

Betweenness Centrality in Large Complex Networks

Marc Barthélemy

CEA, Département de Physique Théorique et Appliquée
BP12 Bruyères-Le-Châtel, France

We analyze the betweenness centrality (BC) of nodes in large complex networks. In general, the BC is increasing with connectivity as a power law with an exponent η . We find that for trees or networks with a small loop density $\eta = 2$ while a larger density of loops leads to $\eta < 2$. For scale-free networks characterized by an exponent γ which describes the connectivity distribution decay, the BC is also distributed according to a power law with a non universal exponent δ . We show that this exponent δ must satisfy the exact bound $\delta \geq (\gamma + 1)/2$. If the scale free network is a tree, then we have the equality $\delta = (\gamma + 1)/2$.

I. INTRODUCTION

In large complex networks, not all nodes are equivalent. For example, the removal of a node can have a very different effect depending on the node. If the node is at a dead-end, its removal will be without any effect in contrast with the case of a cut-vertex (the analog of a bridge for edges) which removal creates new disconnected components [1,2]. This question of the importance of nodes in a network is thus of primary interest since it concerns crucial subjects such as networks resilience to attacks [3–5] and also immunization against epidemics [6]. In social network analysis, this problem of determining the rank—or the “centrality”—of the actors according to their position in the social structure was studied a long time ago [7,8]. Different quantities were then defined in this context of social networks in order to quantify this centrality. The simplest proxy for centrality one could think of is the connectivity. However, the inspection of a simple example such as the one in Fig. 1 shows that centrality is in general not related to connectivity. The reason is that connectivity is a local quantity which does not inform about the importance of the node in the network. Indeed, the node v in Fig. 1 has a small connectivity and the effect of its removal is not determined by its connectivity but by the fact that it links together different parts of the network. A good measure of the centrality of a node has thus to incorporate a more global information such as its role played in the existence of paths between any two given nodes in the network. One is thus naturally led to the definition of the betweenness centrality (BC) which counts the fraction of shortest paths going through a given node. More precisely, the BC of a node v is given by [7,8]

$$g(v) = \sum_{s \neq v \neq t} \frac{\sigma_{st}(v)}{\sigma_{st}} \quad (1)$$

where σ_{st} is the total number of shortest paths from node s to node t and $\sigma_{st}(v)$ is the number of shortest paths

from s to t going through v . In the following we will also use the pair-dependency defined as [9]

$$\mu_{st}(v) = \frac{\sigma_{st}(v)}{\sigma_{st}} \quad (2)$$

The betweenness centrality g scales as the number of pairs of nodes ($s \neq t \neq v$) and some authors rescale it by $(N-1)(N-2)/2$ in order to get a number in the interval $[0, 1]$ (N is the number of nodes in the giant component of the network). A naive algorithm for computing g would lead to a complexity of order $\mathcal{O}(N^3)$ and would thus be prohibitive for large networks. Fortunately a rapid algorithm was recently proposed [9] which reduces the complexity to $\mathcal{O}(N^2)$ allowing the computation of the centrality for large networks.

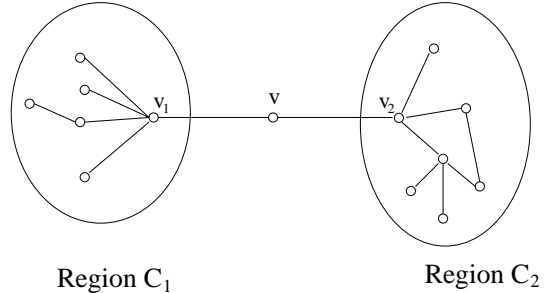


FIG. 1. The node v has a small connectivity (only two neighbors) but all shortest paths from region 1 to region 2 has to go through v which implies a very large centrality. In fact, v is here a cut-vertex; its removal will break the network into two disconnected components.

The definition (1) is indeed a good description of centrality as can be easily seen on the example of figure 1. The BC of the node v is given by

$$g(v) = 2 \sum_{s \in C_1, t \in C_2} \frac{\sigma_{st}(v)}{\sigma_{st}} \quad (3)$$

$$= 2 \sum_{s \in C_1, t \in C_2} 1 \quad (4)$$

$$= 2N_1N_2 \quad (5)$$

where N_1 (N_2) is the number of nodes in region C_1 (C_2). The first equality comes from the fact that the term for which s and t are in the same region does not contribute since in this case $\sigma_{st}(v) = 0$. This result shows that although v has a small connectivity, its BC defined by (1) is large as intuitively expected. This little argument prefigures the more general one about centrality for trees (see below).

High values of the centrality thus indicate that a node can reach the others on short paths or that this vertex lies on many short paths. If one removes a node with large centrality it will lengthen the paths between many pairs of nodes. The extreme case is when the node is a cut-vertex [1,2] and its removal creates new connected components. This was for example used in [10] to determine recursively different communities in large networks.

There are other centrality indices based on shortest paths linking pairs of nodes (stress, closeness, or graph centrality [8,9]). In order to take into account the fact that shortest paths are not always relevant, other definitions were introduced such as the flow betweenness [11] and recently a betweenness centrality based on random walks [12]. This definition (1) differs from the following one which includes the paths endpoints s and t

$$\tilde{g}(v) = \sum_{s \neq t} \frac{\sigma_{st}(v)}{\sigma_{st}} \quad (6)$$

(s or t can be v). It can be easily checked that

$$\tilde{g}(v) = \sum_{s=v \neq t} \mu_{st}(v) + \sum_{s \neq t=v} \mu_{st}(v) + \sum_{s \neq v \neq t} \mu_{st}(v) \quad (7)$$

$$= 2(N-1) + g(v) \quad (8)$$

This additional term $2(N-1)$ is sub-dominant since $g \sim \mathcal{O}(N^2)$ and is thus negligible in the limit of large networks leading to the same results for both definitions (for a typical value of the order $N = 10^4$, the relative difference for large connectivities is negligible—of order 10^{-4} —but could be larger for lower k). In this work, we use the definition (1) and restrict ourselves to non-weighted and non-directed graphs. We will rescale the BC by $(N-1)(N-2)/2$ so that $g \in [0, 1]$. We will keep the same notation g for this normalized centrality.

II. CENTRALITY AND CONNECTIVITY

It has been observed [13] that large networks can be essentially classified in two categories according to the decay of the connectivity distribution $P(k)$. The first category comprises the “exponential” networks with a connectivity distribution decaying faster than any power law (random graph, Poisson graph, etc). In contrast, the

second category is constituted by the “scale-free” networks which have a probability distribution decaying as a power law characterized by an exponent γ

$$P(k) \sim k^{-\gamma} \quad (9)$$

For these networks, there are no typical nodes since the connectivity can vary over a large range of values. In this sense, scale-free networks are very heterogeneous compared to exponential networks for which connectivity fluctuations are small.

In the following, we will investigate the BC for networks which are simple models representative of each class.

A. Scale-Free Networks

In the case of scale-free networks, Goh et al have presented a numerical study of the BC (or “load”) distribution in a static scale-free network model [14]. For this scale-free model, the exponent $\gamma \in]2, \infty[$ is a tunable parameter. They also studied the scale-free model obtained by preferential attachment [15] for which $\gamma = 3$. They showed that the BC is distributed according to a power-law with exponent δ [16]

$$P(g) \sim g^{-\delta} \quad (10)$$

This behavior holds for large g up to a cut-off value which is controlled by finite-size effects. On the basis of their numerical results, they conjectured that the value of $\delta \simeq 2.2$ is “universal” for all values of $\gamma \in]2, 3]$. Universality is usually invoked in physics when different systems show the same behavior [18]. For example many of the observed second order phase transitions have a behavior which depends only on the dimension of the system and the symmetry of the order parameter. In terms of the renormalization group, all these systems are described by the same fixed point of the renormalization group transformation and their critical exponents are then equal. In the case of networks, Goh et al [14,17] measured the exponent δ for different real-world and in silico systems and found only two classes [17]: Either $\delta \simeq 2.2$ (Class I) or $\delta = 2$ (Class II). According to these numerical findings, they claimed that there is “universality” and that networks could be classified according to the value of δ . This means that within a given class, δ is independent of the details of the network such as the mean connectivity $\langle k \rangle = 2m$, or the exponent γ .

The value of δ is however not universal [19] and varies significantly as γ changes in the interval $]2, 3]$ or as m varies. In order to see this non-universality, we first computed the cumulative function $F(g) = \text{Prob}(\text{BC} \geq g)$ for the model proposed in [14] and for the scale-free network obtained by preferential attachment [15]. The results are shown on Fig. 2 and even if the variations are small, the

differences are significant enough to show that δ varies. However, as it can be seen on this Fig. 2 for the BA case, the power law is screened by a cut-off which can be small due to finite-size effects.

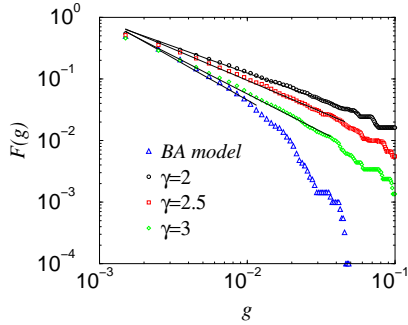


FIG. 2. Cumulative function of the load for different values of $\gamma = 2, 2.5$, and 3 (for $m = 2$). These results were obtained with the same values as in [14] $N = 10^4$ and for 10 configurations. The power law fits (straight lines) give the values $\delta = 1.86, 2.01$ and 2.23 while for the BA model $\delta \simeq 2.3$.

The variations of δ obtained with $F(g)$ are significant enough to claim that it is not a universal exponent but in order to double-check our results we can also use an indirect way of computing δ . We study the relation between the load and the connectivity [14,20] which is of the form

$$g \sim k^\eta \quad (11)$$

where the exponent η depends on the network. This relation (between two random variables) implies that for a given value of k , the corresponding value g_k of the centrality is fixed. Due to noise such as finite-size effects, g_k can however have small fluctuations and we compute the average of g_k at fixed k . The result is shown on Fig. 3 and as can be seen on this plot, the power law (11) holds remarkably for a large range of k and allows an accurate measure of η . In addition, this relation (11) enables us to estimate the cut-off value above which the power-law (10) does not hold. Indeed, the maximum connectivity scales as [21] $k_c \sim N^{1/(\gamma-1)}$ which thus implies that the maximum BC scales as $g_c \sim N^{\eta/(\gamma-1)}$. Finally, we also checked that the value of η does not change significantly for different values of the system size: For $\gamma = 2.5$, we obtain $\eta(N = 10^4) = 1.461 \pm 0.005$, $\eta(N = 2 \cdot 10^4) = 1.467 \pm 0.006$, and $\eta(N = 5 \cdot 10^4) = 1.467 \pm 0.006$ which represents a relative variation due to size less than 1%.

The exponents η and δ are not independent since Eq. (11) implies that

$$P(g) = \int dk P(k) \delta(g - k^\eta) \quad (12)$$

which for large g implies a large k and

$$P(g \gg 1) \sim \int dk k^{-\gamma} \delta(g - k^\eta)$$

$$\sim g^{-1-\frac{\gamma-1}{\eta}} \quad (13)$$

which proves the following equality [20]

$$\eta = \frac{\gamma - 1}{\delta - 1} \quad (14)$$

If the value of $\delta \simeq 2.2$ is universal then η is a linear function of γ with slope $\simeq 1/1.2 \simeq 0.83$.

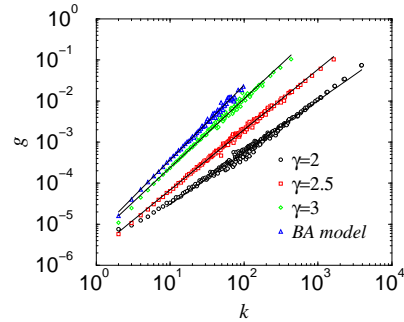


FIG. 3. Log-Log plot of the normalized average load versus connectivity for the same models as in [14] with $m = 2$. The power law fits (straight lines) give $\eta = 1.27 \pm 0.01$ ($N = 3 \cdot 10^4$), 1.467 ± 0.006 ($N = 5 \cdot 10^4$), and 1.68 ± 0.02 ($N = 5 \cdot 10^4$) for $\gamma = 2, 2.5$, and 3 respectively. For the BA model, $\eta = 1.81 \pm 0.02$ ($N = 5 \cdot 10^4$).

In Fig. 4 we plot the measured η versus γ for the different types of networks studied and the corresponding value predicted by universality. This Fig. 4 shows that if for $\gamma \simeq 3$ the value $\delta = 2.2$ seems to be acceptable, the claim of universality for $\gamma \in [2, 3]$ proposed in [14] does not hold (our results do not fit in the other class $\delta = 2.0$ either). In addition, we tested the universality for different values of m and we also obtain variations ruling it out: For $\gamma = 2.5$ and for $N = 2 \cdot 10^4$, we obtain $\eta = 1.477 \pm 0.006, 1.56 \pm 0.006$, and 1.64 ± 0.01 for $m = 2, 4, 6$ respectively. Even if Goh et al have recently shown [22] with a variant of the BA model that for $m \in [1, 2]$, the exponent δ is close to 2.2 for other models supposed to be within the same universality class (BA model, static model, etc.), the exponent δ varies with m or γ and is therefore not universal.

We also note in Figure 4 that for larger values of γ , the exponent η seems to converge to the value $\eta = 2$. This seems to show that for an exponential network, formally characterized by $\gamma = \infty$, the exponent η is equal to two. We will discuss this fact in more details below.

Finally, the case $m = 1$ for the preferential attachment is special in the sense that the obtained scale-free network is a tree. Exact calculations in this case [23,17] show that $\delta = 2 = \eta$. We will see below that the value $\eta = 2$ is in fact expected for any tree and that $\delta = 2$ is the expected value for a scale-free tree only with $\gamma = 3$.

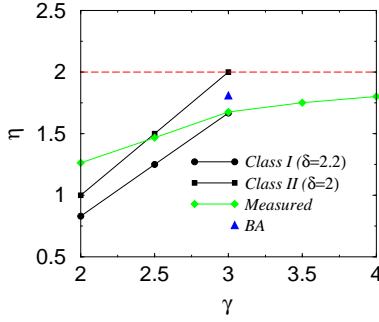


FIG. 4. Exponent η versus γ . If the universality proposed in [14] would be correct, the measured values for $\gamma \in [2, 3)$ should lie on the “universal” straight line corresponding to $\delta = 2.2$ (class I).

B. Random graph

We have seen different examples of scale-free networks in the previous section and we focus now on the random graph [24,25] (often called Erdos-Renyi graph) which is a typical example of exponential networks for which the connectivity distribution is decaying at least as fast as an exponential. This network is constructed as follows. Starting from N nodes, one connects with probability p each pair of nodes. The average final number of edges is thus $E = pN(N-1)/2$ and the average connectivity is $2E/N = p(N-1) \simeq pN$ for large graphs. More generally, the probability that a node has connectivity k is given by the Binomial law

$$p(k) = \binom{N-1}{k} p^k (1-p)^{N-1-k} \quad (15)$$

which converges to a Poisson law of parameter $\langle k \rangle$ for large N and small p such that $\langle k \rangle = pN$ is fixed. We studied the centrality for this network and in Fig. 5 we plot the measured BC versus the connectivity. Even if the connectivity is not varying over a very large range, this plot shows that for large k we have $\eta = 2$. We will discuss this result in more details below but we already note that the random graph has a very small clustering coefficient $C \sim 1/N$ (C counts the average fraction of pairs of connected neighbors [26]) and that this property could possibly be related to the fact that $\eta = 2$.

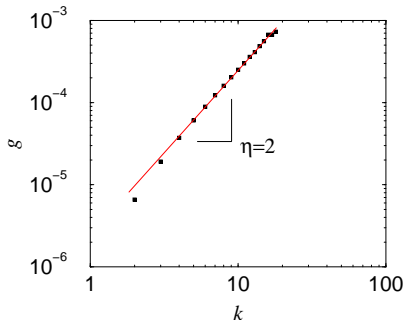


FIG. 5. Log-Log plot of the normalized average load versus connectivity for the random graph model with $N = 5 \cdot 10^4$ and $\langle k \rangle = 6$. The straight line is of slope $\eta = 2$.

III. DISCUSSION AND ANALYSIS OF THE RESULTS

The results obtained above show that the exponents η and δ are not universal and depend on the details of the network. In particular, if the network is scale-free (tree-like or not) δ depends on the exponent γ which describes the power law decay of the connectivity distribution.

The important exponent appears to be η which describes how the betweenness centrality depends on the connectivity. The “optimal” situation which maximizes the BC for a vertex is obtained when all shortest paths are going through it, which happens for a tree structure (ie. a network without loops). To this optimal tree situation corresponds the maximum value of $\eta = 2$. In order to show this, we first define some objects. If a vertex v has connectivity k , we denote by v_i ($i = 1, \dots, k$) its k neighbors. Each neighbor v_i defines a “neighborhood” C_i constituted by nodes which are closer to this neighbor than to any other one. More formally, C_i is defined as follows

$$C_i = \{s \mid d(s, v_i) \leq d(s, v_j) \forall j \neq i\} \quad (16)$$

When the equality of distances $d(s, v_i) = d(s, v_j)$ is obtained for some j then the node s belongs to the two neighborhoods C_i and C_j . The existence of a non empty intersection between different neighborhoods allows for the possibility of paths by-passing the node v .

In the following we denote by N_i the size of each region C_i . In general the shortest paths from $s \in C_i$ to $t \in C_j$ go through v or avoid v by using paths on nodes belonging to $C_i \cap C_j$ [see Fig. 6]. If the two nodes s and t belong to the same neighborhood, say C_l , there is always a shortest path within C_l (in the worst case the shortest path goes through v_l but not through v) and therefore

$$\mu_{st}(v) = 0 \text{ if } s, t \in C_l \quad (17)$$

In terms of these neighborhoods C_i , the BC can be rewritten as

$$g(v) = \sum_{s \neq v \neq t} \mu_{st}(v) \quad (18)$$

$$= \sum_{i \neq j} \sum_{s \in C_i, t \in C_j} \mu_{st}(v) \quad (19)$$

(the term $i = j$ gives zero).

For a tree, these regions C_i are disconnected one from the other and the BC can then be rewritten as

$$g(v) \sim \sum_{i \neq j} N_i N_j \quad (20)$$

If in addition these different parts are of the same order of magnitude $N_i \simeq N_0$ (which is similar to a statistical isotropy condition) we obtain

$$g(v) \sim N_0^2 k(k-1) \quad (21)$$

which for large k behaves as k^2 leading to the value $\eta = 2$. Obviously, the “isotropy” condition $N_i \simeq \text{const.}$ is necessary and if it is not satisfied then the preceding argument does not apply [27]. We note that an exactly solvable model for which this assumption is satisfied is the tree graph obtained with the BA model with $m = 1$ and where one indeed finds $\eta = 2$ [17]. The tree situation maximizes the BC since all shortest paths are going through the node v . In any other cases, the centrality will be less and the maximum possible value of η is 2. More generally, if for a network the density of loops is small enough such that most shortest paths which go from C_i to C_j have to go through v then we obtain $\eta = 2$. This is the case for trees but also for random graphs for which the clustering is small $\sim 1/N$.

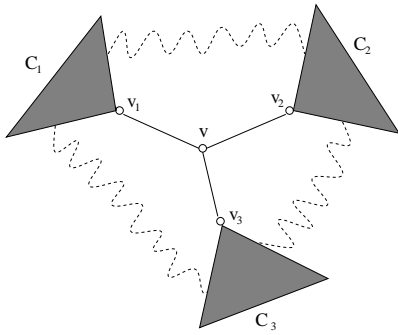


FIG. 6. The node v has here 3 neighbors v_1, v_2, v_3 . These neighbors define three different regions which are disconnected in the case of a tree. When the intersection of these regions is not empty, (shortest) paths between these regions which by-pass v can exist (and are represented by the dotted line between the regions C_i).

If in addition to be a tree, the network is scale-free we can use the relation (14) which together with $\eta = 2$ leads to

$$\eta = 2 \Rightarrow \delta = \frac{\gamma + 1}{2} \quad (22)$$

This relation in particular implies that for the scale-free BA network with $m = 1$ and $\gamma = 3$, we obtain $\delta = (\gamma + 1)/2 = 2$ in agreement with previous results [23,17]. It should be noted that in both these papers [23,17] the authors demonstrate that $\delta = 2$ in the specific case of preferential attachment. However, in [17], the authors claim that their result is valid for any scale-free tree with $\gamma > 2$. This is an incorrect statement since their derivation is only valid for preferential attachment and in general δ depends on γ as predicted by Eq. (22).

On the other hand—and this is the second possible category of networks—if there is a significant fraction of shortest paths which by-pass v then the exponent η will be less than 2. If the network is scale-free then we can use the relation (14) which together with $\eta < 2$ leads to the exact bound

$$\eta < 2 \Rightarrow \delta > \frac{\gamma + 1}{2} \quad (23)$$

The quantity $2 - \eta$ is thus a measure of the density of loops in the network. The fact that $\eta < 2$ indicates that the different parts are also connected by shortest paths which do not pass through the central node. More generally, it would be interesting to understand how η depends on the different parameters of the network such as γ , the clustering coefficient, the loop density, the “anisotropy”, or any other correlation function.

In summary, it seems that concerning the betweenness centrality, we can distinguish two main categories. For the first one which comprises the trees and tree-like networks (clustering almost zero, density of loops very small), we have $\eta = 2$. If in addition, the tree is scale-free with exponent γ , we have the relation $\delta = (\gamma + 1)/2$. The second category comprises the networks for which the density of loops is large enough so that the networks are very different from trees. In this case, the exponents δ, η —when they exist—are not universal and depend on the different details (average connectivity, correlations, etc). If this “clustered” network is scale-free with exponent γ , the exponent δ must obey an exact bound [Eq. (23)]. Although we believe that the present picture is the correct one, further studies are still necessary to understand which are exactly the parameters which control the behavior of η . In this respect, analytical insights would be particularly valuable.

Acknowledgments: I thank the department of physics-INFN in Torino for its warm hospitality during the time this work was started and the Equipe Réseaux, Savoirs & Territoires at the Ecole Normale Supérieure, Paris.

-
- [1] C. Bergé, *Graphs and Hypergraphs* (North-Holland, Amsterdam, 1976) 2nd ed.
 - [2] J. Clark and D.A. Holton, *A first look at graph theory*, World Scientific (1991).
 - [3] R. Albert, H. Jeong, and A.-L. Barabási, *Nature* (London) **401**, 130 (1999).
 - [4] R. Cohen, K. Erez, D. benAvraham, and S. Havlin, *Phys. Rev. Lett.* **86**, 3682 (2001).
 - [5] P. Holme, B.J. Kim, C.N. Yoon, and S.K. Han, *Phys. Rev. E* **65**, 056109 (2002).
 - [6] R. Pastor-Satorras and A. Vespignani, *Phys. Rev. Lett.* **86**, 3200 (2001).
 - [7] L.C. Freeman, *Sociometry* **40**, 35 (1977).

- [8] S. Wasserman and K. Faust, *Social Network Analysis: Methods and applications*, Cambridge University Press (1994).
- [9] U. Brandes, *Journal of Mathematical Sociology*, **25**, 163 (2001).
- [10] D. Wilkinson and B.A. Huberman, condmat/0210147.
- [11] L.C. Freeman, S.P. Borgatti, and D.R. White, *Social Networks* **13**, 141 (1991).
- [12] M.E.J. Newman, condmat/0309045.
- [13] L.A.N. Amaral, A. Scala, M. Barthélemy, and H.E. Stanley, *Proc. Natl. Acad. Sci. USA* **97**, 11149 (2000).
- [14] K.-I. Goh, B. Kahng, and D. Kim, *Phys. Rev. Lett.* **87**, 278701 (2001).
- [15] A.-L. Barabasi and R. Albert, *Science* **286**, 509 (1999).
- [16] For networks with peaked connectivity distributions such as the random graph, the centrality is also peaked and the exponent δ is not defined.
- [17] K.-I. Goh, H. Jeong, B. Kahng, and D. Kim, *Proc. Natl. Acad. Sci. (USA)* **99**, 12583 (2002).
- [18] L.P. Kadanoff, *Statistical Physics: Statics, Dynamics and Renormalization* (World Scientific 2000).
- [19] M. Barthélemy, *Phys. Rev. Lett.* **91**, 189803 (2003).
- [20] A. Vazquez, R. Pastor-Satorras, and A. Vespignani, *Phys. Rev. E* **65**, 066130 (2002).
- [21] S.N. Dorogovtsev and J.F.F. Mendes, *Adv. Phys.* **51**, 1079 (2002).
- [22] K.-I. Goh, B. Kahng, and D. Kim, *Phys. Rev. Lett.* **91**, 189804 (2003).
- [23] G. Szabo, M. Alava, and J. Kertesz, *Phys. Rev. E* **66**, 036101 (2002).
- [24] B. Bollobas, *Random Graph* (Academic Press, New York 1985).
- [25] A. Renyi, *Probability theory* New York, Elsevier, 1980.
- [26] D.J. Watts and D.H. Strogatz, *Collective Dynamics of Small-World Networks*, *Nature* **393**, 440 (1998).
- [27] It would be interesting to quantify for different types of networks the degree of anisotropy—measured by the N_i 's—versus the connectivity.