

UNIVERSITY OF FRIBOURG

MASTER THESIS

Giant connected component in networks

Author:

Benoît RICHARD

Supervisors:

Dr. Guiyuan SHI

Prof. Yi-Cheng ZHANG

Theoretical Interdisciplinary Physics Group
Department of Physics

October 7, 2018

University of Fribourg

Abstract

Faculty of Science
Department of Physics

Master Project

Giant connected component in networks

by Benoît RICHARD

Write the abstract

Contents

Abstract	iii
Contents	v
1 Single layer networks	1
1.1 Introduction	1
1.2 Configuration model	1
1.3 Self-edges and multi-edges	6
1.4 Uniformity of network space sampling	7
1.5 Excess degree distribution	8
1.6 Generating functions	9
1.7 Erdos-Renyi networks	10
1.8 Geometric networks	10
1.9 Scale-free networks	11
1.10 Giant connected component	12
1.10.1 Size of the GCC	12
1.10.2 Algorithm to find the connected components	15
1.10.3 Degree distribution in the GCC	16
1.11 Generating connected networks	16
1.11.1 Algorithm	16
1.11.2 Erdos-Renyi reconstruction	18
1.11.3 Real world networks	19
2 Multiplex networks	21
2.1 Introduction	21
2.2 Giant viable cluster	22
2.2.1 Size of the GVC	22
2.2.2 Algorithm to find the viable clusters	24
2.2.3 Single parameter multiplex networks	24
2.3 Boundary condition	25
2.4 Interval arithmetic	28
2.4.1 Motivation	28
2.4.2 Theoretical foundation	28
2.4.3 Algorithm to test the discontinuity of the phase transition	30
2.4.4 Algorithm to find the critical region	30
2.5 Results	31
A Fixpoint iteration for connected networks generation	35
A.1 Convergence of the fixpoint iteration	35
A.2 Monotonicity of $\mu(z)$	35
Bibliography	37

List of Symbols

Single layer networks

Symbol	Description	Math. definition
c	Expectation value for the degree	$\mathbb{E}[\deg v]$
C	A connected component	
$\deg v$	Degree of vertex v	
$\mathbb{E}[\dots]$	Expectation value	
$g_0(z)$	Generating function for the degree of uniformly chosen nodes	$\sum_{k=0}^{\infty} p_k z^k$
$g_1(z)$	Generating function for the degree of nodes reached by following an edge	$\sum_{k=0}^{\infty} q_k z^k$
k_i	Degree of vertex i	
m	Number of edges in the network	
n	Number of nodes in the network	
$N(v)$	Neighborhood of a vertex v	
p_{ij}	Probability that vertices i and j are connected	
p_k	Probability that a random node has degree k	$P_0(\deg v = k)$
$P_0(\dots)$	Probability starting from a uniformly chosen node	
$P_1(\dots)$	Probability starting from a node reached by following an edge	
q_k	Probability that a node reached by following an edge has degree $k + 1$	$P_1(\deg v = k + 1)$
r_k	Probability that a random node of the GCC has degree k	$P_0(\deg v = k v \in \text{GCC})$
S	Fraction of the network which is part of the GCC in the large n limit	$P_0(v \in \text{GCC})$
u	Probability that a node reached by following an edge is not part of the GCC	$P_1(v \notin \text{GCC})$
v	Random variable representing a vertex chosen in a network, either uniformly or by following an edge depending of the context	
α	Exponent of a power law distribution	
$\zeta(\alpha)$	Riemann zeta function	

Multiplex networks

For multiplex networks, the convention we follow to indicate a value refer to layer i is to add an upper index (i) if the quantity has a lower index and to add a lower index i otherwise.

Symbol	Description	Math. definition
L	Number of layers of the multiplex network	
N	Number of parameters determining the degree distributions of a multiplex network	
$N_i(v)$	Neighborhood of a vertex v in layer i	
$F(\boldsymbol{\lambda}, \mathbf{u})$	Residual function for \mathbf{u}	
$F_{\boldsymbol{\lambda}}(\mathbf{u})$	Residual function for \mathbf{u} with fixed parameters $\boldsymbol{\lambda}$	
$g_0^{(i)}(z)$	Generating function for the degree of uniformly chosen nodes in layer i	$\sum_{k=0}^{\infty} p_k^{(i)} u_i^k$
$g_1^{(i)}(z)$	Generating function for the degree of nodes reached by following an edge in layer i	$\sum_{k=0}^{\infty} q_k^{(i)} u_i^k$
$J(\boldsymbol{\lambda}, \mathbf{u})$	Jacobi matrix of $F_{\boldsymbol{\lambda}}$	
$p_k^{(i)}$	Probability that a uniformly chosen node has degree k in layer i	$P_0^{(i)}(\deg v = k)$
$P_0^{(i)}(\dots)$	Probability for an event in layer i starting from a uniformly chosen node	
$P_1^{(i)}(\dots)$	Probability for an event in layer i starting from a node reached by following an edge in layer i	
$q_k^{(i)}$	Probability that a node reached by following an edge has degree $k + 1$	$P_1^{(i)}(\deg v = k + 1)$
\mathcal{R}	Critical region	
u_i	Probability that a node reached by following an edge in layer i is not part of the GVC	$P_1^{(i)}(v \notin GVC)$
\mathbf{u}	Vector of all u_i	(u_1, u_2, \dots, u_L)
\mathbf{u}_T	Trivial solution for \mathbf{u}	$(1, 1, \dots, 1)$
\mathbf{u}^\dagger	A non trivial solution for \mathbf{u}	
λ_j	One of the N parameters determining the degree distributions of a multiplex network	
$\boldsymbol{\lambda}$	Parameter vector	$(\lambda_1, \lambda_2, \dots, \lambda_N)$
$\boldsymbol{\lambda}^c$	Parameter vector in the critical region	
$\psi(\boldsymbol{\lambda}, \mathbf{u})$	Fixpoint function for \mathbf{u}	

Todo list

Write the abstract	iii
Think about that and the need for that for GCC calculation	6
The following calculation is not correct, nodes with the same degree are indistinguishable as well. Correct and note the disagreement with Newman	7
Missing reference	8
Ask Guiyuan which version to use and if there is a some motivation in using this geom. dist.	10
Add examples	11
Credit the guy the polylog code comes from. Maybe say a bit more about polylog/reimplement polylog using intergation	11
Prove uniqueness of GCC	12
Speak about phase transition	14
Add something about the fact we will refer to u and S thinking about the GCC	14
Add the position of the critical point	16
Do something with with this information	16
Move this to the general introduction ? Put the actual general introduction into the single layer network chapter ?	21
Move introduction of λ parameters in this section	22
Correct the following explanation, it is going nowhere	22
Find a ref ?	22
Whaaat ?	22
Verify ref	24
Write something smart to conclude the section	24
Write. Make sure to introduce small components problem	24
Write some introduction	25
Ref this schemes	27
Missing reference	28
Missing reference	29
ref the equation	29
Introduce Krawczyk operator as it is used below. Maybe in appendix.	29
Some conclusion	30
Speak more about Krawczyk and Newton ?	30
Find back which one exactly	30
Missing reference	30
Choose what multiplex networks should be used and with what parameters and actually produce the results.	31
Do other than ER, put label with number of layers add simulations to plot	32
Write corresponding equation, reference figure somewhere, put it somewhere better	32
Make the plot more readable and less ugly	33
Find why the two methods seems drift away one from the other far from the center.	33

Chapter 1

Single layer networks

layer networks)? 1.1 Introduction

Many systems in real world can be conceptually represented as objects being connected to others. Such representation is called a *network*. For example, a power grid can be represented as stations connected by power lines, as shown in fig. 1.1 for the power grid of western USA. The concept of network does not require the object or the links between them to be physical. As an example, we can represent collaboration as a network: two people are connected if they did collaborate. Such network is shown in fig. 1.2 for collaboration of jazz musicians. This implies that a very different systems can be represented as networks, from cities connected by road (see fig. 1.3 for E-road network) to protein connected if they interact (see fig. 1.4 for a network of human proteins). Insight on fundamental properties of networks may shed light on a very broad range of problems. To gain such insights, theoretical studies of general networks, such as the one presented in this thesis, is required.

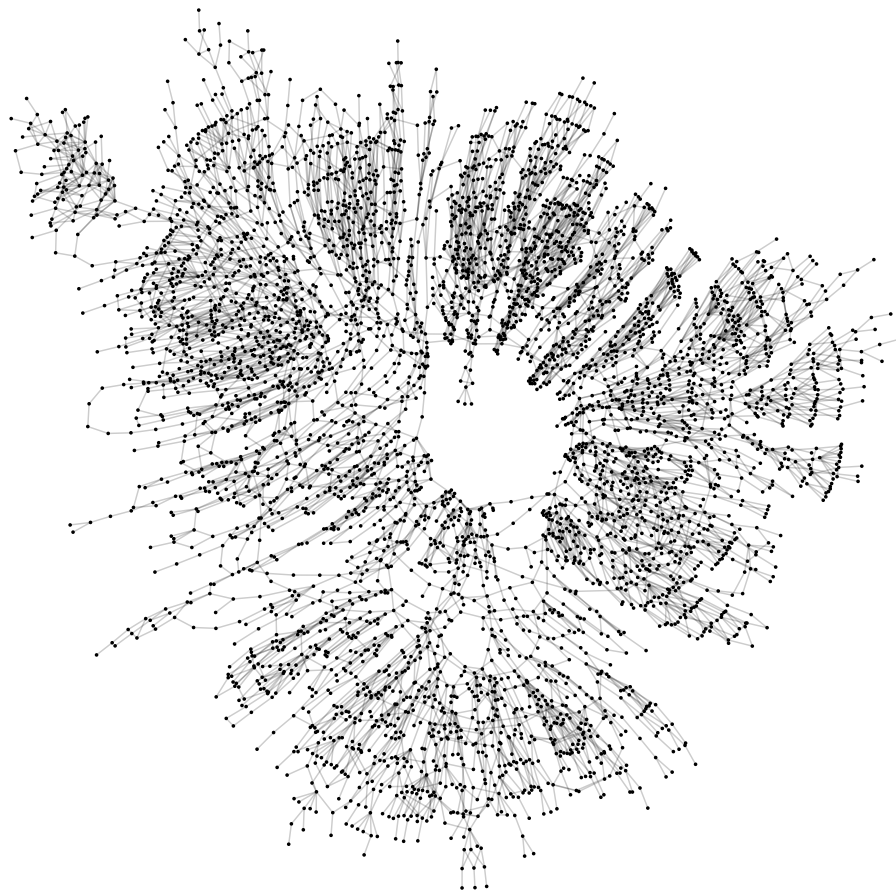
Mathematically, networks are represented as *graphs*. A graph is an object composed of a set V of *vertices* (also referred to as nodes) and a set of *edges* E . An edge is characterized by the fact that it connects two vertices together, which in mathematical terms translate to the fact that an edge can be written as a pair of vertices or equivalently $E \subset \{(v_1, v_2) | v_1, v_2 \in V\}$. Many extensions of this model exist, for example edges may have a direction (*directed graph*), implying that $(v_1, v_2) \neq (v_2, v_1)$, or edges can carry a value (*weighted graph*).

In this first chapter we focus on the standard unweighted and undirected networks, even restricting our theoretical analysis to the case of infinitely large networks as it allows for several convenient simplifications while still capturing important feature networks. In Chapter 2 we study an extension of this network category, multiplex networks, which consists in allowing distinct type of edges between the vertices.

1.2 Configuration model

figuration model) Since we are interested in fundamental properties of networks, we need to abstract from the specificity of one network. To do so, we consider that networks are fully determined by their *degree distribution* $\{p_k\}$ where p_k is the probability for a node chosen randomly and uniformly to have degree k . Since knowing how a network can be constructed is useful both conceptually and to perform computation on properties of the network, we now present an algorithm called *configuration model* [2] that sample uniformly the space of all network with a given degree distribution.

Consider a network with n vertices with a given degree distribution $\{p_k\}$. If we cut every edges in two, every vertex keep a number of *stubs* (half-edges) equal to



western US powergrid)

FIGURE 1.1: Power grid network of the Western States of USA. Nodes represent electrical stations (generator, transformer, substation) and edges represent power supply lines. Data retrieved from the Konect database [1] (Konect code UG).

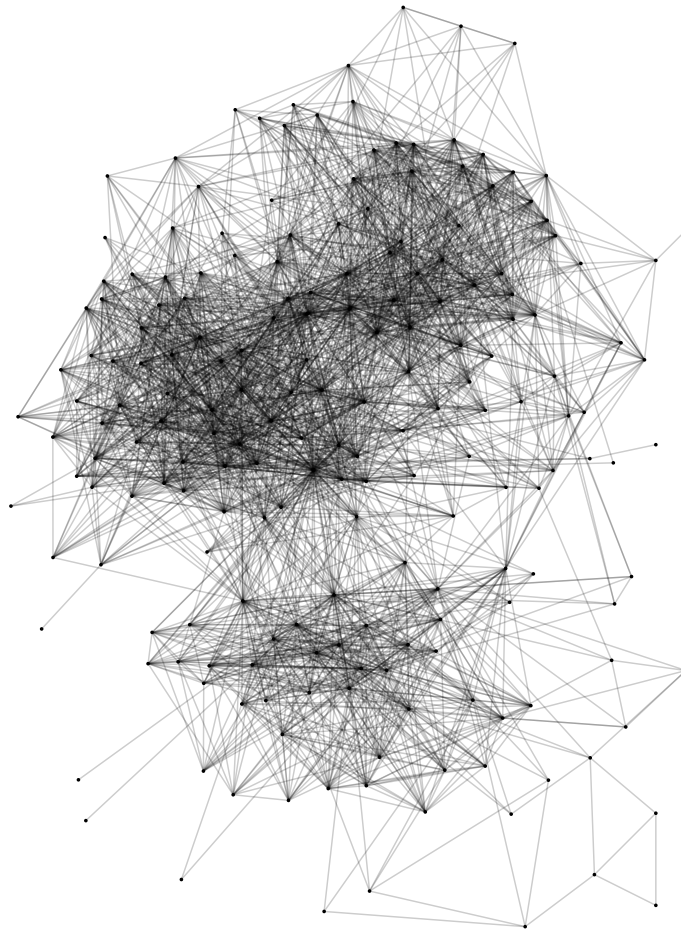
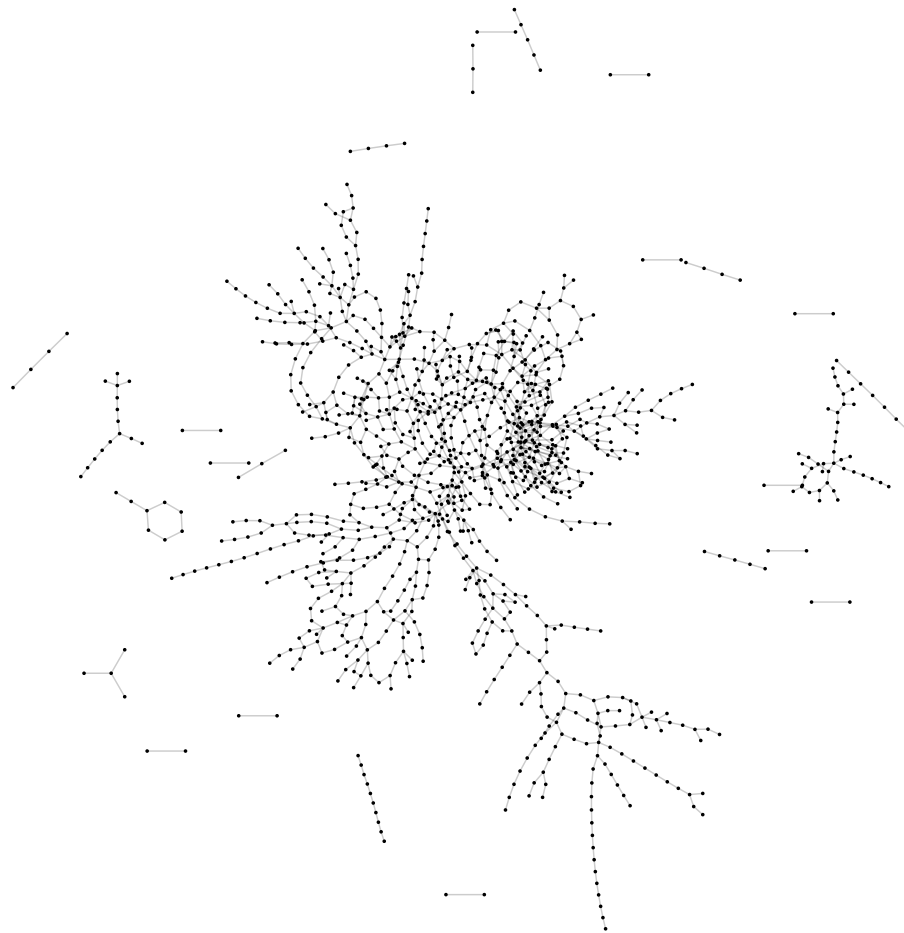


FIGURE 1.2: Collaboration network of jazz musicians as of 2003. Nodes represent musicians and edges represent the fact that the two musicians had played in the same band. Data retrieved from the Konect database [1] (Konect code JZ).



e: Network euroroad)

FIGURE 1.3: International E-road network, a network of road situated mainly in Europe. Nodes represent cities and edges represent E-roads connecting them. Data retrieved from the Konect database [1] (Konect code ET).

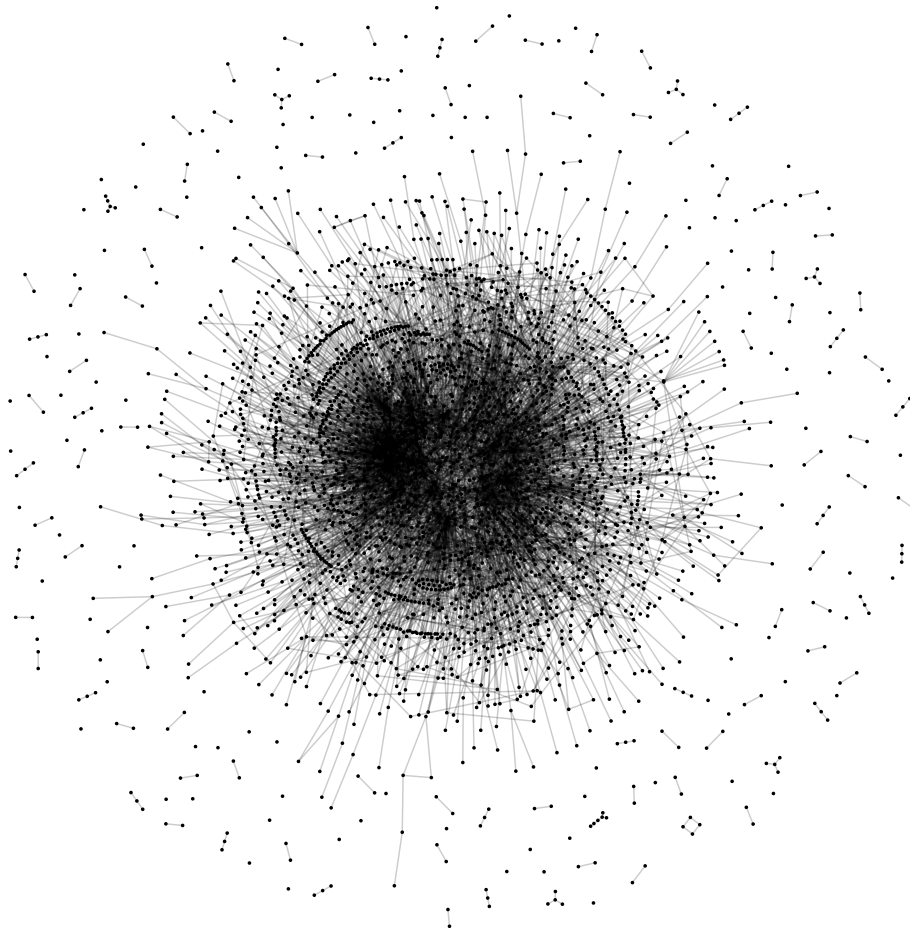


FIGURE 1.4: Network of protein-protein interaction in humans. Nodes represent proteins and edges a binary interaction. Data retrieved from the Konekt database [1] (Konekt code MV).

human proteins)

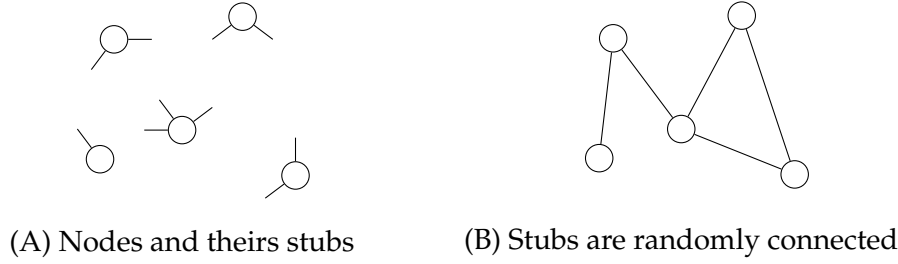


FIGURE 1.5: Schematic representation of the configuration model.

its degree. The resulting set of vertices and stubs is independent of the network structure, but common for all networks with the same degree distribution. The idea of this algorithm is thus to start from this state, a set of nodes with stub degree distribution $\{p_k\}$. Then each stub is bonded to another chosen uniformly amongst other stubs to form all edges of the network. By construction, the produced network has degree distribution $\{p_k\}$.

The fact that stubs are paired uniformly and independently is important in that it implies that the vertex reached by following a random edge does not depend, in probability, on the vertex at the other end of the edge.

Think about that and the need for that for GCC calculation

1.3 Self-edges and multi-edges

Using the insight given by the algorithm, we can compute the probability that a node is connected to itself, thus making a so-called *self-edge*. In a network with m edges, there are $2m$ stubs. The probability p_{ii} that a stub of vertex i connect to another stub of the same vertex is thus

$$p_{ii} = m \frac{\binom{k_i}{2}}{\binom{2m}{2}} = \frac{k_i(k_i - 1)}{2(2m - 1)}, \quad (1.1) \{?\}$$

where k_i is the degree of vertex i . In the limit of large n , the number of vertices with degree k is equal to np_k and as a consequence the total number m of edges is equal to

$$m = \frac{1}{2} \sum_{k=0}^{\infty} np_k k = \frac{n}{2} \mathbb{E}[\deg v], \quad (1.2) \{?\}$$

with $\mathbb{E}[\dots]$ denoting the expectation value. The average number of self-edges is therefore asymptotically constant as n becomes large, so the fraction of vertices having self edges goes to zeros as n grows and we can safely consider the generated network has no self-edges at all.

Similarly, we find that the probability p_{ij} that two vertices i and j are connected is equal to

$$p_{ij} = m \frac{k_i k_j}{\binom{2m-2}{2}} = \frac{k_i k_j}{2m-1}. \quad (1.3) \{?\}$$

The probability to have two or more edges between the vertices i and j is equal to the probability that i and j are connected and that they remain so after we remove one edge between them. The probability for them to be connected with one edge less is the same as p_{ij} but with one edge less in total and one stub less at both i and j , giving

$$\frac{(k_i - 1)(k_j - 1)}{2m - 3}. \quad (1.4) \{?\}$$

In consequence we find the probability to have at least two edges between i and j to be

$$p_{ij} \frac{(k_i - 1)(k_j - 1)}{2m - 3}, \quad (1.5) \{?\}$$

giving the average number of multi-edges

$$\frac{\sum_{ij} k_i k_j (k_j - 1)(k_i - 1)}{2(2m - 1)(2m - 3)} \approx \frac{1}{8m^2} \sum_i k_i (k_i - 1) \sum_j k_j (k_j - 1) \quad (1.6) \{?\}$$

$$= \frac{1}{2} \left[\frac{\mathbb{E}[(\deg v)^2] - \mathbb{E}[\deg v]}{\mathbb{E}[\deg v]} \right]^2. \quad (1.7) \{?\}$$

The approximation arise as in the limit of large n we have $2m - 3 \approx 2m - 3 \approx 2m$ as m scale proportionally to n . As in the case of self-edges, the number of multi-edges is asymptotically constant and we can therefore consider that the generated networks has no multi-edges.

Since $\{p_k\}$ is a probability distribution, it is independent of the number of nodes n of the network. Therefore for any degree distribution, we can consider the limit for large n , which we do as it allows several mathematical simplifications of the problem as outlined below. For sufficiently large networks, the difference between the large n limit and the actual network is small and can thus safely be neglected.

A network having this two properties, absence of self-edges and of multi-edges, is said to be a *simple graph*. Since for n large enough all networks in the context of the configuration model have approximately these two properties, we always consider that the networks are simple graphs in the remaining of this thesis.

1.4 Uniformity of network space sampling

In the previous section we demonstrated that the configuration model asymptotically produces simple graphs in the limit of large n . This is necessary to prove the claims we make in Section 1.2 that the configuration uniformly sample the space of networks with a given degree distribution. To prove that claim is equivalent to prove that each different network appears with the same probability.

The following calculation is not correct, nodes with the same degree are indistinguishable as well. Correct and note the disagreement with Newman

A *matching* of the stubs is a possible way to bond $2m$ distinguishable stubs. By construction the pair of stubs that are bound are chosen uniformly and thus each matching appears with the same probability than any other. However, in practice the stubs of a single vertex are not distinguishable and several matching leads to the same network. Therefore there is $k!$ way of arranging the stubs of a vertex of degree

k , resulting in a total of

$$\prod_{i=1}^n k_i! \quad (1.8) \{?\}$$

matchings corresponding to the same network, where k_i is the degree of vertex i . This quantity only depends on the degree distribution, thus all networks with the same degree distribution have the same number of different matchings. Since each matching appears with the same probability, we can conclude that each network with a given degree distribution appears with the same probability in the configuration model.

In a sense the configuration model is therefore optimal for the point of view we adopt in this thesis. Indeed we consider a network fully determined by its degree distribution and all network with a common degree distribution are equal with regard to the configuration model as it sample them uniformly. We could in fact present the configuration model as a consequence of the requirement that we want a sampling method introducing no unnecessary constrain on the produced networks.

See Section 15.2 of [2] or [?] for discussions of that point of view.

sing ref: maximum entropy paper

1.5 Excess degree distribution

As we will see below, while we consider that a network is fully determined by its degree distribution, considering vertices reached by following an edge gives valuable insights on the network structure. We call such vertex a *first neighbor* vertex and we denote $P_1(\dots)$ the probability associated with a first neighbor, while we denote $P_0(\dots)$ the probability associated with uniformly chosen vertices¹. We can define the *excess degree distribution* $\{q_k\}$ as

$$q_k = P_1(\deg v = k + 1), \quad \forall k \in \mathbb{N}. \quad (1.9) \{?\}$$

The probability q_k correspond to a first neighbor having degree $k + 1$, or equivalently to the probability to have k edges other than the one used to reach the node in the first place, hence the name excess degree distribution.

The excess degree distribution can be computed explicitly by noting that a stub has the same probability to be connected to any of the other $2m - 1$ stubs, thus the probability that this stub is connected to a given node of degree k is $k/(2m - 1)$. Multiplying by the total number of node of degree k , np_k in the large n limit, gives the probability that a given node is attached to a node of degree k as

$$\frac{k}{2m - 1} np_k = \frac{kp_k}{\mathbb{E}[\deg v]}. \quad (1.10) \{?\}$$

Now q_k is the probability that a first neighbor has degree $k + 1$, so we can conclude

$$q_k = \frac{(k + 1)p_{k+1}}{\mathbb{E}[\deg v]}. \quad (1.11) \text{qk as function}$$

¹In principle $P_j(\dots)$ could be defined, corresponding to the probability associated with vertices reached after following j edges.

1.6 Generating functions

ing functions)? A powerful way of representing a degree distribution (or any discrete probability law) is the *generating function* of the distribution. For a degree distribution $\{p_k\}$ it is defined as the function

$$g_0(z) = \sum_{k=0}^{\infty} p_k z^k. \quad (1.12) \quad \text{Definition of } g_0$$

Observe that the derivative with respect to z of $g_0(z)$ is

$$g'_0(z) = \sum_{k=0}^{\infty} p_k k z^{k-1}. \quad (1.13) \quad \text{Derivative of } g_0$$

Comparing with definition of the expectation value we find

$$\mathbb{E}[\deg v] = \sum_{k=0}^{\infty} p_k k = g'_0(1). \quad (1.14) \quad \text{Expectation value as } g'_0(1)$$

Therefore, the generating function is sufficient to know the expectation value of a distribution. This result can be generalized: if we look at the second derivative of the generating function, by differentiating eq. (1.13), we find

$$g''_0(1) = \sum_{k=0}^{\infty} p_k k(k-1) = \sum_{k=0}^{\infty} p_k k^2 - \sum_{k=0}^{\infty} p_k k \quad (1.15) \quad \{?\}$$

$$= \mathbb{E}[(\deg v)^2] - \mathbb{E}[\deg v]. \quad (1.16) \quad \{?\}$$

Thus we can write

$$\mathbb{E}[(\deg v)^2] = g''_0(1) + g'_0(1) = \left. \frac{\partial}{\partial z} (z g'_0(z)) \right|_{z=1}, \quad (1.17) \quad \text{Second moment as derivative}$$

which means that the second moment is fully determined by the generating as well. In fact such formula exists for all moments, but we do not present it here as we only need the first two moments. The second moment is used to compute the variance of the degree distribution as follow

$$\text{var}[\deg v] = \mathbb{E}[(\deg v)^2] - (\mathbb{E}[\deg v])^2 \quad (1.18) \quad \{?\}$$

$$= \left. \frac{\partial}{\partial z} (z g'_0(z)) \right|_{z=1} - (g'_0(1))^2. \quad (1.19) \quad \text{Variance as derivative}$$

All these properties do not depend on the degree distribution p_k to which the generating function is associated and apply to any discrete probability distribution. One other important distribution for network is the excess degree distribution $\{q_k\}$, we define its generating function g_1 as

$$g_1(z) = \sum_{k=0}^{\infty} q_k z^k. \quad (1.20) \quad \text{Definition of } g_1$$

Inserting eq. (1.11) in this definition, we get

$$g_1(z) = \frac{1}{\mathbb{E}[\deg v]} \sum_{k=0}^{\infty} (k+1)p_{k+1}z^k = \frac{g'_0(z)}{g'_0(1)}, \quad (1.21) \quad \boxed{\text{gl as a function}}$$

where we used eq. (1.32) and eq. (1.14) for the last equality.

1.7 Erdos-Renyi networks

In this thesis we focus on three types of theoretical networks, beside some real examples: Erdos-Renyi networks, scale-free networks and geometric networks.

An Erdos-Renyi network is characterized by the fact that it can be grown as follows: for each pair of nodes i and j , add an edge with probability p . To find the degree distribution in such network, first notice that the expected degree, usually denoted c for Erdos-Renyi network, is equal to the number of other vertices multiplied by the probability to be connected to each of them, i.e.

$$c = \mathbb{E}[\deg v] = (n-1)p. \quad (1.22) \quad \{?\}$$

We generally use c as the parameter defining an Erdos-Renyi network, rather than p , since it makes more to keep c constant when n becomes large, rather than p .

The probability for a node to have degree k is

$$p_k = \binom{n-1}{k} p^k (1-p)^{n-1-k}, \quad \forall k \in \mathbb{N}. \quad (1.23) \quad \boxed{\text{Poisson degree}}$$

We recognize a binomial degree distribution for $n-1$ trials with success probability p . In the limit of large n we can approximate such distribution by a Poisson distribution with parameter $c = (n-1)p$

$$p_k \approx \frac{c^k}{k!} e^{-c}. \quad (1.24) \quad \boxed{\text{pk for Erdos-Ren}}$$

The parameter c is the expected degree in the network, it is proportional to $n-1$ rather than n because we only try to bind each vertex with each other, and not with itself, making a total of $n-1$ trials.

Inserting the degree distribution in the definition of the generating function (1.12), we recognize Taylor series representing the exponential function and thus we get

$$g_0(z) = e^{-c} \sum_{k=0}^{\infty} \frac{z^k c^k}{k!} = e^{-c} e^{cz} = e^{c(z-1)}. \quad (1.25) \quad \boxed{\text{g0 for ER network}}$$

Taking the derivative and inserting in eq. (1.20) yields the generating function for the excess degree distribution

$$g_1(z) = e^{c(z-1)}, \quad (1.26) \quad \{?\}$$

which appears to be equal to $g_0(z)$.

1.8 Geometric networks

Geometric networks have a geometric degree distribution

Ask Guiyuan which version to use and if there is a some motivation in using this geom. dist.

1.9 Scale-free networks

The degree distribution of a so-called scale-free network follows a power law with exponent α

$$p_k = \frac{k^{-\alpha}}{\zeta(\alpha)}, \quad \forall k \in \mathbb{N}^*, \quad (1.27) \text{ ?Power law degree di}$$

where $\zeta(\alpha)$ is the Riemann zeta function and $p_0 = 0$. This kind of networks is interesting as many real networks exhibits power law tail in their degree distribution . However, power law distribution are mathematically more complicated than the two previous examples as their generating function can not be represented in term of elementary function. The best we can do is introducing the *polylogarithm* $\text{Li}_\alpha(z)$

Add examples

$$\text{Li}_\alpha(z) = \sum_{k=1}^{\infty} k^{-\alpha} z^k. \quad (1.28) \text{ Definition of polylog}$$

The polylogarithm is a generalization of the Riemann zeta function, as can be seen by the fact that for $z = 1$ we have

$$\text{Li}_\alpha(1) = \sum_{k=1}^{\infty} k^{-\alpha} = \zeta(\alpha). \quad (1.29) \text{ Polylogarithm of 1}$$

Credit the guy the polylog code comes from. Maybe say a bit more about polylog/reimplement polylog using intergration

With that notation, the generating function $g_0(z)$ for scale-free networks can be written

$$g_0(z) = \sum_{k=1}^{\infty} \frac{k^{-\alpha}}{\zeta(\alpha)} z^k = \frac{\text{Li}_\alpha(z)}{\zeta(\alpha)}. \quad (1.30) \text{ ?Generating function}$$

While no simple expression exist for the polylogarithm, its formal definition (1.28) is sufficient to compute its derivative

$$\frac{\partial}{\partial z} \text{Li}_\alpha(z) = \sum_{k=1}^{\infty} k^{-\alpha+1} z^{k-1} = \frac{1}{z} \text{Li}_{\alpha-1}(z). \quad (1.31) \text{ ?Derivative of the p}$$

We therefore find

$$g'_0(z) = \frac{\text{Li}_{\alpha-1}(z)}{z\zeta(\alpha)} \quad (1.32) \text{ Derivative of } g_0 \text{ fo}$$

that we can insert in eq. (1.21) to get

$$g_1(z) = \frac{\text{Li}_{\alpha-1}(z)}{z\zeta(\alpha-1)}. \quad (1.33) \text{ ?}$$

We used eq. (1.29) to perform the simplification $g'_0(1) = \zeta(\alpha - 1)/\zeta(\alpha)$. Thanks to eq. (1.14), we know that $g'_0(1) = \mathbb{E}[\deg v]$. Thus we have

$$\mathbb{E}[\deg v] = \frac{\zeta(\alpha - 1)}{\zeta(\alpha)}. \quad (1.34) \text{ Mean degree in } \boxed{\text{Mean degree in } \dots}$$

Since we are not considering the analytic continuation of the zeta function, but only its real sum description given in eq. (1.29), $\zeta(\alpha)$ diverges for $\alpha < 1$. We can therefore distinguish three cases. First for $\alpha \leq 2$, eq. (1.34) diverges and the mean degree always goes to infinity as the network grows. The name scale free comes from that fact, as there is no *typical scale* for the degrees and greater degrees can always happen with probability too high to be ignored.

In the case $2 < \alpha \leq 3$, the mean degree is finite, but its second moment is not. It can be seen by calculating explicitly using eq. (1.17),

$$\mathbb{E}[(\deg v)^2] = \frac{\partial}{\partial z} (z g'_0(z)) \Big|_{z=1} = \frac{\partial}{\partial z} \left(\frac{\text{Li}_{\alpha-1}(z)}{\zeta(\alpha)} \right) \Big|_{z=1} \quad (1.35) \{?\}$$

$$= \frac{\text{Li}_{\alpha-2}(1)}{\zeta(\alpha)} = \frac{\zeta(\alpha - 2)}{\zeta(\alpha)}. \quad (1.36) \{?\}$$

The divergence of the second moment come from the fact that the factor $\zeta(\alpha - 2)$ diverges for $\alpha \leq 3$. In consequence, the variance of the degree distribution diverges, which means the degree are very widely distributed.

A nicer interpretation arises from a consequence of eq. (1.21), namely

$$g'_1(1) = \frac{g''_0(z)}{g'_0(z)}. \quad (1.37) \{?\}$$

A similar calculation as for the second moment shows that for $\alpha \leq 3$, this diverges too. Hence, for $2 < \alpha \leq 3$ the degree distribution is not strictly speaking scale free, but the excess degree distribution is.

Finally for $\alpha > 3$, the two first moment are finite and only higher moments diverge, which does not implies any peculiar behavior.

1.10 Giant connected component

connected component) 1.10.1 Size of the GCC

An interesting property of a network is the presence and size of *connected component*. A set of nodes is said to be connected if there is a path formed of successive edges from any of its node to any other. All networks can be divided in connected components such that all nodes are element of exactly one component, as is exemplified in fig. 1.6. The connectedness of network is crucial in many real world realisations of networks. In particular any logistic network, such as power grid networks, rail road network or the internet network, is functional only if it is able to transfer goods or services (electricity, passengers or informations) from any node to any other.

Insight on the property of networks can be found by studying the case where the fraction S of the network occupied by a connected component does not vanish in the large n limit. Such component is called a *giant connected component* (GCC). First note that such component is unique.

Prove uniqueness of GCC

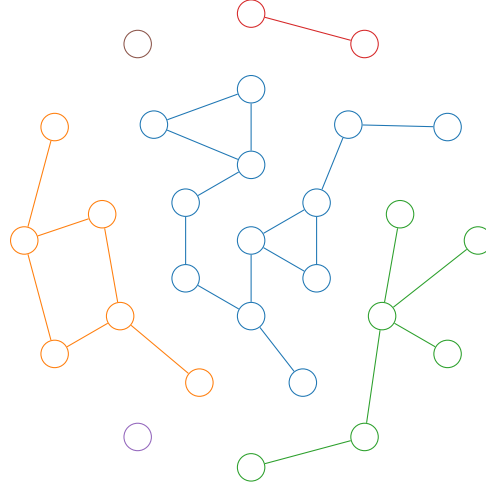


FIGURE 1.6: Scheme of a graph with six connected components, each drawn with a different color.

Then we want to determine the size S of the GCC. To do this consider a connected component C of the network. We define u the probability that a node reached by following an edge is not part of C . We can therefore write the probability S that a randomly chosen vertex is part of C as

$$S = 1 - P_0(w \notin C \forall w \in N(v)) \quad (1.38) \text{ (?)}$$

$$= 1 - \sum_{k=0}^{\infty} P_0(w \notin C \forall w \in N(v) | \deg v = k) P_0(\deg v = k). \quad (1.39) \text{ (?)}$$

The probability $P_0(w \notin C \forall w \in N(v) | \deg v = k)$ is the probability that no neighbors of a node with degree k are part of the component C . This in turn is the probability that by following k independent edges² we find each time a node which is not part of C . This observation allow us to write

$$P_0(w \notin C \forall w \in N(v) | \deg v = k) = [P_1(w \notin C)]^k = u^k. \quad (1.40) \text{ [Probability that a node is not part of } C \text{]} \text{ (?)}$$

In the last equality we introduce the probability u for a first neighbor not to be part of C ,

$$u = P_1(w \notin C). \quad (1.41) \text{ ?Definition of } u \text{ ?}$$

Noting that $P_0(\deg v = k) = p_k$, we find finally

$$S = 1 - \sum_{k=0}^{\infty} u^k p_k \quad (1.42) \text{ (?)}$$

$$= 1 - g_0(u). \quad (1.43) \text{ (?)}$$

We now have a compact expression for S in terms of u and the generating function of the degree distribution g_0 . To determine u we observe that if a vertex is not

²Edges are independant by construction of the configuration model.

part of C , none of its neighbors is either. Thus we can write

$$u = P_1(w \notin C \forall w \in N(v)) \quad (1.44) \{?\}$$

$$= \sum_{k=0}^{\infty} P_1(w \notin C \forall w \in N(v) | \deg v = k) P_1(\deg v = k) \quad (1.45) \{?\}$$

$$= \sum_{k=0}^{\infty} u^k q_k \quad (1.46) \{?\}$$

$$= g_1(u), \quad (1.47) \{?\}$$

where we use eq. (1.40) again together with the fact that by definition of the excess degree distribution $q_k = P_1(\deg v = k)$.

We end up with two equations to describe the size of a connected component

$$S = 1 - g_0(u) \quad (1.48) \text{Single layer } S$$

$$u = g_1(u). \quad (1.49) \text{Single layer } u$$

If we can solve the second one we immediately get the size of the connected component C . However eq. (1.49) only gives u implicitly and its form strongly depends on the degree distribution, therefore no general analytical solutions can be given. From its graphical representation, shown in fig. 1.7, we can see that it has at most two solutions. First the trivial solution $u = 1$ is always present, as by the definition (1.20) of $g_1(z)$, we have $g_1(1) = 1$, implying $S = 0$. The components described by this regime are not giant connected components, as their size S vanish in the large n limit.

Then in some cases, another solution exists with $u < 1$ and $S > 0$. This solution correspond to the GCC. Note that this implies that if a network admit multiple GCC, they must all be of the same relative size S .

Speak about phase transition

Add something about the fact we will refer to u and S thinking about the GCC

To see when a GCC exists, observe that, as can be seen on fig. 1.7, the $g_1(z)$ curve must "go below" the identity curve at $z = 1$ to create a non trivial solution. This requirement means that the slope of $g_1(z)$ at $z = 1$ must be greater than the slope of the identity, which is 1. In term of the derivative of $g_1(z)$ the condition is thus

$$g_1'(z) > 1 \quad (1.50) \text{Boundary condit}$$

to have a non trivial solution to eq. (1.49).

Finally eq. (1.49) can be solved numerically by noticing that the solution u is a fixpoint of the function $g_1(z)$. Since all coefficient in eq. (1.20) are non negative, $g_1(z)$ and all its derivative are positive for $z > 0$. As a consequence starting from $z_0 = 0$, the sequence z_k defined by the iteration

$$z_{k+1} = g_1(z_k) \quad (1.51) \text{Single layer fi}$$

always converges toward the smallest solution of eq. (1.49).

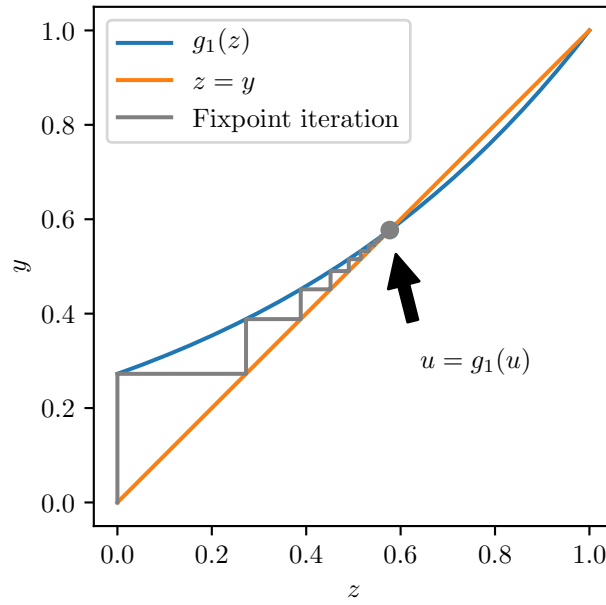


FIGURE 1.7: Graphical representation of eq. (1.49) for an Erdos-Renyi network with $c = 1.3$. Gray solid line represent the successive steps of the fixpoint iteration.

1.10.2 Algorithm to find the connected components

Since we are working in the limit of large n , the networks we are generating and analysing have large n as well. Therefore the algorithm we use must be designed with some care to avoid consuming too much computing time, which would make them impractical to use. This motivates us to present the algorithm we use here.

To find all connected components of a network we proceed as follows

1. Add all nodes to the set $\mathcal{S}_{\text{unprocessed}}$ of unprocessed nodes.
2. Remove one node from $\mathcal{S}_{\text{unprocessed}}$ and add it to the set $\mathcal{S}_{\text{queued}}$ of queued nodes.
3. Start a new component C .
4. Remove one node from $\mathcal{S}_{\text{queued}}$ nodes and name it v .
5. Add v to the current component C .
6. Add all unprocessed neighbors of v to $\mathcal{S}_{\text{queued}}$.
7. If $\mathcal{S}_{\text{queued}}$ is not empty go to 4, else store the component C and continue.
8. If $\mathcal{S}_{\text{unprocessed}}$ is not empty go to 2, else terminate.

This algorithm goes through each node exactly once and is thus of complexity $\mathcal{O}(n)$ which is the optimal complexity since all nodes must be associated to a component.

However, to ensure that the algorithm is fast we must be able to efficiently find all neighbors of a node. To do that we represent the network as an *adjacency list*: each node is given an index i and the adjacency list A_i contains all the neighbors of i . The whole network is thus represented as a list of adjacency list $A = (A_1, \dots, A_n)$.³

³For better performance during the creation of the network, our implementation goes a step further and actually request sorted adjacency lists.

Other representation of networks exist, which are more convenient and efficient for some purposes. However we do not use them in this thesis and we stick to the adjacency list representation.

In fig. 1.8 we compare the fraction of the network occupied by the biggest component found using this algorithm with the numerical solution of eq. (1.48) and (1.49) computed using the fixpoint iteration (1.51). As we can see the agreement between the two is very good for large n , but the drift is significant for small n , especially near the critical point.

Add the position of the critical point

1.10.3 Degree distribution in the GCC

distribution in the GCC)? Per Bayes theorem we have for two random events A and B

$$P(A|B) = P(B|A) \frac{P(A)}{P(B)}. \quad (1.52) \text{ Bayes theorem}$$

We can apply it to compute the probability r_k that a vertex in the GCC has degree k

$$r_k = P(\deg v = k | v \in GCC) \quad (1.53) \{?$$

$$= P(v \in GCC | \deg v = k) \frac{P(\deg v = k)}{P(v \in GCC)} \quad (1.54) \{?$$

$$= \frac{p_k}{S} (1 - P(v \notin GCC | \deg v = k)) = \frac{p_k}{S} (1 - u^k). \quad (1.55) \text{ Degree distribution}$$

To get the final result we used eq. (1.40). This result was previously presented more generally and following a very different path in [3].

Therefore we see that considering a vertex in the GCC biases the probability that it has degree k by a factor $(1 - u^k)/S$ as compared to choosing a vertex uniformly in the network. Since both u and S are smaller than 1, the net effect is to lower the proportion of low degree vertices in the GCC and thus to increase the proportion of high degree vertices. This result can be understood intuitively since in the configuration model the stubs are connected independently. Therefore each edge of a node increases the probability that this node is connected to the GCC, making high degree node over represented in the GCC.

1.11 Generating connected networks

connected networks)? **1.11.1 Algorithm**

to something with
with this information

Algorithm previously in [4]

The knowledge of the degree distribution in the GCC can be used generate a connected component of a given degree distribution r_k as follow: we first determine a degree distribution p_k fulfilling eq. (1.55) for some target degree distribution r_k . Then we generate a network with degree distribution p_k using the configuration model. Finally we take its GCC as our connected network. By construction the vertices in the GCC will have degree distribution r_k . Determining the factors p_k is not immediate however since u is an unknown which is itself a function of p_k . We propose an algorithm to determine it numerically.

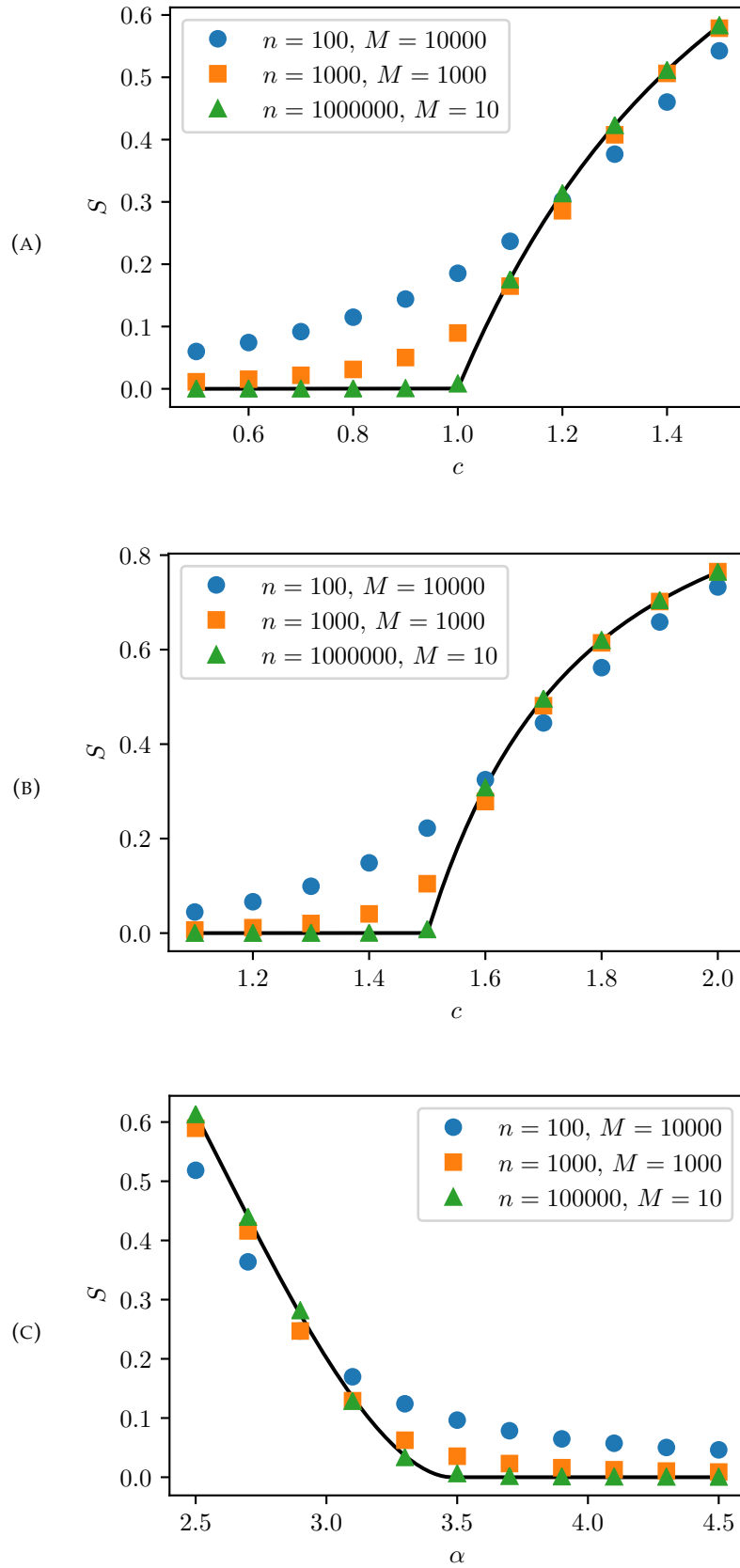


FIGURE 1.8: Numerical solution of eq. (1.48) and (1.49) (solid black line), together with results on simulated networks. Results of simulations are average over M independently generated networks for various values of the number of nodes n in the network. (A) Poisson degree distribution. (B) Geometric degree distribution. (C) Power law degree distribution.

First we isolate p_k from eq. (1.55) to get

$$p_k = S\pi_k(u), \quad \text{with} \quad \pi_k(z) = \frac{r_k}{1 - z^k} \quad (1.56) \{?\}$$

Inserting this in the expression (1.49) for u , we get

$$u = \frac{\sum_{k=1}^{\infty} k\pi_k(u)u^{k-1}}{\sum_{k=1}^{\infty} k\pi_k(u)}. \quad (1.57) \text{Fixpoint equati}$$

Therefore u is a fixpoint of the function

$$\mu(z) = \frac{\sum_{k=1}^{\infty} k\pi_k(z)z^{k-1}}{\sum_{k=1}^{\infty} k\pi_k(z)}, \quad (1.58) \text{Defition of mu}$$

which is fully determined by the GCC degree distribution r_k . Note that for $r_1 = 0$, we have the fixpoint $u = 0$ and $p_k = r_k$ for all k . This is consistent with the fact that small component of a network produced with the configuration model have a probability 0 to have loop [2]. Indeed if $p_1 = 0$ all components must have loops, therefore the probability to have small components is 0 as well.

On the other hand $r_1 > 0$ implies $u > 0$. To approximate its value we define the sequence $u_{j+1} = \mu(u_j)$, with $u_0 = r_1$. This sequence will converge toward u for large j . A proof of this statement is given in Appendix A.1.

In practice we can not deal evaluate infinite sums numerically, thus we need to choose a cutoff index K for the sums such that

$$\sum_{k=K+1}^{\infty} k\pi_k(u) \ll 1. \quad (1.59) \{?\}$$

For scale-free network with exponent smaller than 2 for example, this sum always diverges and thus this method is not applicable.

Once u is approximated, we can compute the first K probabilities p_k , which is sufficient to sample random numbers between 1 and K with relative probability p_k . If K is chosen such that $r_k \ll 1$ for $k > K$, the degree distribution in the GCC closely approximate the distribution r_k .

1.11.2 Erdos-Renyi reconstruction

In order to test the algorithm presented, we choose the target connected degree distribution r_k to be the degree distribution of the GCC of an Erdos-Renyi network. It is then expected that the reconstructed p_k closely approximate a Poisson degree distribution.

The probability p_k to have degree k in an Erdos-Renyi network is given in eq. (1.24). Using eq. (1.49) and (1.48) to find u and S yield everything we need to be able to determine the GCC degree distribution r_k from eq. (1.55). We can therefore use the algorithm on these r_k .

When computing S for the original Poisson distribution, we should however be cautious, as the reconstructed p_0 will always be 0. The expected resulting degree distribution, correctly normalized, is therefore

$$p_k = \frac{c^k}{k!} \frac{1}{e^c - 1}. \quad (1.60) \{?\}$$

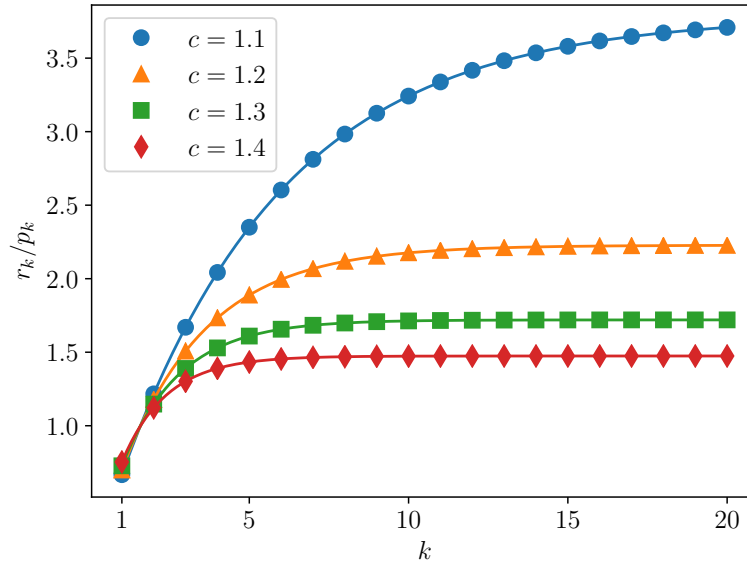


FIGURE 1.9: Bias factor r_k/p_k for r_k being the degree distribution of the GCC of an Erdos-Renyi network with various mean degree c and cutoff constant $K = 10000$. The p_k have been determined using the algorithm presented in the text. Solid line is the expected value $(1 - u^k)/S$ for the bias factor.

The expected bias ratio r_k/p_k is shown for various mean degree c and a cutoff constant $K = 10000$ in fig. 1.9 together with the same value computed from the algorithm presented above. As it can be seen, the agreement is very good. During the computations it has been observed that the closer the mean degree is to the critical value $c = 1$, the slower the fixpoint iteration converges.

1.11.3 Real world networks

As an example of use of the algorithm presented, we apply it to real networks. We choose two by design connected network from the Konect network database [1], the powergrid of the western states of the United States and the road network of the state of California⁴. The power grid network is represented in fig. 1.1, while the Californian road network is not shown due to his high number of vertices ($n \approx 2 \times 10^6$).

The connected degree distribution r_k is taken to be the empirical degree distribution of the real network considered. As a consequence, the cutoff constant K is the maximal degree appearing in the network. Then, to sample the resulting degree distribution p_k , we simply take the number of vertex d_k of degree k to be the closest integer to np_k , where n is the total number of nodes. In the examples presented we use $n = 10^7$.

We compare the degree distribution r_k^{gen} of the GCC of the newly generated network with the target distribution r_k by looking at the ratio r_k^{gen}/r_k . Resulting ratios are shown in fig. 1.10, together with the results obtained by taking the GCC of the reshuffled network.

⁴The codes of the networks in the Konect database are respectively UG and RO.

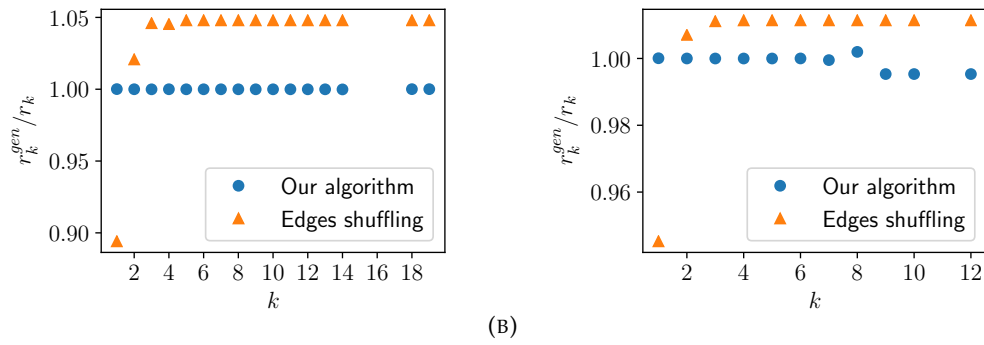


FIGURE 1.10: Plot of the ratio of the generated connected degree distribution r_k^{gen} to the target degree distribution r_k taken from a real network for our algorithm and the edges shuffling method. Missing values correspond to degree with probability 0 to appear. (A) Western US power-grid network. (B) California road network.

Chapter 2

Multiplex networks

2.1 Introduction

Move this to the general introduction ? Put the actual general introduction into the single layer network chapter ?

The concept of network can be extended by allowing different types of edges between the nodes of a network. Such generalized networks are called *multiplex networks* or *multi layers networks*. Real world example of such multiplex networks are numerous, for example there are many different types of connections between cities: transport connections (road, railroads, airlines), information connections (phone lines, internet connections) or supply connections (electricity, water).

Mathematically, it is convenient to represent multiplex networks in a slightly different but equivalent way. We consider each type of edge to generate its own network. If we started with L different type of edges, we therefore get L networks, all sharing the same set of nodes. Conceptually it is similar to stack L networks on top of each other, each network forming a *layer* of the composed multi layers network (hence the name). An illustration of both representation of multiplex networks is shown in fig. 2.1 for a two layers network.

Thinking in term of layers allows to reuse all quantities define in the context of single layer networks, such as the generating functions. We the extended quantity the same way as the original one, adding a subscript i to its symbol to denote the fact that it applies to layer i , except if the symbol already has a subscript. In this case we add a superscript (i) .

Finally, we only consider that each layer of a multiplex network is independant to each other. This is most likely not true for real networks, but we restrict ourself to this case for simplicity.

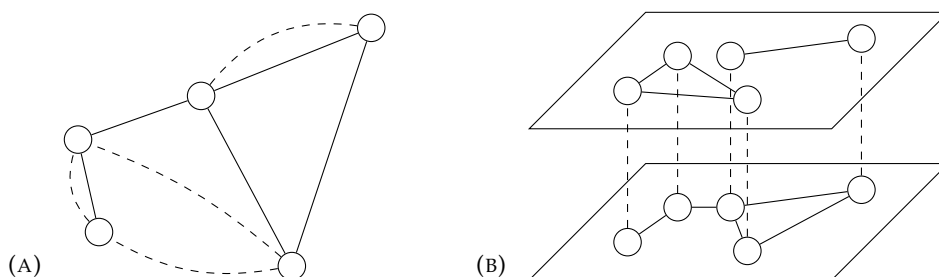


FIGURE 2.1: Possible representations of a two layers network. (A) A network with two kinds of edges, straight solid edges and curved dashed ones. (B) Two networks in separate layers with corresponding nodes.

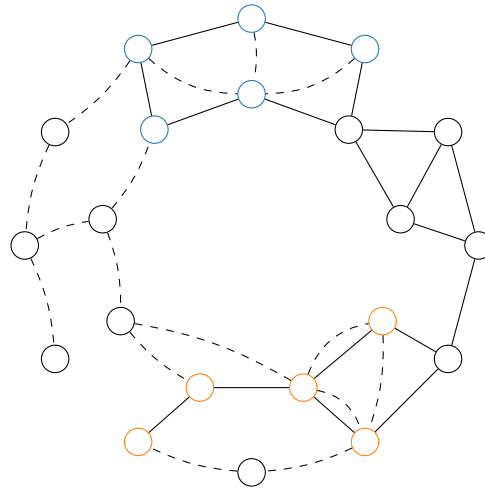


FIGURE 2.2: Intersection of connected components of two layers of a multiplex networks, represented respectively with solid straight line and with curved dashed line. The intersection is composed of the colored nodes, however there is no path in the intersection connecting the blue and orange nodes. Therefore the intersection results in two distinct viable clusters.

2.2 Giant viable cluster

2.2.1 Size of the GVC

Move introduction of lambda parameters in this section

For single layer network, we are interested in the GCC as it is common for a system to be operational only if it is connected to the whole corresponding network. In the context of multiplex networks, the concept of connected component is extended by introducing *viable clusters*. A viable cluster is defined as a subset of the multiplex network in which there is a path between each pair of vertices in all layers. Note that the path between two vertices must be totally included in the viable cluster. As a consequence, the intersection of connected components is not in general a viable clusters as it is illustrated in fig. 2.2.

Correct the following explanation, it is going nowhere

The rationale behind this definition can be clarified by the introduction of the *giant viable cluster* (GVC) which is the extension of the single layer concept of GCC, i.e. a viable cluster which occupy a non zero fraction of the whole network in the limit of large n . Let consider that a node is functional only if it is connected to the GVC. One example of such requirement is the network of powergrid regulation stations. Such stations must be connected to both the powergrid and an internal

Let $g_0^{(i)}$ and $g_1^{(i)}$ be the generating functions of respectively the degree and the excess degree distribution in layer i . Moreover define u_i as the probability that a vertex reached after following an edge in layer i is not part of the giant viable cluster. Be aware that u_i is not the direct extension of the single layer u quantity, which correspond to the probability not to be part of the GCC.

Then if we pick a vertex v at random the probability S that it is part of the giant viable cluster can be written as

$$S = P_0 \left(\bigcap_{i=1}^L \exists w \in N_i(v) \ w \in GVC \right). \quad (2.1) \{?\}$$

By requiring that the layers are independent from one others, we can rewrite S as a product

$$S = \prod_{i=1}^L P_0^{(i)} (\exists w \in N_i(v) \ w \in GVC) \quad (2.2) \{?\}$$

$$= \prod_{i=1}^L \left[1 - P_0^{(i)} (w \notin GVC \ \forall w \in N_i(v)) \right] \quad (2.3) \{?\}$$

$$= \prod_{i=1}^L \left[1 - \sum_{k=0}^{\infty} P_0^{(i)} (w \notin GVC \ \forall w \in N_i(v) | \deg v = k) p_k^{(i)} \right] \quad (2.4) \{?\}$$

$$= \prod_{i=1}^L \left[1 - \sum_{k=0}^{\infty} P_1^{(i)} (v \notin GVC | \deg v = k) p_k^{(i)} \right] \quad (2.5) \{?\}$$

$$= \prod_{i=1}^L \left[1 - \sum_{k=0}^{\infty} u_i^k p_k^{(i)} \right] \quad (2.6) \{?\}$$

$$= \prod_{i=1}^L \left[1 - g_0^{(i)}(u_i) \right]. \quad (2.7) \text{Multiplex GCC size}$$

We can find u_j by a similar reasoning. First note that $1 - u_j$ is the probability that a vertex reached by following an edge in layer j is in the giant viable cluster. Which as before can be written in the form

$$1 - u_j = P_1^{(j)} \left(\bigcap_{i=1}^L \exists w \in N_i(v) \ w \in GVC \right) \quad (2.8) \{?\}$$

$$= \prod_{i=1}^L P_1^{(j)} (\exists w \in N_i(v) \ w \in GVC). \quad (2.9) \{?\}$$

Since the layers are independent, the fact that we reached v by following an edge in layer j to reach vertex v is irrelevant in all other layers. However in layer j this means that the degree distribution follows the distribution $q_k^{(j)}$ rather than $p_k^{(j)}$. Putting this together we get

$$1 - u_j = \left[1 - \sum_{k=0}^{\infty} u_j^k q_k^{(j)} \right] \prod_{\substack{i=1 \\ i \neq j}}^L \left[1 - \sum_{k=0}^{\infty} u_i^k p_k^{(i)} \right] \quad (2.10) \{?\}$$

$$= \left[1 - g_1^{(j)}(u_j) \right] \prod_{\substack{i=1 \\ i \neq j}}^L \left[1 - g_0^{(i)}(u_i) \right]. \quad (2.11) \text{Multiplex u final}$$

As in the single layer case, it is useful to express the u_j as fixpoint of some function. In addition to providing a simple way of numerically solving eq. (2.11), it can

provide valuable information on the presence of solution as presented in Section 2.4.

To make the fixpoint nature of the u_j manifest, we define the functions

$$\psi_j(\mathbf{z}) = 1 - \left[1 - g_1^{(j)}(z_j)\right] \prod_{\substack{i=1 \\ i \neq j}}^L \left[1 - g_0^{(i)}(z_i)\right], \quad (2.12) \text{Definition of } \psi_j$$

where \mathbf{z} is the L dimensional vector

$$\mathbf{z} = (z_1, z_2, \dots, z_L). \quad (2.13) \{?\}$$

Introducing a similar L dimensional vector \mathbf{u} and a L dimensional function ψ as respectively

$$\mathbf{u} = (u_1, u_2, \dots, u_L) \quad (2.14) \{?\}$$

$$\psi(\mathbf{z}) = (\psi_1(\mathbf{z}), \psi_2(\mathbf{z}), \dots, \psi_L(\mathbf{z})), \quad (2.15) \{?\}$$

we can compactly rewrite eq. (2.11) as

$$\mathbf{u} = \psi(\mathbf{u}), \quad (2.16) \text{uvec as a fixpo}$$

which explicitly show that \mathbf{u} is a fixpoint of the function ψ .

Verify ref

Moreover, per definition of the generating functions (1.12) and (1.20), we have $g_0^{(i)}(1) = 1$ and $g_1^{(i)}(1) = 1$ for all layers. As a consequence, eq. (2.11) always admits the trivial solution

$$\mathbf{u}_T = (1, 1, \dots, 1). \quad (2.17) \text{Definition of m}$$

Write something smart to conclude the section

2.2.2 Algorithm to find the viable clusters

Write. Make sure to introduce small components problem

2.2.3 Single parameter multiplex networks

Graphical representation of solutions of multi dimensional equations such as eq. (2.11) is challenging. To avoid this problem we introduce single parameter networks. We define them as multiplex networks in which each layer has the same degree distribution. By letting the degree distribution depend on only one parameter λ , the dimension of the problem reduces to one, as in the single layer case, while still providing important insight in the difference brought by allowing multiple layers.

Mathematically, imposing constant degree distribution across layers implies that the generating functions are constant as well,

$$g_0^{(i)}(z) = g_0(z), \quad (2.18) \{?\}$$

$$g_1^{(i)}(z) = g_1(z), \quad \forall i = 1, \dots, L. \quad (2.19) \{?\}$$

Moreover by definition u_j is the probability that, in the large n limit, the probability that a vertex reached by following an edge is not part of the GVC. Having the same

degree distribution in each layer, thus implies that the $u = u_j$ is the same for each layer.

Plugging these two observations into eq. (2.11) yields the single one dimensional equation for u

$$u = 1 - [1 - g_1(u)] \prod_{\substack{i=1 \\ i \neq j}}^L [1 - g_0(u_i)] \quad (2.20) \{?\}$$

$$= 1 - [1 - g_1(u)] [1 - g_0(u)]^{L-1}. \quad (2.21) \{?\}$$

Similarly using eq. (2.7) we find for S

$$S = [1 - g_0(u)]^L. \quad (2.22) \{?\}$$

2.3 Boundary condition

for multilayer)?

Write some introduction

Up to now, we have considered the multiplex generated to be determined via the degree distributions of each of its layer. However a degree distribution has an infinite number of degrees of freedom, therefore it is more practical to let the degree distributions depend on a finite vector of parameters $\boldsymbol{\lambda} = (\lambda_1, \lambda_2, \dots, \lambda_N)$ and express the behaviour of the network in term of them. Note that the number of parameters N does not need to match the number of layers L .

In order to make our main statement about the critical region for a multiplex network, we need to introduce the functions

$$f_j(\boldsymbol{\lambda}, \mathbf{z}) = 1 - z_j - [1 - g_1^{(j)}(z_j)] \prod_{\substack{i=1 \\ i \neq j}}^L [1 - g_0^{(i)}(z_i)] \quad (2.23) \text{?Definition fj?}$$

$$= \psi_j(\boldsymbol{\lambda}, \mathbf{z}) - z_j. \quad (2.24) \{?\}$$

The functions $\psi_j(\boldsymbol{\lambda}, \mathbf{z})$ are the ones introduced in eq. (2.12) with the dependence on the parameter vector $\boldsymbol{\lambda}$ written explicitly for clarity. We do not write the $\boldsymbol{\lambda}$ dependency for the generating function in order to keep the notation reasonably light.

To get a compact representation of our problem, we gather all functions f_j in a high dimensional function

$$F : \mathbb{R}^N \times \mathcal{I}^L \rightarrow \mathbb{R}^L \quad (2.25) \{?\}$$

$$(\boldsymbol{\lambda}, \mathbf{z}) \mapsto F(\boldsymbol{\lambda}, \mathbf{z}) = (f_1(\boldsymbol{\lambda}, \mathbf{z}), f_2(\boldsymbol{\lambda}, \mathbf{z}), \dots, f_L(\boldsymbol{\lambda}, \mathbf{z})) \quad (2.26) \text{?Definition F?}$$

$$= \psi(\boldsymbol{\lambda}, \mathbf{z}) - \mathbf{z}, \quad (2.27) \{?\}$$

where $\mathcal{I} = [0, 1]$.

Then, since the functions $g_0^{(i)}$ and $g_1^{(i)}$ are analytic with respect to the z_i , the function

$$F_{\boldsymbol{\lambda}} : \mathcal{I}^L \rightarrow \mathbb{R}^L \quad (2.28) \{?\}$$

$$\mathbf{u} \mapsto F_{\boldsymbol{\lambda}}(\mathbf{u}) = F(\boldsymbol{\lambda}, \mathbf{u}), \quad (2.29) \{?\}$$

is continuously differentiable for all parameters λ . Therefore we can define Jacobi matrix $J_\lambda(\mathbf{u})$ of F_λ as having coefficients

$$[J_\lambda(\mathbf{u})]_{ij} = \frac{\partial f_i(\lambda, \mathbf{u})}{\partial u_j}. \quad (2.30) \{?\}$$

With the help of the notation introduced, we can now express solving eq. (2.11) as being equivalent to finding $\mathbf{u} \in \mathcal{I}^L$ such that

$$F(\lambda, \mathbf{u}) = 0. \quad (2.31) \text{Implicit equation}$$

If this equation only admits the trivial solution $\mathbf{u}_T = (1, 1, \dots, 1)$, the parameters λ correspond to a state without GVC. On the other hand if multiple solutions \mathbf{u} exist, a GVC must exist as well. To determine the boundary between these two regions, that we call the *critical region* \mathcal{R} , we use the implicit function theorem.

First, we assume that F (and not only F_λ) is continuously differentiable and that we know a solution \mathbf{u} of (2.31) for some parameter vector λ . With that assumption the implicit function theorem can be state as follow:

If $\det [J_\lambda(\mathbf{u})] \neq 0$ then there is an open neighbourhood $U \subset \mathbb{R}^L$ of λ such that there is an unique continuously differentiable function $h : U \rightarrow \mathcal{I}^L$ with

$$h(\lambda) = \mathbf{u} \quad (2.32) \{?\}$$

$$F(\lambda, h(\lambda)) = 0, \quad \forall \lambda \in U. \quad (2.33) \text{Implicit solution}$$

The result in which we are interested here comes from the contrapositive of this statement, namely that if for all neighbourhoods U we can not find a uniquely defined continuous function h , then the determinant of the Jacobi matrix $J_\lambda(\mathbf{u})$ must be zero,

$$\det [J_\lambda(\mathbf{u})] = 0. \quad (2.34) \text{Boundary condition}$$

This condition has previously been outlined, without proof in [5]. We now prove that such situations arise if λ is a critical point λ^c of the phase transition between the absence and existence of a GVC, and therefore that eq. (2.34) is a sufficient condition to find the critical region of such phase transition.

First notice that in the context of multiplex network a phase transition appears between the trivial solution $\mathbf{u}_T = (1, 1, \dots, 1)$ (where $S = 0$) and non trivial solutions ($S > 0$). However, the trivial solution \mathbf{u}_T solves eq. (2.11) for any generating function and thus for any parameter vector λ . Moreover, on one side of the phase transition occurring at λ^c one solution exists, while on the other at least two do. Therefore for any U open containing λ^c we can define two distinct functions on U that fulfil eq. (2.33), the trivial $h_T(\lambda) = \mathbf{u}_T$ and another function h corresponding to the non trivial solutions, with $h(\lambda^c) = h_T(\lambda^c) = \mathbf{u}_T$. So the function h of the implicit function theorem is not uniquely defined and thus $\det [J_\lambda(\mathbf{u}_T)] = 0$.

On the other hand, let consider a discontinuous phase transition at λ^c . For any neighbourhood U of λ^c there are two sequences $(\lambda_n, \mathbf{u}_n)$ and (η_m, \mathbf{v}_m) with $\lambda_n, \eta_m \in$

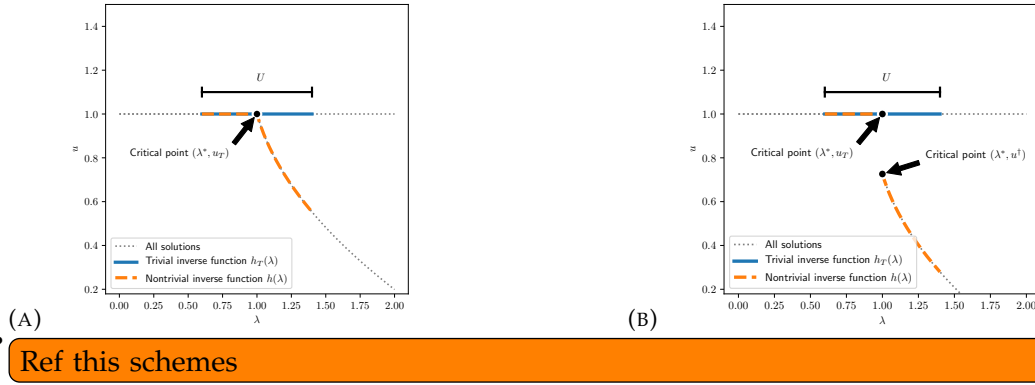


FIGURE 2.3: (a) Scheme of a continuous phase transition. (b) Scheme of a discontinuous phase transition.

U for all $m, n \in \mathbb{N}$ such that

$$\lim_{n \rightarrow \infty} (\lambda_n, \mathbf{u}_n) = (\lambda^c, \mathbf{u}^\dagger) \quad \text{with } \mathbf{u}^\dagger \neq \mathbf{u}_T \quad (2.35) \{?\}$$

$$\lim_{m \rightarrow \infty} (\eta_m, \mathbf{v}_m) = (\lambda^c, \mathbf{u}_T) \quad (2.36) \{?\}$$

$$F(\lambda_n, \mathbf{u}_n) = 0 \quad \forall n \quad (2.37) \{?\}$$

$$F(\eta_m, \mathbf{v}_m) = 0 \quad \forall m. \quad (2.38) \{?\}$$

If we assume that an unique continuous function h solving eq. (2.33) exists, we would have

$$h(\lambda_n) = \mathbf{u}_n \quad \forall n \quad (2.39) \{?\}$$

$$h(\eta_m) = \mathbf{v}_m \quad \forall m. \quad (2.40) \{?\}$$

The continuity of h would furthermore imply

$$h(\lambda^c) = \lim_{n \rightarrow \infty} h(\lambda_n) = \lim_{n \rightarrow \infty} \mathbf{u}_n = \mathbf{u}^\dagger, \quad (2.41) \{?\}$$

but also

$$h(\lambda^c) = \lim_{m \rightarrow \infty} h(\eta_m) = \lim_{m \rightarrow \infty} \mathbf{v}_m = \mathbf{u}_T. \quad (2.42) \{?\}$$

Since $\mathbf{u}_T \neq \mathbf{u}^\dagger$, this gives raise to the contradiction $h(\lambda^c) \neq h(\lambda^c)$. Therefore our assumption must be false and no continuous function h can be defined to solve eq. (2.33). So finally, we have $\det [J(\lambda^c, \mathbf{u})] = 0$, \mathbf{u} being either \mathbf{u}_T or \mathbf{u}^\dagger .

If $L = N = 1$, the problem simplifies to the single layer case 1.10 in which the degree distribution is determined by a single parameter λ . In that case the Jacobi matrix J reduces to the scalar quantity

$$J_\lambda(u) = \frac{\partial}{\partial u} (g_1(u) - u) = \frac{\partial g_1(u)}{\partial u} - 1. \quad (2.43) \{?\}$$

Therefore the condition for the boundary $\det [J_\lambda(u)] = 0$ becomes

$$\frac{\partial g_1(u)}{\partial u} = 1. \quad (2.44) \{?\}$$

As expected this condition is the same as the one given previously in eq. (1.50).

2.4 Interval arithmetic

2.4.1 Motivation

In order to verify that eq. (2.34) indeed gives the critical region \mathcal{R} for a multiplex network, we need to introduce an independent numerical method to approximate it. This comes down to find the number of solutions of eq. (2.11) in the whole parameter space and then draw the boundary of the regions with a constant number of solution. However this leads to two problems.

First, standard algorithms to find zeros of functions are not guaranteed to find all zeros of a function. These algorithms usually find one solution at a time, the one being found depending on the initial guess solution provided by the user, no global information about the solution space being computed during the process. This causes immediate problem for our purpose, since missing a solution of eq. (2.11) may change our estimation of the boundary region by leading to the conclusion that some regions of the parameter space are part of trivial domain while they are not.

Secondly, the standard algorithms can only deal with one point of the parameter space at a time. This forces us to restrict our search on a discrete set of parameters. If this set is ill chosen, we may miss important features of the critical region. This second point is rather minor compared to the first. We indeed expect the critical region to be smooth and therefore restricting our analysis on any lattice should reasonably approximate it. We still mention this caveat however, since it can be solved using the same framework needed to solve the first issue.

Moreover, in Section [?] we present cases where the phase transition displayed by eq. (2.11) is discontinuous. However, to prove this fact, we need to find some region between the trivial and non trivial solutions that can be proved to be exempt of any solution. As outlined above, standard zero finding algorithm are not suitable for such purpose.

To solve all these problems, we therefore introduce the field of *interval arithmetic* and methods thereof that allows to rigorously guarantee presence or absence of solutions of equations in whole region of space.

2.4.2 Theoretical foundation

First of all an *interval* I is defined a set of the form

$$I = [a, b] = \{x \in \mathbb{R} | a \leq x \leq b\}. \quad (2.45) \{?\}$$

The set of all intervals is denoted as \mathbb{IR} . The N -dimensional equivalent of an interval is an *interval box* X , defined as the Cartesian product of N intervals,

$$X = I_1 \times I_2 \times \cdots \times I_N, \quad I_k \in \mathbb{IR} \quad \forall k = 1, \dots, N \quad (2.46) \{?\}$$

The set of all N -dimensional interval boxes is denoted \mathbb{IR}^N .

Given a function $\phi : \mathbb{R}^M \rightarrow \mathbb{R}^N$ we define a new interval valued function Φ such that

$$\Phi : \mathbb{IR}^M \rightarrow \mathbb{IR}^N \quad \text{and} \quad x \in X \quad \Rightarrow \quad \phi(x) \in \Phi(X). \quad (2.47) \text{Definition inte}$$

A function Φ with this property is called an *interval extension* of ϕ . In other words $\Phi(X)$ is guaranteed to contains all possible image of X under ϕ . Note that an interval extension is not unique, and we would like to use the *tightest* possible extension, meaning that we want $\Phi(X)$ to approximate $\phi(X)$ (i.e. the image of the set X by ϕ) as closely as possible. A perfect estimation is however not possible in general since the image of an interval box is not always itself an interval box.

Note that as oppose to ϕ , its interval extension Φ is exactly representable numerically. This can be done by requiring the fact that the implementation of Φ includes numerical inaccuracy in such a way that the eq. (2.47) holds for the finite precision numerical result $\Phi(X)$. Hopefully good implementations, that respect this condition and are reasonably tight, can be found in existing libraries. In this thesis we use the implementation provided by the Julia package `ValidatedNumerics.jl`[?].

Missing ref

The properties of the interval extension allows solve equations in a guaranteed way. To see that first consider an equation expressed as the fixpoint problem

$$\phi(x) = x, \quad (2.48) \quad \text{Interval arithmetic}$$

and Φ an interval extension of ϕ . There are many ways of defining ϕ to make an arbitrary equation equivalent to a fixpoint problem. For interval arithmetic the most common are Newton and Krawczyk operators [6, 7]. However we will also use the more natural fixpoint displayed by eq. .

ref the equation

Introduce Krawczyk operator as it is used below. Maybe in appendix.

Now that we have expressed a problem in term of a fixpoint, let the interval box X contain a solution x^* of eq. (2.48). By the definition 2.47 of an interval extension and the fact that x^* is a fixpoint of ϕ , we have

$$x^* = \phi(x^*) \in \Phi(X) \quad (2.49) \quad \{?\}$$

Since $x^* \in X$ we can even say

$$x^* \in X \cap \Phi(X). \quad (2.50) \quad \text{Generic solution re.}$$

Therefore in general the relation between X and $\Phi(X)$ can give information about the presence of solution in X . First if $X \cap \Phi(X) = \emptyset$, there can not be any solution of eq. (2.48) in X .

On the other hand, if $X \subset \Phi(X)$, we can use the fact that both X and $\Phi(X)$ are interval boxes and thus are compact convex subspaces of \mathcal{R}^N . Brouwer fixpoint theorem can be used, which states that under such conditions, $\Phi(X)$ must have a fixpoint. As a consequence, X must contain at least one solution. It has also been proven that for Newton and Krawczyk operators if X is a subset of the interior of $\Phi(X)$, then the solution in X is unique [6, 7]. We do not discuss the uniqueness further in this thesis though, since in our problem the solutions are not discrete and therefore are never unique inside an interval box.

Finally, in the case where X is neither outside of or included in $\Phi(X)$, we can not conclude anything about the number of solutions of eq. (2.48) in X . However, in this case X could be bisected in the hope to be able to get more conclusive result on finer interval box.

To conclude this section, consider an interval box X_0 and the sequence X_k of interval boxes recursively defined as

$$X_{k+1} = X_k \cap \Phi(X_k), \quad \forall k \in \mathbb{N}. \quad (2.51) \quad \text{Generic refinement}$$

Equation (2.50) tells that all zeros present in the initial interval box X_0 are also present in X_k for all k . Moreover by definition $X_{k+1} \subseteq X_k$, so iterating in the sequence for large k can only yield a closer enclosure of the solutions in X_0 . In fact, if X_0 contains one isolated solution x^* in its interior and x^* is an attractive fixpoint of ϕ , then X_k converges to $\{x^*\}$ for large k .

Some conclusion

Speak more about Krawczyk and Newton ?

2.4.3 Algorithm to test the discontinuity of the phase transition

We now go back to

2.4.4 Algorithm to find the critical region

We would like to approximate the critical region \mathcal{R} by an union of interval boxes in the

Now, let $\mathbf{u}^* \in U_0$ be a solution of (2.31) for some $\lambda \in \Lambda$. By the definition of ψ and eq. (2.47),

$$F(\lambda^*, \mathbf{u}^*) = 0 \quad \Rightarrow \quad \mathbf{u}^* = \psi(\lambda^*, \mathbf{u}^*) \in \Psi(\Lambda, U_0). \quad (2.52) \{?\}$$

Therefore if we apply $\Psi(\Lambda, \cdot)$ to an interval box U_0 containing a solution, the resulting interval box contains the solution as well. We know that all solutions \mathbf{u}^* are contained in \mathcal{I} by definition of \mathbf{u} and thus by iterating the previous argument from $U_0 = \mathcal{I}$, all solutions are elements of the interval boxes $U_k(\Lambda)$ recursively defined by

$$U_{k+1}(\Lambda) = \Psi(\Lambda, U_k(\Lambda)). \quad (2.53) \text{Recursion relat}$$

Therefore if we find k such that $U_k = \{\mathbf{u}_T\}$, we know that the system for any $\lambda \in \Lambda$ only admits the trivial solution.

In practice however, the sequence U_k never converges to exactly the set $\{\mathbf{u}_T\}$, we therefore consider the criterion to be met if

$$U_k \subset [1 - \varepsilon, 1]^L, \quad (2.54) \text{Criterion for t}$$

for some small tolerance ε .

Furthermore it is possible in some cases to guarantee the presence of non trivial solutions, allowing to conclude that a GVC emerges. Indeed, if we can find some interval boxes Λ and U such that

$$\Psi(\Lambda, U) \subset U \quad \text{and} \quad \mathbf{u}_T \notin U, \quad (2.55) \text{Criterion for n}$$

then by definition of the interval extension (eq. (2.47)), we have

$$\psi(\lambda, U) \subset U, \quad \forall \lambda \in \Lambda. \quad (2.56) \{?\}$$

and back which one
actly

sing ref

Since U is closed and simply connected, the fixpoint theorem [?] applies, implying that for each $\lambda \in \Lambda$ there must be at least one \mathbf{u}^* in U such that \mathbf{u}^* is a fixpoint, or in other words such that $\mathbf{u}^* = \psi(\lambda, \mathbf{u}^*)$. Therefore eq. (2.55) is a sufficient condition to prove the existence of at least one solution. Moreover since we imposed $\mathbf{u}_T \notin U$, the solution present can not be the trivial one.

Equations (2.54) and (2.55) give guaranteed criteria for respectively the absence and presence of a non trivial solution in the region Λ considered. This is sufficient to propose an algorithm to estimate the critical region \mathcal{R} .

1. Choose an initial parameter region and store it in the set $\mathcal{S}_{\text{working}}$ of regions yet to be processed.
2. If $\mathcal{S}_{\text{working}}$ is empty terminate, otherwise retrieve the next parameter region from $\mathcal{S}_{\text{working}}$, and call it Λ .
3. If the radius of Λ is smaller than some tolerance δ , store it in the set $\mathcal{S}_{\text{unknown}}$ of regions for which the algorithm is unable to conclude using the tolerance δ .
4. Compute $U_k(\Lambda)$ for k big, using eq. (2.53).
5. If $U_k(\Lambda)$ fulfil eq. (2.54), store Λ in the set $\mathcal{S}_{\text{trivial}}$ of trivial regions and go to 2.
6. Take a subset V of $U_k(\Lambda)$ such that $\mathbf{u}_T \notin V$.
7. If V fulfil eq. (2.55), store Λ in the set \mathcal{S}_{GVC} of non trivial regions and go to 2.
8. Bisect Λ in two sub regions and add both to $\mathcal{S}_{\text{working}}$. Go to 2.

By construction, at any step the critical region is contain in the union of the intervals with unknown status, i.e.

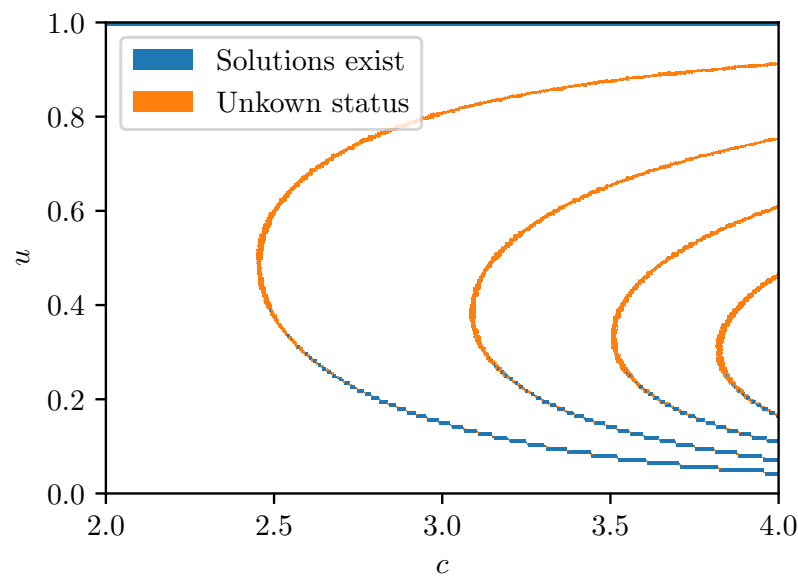
$$\mathcal{R} \subset \bigcup_{U \in \mathcal{S}_{\text{unknown}}} U. \quad (2.57) \{?\}$$

Thus when the algorithm terminates, this union gives an estimation of C .

2.5 Results

We apply the two methods presented to different cases of multiplex network in order to test it. The algorithm using interval arithmetic is compared to the numerical solution of the system composed by eq. (2.11) and (2.34), computed using the `NLSolve.jl` Julia package.

Choose what multiplex networks should be used and with what parameters and actually produce the results.

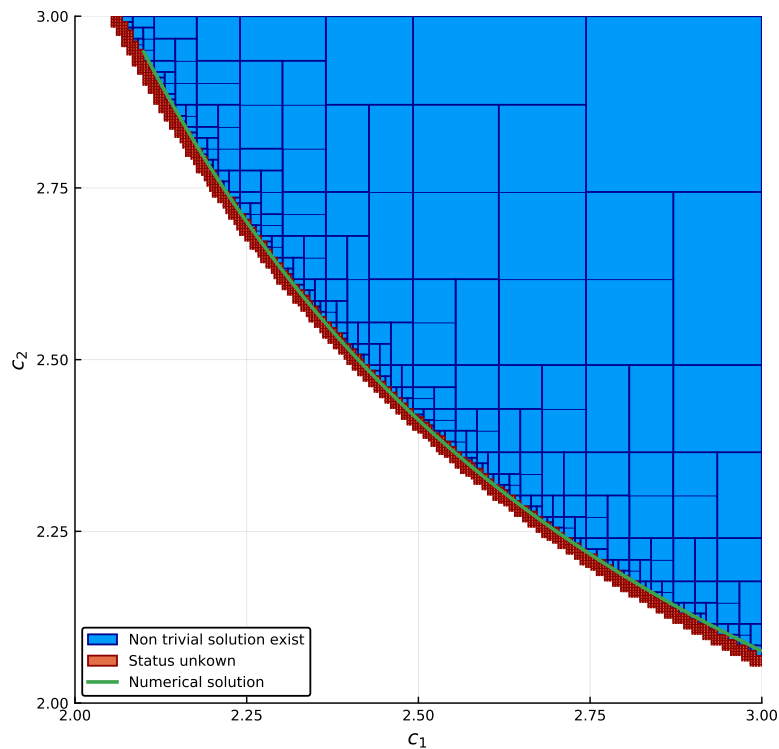


single parameter ER)?

Do other than ER, put label with number of layers add simulations to plot

Write corresponding equation, reference figure somewhere, put it somewhere better

FIGURE 2.4: Solution of eq. ?? for various number of layers. At least one solution is guaranteed to exist in every blue box, while uncolored part of the plot are guaranteed to contain no solutions. No conclusion can be made about existence of solution in the orange regions. This plot is sufficient to prove that the phase transition is discontinuous for the presented number of layers.



s and boundary)?

Make the plot more readable and less ugly

Find why the two methods seems drift away one from the other far from the center.

FIGURE 2.5: Phase diagram for a multiplex network composed of two Erdos-Renyi layer with mean degree c_1 and c_2 . In the blue region a non trivial solution for \mathbf{u} has been found, in the uncolored region only the trivial solution \mathbf{u}_T exists and in the red region the algorithm was unable to conclude in favor of either case. The solid green line is the numerical solution to eq. (2.11) and (2.34).

Appendix A

Fixpoint iteration for connected networks generation

A.1 Convergence of the fixpoint iteration

int convergence) First notice that the case $r_1 = 0$ is trivial, as described in the main text. We will therefore assume in this Appendix that $r_1 > 0$, immediately giving $\mu(0) = r_1 > 0$. Second, note that (1.58) tells us that $z < 1$ implies $\mu(z) < 1$. From there we separate two cases:

If $u = 1$ is the unique solution of eq. (1.57) then $\mu(z)$ must be continuous for $z \in [0, 1]$ and $\mu'(1) < 1$, making $u = 1$ an attractive fixpoint. On the other hand if there is another solution u^* to eq. (1.57), it is the unique solution with $0 \leq u^* < 1$ since $\mu(z)$ is an increasing function of z , as it is demonstrated in Appendix A.2. Moreover, since $\mu(0) > 0$ we have $\mu'(u^*) < 1$, which makes it an attracting fixpoint and makes $u = 1$ a repulsive one.

We can thus conclude that the fixpoint iteration proposed always converges and converges to the degenerate case $u = 1$ only if it is the unique possibility.

A.2 Monotonicity of $\mu(z)$

x: Monotonicity) To prove that $\mu(z)$ is an increasing function, we compute its derivative with respect to z , which yields

$$\mu'(z) = \left[\sum_{k=1}^{\infty} k \pi_k(z) \right]^{-2} (s_1(z) + s_2(z)) \quad (\text{A.1}) \{?\}$$

$$s_1(z) = \sum_{j,k} k j \pi'_k(z) \pi_j(z) (z^{k-1} - z^{j-1}) \quad (\text{A.2}) \{?\}$$

$$s_2(z) = \sum_{j,k} k(k-1) j \pi_k(z) \pi_j(z) z^{k-2}. \quad (\text{A.3}) \{?\}$$

The sum $s_1(z)$ can be rewritten as

$$s_1(z) = \sum_{j>k} k j (\pi'_k(z) \pi_j(z) - \pi'_j(z) \pi_k(z)) (z^{k-1} - z^{j-1}) \quad (\text{A.4}) \{?\}$$

$$= \sum_{j>k} \frac{k r_k}{1 - z^k} \frac{j r_j}{1 - z^j} \frac{z^k - z^j}{z^2} \left(\frac{k}{z^{-k} - 1} - \frac{j}{z^{-j} - 1} \right) \quad (\text{A.5}) \{?\}$$

$$= \sum_{j>k} k j \pi_k(z) \pi_j(z) \frac{z^k - z^j}{z^2} \left(\frac{k}{z^{-k} - 1} - \frac{j}{z^{-j} - 1} \right). \quad (\text{A.6}) \{?\}$$

Using the fact that the function

$$f_z(\lambda) = \frac{\lambda}{z^{-\lambda} - 1} \quad (\text{A.7}) \{?\}$$

is a decreasing function of λ we can see that for $z \in [0, 1)$ and $j > k$ we have

$$z^k - z^j \geq 0 \quad (\text{A.8}) \{?\}$$

$$\frac{k}{z^{-k} - 1} - \frac{j}{z^{-j} - 1} \geq 0, \quad (\text{A.9}) \{?\}$$

and thus $s_1(z) \geq 0$. Moreover each terms in $s_2(z)$ is non-negative, so we have $s_2(z) \geq 0$. We can therefore conclude that $\mu'(z) \geq 0$ and thus that $\mu(z)$ is an increasing function of z .

Bibliography

- [1] J. Kunegis. “Konekt: the koblenz network collection”. In: *Proceedings of the 22nd International Conference on World Wide Web*. ACM. 2013, pp. 1343–1350.
- [2] M. Newman. *Networks: an introduction*. 2010.
- [3] M. Bauer and D. Bernard. “Maximal entropy random networks with given degree distribution”. In: *arXiv preprint cond-mat/0206150* (2002).
- [4] P. Bialas and A. K. Oleś. “Correlations in connected random graphs”. In: *Physical Review E* 77.3 (2008), p. 036124.
- [5] G. Baxter et al. “Avalanche collapse of interdependent networks”. In: *Physical review letters* 109.24 (2012), p. 248701.
- [6] R. E. Moore, R. B. Kearfott, and M. J. Cloud. *Introduction to interval analysis*. Vol. 110. Siam, 2009.
- [7] W. Tucker. *Validated numerics: a short introduction to rigorous computations*. Princeton University Press, 2011.