, which are highly important for our social behavior.

Recently, a wave of studies focused the attention on the effect of the community structure in the modelling of epidemic spreading [98–100]. However, the focus was only set towards the interaction between the viruses' features and the topology, without considering the important relation between cognitive strategies used by subjects and the structure of their (local) community/neighborhood.

Considering this scenario, an important challenge is the comprehension of the structure of real-world networks [31, 45, 101]. Given a graph, a community is a group of vertices that is “more linked” within the group than with the rest of the graph. This is clearly a poor definition, and indeed, in a connected graph, there is not a clear distinction between a community and a rest of the graph. In general, there is a continuum of nested communities whose boundaries are somewhat arbitrary: the structure of communities can be seen as a hierarchical dendrogram [30].

It is generally accepted that the presence of a community structure plays a crucial role in the dynamics of complex networks; for this reason, lots of energy has been invested to develop algorithms for the detection of communities in networks [43, 90, 102]. However, in complex networks, and in particular in social networks, it is very difficult to give a clear definition of a community: nodes often belong to more than just one cluster or module. The problem of overlapping communities was exposed in [32] and recently analyzed in [33]. People usually belong to different communities at the same time, depending on their families, friends, colleagues, etc. For instance, if we want to analyze the spreading of sexual diseases in a social environment, it is important to understand the mechanism that leads people to interact with each other. We can surely detect two distinct groups of people (*i.e.*, communities): heterosexual and homosexual, with bisexual people that act as overlapping vertexes between the two principal communities [99, 103, 104].

The strategies used to face the infection spreading in a community is itself a complex process (*i.e.*, social problem solving) in which strategies spread (as the epidemics) along the community, and are negotiated and assumed or discarded depending on their social success.

Several factors can affect the social problem solving which is represented by the adoption of a behaviour to reduce the infection risk. Of course, personality factors, previous

experiences and the social and economical states of a subject can be considered as influencing variables. Another important variable is represented by the structure of the environment in which the social communities live, because it determines at the same time the speed of the epidemic diffusion and the strategy of the negotiation process; in particular large and more connected communities are often characterized by conservative strategies while small and isolated communities allow more relaxed strategies.

The same strategy can be more or less effective depending on the strategies adopted by the neighbours (community) of the subject. For instance, a subject in a conservative community can adopt a more risky (and presumably profitable) attitude with a certain confidence since he would be protected from the infection because of their neighbours' behaviours. This "parasitic" behavior (like refusing vaccinations) can be tolerated up to a certain level without lowering the community's fitness.

Not only the neighbor's behaviours affect the evolution of the cognitive strategies of a subject, but also the position he has in the network should be a relevant factor. A hub, or a subject with a great social betweenness, is usually more exposed to the infection than a leaf, and as a consequence, the best strategy for him has to be different. In the same way, since the topology of the network (*e.g.*, small world, random) determines variables such as the speed of the spreading, or its pervasiveness, it should also affect the development of the "best strategy".

Moreover, while the negotiation process evolves, the cognitive strategies usually develop within the most intimate community of a subject, thus the behaviour adopted by subjects could be an interesting feature for the community detection problem as well.

The understanding of the effects of the community structure on the epidemic spreading in networks is still an open task. In this section we investigate the role of risk perception in artificial networks, generated in order to reproduce several types of overlapping community structures.

The rest of this section is organized as follows: we start by describing a mechanism for generating networks with overlapping community structures in section 4.1.2. In section 4.1.3, we describe the SIS model adopted to model the risk perception of subjects in those networks. Finally, section 4.1.4 contains simulation results from our model with a throughout discussion and future work proposals.

#### 4.1.2 The networks model

There are  $n_c$  different communities with  $n_v$  vertices (here we consider only undirected and unweighted graphs); we assume that the probability to have a link between the

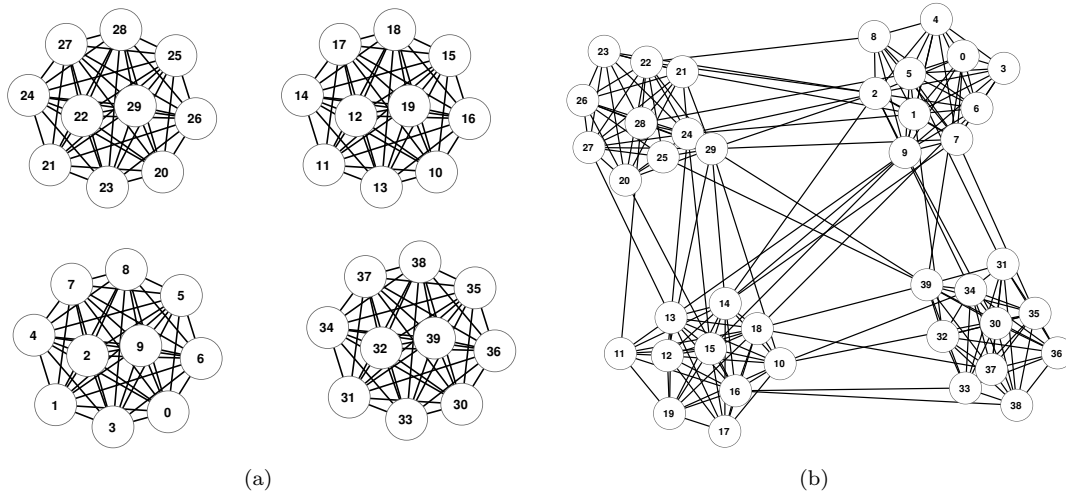


FIGURE 4.1: Schematic representation of our network-generation model.

(a) An example network with 4 different communities composed by 10 vertices: in this case, considering  $p_1 = 1$  and  $p_2 = 0$ , we generate 4 non-interconnected fully connected networks. (b) The same 4 communities with parameters  $p_1 = 0.95$  and  $p_2 = 0.05$ .

vertexes in the same community is  $p_1$ , while  $p_2$  is the probability to have a link between two nodes belonging to different communities. For instance, with  $p_1 = 1$  and  $p_2 = 0$ , we generate  $n_c$  fully connected graphs, with no connections among them as shown in the Figure 4.1(a). It is possible to use the parameters  $p_1$  and  $p_2$  to control the interaction among different communities, as shown in Figure 4.1(b). The algorithm for generating this kind of networks can be summarized as:

1. Define  $s_1$  as number of vertexes in the communities;
2. Define  $n_c$  as number of communities;
3. For all the  $n_c$  communities create a link between the vertexes on them with probability  $p_1$ ;
4. For all the vertexes  $N = s_1 n_c$  create a link between them and a random vertex of other communities with probability  $p_2$ ;

Constraining the condition  $p_1 = 1 - p_2$ , we can reduce the free parameters to just one. The connectivity degree itself depends on the size of the network and on the probabilities  $p_1$  and  $p_2$ . In particular, the connectivity function  $f(k)$  has a normal distribution from which we could define the mean connectivity  $\langle k \rangle$  as

$$\langle k \rangle = (s_1 - 1)p_1 + (n_c - 1)s_1 p_2 \quad (4.1)$$

with standard deviation  $\sigma^2(k) = (s_1 - 1)p_1(1 - p_1) + (n_c - 1)s_1 p_2(1 - p_2)$ .



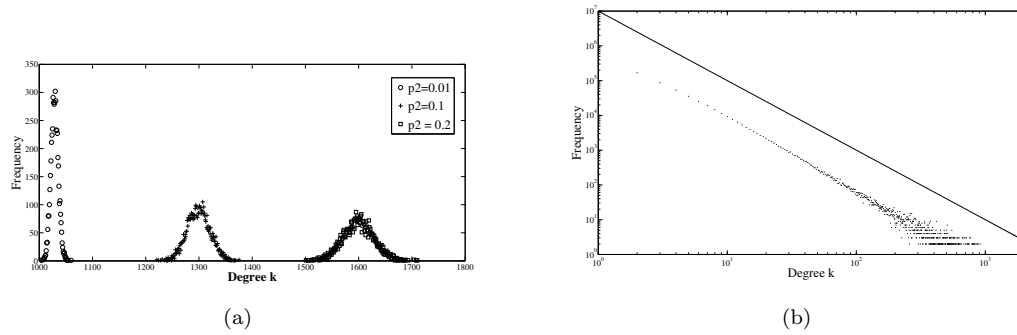


FIGURE 4.2: Connectivity degree distribution of generated networks.

(a) Random networks: in this figure, we show the frequency distribution of the connectivity degree changing the value of the parameter  $p_2$ . The circles represent the values for  $p_2 = 0.01$ , crosses for  $p_2 = 0.1$  and eventually squares for  $p_2 = 0.2$ . Here,  $s_1 = 1000$  and  $n_c = 5$ , thus we have generated networks with 5 communities of 1000 nodes for each. (b) Distribution of connectivity degree for the scale-free network generated with the mechanism described above (dots). The straight line is a power law curve with exponent  $\gamma = 2.5$ .

In Figure 4.2(a) we show the frequency distribution of the connectivity degree of nodes varying the value of the parameter  $p_2$  for a network composed by  $N = 5000$  nodes and  $n_c = 5$  communities.

It is widely accepted that real-world networks from social networks to computer networks are scale-free networks, whose degree distribution follows a power law, at least asymptotically. In this network, the probability distribution of contacts often exhibits a power-law behavior:

$$P(k) \propto ck^{-\gamma}, \quad (4.2)$$

with an exponent  $\gamma$  between 2 and 3 [105, 106]. For generating networks with this kind of characteristics, we adopt the following mechanism:

1. Start with a fully connected network of  $m$  nodes;
2. Add  $N - m$  nodes;
3. For each new node add  $m$  links;
4. For each of these links choose a node at random from the ones already belonging to the network and attach the link to one of the neighbors of that node, if not already attached.

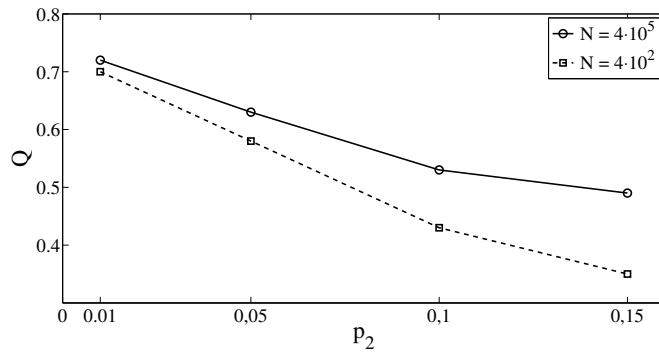


FIGURE 4.3: Modularity values of generated networks.

Different values of modularity ( $Q$ ) by increasing the mixing parameter  $p_2$  for two networks with  $N = 4 \cdot 10^5$  nodes and  $N = 4 \cdot 10^2$  nodes.

Through this mechanism we are able to generate scale-free networks with an exponent  $\gamma = 2.5$  as shown in Figure 4.2(b). There, we show the frequency distribution of the connectivity degree for a network of  $10^6$  nodes. To generate a community structure with a realistic distribution, we first generate  $n_c$  scale-free networks as explained above. Then, for all nodes and all outgoing links, we replace the link pointing inside the community with that connecting a neighbor of a random node in a random community with a probability of  $p_2 = 1 - p_1$ . Thus, the algorithm can be summarized as:

1. Generate  $n_c$  communities as scale-free networks with  $n_v$  vertices;
2. For all the vertices, with a probability  $p_2 = 1 - p_1$ ;
  - Delete a random link;
  - Select a random node of another community and create a link with one of its adjacent vertex;
3. End.

In this way, we are able to generate scale-free networks with a well defined community structure. A good measure for the estimation of the strength of the community structure is the so-called modularity [45]. The modularity  $Q$  is defined to be:

$$Q = \frac{1}{2} \sum_{vw} \left[ A_{vw} - \frac{K_v K_w}{2m} \right] \delta(c_v, c_w), \quad (4.3)$$

where  $A$  is the adjacency matrix in which  $A_{vw} = 1$  if  $w$  and  $v$  are connected and 0 otherwise.  $m = \frac{1}{2} \sum_{vw} A_{vw}$  is the number of edges in the graph,  $K_i$  is the connectivity degree of node  $i$  and  $(K_v K_w)/(2m)$  represents the probability of an edge existing between

vertices  $v$  and  $w$  if connections are made at random but with respecting vertex degrees.  $\delta(c_v, c_w)$  is defined as follows:

$$\delta(c_v, c_w) = \sum_r^{n_c} \hat{c}_{vr} \hat{c}_{wr} \quad (4.4)$$

where  $\hat{c}_{ir}$  is 1 if vertex  $i$  belongs to group  $r$ , and 0 otherwise.

In Figure 4.3 we show the values of modularity for two networks that were generated with the same algorithm, but with different sizes. Here, we consider a network with 4 communities: in the first case  $s_1 = 10^5$ , while in the second case  $s_1 = 10^2$ . What one can observe in Figure 4.3 is that the modularity's behaviour does not change significantly for different network sizes with the same number of communities.

In the case of scale-free networks, the mean connectivity degree  $\langle k \rangle$  is fixed a priori when we choose the number of links the new nodes create. In the case of random networks the mean connectivity is given by Eq. 4.1.

### 4.1.3 The risk perception model

We use the susceptible-infected-susceptible model (SIS) [92, 104] for describing an infectious process. In the SIS model, nodes can be in two distinct states: healthy and ill. Let us denote by  $\tau$  the probability that the infection can spread along a single link. Thus, if node  $i$  is susceptible and it has  $k_i$  neighbors of which  $s_n$  are infected, then, at each time step, node  $i$  will become infected with the probability:

$$p(s, k) = [1 - (1 - \tau)^{s_n}]. \quad (4.5)$$

We model the effect of risk perception considering the global information of the infection level for the whole network, the information about the infected neighbors and the information about the average state of the community. Thus, the risk perception for the individual  $i$  is given by:

$$I_i = \exp \left\{ -H + J_1 \left( \frac{s_{ni}}{k_i} \right) + J_2 \left( \frac{s_{ci}}{n_{ci}} \right) \right\}, \quad (4.6)$$

where  $H = J(s/N)$  is the perception about the global network on which  $s$  is the total number of infected agents while  $N$  is the number of agents in the network. The second term of the Eq. 4.6 represents the perception about the neighborhood, while the third term represents the perception about the local community of the agent  $i$ .

In this model, we assume that people receive information about the network's state through examination of people in the neighborhood. The global information could refer to entities like media while the information about the community could be assumed as *word of mouth*. Here, we don't consider the cost that people should pay in order to get these information, but it is clearly an important constraint to consider in future works.

The risk perception  $I_i$ , defined in Eq. 4.6, is assumed to determine the probability that the agents meet someone in its neighbourhood. The algorithm is given by:

1. For all nodes  $i = 1, 2, \dots, N$ ;
2. For all its neighbors  $j = 1, 2, \dots, k_i$ ;
3. If  $I_i > rand$ ;

  - $i$  meets  $j$ ;
  - If  $j$  at time  $t - 1$  was infected then  $i$  becomes ill with probability  $\tau$ ;

4. End.

Then, we propose a gain function defined as the number of meetings in time considering different values of  $j = J, J_1, J_2$  and different kind of scale-free and random networks; the gain function  $G(j)$  is given by:

$$G(j) = \frac{\sum_{t=1}^{T_e} M_t}{T_e}, \quad (4.7)$$

in which  $T_e$  is the time for the extinction, while  $M_t$  is the number of meetings during time. Based on that, we can eventually define a fitness function that considers the probability to extinct the epidemic in the given time. Thus, the fitness function is given by:

$$F_j^T = G(j)P_e(j) \quad (4.8)$$

It is possible to make a mean-field approximation of this model. Pastor-Satorras and Vespignani defined the mean-field equation for scale-free networks in [92]. In 2010, Kitchovitch and Liò [107] modeled the mean number of infected neighbors  $g(k)$  for individuals  $i$  with connectivity degree  $k$ . In fact, given the probability of receiving an infection by at least one of the infected neighbors (Eq. 4.5), it is possible to define the rate of change of the fraction of individuals  $i$  with degree  $k$  at time  $t$  by the following:

$$\frac{di_k}{t} = -\gamma + (1 - i_k)g(k), \quad (4.9)$$

on which  $\gamma$  is the rate of recovery (in our simulations we set  $\gamma = 1$ ).

Then, as shown by Boccaletti et al. [108], for any node, the degree distribution of any of its neighbors is,

$$q_k = \frac{kP(k)}{\langle k \rangle}, \quad (4.10)$$

hence, it is possible to define the number of infected neighbors as:

$$i_n = \sum_{K_{min}}^{k_{max}} q_k i_k, \quad (4.11)$$

and it allows to give a definition of  $g(k)$  as:

$$g(k) = \sum_{s=0}^k \binom{k}{s} p(s, k) i_n^s (1 - i_n)^{k-s}, \quad (4.12)$$

where  $s = s_n$  is the number of infected neighbors.

The temporal behavior of the mean fraction  $c$  of infected individuals in the case of a network with fixed connectivity is given by:

$$c' = \sum_{s_n=1}^k \binom{k}{s_n} c^{s_n} (1 - c)^{k-s_n} [1 - (1 - \tau)^{s_n}], \quad (4.13)$$

where  $c \equiv c(t)$ ,  $c' \equiv c(t + 1)$  and the sum runs over the number  $k_{inf}$  of infected individuals.

#### 4.1.4 Results and Discussion

We studied the behavior of our model for different scenarios. In Figure 4.5, we show results considering a network of 500 nodes and 5 communities where the initial number of infected agents is  $\approx 10\%$  of all agents in the network. We focus on the information about the community (parameter  $J_2$ ), while we kept  $J = J_1 = 1$  fixed. It is very interesting to observe the time necessary for the extinction of the epidemics, with the probability of being infected  $\tau = 0.5$  and changing the community structure of the network.

The effects of the parameters  $J_2$ ,  $J_1$  and  $J$  on the fitness function  $F$  considering different scale-free and random networks with different values of modularity are shown in Figure 4.4. The results were averaged over 100 simulations for each value of  $J$ ,  $J_1$  and  $J_2$ .

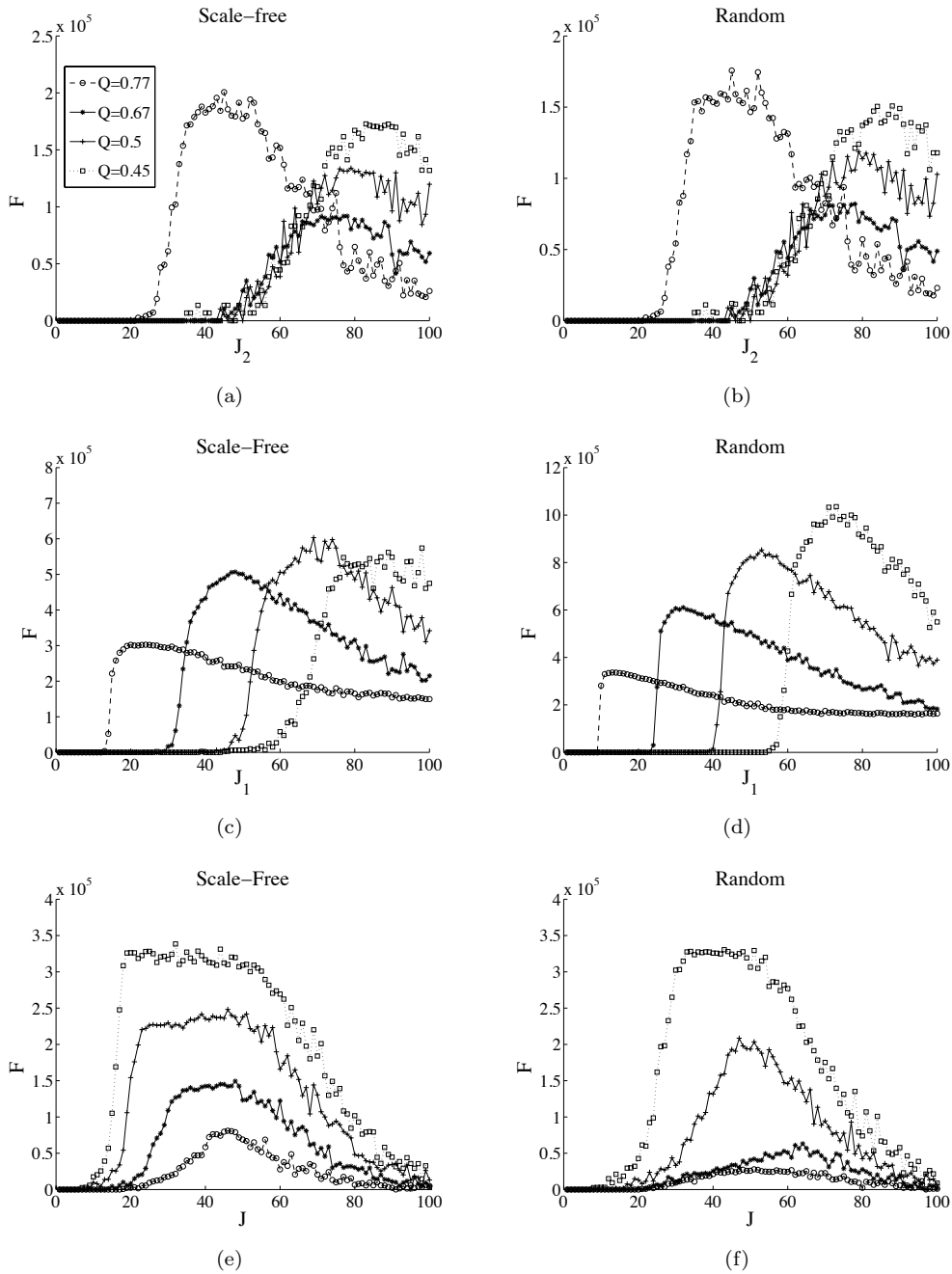


FIGURE 4.4: Effect of the precaution parameters in different networks.

$J_2$  (a),  $J_1$  (b) and  $J$  (c) ( $x$ -axis) versus the fitness function  $F$  ( $y$ -axis) considering different scale-free and random networks with different values of modularity. Results are averaged over 100 simulations for each value of  $J$ ,  $J_1$  and  $J_2$ .

On the left side of Figure 4.5 we show the temporal evolution of the percentage of infectious agents for different kind of networks and different values of  $J_2$ . We can observe that the extinction time increases when the modularity of network decreases, even if we use higher values of  $J_2$ . On the right side of Figure 4.5, we show the effect of the precaution on the extinction time. The straight line corresponds to different values

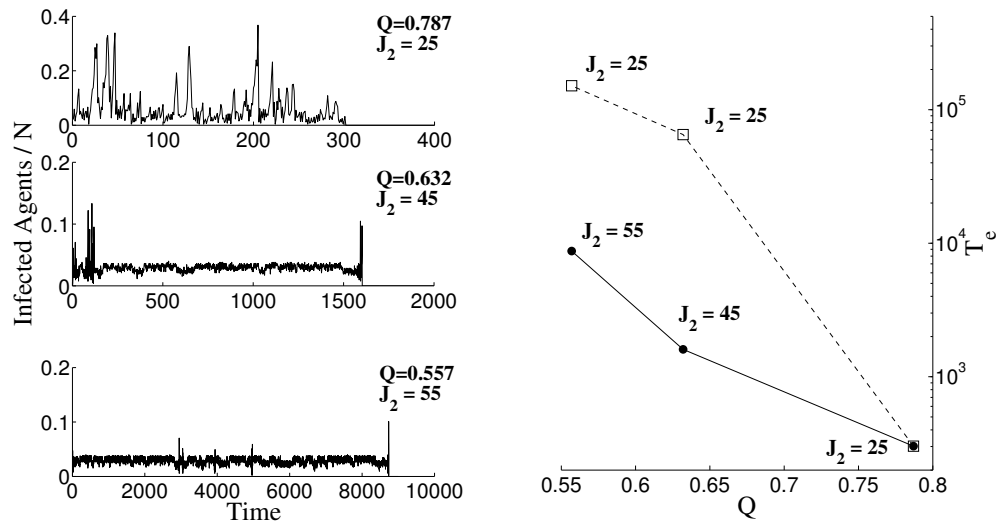


FIGURE 4.5: Infected individuals and the effect of precaution parameter in community-structured networks.

On the left side of the figure we show the temporal evolution of infected individuals by varying the mixing parameter  $p_2$ . The time necessary for the epidemic extinction increases as the modularity  $Q$  decreases. On the right side, we show the effects of the precaution parameter  $J_2$  on the extinction time by varying the modularity  $Q$ . The straight line represents the results for different value of  $J_2$ , while the dashed line represents the results for a constant value of  $J_2$ .

of  $J_2$ , while the dotted line corresponds to the same value of  $J_2$  in different kind of networks. It is also possible to observe that when a network becomes *less clustered*, the information about the community becomes less important.

In the case of scale-free networks, the mean connectivity degree  $\langle k \rangle$  is related to the number  $m$  of links the new nodes create. In the above example, considering  $m = 5$ , we obtained a mean connectivity degree  $\langle k \rangle = 7.8$ .

For comparisons, we generated random networks with a mean connectivity degree  $\langle k \rangle \in (7, 8)$ . The first result that we obtained is that the extinction time is larger than in the scale-free case. In Figure 4.6, we show the temporal evolution of the infected agents for a random network with modularity  $Q = 0.78$  considering  $J_2 = 25$  as in the upper plot on the left side of Figure 4.5. For the scale-free network the time necessary for the extinction is  $T_e \simeq 3 \cdot 10^2$  while for the random one it is  $T_e \simeq 3 \cdot 10^3$ .

Regarding the effects of the global and local (neighborhood) information, we investigated scale-free networks composed by 500 nodes and 5 communities with a fixed maximum threshold time  $T_{max}$ , necessary for the extinction of the epidemics. We assume  $T_{max} = 1000$  and separately measure critical values of  $J$ ,  $J_1$  and  $J_2$ . In table 4.1, we show the

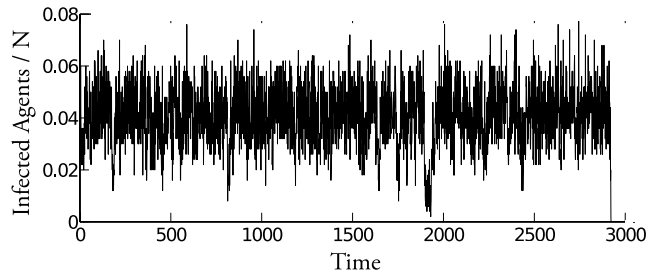


FIGURE 4.6: Percentage of infected agents for a random network.

Here  $N = 100$  nodes and 5 communities with modularity  $Q = 0.78$ . Adopting the same parameter used for the simulation reported in Figure 4.5, we show how the time for the extinction (approx. 2900 units) of the epidemic is greater than for the scale-free case (*i.e.*, upper plot on the left side of Figure 4.5).

$Q$ (modularity)	Critical Values		
	$J$	$J_1$	$J_2$
0.78	45	15	25
0.64	40	15	45
0.35	40	20	55

TABLE 4.1: Critical values for the extinction of the epidemic on case of scale-free networks of 500 nodes and 5 communities considering a maximum threshold time  $T_{max} = 1000$ , necessary for the extinction of the epidemics.

critical values of the three parameters by changing the modularity  $Q$ . As we can observe, the most variable parameter is  $J_2$  while the other two parameters do not appear to change. From this figure, we observe that the information about the fraction of infected neighbors is the most effective for stopping the disease. However, in order to get this piece of information, each node needs to check the status of all its neighbours, a task that can be quite hard and possibly conflicting with privacy. On the other hand, the information on the average infectious level in the community or in the whole population is more easily obtained. Therefore, one needs to add the cost of information into the model in order to decide what the most effective solution for risk perception is.

Summarizing, we have studied the progression and extinction of a disease in a SIS model over modular networks, formed by a certain number of random and scale-free communities. The infection probability is modulated by a risk perception term (modeling the probability of an encounter). This term depends on the global, local and community infection level. We found that in scale-free networks the progression is slower than in random ones with the same average connectivity. For what concerns the role of perception, we found that the local one (information about infected neighbours) is the most effective for stopping the spreading of the disease. However, it is also the piece of information that requires most efforts to be gathered, and therefore it may result a high cost/efficacy ratio.



The main element of originality of this model is that we introduced a network model based on communities, which still retains the scale-free structure with the possibility of changing the modularity, and we think that this structure (albeit being quite theoretical) is more realistic than standard scale-free networks. The fact the knowledge about own community is more effective than other indicators is surely trivial (and we expected to get this result), but it is also the information that is more expensive to get, at least for the standard data gathering existing today. We would like to quantify the advantage in using this indicator in order to compare its efficiency with respect to its cost (and for doing it we need to include a cost model, that will be done in a future work) and also point to the necessity of gathering this kind of local information, that in a real case may also present problems related to the privacy, but might be of great importance in the case of a pandemic.

We plan to extend the model by inserting a cost function, taking into account what are the best strategies to avoid the spreading of epidemics in different environments considering agents as intelligent entities capable to change or select the best strategies dynamically in order to minimize the risk and to maximize the economy of the system. In combination to this, we plan to add a more complex model such as the SIR eventually with vaccinations, for which there are important factors like the penetration and the possibility that the modular structure may be exploited to “shield” a community that may remain not exposed – similarly like people that refuse vaccinations, but are “shielded” by a surrounding community that vaccinates.

## 4.2 Multiplex Networks

### 4.2.1 Introduction

Recently, the Health magazine reported “Although H1N1 influenza killed more than 4,000 people in the United States in 2009-2010, this outbreak was relatively mild compared to some flu pandemics” [109].

Indeed, the twentieth century was characterized by a series of more serious events. During the 1918-19 the world assisted to the so-called *Spanish Flu*. Starting three different places: Brest (France), Boston (Massachusetts) and Freetown (Sierra Leone), the disease spread worldwide, killing 25 million people in 6 months (about 17 million in India, 500,000 in the United States and 200,000 in the United Kingdom). According to recent studies on the bodies of some American soldiers who died from the flu, recovered in the ice of Alaska [110, 111], it has been possible to reconstruct the 1918 virus through the synthesis of all of its eight subunits. The sequencing has put fully in the light that it

was a much more lethal virus than “normal” flu strains, belonging to the H1N1 subtype. The H1N1 is a bird flu virus that seems to have jumped directly from birds to humans in 1918.

In 1957, another pandemic originated in China and spread rapidly in Southeast Asia, taking hence the name of *Asian*. The virus responsible was identified in the subtype H2N2, new to humans, resulting from previous human H1N1 virus that was remixed with a duck virus from which it received the genes encoding the H2 and N2. This pandemic took eight months to travel worldwide and caused one to two million victims.

The 1968 pandemic was the mildest of the twentieth century and began once again to China in July 1968. From there it spread to Hong Kong, where more than half a million people fell ill, and in the same year reached the United States and the rest of the world. The virus was identified in subtype H3N2, a mutation of the H2N2 virus in 1957. Estimates of casualties range from 750,000 to 2 million people worldwide (34,000 in the U.S.) in the two years (1968-69) of its activity.

Given these facts (and the whole records of pandemics in history [112]), it is not surprising that the public health organizations are concerned about the appearance of a new deadly pandemics.

However, in recent decades there have been many cases of false or exaggerated information about epidemics. One example is the *Swine flu* of 1976, or the *Avian flu* in 1997 where a United Nations health official warned that the virus could kill up to 150 million worldwide [109] or the more recent 2009 *H1N1 flu*, during whose outbreak the U.K. Department of Health warned about 65000 possible deaths as reported by the Daily Mail in 2010 [113]. Fortunately these fears did not realized.

These catastrophic scenarios and the extent of their impact on the economic and social contexts induced a reflection on the method used to forecast the evolution of a disease in real world. It is well known that in deeply connected networks (and in particular in scale-free ones without strong compartmentalization) the epidemic threshold of standard epidemic modeling is vanishing [114–118]. Indeed, *lazzarettos* [119] experience first in Venice and then in many ports and cities was so successful in the absence of effective treatments because it was able to break the contact network. Indeed, the pest was last observed in Venice in 1630, whereas in southeastern Europe, it was present until the 19th century [120].

However, the last deadly pandemics of pest in Europe happened in 1820 [112], and worldwide in Vietnam in the 60’s; the last pandemic influenza, Hong Kong flu, in 1968-1969. In other words, the last rapid deadly pandemics happened well before the appearance

of the highly-connected human networks. We are here interested in diseases for which no vaccination is possible.

Is the absence of “modern” pandemics related to higher prevention efforts, more accurate hygiene and better health services or is it related to the influence of information on the effective infection rate of diseases?

Clearly, the public health systems put a lot of efforts in trying to make people aware of the dangers connected to hygiene, dangerous sexual habits and so on. However, the current worldwide diffusion of large-scale diseases (HIV, seasonal influenza, cold, papilloma virus, herpes virus, viral hepatitis among others) is deeply related to their silent (and slow) progression or to the assumption (possibly erroneous) or their harmlessness.

Therefore, in order to accurately modeling the spreading of a diseases in human societies, we need to take into account the perception of its impacts and the consequent precautions that people take when they become aware of an epidemic. These precautions consist in changes in personal habits, vaccination or modifications of the contact network.

In a previous work [121], some of us investigated the influence of the risk perception in epidemic spreading in the case in which the network of contacts is not affected and no vaccination is possible, so the only modification is through a change in personal habits. We assumed that the knowledge about the diffusion of the disease among neighbors (without knowing who is actually infected) effectively lowers its infectivity. We studied the worst case of an infection on a scale-free network with exponent  $\gamma = 2$  and we showed that in this case no degree of prevention is able to stop the infection and one has to take additional precaution for hubs (such as public officers and physicians). We extend here the investigation to different network structures, on order to obtain a complete reference frame. We investigated also the case in which the risk perception of being infected was function of the global information (e.g. media) and the fraction of neighbors infected nodes [121]. In this model there are always a finite level of precaution parameter for which the epidemics go extinct.

Here, we consider an important factor of modern society: the fact that our information no more comes mainly from physical contacts nor from broadcast media, but rather from our “virtual” contact network [122–124]. A recent study, *State of the news media* for the United States [125], highlights this phenomena. It shows the extent of the influence of social networks, for what concerns subscribers who can read news published by newspapers. The 9% of the population claims to inquire “very often” through Facebook and Twitter and the seven out of ten members is addressed to articles (from newspapers and other sources) by friends and family members.

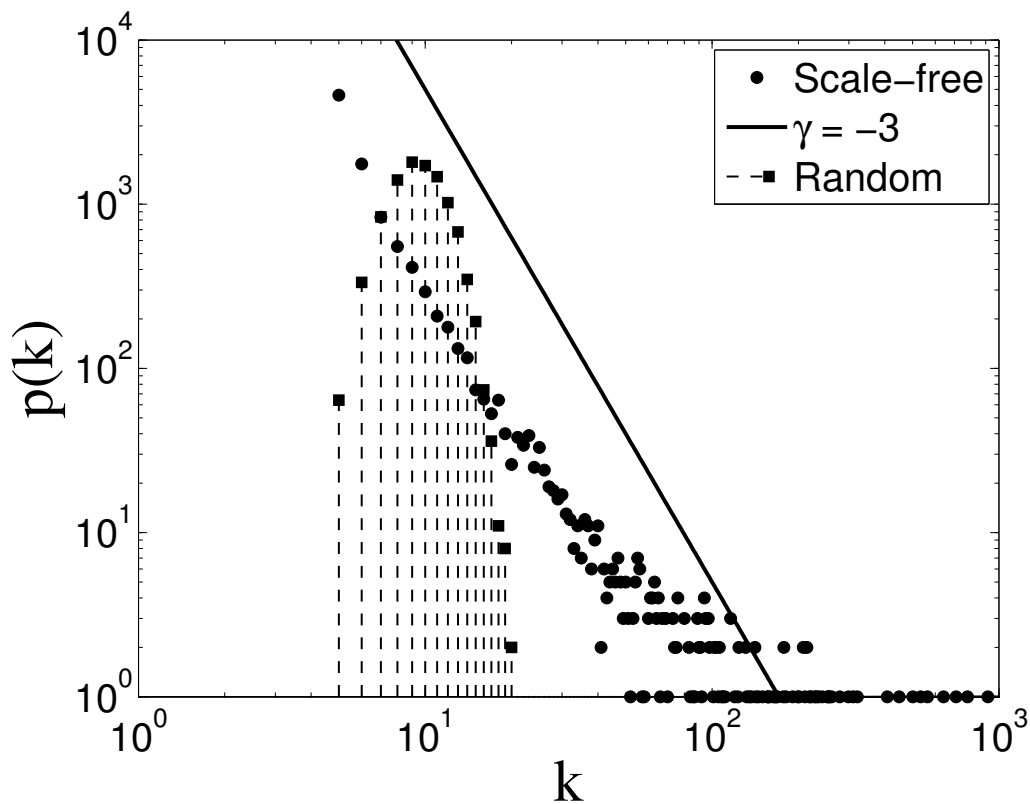


FIGURE 4.7: Distribution of connectivity degree of scale-free and random networks. Scale-free (points) and random (stems) networks with  $N = 10000$ ,  $m = 5$ ,  $\langle k \rangle = 2m = 10$  (the solid line corresponds to a power-law with exponent  $\gamma = -3$ ).

We are therefore confronted with news coming mainly from a virtual network. On the other hand the real network of contacts is the environment where physical and biological reactions take place. We extend our model to the case in which the source of information (the virtual contacts) does not coincide with the actual source of infection (the real contacts). This system is well-represented as multiplex-networks [126–130], i.e., graphs composed by  $M$  different layers in which the same set of  $N$  nodes can be connected to each other by means of links belonging to  $M$  different classes or types, which represents a specific case of Interdependent Networks [131, 132]. Recently, Granell et al. [133] have pointed out the attention to an interesting scenario where the multiplex corresponds to a two-layers network, one where the dynamics of the awareness about the disease (the information dynamics) evolves and another where the epidemic process spreads.

The first layer represents the information network where people become aware of the epidemic thanks to news coming from virtual and real contacts in various proportions. The second layer represents the real contact network where the epidemic spreading takes place.

In this paper we want to model the effect of the virtual information for simulating the awareness of the agents in the real-world network contacts. We study how the percolation threshold of a susceptible-infected-susceptible (SIS) dynamics depends on the risk perception (that affects the infectivity probability) when this information comes from the same contact network of the disease or from a different network. In other words, we study the interplay between risk perception and disease spreading on multiplex networks.

We are interested in the epidemic threshold, which is a quantity that it is not easy to obtain automatically (for different values of the parameters) using numerical simulations. We extend a self-organized formulation of percolation phenomena [134] that allows to obtain this threshold in just one simulation (for a sufficiently large system).

In Section 2 we present the multiplex network system. In Section 2 the self-organized percolation method and in Section 3 and 4 we apply the method to the standard risk perception problem and to the multiplex version, respectively. Conclusions are drawn in the last section.

#### 4.2.2 The network model

In this section we show our method for generating multiplex networks. First of all we describe three different mechanisms for generating three different kind of undirected unweighted networks, namely *regular*, *random* and *scale-free* networks. Let us denote by  $a_{ij} = 0, 1$  the adjacency (symmetric) matrix of the problem,  $a_{ij} = a_{ji} = 1$  if there is a link from  $j$  to  $i$  and zero otherwise. We shall denote by  $k_i = \sum_j a_{ij}$  the connectivity of site  $i$  and by  $j_1^{(i)}, j_2^{(i)}, \dots, j_{k_i}^{(i)}$  its neighbors ( $a_{i,n_j} = 1$ ). For all these mechanisms we define the number of links  $m$  of the incoming node. Here  $m$  is a parameter controlling the average connectivity degree  $\langle k \rangle = 2m$ . Then at time  $t = 0$  we have fully connected networks composed by  $N(0) = m + 1$  nodes.

- Regular: for generating regular networks of  $N$  nodes at time  $t + 1$  the incoming nodes will link to the closest  $m$  nodes. For instance if we define  $m = 2$ , at time  $t = 0$  the network is composed by nodes 1, 2, 3 fully connected, then at time  $t + 1$  the new node 4 will link to nodes 2 and 3 and so on until the closing of the circle where the node 1 will receive two more links from nodes  $N$  and  $N - 1$ . In this way all the nodes will have a connectivity  $\langle k \rangle = 2m$ .
- Random: starting from  $N(0) = m + 1$  nodes at  $t + 1$  the incoming node will like with an uniform probability  $p$  to  $m$  different nodes: i.e. the algorithm is developed

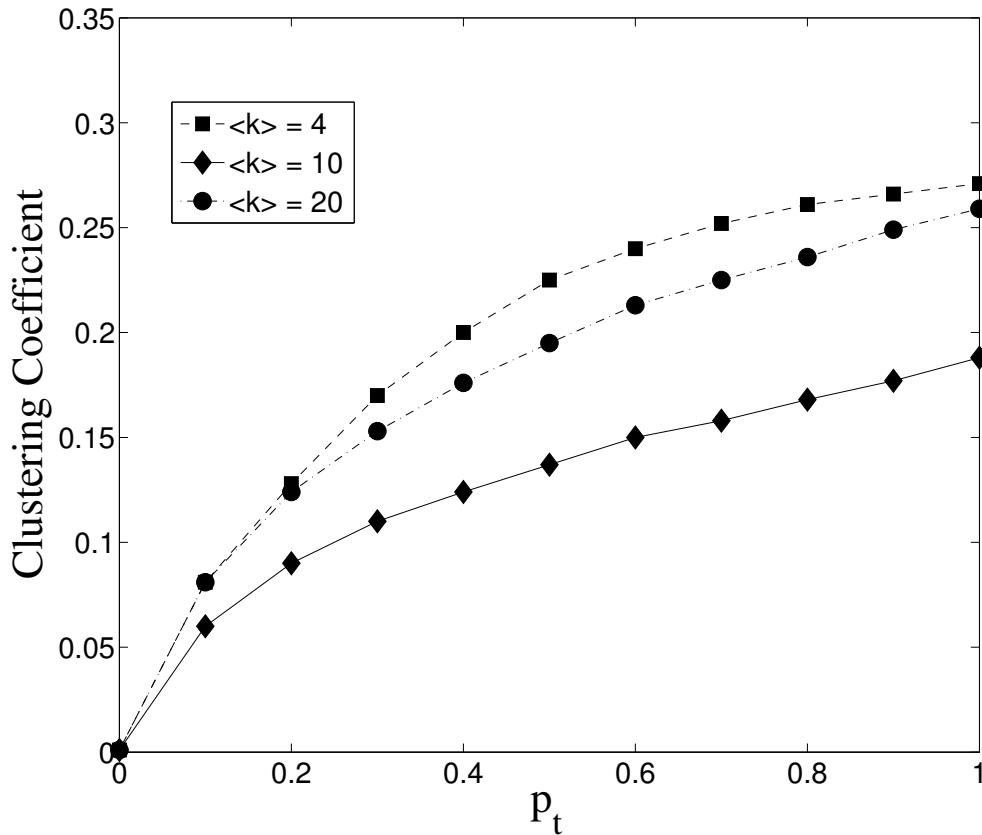


FIGURE 4.8: Clustering coefficient (y-axis) as function of *transitivity parameter*  $p_t$ . The three different lines correspond to three random networks composed by  $N = 10000$  nodes and average connectivity degree  $\langle k \rangle = 4, 10, 20$  as shown in the legend.

for avoiding self-loops and multiple links. So the total number of links in the networks will be  $L = 2 \cdot m \cdot N$  and the average connectivity degree is defined by:

$$\langle k \rangle = \frac{1}{N} \sum_i^N k_i = \frac{1}{N} 2mN = 2m. \quad (4.14)$$

The probability distribution of random networks is Poissonian,  $P(k) = \frac{z^k e^{-z}}{k!}$ , where  $z = \langle k \rangle$ .

- Scale-Free: they are generated following the so-called Barabási–Albert model [135]. Also here we start with  $m$  fully connected nodes, then at each time step each node establishes  $m$  links according to the following mechanism:

1. pick a random node in a network ( $r_n$ );
2. link to a random adjacent node of  $r_n$  (if not already connected);
3. repeat 1, 2 until it makes  $m$  links.

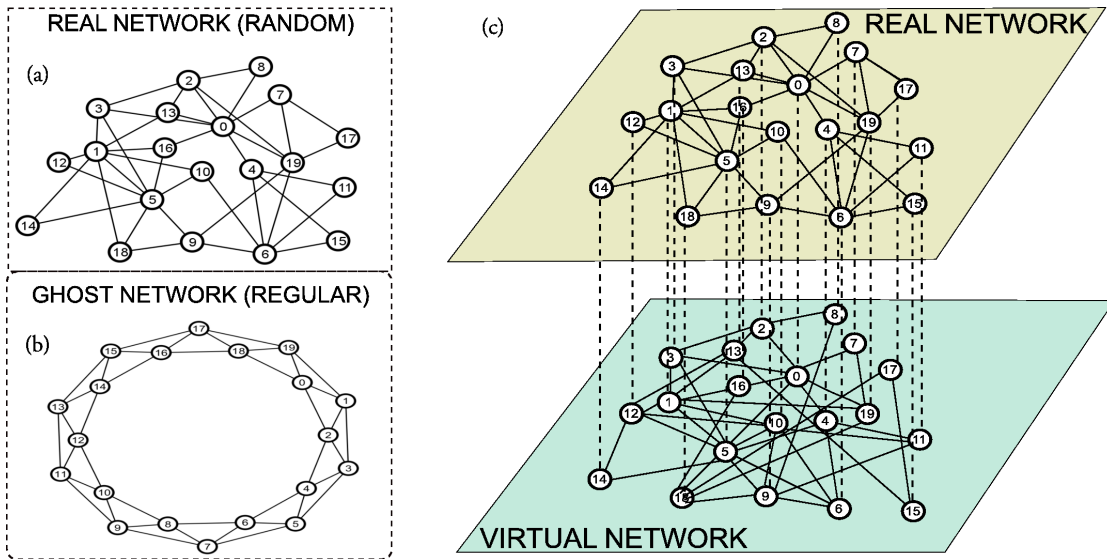


FIGURE 4.9: Example of multiplexes generated with our method.

(a)  $R_{net}$  and (b)  $G_{net}$  are both networks of  $N = 20$  nodes and  $\langle k \rangle = 4$ . (c) Multiplex networks where the *Virtual* network ( $V_{net}$ ) is given by the mixing of (a) and (b) with  $p_r = 0.5$ .

this mechanism allow us to generate scale-free networks ( $P \sim k^\gamma$ ) with exponent  $\gamma = -3$ . Also in this case, the total number of links is  $L = 2 \cdot m \cdot N$  and so  $\langle k \rangle = 2m$ .

Now we can turn the problem of generating networks by controlling the clustering coefficient. We remember that the clustering coefficient is a function of the number of local triples. Following the definition of Barrat et al. [7], the clustering coefficient  $C$  is defined as the average of the local clustering coefficients  $c_i$ :

$$C = \frac{1}{N} \sum_i c_i = \frac{1}{N} \sum_i \left( \frac{1}{k_i(k_i - 1)} \sum_{j,h} a_{ij} a_{ih} a_{jh} \right), \quad (4.15)$$

where  $k_i$  is the connectivity degree of node  $i$ . Then in order to increase the number of triples in our model we simply define a *transitivity parameter*  $p_t$  which is the probability to link to the adjacent nodes of the first selected nodes instead of another random node. In the Figure 4.8 we show that the clustering coefficient increases decreasing the average connectivity degree  $\langle k \rangle$ .

### 4.2.3 Multiplex network model

Now we can turn the problem to generate multiplex networks. In particular we are going to generate systems where the multiplex is given by two networks that we call *Real*

( $R_{net}$ ) and *Virtual* ( $V_{net}$ ) networks. For doing that we first generate the real network by choosing one from regular, random or scale-free. Then we generate a *Ghost* network ( $G_{net}$ ) chosen also in this case from our three benchmark networks with same average connectivity  $\langle k \rangle = 2m$ . We define a probability  $p_r$  for which each node  $i_v$  in  $V_{net}$  links with probability  $1 - p_r$  to all the adjacent nodes of  $i_r$  in  $R_{net}$  and with probability  $p_r$  to all the adjacent vertices of  $i_g$  in  $G_{net}$ . In this way for  $p_r = 0$ ,  $a_v(i, j) = a_r(i, j)$ ; viceversa for  $p_r = 1$ ,  $a_v(i, j) = a_g(i, j)$ ; where  $a_r(i, j)$ ,  $a_g(i, j)$  and  $a_v(i, j)$  are respectively the adjacency matrices of  $R_{net}$ ,  $G_{net}$  and  $V_{net}$ . This procedure allows us to study the effect of the difference between the  $R_{net}$  where the epidemic spreading takes place and  $V_{net}$  which is the information network where actors become aware of the disease (i.e. they evaluate the perception of the risk of being infected). An example of multiplexes is reported in the Figure 4.9.

#### 4.2.4 The self-organized percolation method

Here we show a self-organized percolation method that allows to obtain the critical value of the percolation parameter in a single run. We consider a SIS problem and the (directed) percolation direction is time (bond percolation). The dynamics of the SIS model is described by the following differential equations:

$$\frac{ds}{dt} = -\tau is + \gamma i, \quad \frac{di}{dt} = \tau is - \gamma i \quad (4.16)$$

where  $i$  and  $s$  correspond to the fraction of infected and susceptible individuals while  $\tau$  and  $\gamma$  are the infection and the recovery rates. Let us assume  $\gamma = 1$ . Let us denote by  $x_i(t) = 0, 1$  (0 = healthy, 1 = infected), the percolating variable and by  $p$  the control parameter (percolation probability). Considering  $p$  is fixed, the stochastic evolution process for the network is defined as

$$x_i(t+1) = \bigvee_{j=j_1^{(i)}, \dots, j_{k_i}^{(i)}} [p > r_{ij}(t)] x_j(t) \quad (4.17)$$

where  $\vee$  represents the OR operator and the multiplication represents the AND. The square bracket represents the truth function,  $[\cdot] = 1$  if “.” is true, and zero otherwise. The quantity  $r_{ij}(t)$  is a random number between 0 and 1 that varies with  $i, j$  and  $t$ . We want to derive an equation for  $p_i(t)$ , which is the minimum value of  $p$  for which  $x_i(t)$  is infected. We can replace  $x_i(t)$  by  $[p > p_i(t)]$ . (4.17) becomes

$$[p > p_i(t+1)] = \bigvee_{j=j_1^{(i)}, \dots, j_{k_i}^{(i)}} [p > r_{ij}(t)] [p > p_j(t)]. \quad (4.18)$$



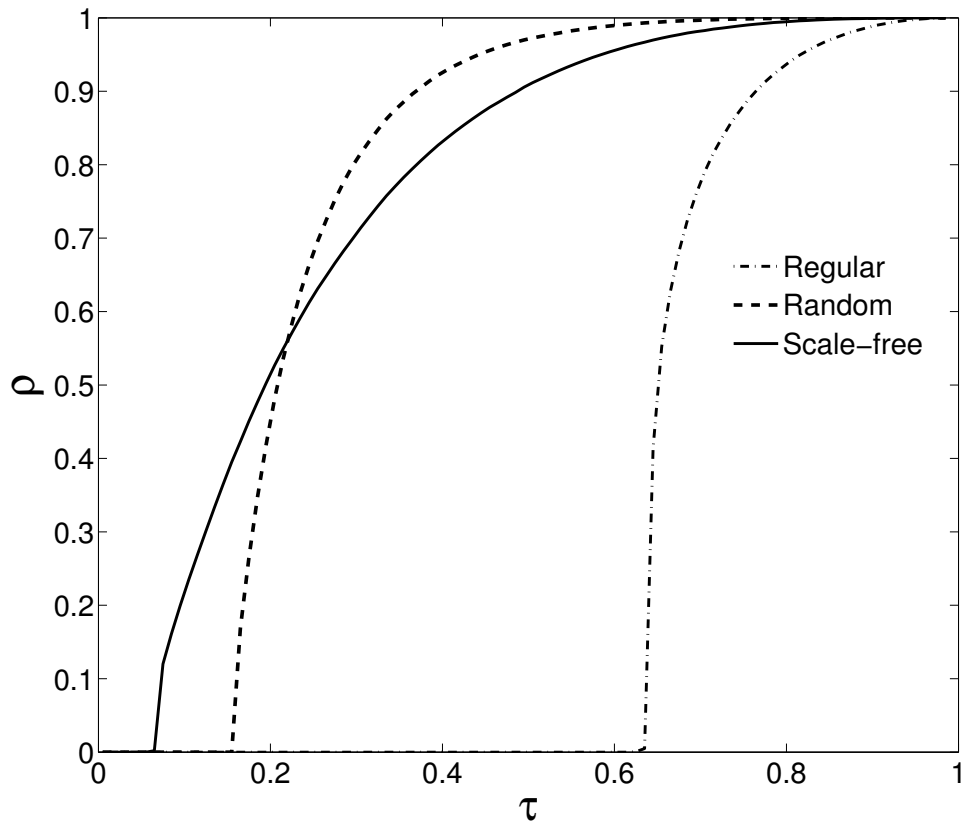


FIGURE 4.10: Cumulative frequency distribution of the infected individuals for the SIS dynamics achieved with our self-organized percolation method for different networks. In x-axis there is the probability  $\tau$  to be infected while on the y-axis the fraction  $\rho$  of infected individuals (realizations are averaged over 100 simulations).

Now  $[p > a][p > b]$  is equal to  $[p > \max(a, b)]$  and  $[p > a] \vee [p > b]$  is equal to  $[p > \min(a, b)]$ , therefore Eq.4.18 becomes

$$[p > p_i(t+1)] = \left[ p > \left( \text{MIN}_{j=j_1^{(i)}, \dots, j_{k_i}^{(i)}} \max(r_{ij}(t), p_j(t)) \right) \right], \quad (4.19)$$

and therefore we get the equations for the  $p_i$ 's

$$p_i(t+1) = \text{MIN}_{j=j_1^{(i)}, \dots, j_{k_i}^{(i)}} \max(r_{ij}(t), p_j(t)). \quad (4.20)$$

Let assume that at time  $t = 0$  all sites are infected, so that  $x_i(0) = 1 \forall p$ . We can alternatively write  $p_i(0) = 1$  (since the minimum value of  $p$  for which  $x_i(0) = 1$  is one). We can therefore iterate Eq. 4.20 and get the asymptotic distribution of  $p_i$ . The minimum of this distribution gives the critical value  $p_c$  for which there is at least one percolating cluster with at least one "infected" site at large times. As usual,  $t$  cannot be

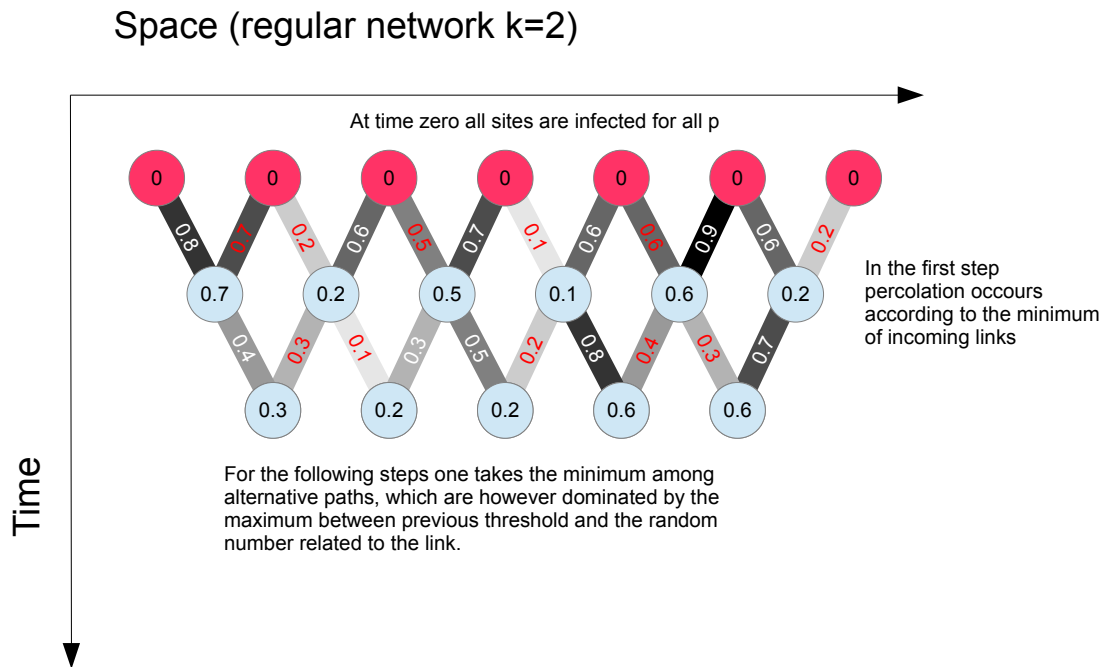


FIGURE 4.11: Evolutionary process of the direct percolation. The case of the regular network with  $\langle k \rangle = 2$ .

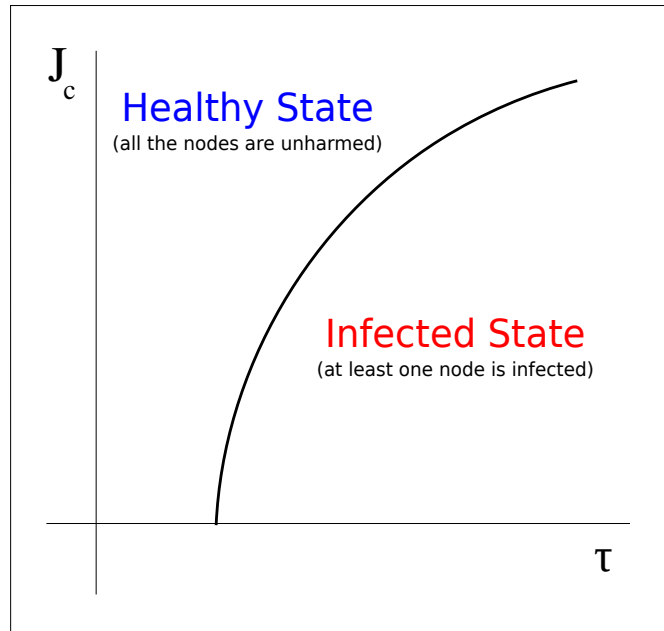
infinitely large for finite  $N$  otherwise there will be surely a fluctuation that will bring the system into the absorbing (healthy  $x_i = 0$ ) configuration. A schematic representation of this modus operandi is illustrated in the Figure 4.11. In Figure 4.10 we show the results of the self-organized percolation method for the SIS dynamics over different networks: in particular for regular networks with  $k = 2$  ( $p_c = 0.644$ , compatible with the results of the bond percolation transition in the Domany-Kinzel model [136]), random networks with mean connectivity degree  $\langle k \rangle = 6$  and  $\langle k^2 \rangle = 39$  ( $p_c = 0.165$ ) and for scale-free networks ( $\langle k \rangle = 6$ ,  $k^2 = 200$ ) with  $\gamma = -3$  ( $p_c = 0.0715$ ) which can be compared with other well-known results [114–118].

#### 4.2.4.1 Risk perception

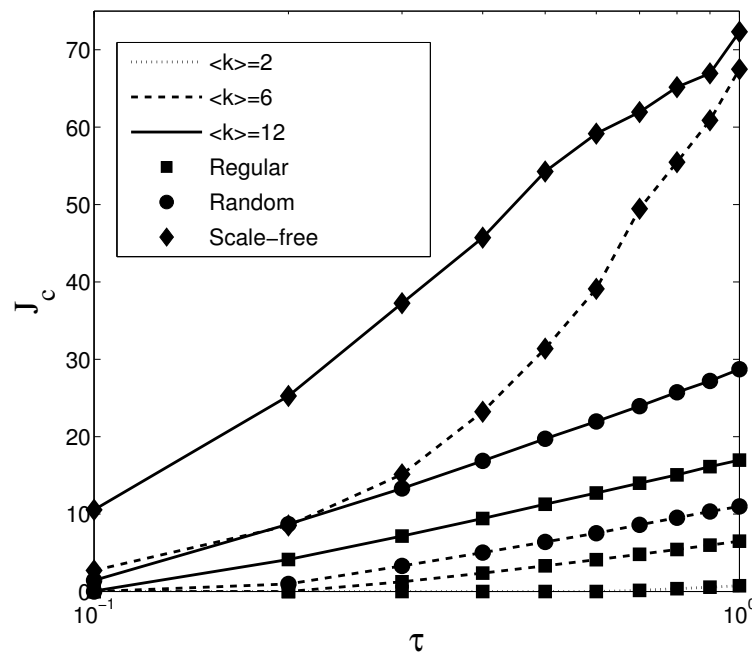
Now, let us apply the method to a more difficult problem, for which the percolation probability depends on the fraction of infected sites in the neighborhood (risk perception). We define the infection probability  $q$  as

$$q = \tau \exp\left(-J \cdot \frac{s}{k}\right) \quad (4.21)$$

where  $\tau$  is the bare infection probability,  $s$  is the number of infected neighbors and  $k$  is the node connectivity. In this case we want to find the minimum value of the



(a)



(b)

FIGURE 4.12: Schematic phase diagram and results for the risk perception.

(a) The perception threshold  $J_c$  (y-axis) separates two regions: one in which for values of  $J \geq J_c$  it is possible to stop the spread of the disease and one for which for values of  $J < J_c$  the system is not able to avoid the epidemic spreading. (b) (log x) Critical values of the precaution level  $J_c$  (y-axis) increasing the bare infection  $\tau$  (x-axis) for different networks (■ regular, ● random and ◆ scale-free). Results are averaged over 10 simulations.

parameter  $J$  for which there is no spreading of the infection at large times. The quantity  $[q > r] = [\tau \exp(-Js/k) > r]$  is equivalent to  $[J < -(k/s) \ln(r/\tau)]$ . Therefore Eq. 4.18 is replaced by

$$\begin{aligned} [J < J_i(t+1)] &= \\ &= \bigvee_{j=j_1^{(i)}, \dots, j_{k_i}^{(i)}} \left[ J < -\frac{k_i}{s_i} \ln \left( \frac{r_{ij}(t)}{\tau} \right) \right] [J < J_j(t)] \end{aligned} \quad (4.22)$$

where

$$s_i \equiv s_i(J) = \sum_{j=j_1^{(i)}, \dots, j_{k_i}^{(i)}} x_j = \sum_{j=j_1^{(i)}, \dots, j_{k_i}^{(i)}} [J_j(t) \geq J]. \quad (4.23)$$

So

$$\begin{aligned} [J < J_i(t+1)] &= \\ &= \bigvee_{j=j_1^{(i)}, \dots, j_{k_i}^{(i)}} \left[ J < -\frac{k_i}{s_i(J_j(t))} \ln \left( \frac{r_{ij}(t)}{\tau} \right) \right] [J < J_j(t)] \end{aligned} \quad (4.24)$$

and therefore

$$\begin{aligned} J_i(t+1) &= \\ &= \text{MAX}_{j=j_1^{(i)}, \dots, j_{k_i}^{(i)}} \min \left( -\frac{k_i}{s_i(J_j(t))} \ln \left( \frac{r_{ij}(t)}{\tau} \right), J_j(t) \right) . \end{aligned} \quad (4.25)$$

Analogously to the previous case, the critical value of  $J_c$  is obtained by taking the maximum value of the  $J_i(t)$  for some large (but finite) value of  $t$ . The results of the method are reported in Figure 4.12 for different networks, increasing the probability of infection  $\tau$ . Results are quite interesting if compared with the simple SIS dynamics (Figure 4.10) here we are able to stop the epidemic increasing the bare infection  $\tau$  until  $\tau = 1$ . Considering for instance the case of random networks with  $\langle k \rangle = 6$  in which we found a critical value of  $\tau_c = 0.165$  (Figure 4.10) while in Figure 4.12 we observe that after that value of  $\tau_c$  we need to adopt a precaution level  $J > 0$  in order to stop the spreading of the disease. The same consideration can be done also for the other scenarios.

#### 4.2.4.2 Virtual risk perception

We can now turn to the problem of computing the critical value  $J_c$  if the perception is computed on the information network that is different from that of the real infection. Here the perception  $s_i$  is computed on the neighbors  $\check{j}^{(i)}$  on the information network. The perceived number of infected neighbors depends on how many, in the information

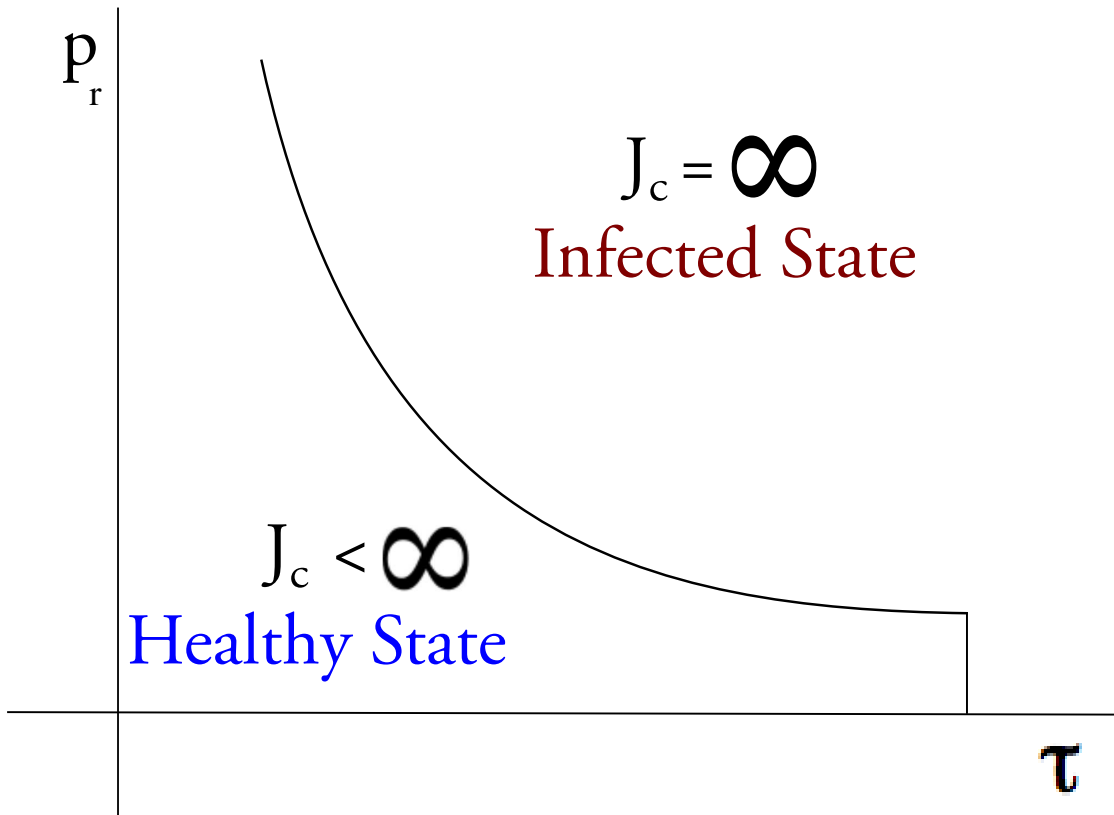


FIGURE 4.13: Schematic phase diagram for the virtual risk perception.

The perception threshold  $J_c$  separates two regions: one in which for finite values of  $J$  it is possible to stop the spread of an epidemic and one for which there are no values of  $J$  able to avoid the epidemic spreading. On the y-axis there is the difference  $p_r$  between the real contacts network and the virtual one; while on the x-axis there is the bare infection probability  $\tau$ .

network, own  $J_j$  larger than that computed in the real network, i.e.,

$$J_i(t+1) = \text{MAX}_{j=j_1^{(i)}, \dots, j_{k_i}^{(i)}} \min \left( -\frac{k_i}{\check{s}_i(J_j)} \ln \left( \frac{r_{ij}(t)}{\tau} \right), J_j(t) \right), \quad (4.26)$$

where

$$\check{s}_i(J) = \sum_{\check{j}=j_1^{(i)}, \dots, j_{k_i}^{(i)}} [J_{\check{j}} \geq J]. \quad (4.27)$$

In other words: for any value of  $J$  in the real neighborhood one computes how many of the *information* neighbors  $\check{j}$  have  $J_{\check{j}} \geq J$ . This is the perceived value of the risk. A schematic phase diagram for the risk-perception in modelling SIS dynamics in multiplex networks is reported in the Figure 4.13 where it is possible to catch the most interesting point of this paper. In particular, by increasing the difference  $d$  between the information network and the real one (*e.g.*, by increasing the weight of the virtual network on the information dynamics) the stopping of the infection spreading becomes harder to achieve

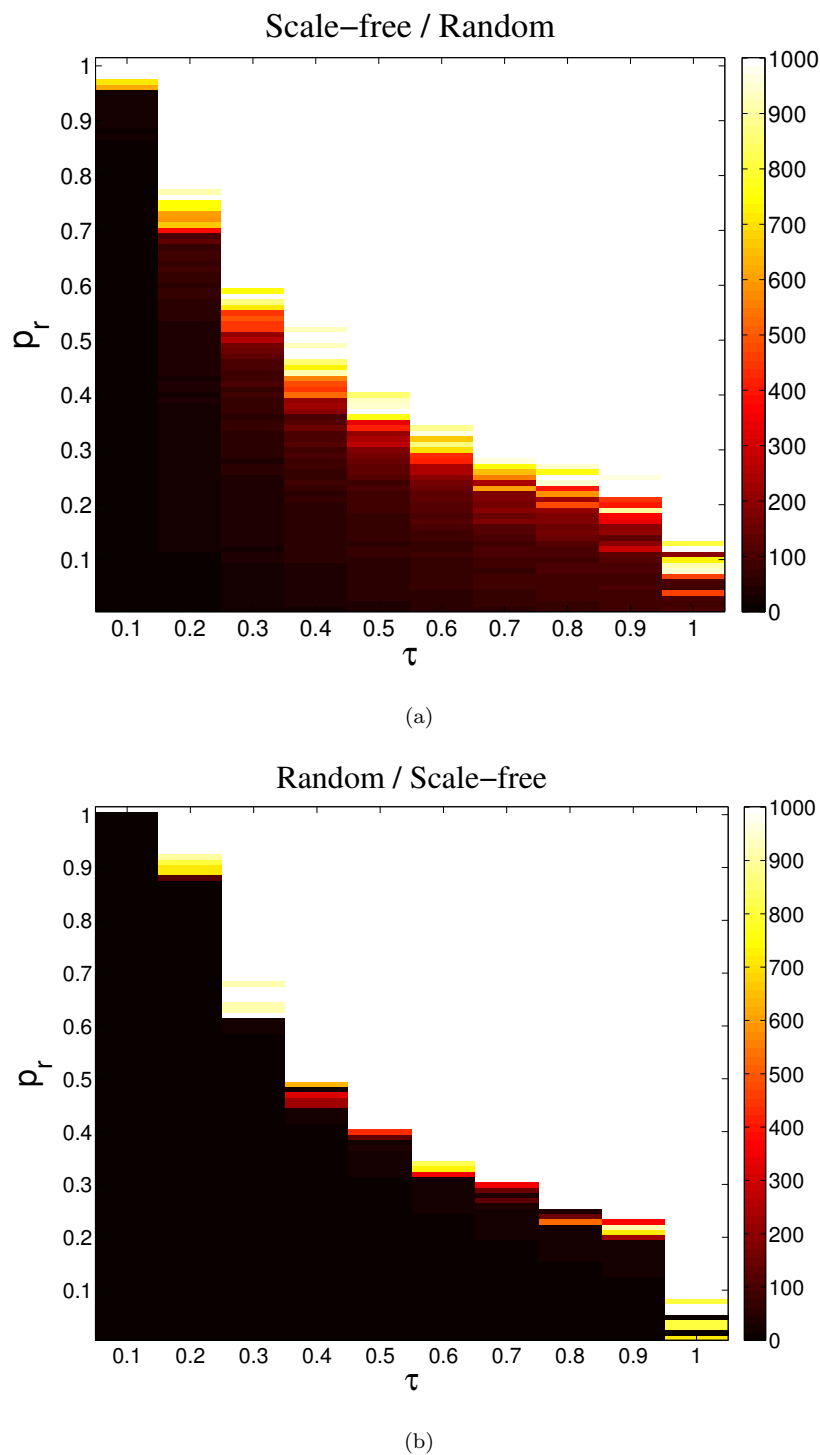


FIGURE 4.14: Risk Perception in multiplex networks.

(a)  $R_{net}$  scale-free and  $G_{net}$  random networks with  $\langle k \rangle = 6$ . (b)  $R_{net}$  random and  $G_{net}$  scale-free networks with  $\langle k \rangle = 6$ . Here 1000 corresponds to  $\infty$ . In both the simulated scenarios we can observe a drastic phase changing between white area in which there are no values of  $J_c$  able to stop the epidemic and the color area where for finite values of  $J_c$  it is possible to avoid the system's percolation.

(and finally impossible). The effective results of the self-organized percolation method on two different scenarios are reported in Figure 4.14. We can notice that the general trend is that, with the increasing of the difference  $d$  among the information network and the real one it becomes harder to stop an epidemics. However, it seems that, unless the case in which the real contact network is scale-free, there is always a finite value of the bare infection probability  $\tau$  for which the infection is stoppable, which confirms the evident fact that a sufficient level of hygiene, more effective drugs and better health services (that contribute to lowering  $\tau$ ) are fundamental elements also in our web 2.0 society.

#### 4.2.5 Final remarks

In conclusion, we have studied the effects of virtual risk perception in a SIS model in multiplex networks. The purpose of this work was to demonstrate how the virtual information can be counterproductive in the case of epidemic spreading in the real world networks. These effects are modulated by two parameters,  $J$  and  $p_r$  that respectively represent the precaution level and the mixing parameter between the real networks and a *ghost* network for generating the multiplex environment. We found that increasing the value of  $p_r$  there are not finite values of  $J_c$  that make the epidemics go extinct. Finally, we proposed a self-organized percolation method for detecting the critical values of parameters in the case of both single-layer and multiplex networks.





## Chapter 5

# Conclusions

In this thesis we have proposed pioneer methods for analysing both the structure and dynamics of complex networks.

### Community detection

We described an algorithm to extract communities structures in a network from a local point of view using algorithms inspired by a competitive interaction. The nonlinear part of this method, responsible for the actual elaboration of information, is inspired by a chemical/ecological competition model [36]. There is not a unique definition of a community, so an exploratory algorithm, like the one that humans have presumably developed during their evolution, should present different clustering for different values of the parameters, or for different iterations.

In a first implementation we adopted a frequency-based approach and an unbounded memory at the level of nodes. Unbounded memory means that the node's state vector  $S_i$  is not been limited and it could potentially reach a size equal to the network size  $N$ . Despite this, the expansion and normalization phases are sufficient to avoid this problem. Nevertheless it is very important to explicitly limit the computational resources of the node, as suggested by Simon in 1955 [73], increasing both the ecological plausibility of the model and the insights which drive the algorithm design. The first results that we obtained are promising. The method provides a natural “scanning” of the various clustering levels. Moreover, our method can be naturally applied to weighted graphs. We have demonstrated, trough the use of information entropy, that our algorithm efficiently discovers all cluster levels for general networks. We believe that the local algorithm procedure will not only allow to study much larger networks but also to mimic single human behavior in social noetwork trough specific and simple heuristics decision rules.

These results suggest that cognitive heuristics could be defined as those mechanism that allow humans to optimize those parameters in order to maximize the gathered information from the environment. Following this assumption, the future works will investigate what kind of computational procedures could be used to mimic this human behavior.

We also applied this method to the problem of discovering communities also in dynamic networks. Given the growing interactions between mobile devices and humans we focused our attention on the importance of the spreading and elaboration of the information which has a crucial role in the so-called Cyber-Physical-World [47]. We evaluated it on different synthetic human mobility scenarios and we found that our method is capable to detect not only the right communities from an individual viewpoint but also to spontaneously reveal the role of each nodes inside the network (travellers and normal agents). In the future, we would like to evaluate the scaling of our algorithms with the system size and apply it to more realistic scenarios. In particular we plan to compare our algorithm with others targeted to pocket switched networks (that use also global information). We would also like to combine the geographic proximity with additional social information so as to better catch the complex association between the real and the virtual world.

Moreover, we proposed a computational scheme inspired by the workings of human cognition. We embedded some fundamental aspects of the human cognitive system into this scheme in order to obtain a minimization of computational resources and the evolution of a dynamic knowledge network over time, and apply it to computer networks. Such algorithm is capable of generating suitable strategies to explore huge graphs like the Internet that are too large and too dynamic to be ever perfectly known. The algorithm equips each node with a local information about possible hubs which are present in its environment. Such information can be used by a node to change its connections whenever its fitness is not satisfying some given requirements. Eventually, we have compared our algorithm with a randomized approach within an ecological scenario for the ICT domain, where a network of nodes carries a certain set of objects, and each node retrieves a subset at a certain time, constrained with limited resources in terms of energy and bandwidth. Finally we have shown that a cognitive-inspired approach improves the overall networks topology better than a randomized algorithm.

In addition we performed a more complex heuristic for testing this method in real biological networks. Previous have analyzed the cluster organization of the cat cortical network using both traditional multidimensional scaling methods and evolutionary optimization algorithms. Interestingly, the evolutionary optimization principle of previous works was based on the modularity measure used to find communities in network with

global algorithms. In this thesis, we deepened this point taking into account different community-detection algorithms. We compare the performances of *NetExplorer*, a local information dynamics algorithm for detecting communities in networks, with six well-known community detection algorithms: *Infomap*, *Hierarchical Infomap*, *Louvain*, *Modularity Optimization*, *Label Propagation* and *Oslom*. The results indicate that our method (*NetExplorer*) is able to detect the four functional clusters where misattributions of some areas are explained by their multimodal function. Results have been discussed in terms of misattributions of brain areas to the different clusters emphasizing connections which are explainable (or not) by a cognitive point of view.

The proposed method is able to identify communities in different networks, and we showed that it is capable to identify communities also in dynamic networks. Moreover, we demonstrated that our algorithm is a multi-levels solution that can be used to capture the hierarchical community from a local viewpoint at any resolution. Experimental results on the real-world and synthetic datasets shown that our algorithm achieves good performance. Concluding we have shown that the adaptation of two heuristics, namely *IDA+LTE* and *Double Pruning* are very effective to allow the method to be independent of the two parameters  $m$  and  $\alpha$  showing that it is competitive with respect to other algorithms (see the review [90]).

## Growing networks

We presented a new model of the evolutionary dynamics and the growth of on-line social networks. The model emulates people's strategies for acquiring information in social networks, emphasising the local subjective view of an individual and what kind of information the individual can acquire when arriving in a new social context. The model proceeds through two phases: (a) a discovery phase, in which the individual becomes aware of the surrounding world and (b) an elaboration phase, in which the individual elaborates locally the information through a cognitive-inspired algorithm. We showed that model generated networks reproduce main features of both theoretical and real-world networks, such as high clustering coefficient, low characteristic path length, strong division in communities, and variability of degree distributions.

## Epidemic Spreading

We studied the influence of global, local and community-level risk perception on the extinction probability of a disease in several models of social networks. In particular, we studied the infection progression as a susceptible-infected-susceptible (SIS) model on several modular networks, formed by a certain number of random and scale-free

---

communities. We found that in the scale-free networks the progression is faster than in random ones with the same average connectivity degree. For what concerns the role of perception, we found that the knowledge of the infection level in one's own neighborhood is the most effective property in stopping the spreading of a disease, but at the same time the more expensive one in terms of the quantity of required information, thus the cost/effectiveness optimum is a tradeoff between several parameters. Moreover we have developed a self-organized percolation method for in both single-layer and multiplex networks for the SIS model considering an important factor of the modern society: the fact that our information no more comes mainly from physical contacts nor from broadcast media, but rather from our "virtual" contact network. Here we demonstrated that when the differences between the virtual and real networks is too elevated there are no finite values of precaution parameter  $J_c$  in order to avoid the spreading of a certain disease.

# Appendix A

## Algorithms pseudocode

---

**Algorithm 1** pseudocode for dynamic networks

---

```
1:  $m = 0.25, \alpha = 1.25$ 
2: for  $T = 1$  to  $T_{max}$  do
3:    $A(i, j, t) \rightarrow M(i, j, t) \rightarrow S(i, j, t)$  From Adjacency index to state network
4:   for  $i = 1$  to  $N$  do
5:     for  $j = 1$  to  $N$  do
6:       if  $S(i, j, t) > 0$  then
7:          $C(i, j, t)$   $i$  belong to community  $j$  with probability  $S(i, j, t)$ 
8:          $p(i, j) = [S(i, j, t) + S(i, j, t - 1)]$ 
9:       end if
10:    end for
11:  end for
12:   $A(i, j, t)_{weight} = A(i, j, t) + A(i, j, t - 1)$ 
13: end for
14:  $A_{weight}/T_{max}$  Weighted network
15:  $p/T_{max}$  Community matrix
```

---

---

**Algorithm 2** pseudocode for double pruning heuristic

---

```

1:  $T, T_{max}, t_b = 0, t_1 = 0, c = 0, a = 0$ 
2:  $m = 0.2, \alpha = 1.4$ 
3: while  $t < T$  do
4:    $t = t + 1; b = b + 1;$ 
5:   if  $b < 2$  then
6:      $t_{max} = b * 50$ 
7:   else
8:      $t_{max} = b * 10$ 
9:   end if
10:  for  $t_1 = 1 : t_{max}$  do
11:     $a = a + 1;$ 
12:    if  $a = b$  then
13:      we evaluate the state vector of each node, then  $S^1$  is vector state of node 1,  $S^2$  of node 2 and so
        on...
14:      for  $i = 1 : N$  do
15:         $dx^i = (S_{max}^i - S_{min}^i)/3$ 
16:        for  $j = 1 : N$  do
17:          if  $S_{i,j} < S_{min}^i + 2 * dx^i$  then
18:             $S_{i,j} = 0$ 
19:          end if
20:        end for
21:      end for
22:      we normalize  $S$ 
23:    else
24:       $S_1 = mS_0 + (1 - m)M'S_0, S_2 = S_1^\alpha$ 
25:      we normalize  $S$ 
26:       $a=0$ 
27:    end if
28:  end for
29:   $S = I(N)$ 
30: end while

```

---

# Bibliography

- [1] W. W. Zachary, “An information flow model for conflict and fission in small groups.,” *Journal of Anthropological Research*, no. 33, p. 452–473, 1977.
- [2] L. Euler, “Solutio problematis ad geometriam situs pertinentis,” *Commentarii academiae scientiarum Petropolitanae*, vol. 8, pp. 128–140, 1741.
- [3] D. König, “Theorie der endlichen und unendlichen Graphen,” Leipzig, 1936.
- [4] B. Bollobas, *Random Graphs*. Cambridge University Press, 2001.
- [5] D. J. Watts and S. H. Strogatz, “Collective dynamics of small-world networks,” *Nature*, vol. 393, no. 6684, pp. 440–442, 1998.
- [6] A. L. Barabasi and R. Albert, “Emergence of scaling in random networks,” *Science*, vol. 286, no. 5439, pp. 509–512, 1999.
- [7] A. Barrat, M. Barthlemy, and A. Vespignani, *Dynamical Processes on Complex Networks*. New York, NY, USA: Cambridge University Press, 2008.
- [8] R. Albert, H. Jeong, and A. L. Barabasi, “The diameter of the world wide web,” *Nature*, vol. 401, pp. 130–131, 1999.
- [9] H. Ebel, L. I. Mielsch, and S. Bornholdt, “Scale-free topology of e-mail networks,” *Physical Review E*, vol. 66, pp. 035103+, 2002.
- [10] S. Valverde, R. F. Cancho, and R. V. Solé, “Scale-free networks from optimal design,” *Europhys. Lett.*, vol. 60, no. 4, pp. 512–517, 2002.
- [11] R. Ferrer i Cancho, C. Janssen, and R. V. Solé, “The topology of technology graphs: Small world patterns in electronic circuits,” *Physical Review E*, vol. 64, no. 4, pp. 046119+, 2001.
- [12] R. Ferrer i Cancho and R. V. Solé, “The small world of human language,” *Proceedings of The Royal Society of London. Series B, Biological Sciences*, vol. 268, pp. 2261–2266, 2001.

- 
- [13] M. Newman, “Scientific collaboration networks. i. network construction and fundamental results,” *Physical Review E*, vol. 64, p. 016131, June 2001.
- [14] J. M. Montoya and R. V. Solé, “Small world patterns in food webs,” *Journal of Theoretical Biology*, vol. 214, no. 3, pp. 405–412, 2002.
- [15] H. Jeong, B. Tombor, R. Albert, Z. N. Oltvai, and A. L. Barabási, “The large-scale organization of metabolic networks.,” *Nature*, vol. 407, no. 6804, pp. 651–654, 2000.
- [16] P. Erdős and A. Rényi, “On the evolution of random graphs,” in *Publication of the Mathematical Institute of the Hungarian Academy of Sciences*, pp. 17–61, 1960.
- [17] M. E. J. Newman, “The structure and function of complex networks,” *SIAM Review*, vol. 45, pp. 167–256, 2003.
- [18] F. Karinthy, *Chains. Everything is different*. Budapest, 1929.
- [19] S. Milgram, “The Small World Problem,” *Psychology Today*, vol. 2, pp. 60–67, 1967.
- [20] R. Albert and A. L. Barabási, “Statistical mechanics of complex networks,” *Rev. Mod. Phys.*, no. 74, p. 47, 2002.
- [21] E. Ravasz and A. L. Barabási, “Hierarchical organization in complex networks,” *Physical Review E*, vol. 67, no. 2, p. 026112, 2003.
- [22] R. Dunbar, “Neocortex size as a constraint on group size in primates,”
- [23] S. Wasserman and K. Faust, *Social Networks Analysis*. University Press, Cambridge, England, 1994.
- [24] J. Scott, *Social Networks Analysis: A Handbook*. London: Sage, 2nd, ed., 2000.
- [25] S. N. Dorogovtsev and J. F. F. Mendes, *Evolution of Networks*. Oxford University Press, Oxford, 2003.
- [26] S. H. Strogatz, “Exploring complex networks,” *Nature*, vol. 410, pp. 268–276, Mar. 2001.
- [27] R. Albert and A. L. Barabási, “Statistical mechanics of complex networks,” *Rev. Mod. Phys.*, no. 74, p. 47, 2002.
- [28] U. Brandes, D. Delling, M. Gaertler, R. Gorke, M. Hoefer, Z. Nikoloski, and D. Wagner, “On finding graph clusterings with maximum modularity,” in *Proceedings of the 33rd International Workshop on Graph-Theoretical Concepts in Computer Science (WG’07)*, 2006.



- [29] M. Blatt, S. Wiseman, and E. Domany, “Superparamagnetic clustering of data.,” *Phys. Rev. E*, no. 76, p. 3251–3254, 1996.
- [30] M. E. J. Newman and M. Girvan, “Finding and evaluating community structure in networks.,” *Phys. Rev. E*, no. 69, p. 026113, 2004.
- [31] M. Newman, “Detecting community structure in networks.,” *Europ. Phys. J. B*, no. 38, pp. 331–330, 2004.
- [32] G. Palla, I. Derényi, I. Farkas, and T. Vicsek, “Uncovering the overlapping community structure of complex networks in nature and society,” *Nature*, vol. 435, pp. 814–818, June 2005.
- [33] A. Lancichinetti, S. Fortunato, and J. Kertész, “Detecting the overlapping and hierarchical community structure of complex networks.,” *New Journal of Physics*, vol. 033015, no. 11, 2009.
- [34] G. Gigerenzer and G. Goldstein, “Models of ecological rationality: The recognition heuristic.,” *Psyc. Rev.*, vol. 109, no. 1, p. 75–90, 2002.
- [35] G. Gigerenzer and W. Gaissmaier, “Heuristic decision making,” *Ann. Rev. of Psyc.*, no. 62, pp. 451–482, 2011.
- [36] V. Nicosia, F. Bagnoli, and V. Latora, “Impact of network structure on a model of diffusion and competitive interaction,” *EPL*, vol. 94, no. 68009, 2011.
- [37] J. Murray, *Mathematical biology*. No. v. 1 in Interdisciplinary applied mathematics, Springer, 2002.
- [38] J. E. de Freitas, L. Lucena, L. da Silva, and H. Hilhorst, “Critical behavior of a two-species reaction-diffusion problem,” *Phys. Rev. E.*, no. 61, p. 6330–6336, 2000.
- [39] E. Tulving, D. L. Schacter, and H. A. Stark, “Priming effects in word fragment completion are independent of recognition memory,” *Journ. Exp. Psyc.: Learning Memory and Cognition*, vol. 8, no. 4, 1982.
- [40] K. I. Forster and C. Davis, “Repetition priming and frequency attenuation,” *Journ. Exp. Psyc.: Learning Memory and Cognition*, vol. 10, no. 4, 1984.
- [41] D. Lusseau, K. Schneider, O. J. Boisseau, P. Haase, E. Slooten, and S. M. Dawson *Behavioral Ecology and Sociobiology*, no. 54, p. 396–405, 2003.
- [42] M. Rosvall and C. Bergstrom, “An information-theoretic framework for resolving community structure in complex networks,” *Proceedings of the National Academy of Sciences*, vol. 104, no. 18, p. 7327, 2007.

- 
- [43] S. V. Dongen, *Graph Clustering by Flow Simulation*. PhD thesis, University of Utrecht, 2000.
- [44] M. Rosvall and C. T. Bergstrom, “Maps of random walks on complex networks reveal community structure,” *PNAS*, vol. 105, no. 4, pp. 1118–1123, 2008.
- [45] M. Girvan and M. Newman, “Community structure in social and biological networks.,” *PNAS*, no. 99, pp. 7821–7826, 2002.
- [46] A. Lancichinetti and S. Fortunato, “Community detection algorithms: a comparative analysis,” 2009. cite arxiv:0908.1062Comment: 12 pages, 8 figures. The software to compute the values of our general normalized mutual information is available at <http://santo.fortunato.googlepages.com/inthepress2>.
- [47] M. Conti, S. K. Das, B. C, M. Kumar, L. M. Ni, A. Passarella, G. Roussos, G. Troster, G. Tsudik, and Z. F, “Looking ahead in pervasive computing: Challenges and opportunities in the era of cyber–physical convergence,” *Pervasive and Mobile Computing*, vol. 8, no. 1, pp. 2 – 21, 2012.
- [48] L. Pelusi, A. Passarella, and M. Conti, “Opportunistic networking: data forwarding in disconnected mobile ad hoc networks,” *IEEE Communications Magazine*, vol. 44, no. 11, pp. 134–141, 2006.
- [49] S. Fortunato, “Community detection in graphs,” *Physics Reports*, vol. 486, no. 3–5, pp. 75 – 174, 2010.
- [50] A. Lancichinetti, F. Radicchi, J. J. Ramasco, and S. Fortunato, “Finding statistically significant communities in networks,” *PLoS ONE*, vol. 6, no. 4, pp. e18961+, 2011.
- [51] M. Rosvall and C. T. Bergstrom, “Multilevel compression of random walks on networks reveals hierarchical organization in large integrated systems,” 2011.
- [52] M. Sales-Pardo, R. Guimerà, A. A. Moreira, and L. A. N. Amaral, “Extracting the hierarchical organization of complex systems,” *PNAS*, vol. 104, no. 39, pp. 15224–15229, 2007.
- [53] V. D. Blondel, J. Guillaume, R. Lambiotte, and E. Lefebvre, “Fast unfolding of communities in large networks,” *JSTAT*, vol. 2008, no. 10, pp. P10008+, 2008.
- [54] U. N. Raghavan, R. Albert, and S. Kumara, “Near linear time algorithm to detect community structures in large-scale networks,” *Phys. Rev. E*, vol. 76, no. 3, pp. 036106+, 2007.

- [55] P. Hui, E. Yoneki, S. Chan, and J. Crowcroft, "Distributed community detection in delay tolerant networks," in *MobiArch*, 2007.
- [56] T. Hossmann, T. Spyropoulos, and F. Legendre, "Know thy neighbor: Towards optimal mapping of contacts to social graphs for dtn routing," in *INFOCOM, 2010 Proceedings IEEE*, pp. 1–9, 2010.
- [57] E. Borgia, M. Conti, and A. Passarella, "Autonomic detection of dynamic social communities in opportunistic networks," in *The 10th IFIP Annual Mediterranean Ad Hoc Networking Workshop*, 2011.
- [58] E. Massaro, F. Bagnoli, A. Guazzini, and P. Lió, "Information dynamics algorithm for detecting communities in networks," *CNSNS*, vol. 17, no. 11, pp. 4294 – 4303, 2012.
- [59] F. Bagnoli, E. Massaro, and A. Guazzini, "Community-detection cellular automata with local and long-range connectivity," in *ACRI2012: Conference of Cellular Automata for Research and Industry, Santorini, Greece. To appear in Springer-Verlag in the Lecture Notes in Computer Science (LNCS)*, 2012.
- [60] D. Borkmann, A. Guazzini, E. Massaro, and S. Rudolph, "A cognitive-inspired model for self-organizing networks," in *IEEE Sixth International Conference on Self-Adaptive and Self-Organizing Systems Workshops (SASOW)*, pp. 229–234, 2012.
- [61] C. Boldrini and A. Passarella, "Hcmm: Modelling spatial and temporal properties of human mobility driven by users social relationships," *Computer Communication*, vol. 33, no. 9, pp. 1056–1074, 2010.
- [62] S. Allen, M. Chorley, G. Colombo, and R. Whitaker, "Opportunistic social dissemination of micro-blogs," *Ad Hoc Networks*, vol. 10, no. 8, pp. 1570 – 1585, 2012.
- [63] A. Picu, T. Spyropoulos, and T. Hossmann, "An analysis of the information spreading delay in heterogeneous mobility dtns," in *World of Wireless, Mobile and Multimedia Networks (WoWMoM), 2012 IEEE International Symposium on a*, pp. 1–10, 2012.
- [64] O. Sporns, *Network of the brain*. The MIT Press, 2011.
- [65] O. Sporns, *Discovering the Human Connectome*. The MIT Press, 2012.
- [66] J. W. Scannell, C. Blakemore, and M. P. Young, "Analysis of connectivity in the cat cerebral cortex.," *PNAS*, vol. USA 99, p. 7821–7826, 2002.

- [67] J. W. Scannell, G. A. P. C. Burns, C. Hilgetag, M. O'Neill, and M. P. Young, "The connectional organization of the cortico-thalamic system of the cat.," *Cereb. Cortex*, vol. 9, pp. 277–299, 1999.
- [68] L. Zemanova, C. Zhou, and J. Kurths, "Structural and functional clusters of complex brain networks," *Physica D*, vol. 224, pp. 202–212, 2006.
- [69] C.-C. Hilgetag, G. A. P. C. Burns, M. A. O'Neill, J. W. Scannell, and M. P. Young., "Anatomical connectivity defines the organization of clusters of cortical areas in the macaque monkey and the cat.," *Phil Trans. R. Soc. Lond. B*, vol. 335, pp. 791–110, 2000.
- [70] C. Hilgetag and M. Kaiser, "Anatomical connectivity defines the organization of clusters of cortical areas in the macaque monkey and the cat.," *Neuroinformatics*, vol. 2, pp. 353–360, 2004.
- [71] S. Dehaene and J.-P. Changeoux, "Neural mechanism for access to consciousness," in *The Cognitive Neuroscience III* (M. S. Gazzaniga, ed.), Cambridge: The MIT Press, 2004.
- [72] L. R. Hoffman, *The Group problem solving process : studies of a valence model / edited by L. Richard Hoffman*. Praeger New York, 1979.
- [73] H. Simon, "A behavioral model of rational choice," *The Quarterly Journal of Economics*, vol. 69, no. 1, pp. 99–118, 1955.
- [74] L. Festinger, *A Theory of Cognitive Dissonance*. Stanford University Press, June 1957.
- [75] V. D. Blondel, J. L. Guillaume, R. Lambiotte, and E. Lefebvre, "Fast unfolding of communities in large networks," *Journal of Statistical Mechanics: Theory and Experiment*, vol. 2008, pp. P10008+, Oct. 2008.
- [76] J. Huang, H. Sun, Y. Liu, Q. Song, and T. Wenginger, "Towards online multiresolution community detection in large-scale networks," *PLoS ONE*, vol. 6, no. 8, p. e23829, 2011.
- [77] F. Radicchi, C. Castellano, F. Cecconi, V. Loreto, and D. Parisi, "Defining and identifying communities in networks," *PNAS*, vol. 101, p. 2658, 2004.
- [78] D. E. Knuth, *The Stanford GraphBase: A Platform for Combinatorial Computing*. Addison-Wesley, Reading, MA, 1993.
- [79] A. Lancichinetti, S. Fortunato, and J. Kertész, "Detecting the overlapping and hierarchical community structure in complex networks," *New Journal of Physics*, vol. 11, no. 3, p. 033015, 2009.

- [80] A. Lancichinetti, S. Fortunato, and F. Radicchi, “Benchmark graphs for testing community detection algorithms,” *Phys. Rev. E*, vol. 78, p. 046110, Oct. 2008.
- [81] “Internet world stats.” [www.internetworldstats.com/facebook.htm](http://www.internetworldstats.com/facebook.htm). Accessed: 1/2/2013.
- [82] A. L. Barabási and R. Albert, “Emergence of Scaling in Random Networks,” *Science*, vol. 286, pp. 509–512, Oct. 1999.
- [83] F. Papadopoulos, M. Kitsak, M. Serrano, M. Boguñá, and D. Krioukov, “Popularity versus Similarity in Growing Networks,” *Nature*, vol. 489, pp. 537–540, Sep 2012.
- [84] E. M. Jin, M. Girvan, and M. E. J. Newman, “Structure of growing social networks,” *Physical Review E*, vol. 64, no. 4, 2001.
- [85] A. Buscarino, L. Fortuna, M. Frasca, and A. S. Fiore, “A new model for growing social networks,” *Engineering. Complexity in*, vol. 0, pp. 103–105, 2010.
- [86] D. J. Watts, P. S. Dodds, and M. E. J. Newman, “Identity and search in social networks,” *Science*, vol. 296, pp. 1302–1305, 2002.
- [87] T. G. Brown, “Cognitive pruning in foreign language teaching,” *The Modern Language Journal*, vol. 56, no. 4, pp. 222–227, 1972.
- [88] L. A. Adamic and N. Glance, “The political blogosphere and the 2004 u.s. election: divided they blog,” in *Proceedings of the 3rd international workshop on Link discovery*, LinkKDD '05, (New York, NY, USA), pp. 36–43, ACM, 2005.
- [89] “The koblenz network collection.” [//konect.uni-koblenz.de/](http://konect.uni-koblenz.de/). Accessed: 05/05/2012.
- [90] S. Fortunato, “Community detection in graphs,” *Physics Reports*, vol. 486, no. 3-5, pp. 75 – 174, 2010.
- [91] M. Newman, *Networks: An Introduction*. New York, NY, USA: Oxford University Press, Inc., 2010.
- [92] R. Pastor-Satorras and A. Vespignani, “Epidemic dynamics and endemic states in complex networks,” *Phys. Rev. E*, vol. 63, p. 066117, May 2001.
- [93] M. Boguñá, R. Pastor-Satorras, and A. Vespignani, “Absence of epidemic threshold in scale-free networks with degree correlations,” *Phys. Rev. Lett.*, vol. 90, p. 028701, 2003.

- [94] F. Bagnoli, P. Liò, and L. Sguanci, “Risk perception in epidemic modeling,” *Phys. Rev. E*, vol. 76, p. 061904, Dec 2007.
- [95] R. Cohen, D. ben Avraham, and S. Havlin, “Percolation critical exponents in scale-free networks,” *Phys. Rev. E*, vol. 66, p. 036113, Sep 2002.
- [96] C. Moore and M. E. J. Newman, “Epidemics and percolation in small-world networks,” *Phys. Rev. E*, vol. 61, pp. 5678–5682, May 2000.
- [97] Z. Dezső and A.-L. Barabási, “Halting viruses in scale-free networks,” *Physical Review E*, vol. 65, pp. 055103+, May 2002.
- [98] M. Salathé and J. H. Jones, “Dynamics and control of diseases in networks with community structure,” *PLoS Comput Biol*, vol. 6, p. e1000736, 04 2010.
- [99] J. Chen, H. Zhang, Z.-H. Guan, and T. Li, “Epidemic spreading on networks with overlapping community structure,” *Physica A: Statistical Mechanics and its Applications*, vol. 391, no. 4, pp. 1848 – 1854, 2012.
- [100] A. Saumell-Mendiola, M. A. Serrano, and M. B. na, “Epidemic spreading on interconnected networks,” *arXiv:1202.4087*, 2012.
- [101] S. Fortunato and C. Castellano, “Community Structure in Graphs,” Dec. 2007.
- [102] C. Dorso and A. Medus, “Community detection in networks.,” *International Journal of Bifurcation and Chaos*, vol. 20, pp. 361–367.
- [103] F. Liljeros, C. R. Edling, L. A. Amaral, E. H. Stanley, and Y. Åberg, “The web of human sexual contacts,” *Nature*, vol. 411, pp. 907–908, June 2001.
- [104] R. M. Anderson and R. M. May, *Infectious Diseases of Humans Dynamics and Control*. Oxford University Press, 1992.
- [105] A. L. Barabási and R. Albert, “Emergence of scaling in random networks,” *Science*, vol. 286, pp. 509–512, 1999.
- [106] S. N. Dorogovtsev, J. F. F. Mendes, and A. N. Samukhin, “Structure of growing networks with preferential linking,” *Phys. Rev. Lett.*, vol. 85, pp. 4633–4636, Nov 2000.
- [107] S. Kitchovitch and P. Liò, “Risk perception and disease spread on social networks,” *Procedia Computer Science*, vol. 1, pp. 2339–2348, May 2010.
- [108] S. Boccaletti, V. Latora, Y. Moreno, M. Chavez, and D.-U. Hwang, “Complex networks : Structure and dynamics,” *Phys. Rep.*, vol. 424, pp. 175–308, Fervier 2006.

- [109] “[Pandemic Scares Throughout History](#),” *Health Magazine*, 2013.
- [110] S. Report, “The 1918 flu virus is resurrected,” *Nature*, vol. 437, pp. 794–795, 2005.
- [111] J. Kaiser, “Resurrected influenza virus yields secrets of deadly 1918 pandemic,” *Science*, vol. 310, no. 5745, pp. 28–29, 2005.
- [112] Wikipedia, “[Pandemic](#) — Wikipedia, the free encyclopedia,” 2013.
- [113] “[The 'false' pandemic: Drug firms cashed in on scare over swine flu, claims Euro health chief](#),” *Daily Mail*, 2010.
- [114] C. Moore and M. E. J. Newman, “Epidemics and percolation in small-world networks,” *Phys. Rev. E*, vol. 61, pp. 5678–5682, 2000.
- [115] R. Pastor-Satorras and A. Vespignani, “Epidemic spreading in scale-free networks,” *Phys. Rev. Lett.*, vol. 86, pp. 3200–3203, 2001.
- [116] M. E. J. Newman, “Exact solutions of epidemic models on networks,” Working Papers 01-12-073, Santa Fe Institute, Dec. 2001.
- [117] R. M. May and A. L. Lloyd, “Infection dynamics on scale-free networks,” *Phys. Rev. E*, vol. 64, p. 066112, 2001.
- [118] R. Pastor-Satorras and A. Vespignani, “Immunization of complex networks,” *Phys. Rev. E*, vol. 65, p. 036104, 2002.
- [119] Wikipedia, “[Lazaretto](#) — Wikipedia, the free encyclopedia,” 2013.
- [120] R. Palmer, *L'azione della repubblica di Venezia nel controllo della peste. Lo sviluppo della politica governativa, Venezia e la peste 1348–1797*. Venice (Italy): Marsilio Editori, 1979.
- [121] F. Bagnoli, P. Liò, and L. Sguanci, “Risk perception in epidemic modeling,” *Phys. Rev. E*, vol. 76, p. 061904, 2007.
- [122] J. Ginsberg, M. Mohebbi, R. Patel, L. Brammer, M. Smolinski, and L. Brilliant, “Detecting influenza epidemics using search engine query data,” *Nature*, vol. 457, pp. 1012–1014, 2009.
- [123] D. Scanfeld, V. Scanfeld, and E. L. Larson, “Dissemination of health information through social networks: Twitter and antibiotics,” *American Journal of Infection Control*, vol. 38, no. 3, pp. 182 – 188, 2010.
- [124] C. Chew and G. Eysenbach, “Pandemics in the age of twitter: Content analysis of tweets during the 2009 h1n1 outbreak,” *PLoS ONE*, vol. 5, p. e14118, 11 2010.

- 
- [125] “The State of the News Media,” *The Pew Research Center’s project for Excellence in Journalism*, 2010.
- [126] M. Kurant and P. Thiran, “Layered complex networks,” *Phys. Rev. Lett.*, vol. 96, p. 138701, 2006.
- [127] P. Mucha, T. Richardson, K. Macon, M. Porter, and J.-P. Onnela, “Community structure in time-dependent, multiscale, and multiplex networks,” *Science*, vol. 328, no. 5980, pp. 876–878, 2010.
- [128] M. Szell, R. Lambiotte, and S. Thurner, “Multirelational Organization of Large-scale Social Networks in an Online World,” 2010.
- [129] “Evolution of cooperation in multiplex networks.,” *Scientific reports*, vol. 2, 2012.
- [130] G. Bianconi, “Statistical mechanics of multiplex networks: Entropy and overlap,” *Phys. Rev. E*, vol. 87, p. 062806, 2013.
- [131] “Catastrophic cascade of failures in interdependent networks,” *Nature*, vol. 464, no. 7291, pp. 1025–1028, 2010.
- [132] J. Gao, S. V. Buldyrev, H. E. Stanley, and S. Havlin, “Networks formed from interdependent networks,” *Nat Phys*, vol. 8, no. 1, pp. 40–48, 2012.
- [133] C. Granell, S. Gómez, and A. Arenas, “Dynamical interplay between awareness and epidemic spreading in multiplex networks,” *Phys. Rev. Lett.*, vol. 111, p. 128701, Sep 2013.
- [134] F. Bagnoli, P. Palmerini, and R. Rechtman, “Algorithmic mapping from criticality to self-organized criticality,” *Phys. Rev. E*, vol. 55, pp. 3970–3976, Apr 1997.
- [135] A. Reka and A. L. Barabási, “Statistical mechanics of complex networks,” *Rev. Mod. Phys.*, vol. 74, pp. 47–97, June 2002.
- [136] E. Domany and W. Kinzel, “Equivalence of cellular automata to ising models and directed percolation,” *Phys. Rev. Lett.*, vol. 53, pp. 311–314, 1984.