

RANDOM NETWORKS WITH TUNABLE DEGREE  
DISTRIBUTION AND CLUSTERING

A Thesis

Presented to the Faculty of the Graduate School

of Cornell University

in Partial Fulfillment of the Requirements for the Degree of

Master of Arts

by

Erik McCullough Volz

January 2005

© 2005 Erik McCullough Volz

ALL RIGHTS RESERVED

## ABSTRACT

We present an algorithm for generating random networks with arbitrary degree distribution and clustering (frequency of triadic closure). We use this algorithm to generate networks with exponential, power law, and poisson degree distributions with variable levels of clustering. Such networks may be used as models of social networks and as a testable null hypothesis about network structure. Finally, we explore the effects of clustering on the point of the phase transition where a giant component forms in a random network, and on the size of the giant component. Some analysis of these effects is presented.

## BIOGRAPHICAL SKETCH

Erik Volz received his BA in Mathematics and Russian language from the University of Rochester in 2002. In the Fall of 2002, he began PhD studies in the Department of Sociology, Cornell University.

## ACKNOWLEDGEMENTS

The author thanks Steve Strogatz, Douglas Heckathorn, and Steve Ellner for valuable suggestions and criticism. The author was funded by the NSF (IGERT-0333366) while conducting this research.

## TABLE OF CONTENTS

|          |   |           |
|----------|---|-----------|
| <b>1</b> | <b>Introduction</b>   | <b>1</b>  |
| <b>2</b> | <b>Random network model</b>   | <b>4</b>  |
| <b>3</b> | <b>Results</b>  | <b>13</b> |
| <b>4</b> | <b>Variations on the algorithm</b>  | <b>22</b> |
| 4.1      | Methods for generating degree assortativity . . . . .                                     | 22        |
| 4.2      | Methods for generating lists of potential triads . . . . .                                | 24        |
| <b>5</b> | <b>Phase transitions</b>  | <b>26</b> |
| <b>6</b> | <b>Finite size effects</b>  | <b>29</b> |
| <b>7</b> | <b>Dependence of the clustering coefficient on input parameter <math>C_{input}</math></b> | <b>32</b> |
| <b>8</b> | <b>Implications for sociology</b>   | <b>34</b> |
| <b>9</b> | <b>Discussion</b>   | <b>36</b> |
|          | <b>Bibliography</b>   | <b>38</b> |

## LIST OF TABLES

|     |  |   |
|-----|--|---|
| 2.1 | Detailed description of the clustering method. . . . .           | 7 |
| 2.2 | Detailed description of the clustering method continued. . . . . | 8 |

## LIST OF FIGURES

|     |  |    |
|-----|--|----|
| 2.1 | Overview of the network construction process. The first node (far left) is chosen at random. Then neighbors for that node are chosen as described in the text. Subsequently, neighbors are chosen for the new nodes, but now we have new connections formed with nodes two steps away with probability $C_{input}$ . Triadic connections are indicated with dotted lines. This process continues until the waves die out, and a new component is formed, or all nodes are exhausted. | 6  |
| 2.2 | Two examples of networks generated with the algorithm. Left: Random network with power law degree distribution, $\kappa = 15$ , $\gamma = 2$ , $C = 0.15$ . Right: Random network with poisson degree distribution, $z = 4$ , $C = 0.40$ . [40] Note that these are abstract representations of random networks. The spatial embedding of the network does not have any meaning.   | 9  |
| 2.3 | Random graphs were generated with an exponential degree distribution ( $\lambda = 1.4$ ) with two algorithms: 1. The clustering algorithm described in this text with $C = 0$ 2. A "stub-matching" algorithm as in [2], known to produce true random graphs with specified degree distributions. The frequency of component sizes is illustrated above.  | 12 |
| 3.1 | Random network on 1500 nodes, poisson degree distribution ( $z = 4$ ), $C = 0.00$ . Compare with figures 2.2(right) and 3.5.   | 14 |
| 3.2 | Random network on 1500 nodes, poisson degree distribution ( $z = 4$ ), $C = 0.30$  | 15 |
| 3.3 | Random network on 1500 nodes, poisson degree distribution ( $z = 4$ ), $C = 0.40$ .  | 16 |
| 3.4 | Random network on 1500 nodes, poisson degree distribution ( $z = 4$ ), $C = 0.60$ . The image is zoomed on several of the largest components.  | 17 |
| 3.5 | Random network on 1500 nodes, poisson degree distribution ( $z = 4$ ), $C = 0.97$  | 18 |
| 3.6 | Size of the giant component versus the clustering coefficient in a poisson random network, $z = 3$ . Each point represents the average of 40 trials.   | 18 |
| 3.7 | N=5,000 nodes. Power law with parameters $\kappa = 10$ and $\gamma = 2$ . Each point represents the average of 40 trials. Compare this with 3.6. The phase transition is much less sharp than for the poisson random networks.   | 19 |
| 3.8 | Two random networks are compared over a range of parameter values for the power law degree distribution with parameters $\kappa$ and $\gamma = 2$ . Each point represents the average of 40 trials.  | 20 |

|     |  |    |
|-----|--|----|
| 4.1 | The size of the giant component is shown versus the input clustering parameter $C_{input}$ . The network is Exponential(4), $n = 20000$ . . .  | 23 |
| 5.1 | The size of the giant component is shown vs. $z$ , the parameter of the poisson degree distribution, for four levels of clustering ( $C = 0.0, C = 0.15, C = 0.30, C = 0.40$ ). The vertical lines indicate the point of the phase transition for each level of clustering predicted by equation 5.4 . . . . .   | 28 |
| 6.1 | The percentage reduction in the number of "stubs" is shown versus the Clustering Coefficient for two networks: (i) Poisson degree distribution with parameter = 4, (ii) Exponential degree distribution with parameter = 2. $N=5000$ for both networks. Each point is based on the average of 20 trials. . . . . | 30 |
| 6.2 | The percentage reduction in the number of "stubs" is shown versus the network size. The network has a Poisson degree distribution with parameter = 4, $C = 0.80$ . Each point is based on the average of 20 trial networks. . . . .  | 31 |
| 7.1 | The clustering realized versus the input clustering parameter $C_{input}$ . The random network has a poisson degree distribution with $z = 8$ . $N = 2500$ . . . . .   | 33 |

# Chapter 1

## Introduction

Many random network models have been proposed to replicate important aspects of the topology of real-world networks [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14]. In particular, much attention has been paid to the degree distribution and the clustering coefficient. A great deal of progress has been made on network models which combine certain degree distributions with some level of clustering [15, 16, 13, 17, 18, 19]. It has been an open problem to combine these two topologies in the most general way. Is it possible to have a network model which is flexible enough to accommodate any combination of degree distribution and clustering? In this article we propose such a model and demonstrate its effectiveness by generating networks over a wide range of parameters.

Random network models have fallen in several broad categories. Some models have focused on Monte Carlo techniques to reproduce a specific topology [1, 2, 20]. Other models have specific topologies built into them (e.g. regular lattices) in order to explicate the so-called "small-world" problem [8, 9]. Yet other models have focused on plausible mechanisms for how networks form, such as a growth process with preferential attachment [15, 10, 11]. In common with most mechanism-based models, we produce our networks by growing them from one initial node. We find that being able to construct a network one node at a time also offers sufficient flexibility to combine arbitrary degree distributions and clustering.

Once we have a network model which can combine arbitrary degree distributions and clustering, it is of interest to explore the effects of these parameters on the size of the giant component and the point of the phase transition where a

giant component forms. This is true with regard to clustering in particular, as so far models capable of interpolating between extremes of this parameter have been lacking. In section 3 we explore the effects of clustering on the size of the giant component and point of the phase transition. In section 5 we present some analysis.

Throughout this article we will rely on the following definitions: The *degree distribution* of a network describes how many neighbors a node in a network has. The probability of a node having degree  $k$  in a network is described by the degree distribution  $p_k$ , where  $p_k$  can take the form of any well defined discrete density function over the positive integers. Examples frequently employed in the literature are

- Poisson:  $p_k = \frac{z^k e^{-z}}{k!}, k \geq 0$
- Power-law. For our experiments, we use power-laws with finite cutoffs  $\kappa$ :  
 $p_k = \frac{k^{-\gamma} e^{-k/\kappa}}{Li_\gamma(e^{-1/\kappa})}, k \geq 1$  where  $Li_n(x)$  is the  $n$ th polylogarithm of  $x$ .
- Exponential:  $p_k = (1 - e^{-1/\lambda})e^{-k/\lambda}, k \geq 0$
- Empirical: The degree distribution is estimated from a network sample.
- Gaussian: The ordinary Gaussian must be modified to be positive and discrete.

The *clustering coefficient*  $C$  describes the proportion of triads in a network out of the total number of possible triads. The clustering coefficient is defined:

$$C = \frac{3N_\Delta}{N_3}$$

where  $N_\Delta$  is the number of triads in the network and  $N_3$  is the number of connected triples of nodes. Note that in every triad there are three connected triples.

There is also a measure of *local Clustering* given by

$$C_i = \frac{N_{\Delta}(i)}{\binom{\delta(i)}{2}}$$

where  $N_{\Delta}(k)$  is the number of triads connected to node  $i$ ,  $\delta(i)$  is the degree of node  $i$ , and  $\binom{\delta(i)}{2}$  is the number of potential triads connected to a node of degree  $\delta(i)$ .

The average value of local clustering (i.e. "Watts-Strogatz Clustering" [8]) is also of interest:

$$\frac{\sum C_i}{N}$$

where  $N$  is the number of nodes in the network. This value is frequently close to the clustering coefficient, and will be equal to the clustering coefficient if local clustering is constant throughout the network.

## Chapter 2

### Random network model

Introducing clustering into a network with a specified degree distribution is a non-trivial problem. Any method aspiring to introduce an arbitrary amount of clustering into a network must interpolate between two extremely different topologies. When clustering is 0%, the method must reproduce pure random networks with specified degree distributions. When clustering is 100%, there is only one configuration a network may have: each node must be connected to a small clique where every node has the same degree, and all of a node's neighbors are connected with one another. This challenge is made all the more difficult by trying to make the model networks general enough to accommodate any desired degree distribution.

The most obvious way of introducing triads is to simply define a *rewiring rule* whereby links are swapped between nodes so as to introduce triads while leaving the degree distribution the same. Such rewiring schemes quickly run into problems, as it is impossible to define a rule such that the number of triads is strictly increasing and the number of triads introduced does not max out. The problem is that when links are "swapped" among nodes, triads are not only created but can be destroyed. For example, in our simulations we have found that such schemes are effective only for introducing about 15% clustering into a poisson random network.

Rewiring algorithms have proven effective at the related challenge of adjusting the *average local clustering*. Kim [12] has recently used rewiring algorithms to introduce large amounts of *local clustering* into networks. Using a MC simulations at zero-temperature (i.e. a triad is never destroyed in the rewiring process) and a Hamiltonian of  $\sum -C_k$ , Kim was able to modify various networks with diverse

degree distributions to exhibit average local clustering ( $\sum C_k/N$ ) ranging from 0% to 70%.

Newman [22] and Guillaume et al. [19] have had some success with another approach. These authors define a bipartite network of individuals and affiliations. Then they project the bipartite network onto a unipartite network of only nodes and no affiliations by connecting two nodes if they share a common affiliation. The distributions of affiliation size and the affiliation-degree distribution of the nodes is chosen in such a way as to produce a desired level of clustering. Tuning the degree distribution simultaneously has proven more challenging, however. While the bipartite projection method may actually have the potential to generate pure random networks with tunable degree distributions and clustering, so far it's efficacy has only been shown for exponential and power-law random networks. It remains an open problem to implement it for arbitrary degree distributions.

Our method works by growing networks. The algorithm first initializes all nodes with a degree drawn i.i.d. from the desired degree distribution. Then the random network is constructed by an iterative procedure similar to a branching process. The premise is to start from a single node and then assign new connections entirely at random under the constraint that a certain amount of clustering must exist. The algorithm is described in detail in table 2.1, and is schematized in figure 2.1. Two example networks are shown in figure 2.2.

Our model has similarities and differences with other models proposed in the literature. Like the algorithm of Milo et al. [20], each node is assigned a unique degree prior to any edges being formed between nodes. But like the model networks of Barabasi [3], Dorogovtsev et al. [21] among others, the network is constructed via a growth process. The first node is chosen at random, and subsequently nodes

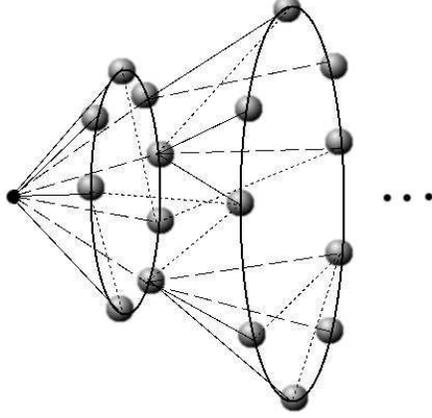


Figure 2.1: Overview of the network construction process. The first node (far left) is chosen at random. Then neighbors for that node are chosen as described in the text. Subsequently, neighbors are chosen for the new nodes, but now we have new connections formed with nodes two steps away with probability  $C_{input}$ . Triadic connections are indicated with dotted lines. This process continues until the waves die out, and a new component is formed, or all nodes are exhausted.

are added to the graph by attaching them to nodes which still have stubs that have not been matched. When the new node forms its own connections, it first forms a list of all nodes which are two steps away. Then with probability  $C_{input}$ , that node is selected as the next neighbor.

One complicated feature of this algorithm concerns the probability of selecting a new neighbor from the stub list. In fact, new neighbors cannot be selected uniformly at random from the stub list, as clustering implies a certain amount of degree assortativity among the nodes in the network. For example, a node connected to a degree  $k$  node has  $k - 1$  potential triads in common with that node, and on average will have  $C(k - 1)$  common triads. This implies that the node must have on average a degree at least equal to  $C(k - 1)$ .

Because triads are distributed uniformly throughout the network, the number

Table 2.1: Detailed description of the clustering method.

1. Initialize all nodes with a degree drawn i.i.d. from the degree distribution
2. Form a list of "stubs" – connections of nodes which have not yet been matched with neighbors. Call this list StubList.
3. Pick a starting node,  $v_0$ , uniformly at random from all nodes.
4. For each of  $v_0$ 's stubs, choose a new neighbor by picking an element  $v_1$  from the stublist with probability  $p_{v_1|d(v_0)}$  as described in the text. If the new neighbor is not
  - the same node as  $v_0$
  - already connected to  $v_0$

then form the connection. Otherwise, repeat the process until a valid neighbor is found. Add all of the new neighbors from this process to a list called NextWave.

5. Copy all elements of NextWave to a list called CurrentWave. Remove all elements from NextWave. For all elements in CurrentWave:

This is continued in table 2.2.

Table 2.2: Detailed description of the clustering method continued.

- (a) Form a list of all nodes 2 steps away. If a node does not have any stubs left in StubList, throw it out. Call this list PotentialTriads
- (b) For all stubs which have not been assigned neighbors
  - i. Scan through PotentialTriads. With probability  $C^{input}$ , connect to node  $v_3 \in \text{StubList}$ . Remove element  $v_3$  from PotentialTriads regardless of whether it was selected. If it was selected, also remove an instance of  $v_3$  from the StubList.
  - ii. If no neighbors were selected from PotentialTriads, select a new neighbor by choosing from StubList as above. If the new neighbor is not in CurrentWave, and if the new neighbor is not already in NextWave, add them to NextWave.

Repeat the last step until NextWave is empty following an iteration. Then, if StubList is empty, the process is complete— all connections have been formed. Otherwise, start a new component by choosing a new starting node uniformly at random from those not yet in the network.

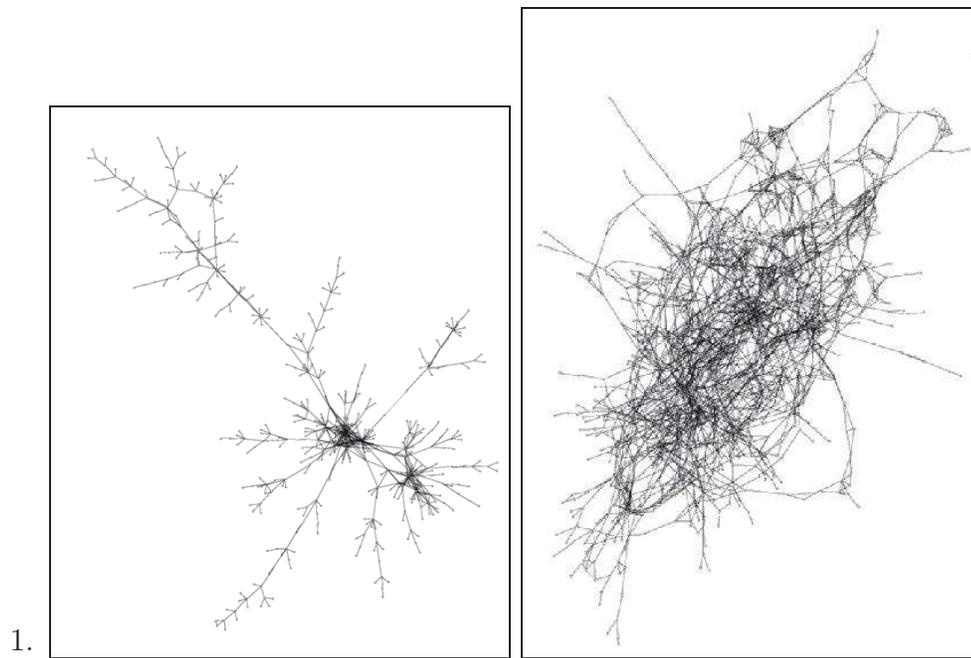


Figure 2.2: Two examples of networks generated with the algorithm. Left: Random network with power law degree distribution,  $\kappa = 15$ ,  $\gamma = 2$ ,  $C = 0.15$ . Right: Random network with poisson degree distribution,  $z = 4$ ,  $C = 0.40$ . [40] Note that these are abstract representations of random networks. The spatial embedding of the network does not have any meaning.

of triads connected to a node of degree  $k$  is distributed  $binomial(\binom{k}{2}, C)$ . As noted above the number of common triads with a neighbor of degree  $k$  is distributed  $binomial(k - 1, C)$ . Let  $\tau_{ij}$  denote the number of triads node  $i$  has in common with node  $j$ , and  $\tau_{ji}$  denote the number of triads  $j$  has in common with  $i$ . Of course these two random variables should be equal. We can calculate the probability of these two potential neighbors as having an equal number of common triads as:

$$p_{ij}^c = \sum_{x=0}^{\min\{d(i), d(j)\}} p(\tau_{ij} = x)p(\tau_{ji} = x)$$

Let  $q_j$  denote the probability of selecting node  $j$  from the stub list. Then the correct probability for selecting node  $j$  as a neighbor is:

$$q_{ij} = \frac{q_j p_{ij}^c}{\sum_{\alpha} p_{i\alpha}^c}$$

which is just  $q_j$  weighted by the probability of the two neighbors having a compatible number of triads in common.

In order to sample from this distribution, we use Markov Chain Monte Carlo techniques. For a large number of iterations we select a new node  $\beta$  from the stub list, then with probability  $a_{\alpha\beta}$  we accept this new neighbor, where  $\alpha$  is the currently selected node in the markov process, and

$$a_{ij} = \frac{p_{i\mu}^c}{p_{i\alpha}^c}$$

If  $\beta$  is not accepted, we keep  $\alpha$  for the next iteration. The final neighbor is the node selected at the last iteration.

It is desirable that our algorithm selects networks as uniformly as possible from the ensemble of all networks which realize a given degree distribution and clustering coefficient. It is difficult to prove that our algorithm is truly unbiased in this sense, though our networks do have many of the properties of an unbiased random

network. The algorithm can be tuned to produce exactly the right proportion of triads to triples in the limit of large graph size. Furthermore, the degree of the nodes were chosen as i.i.d. random variables, so in the limit of large graph size, the degree distribution is unbiased too. Triads are uniformly distributed throughout the network as reflected by the fact that the local clustering is independent of degree. Lastly, when this algorithm is used to produce networks with no clustering at all, it produces networks with the same statistical properties as true random graphs with a specified degree distribution. As shown in figure 2.3, the distribution of component sizes for networks made with this algorithm is identical to true random graphs with specified degree distribution without clustering.

It is worth noting that many real-world networks, particularly in the biological realm, have local clustering which scales as  $1/k$  [23]. Our model in contrast produces constant local clustering, though it may be possible to generalize our method to create networks with any desired schedule of local clustering.

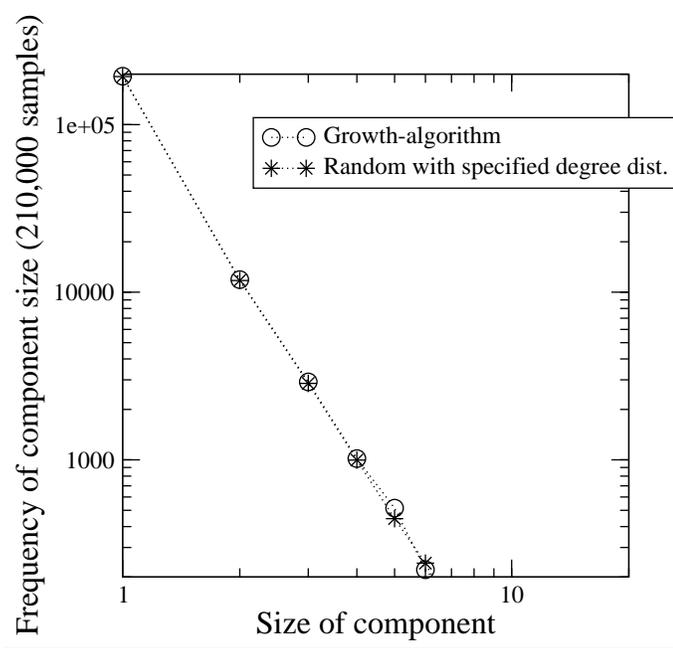


Figure 2.3: Random graphs were generated with an exponential degree distribution ( $\lambda = 1.4$ ) with two algorithms: 1. The clustering algorithm described in this text with  $C = 0$  2. A "stub-matching" algorithm as in [2], known to produce true random graphs with specified degree distributions. The frequency of component sizes is illustrated above.

# Chapter 3

## Results

We have explored the effects of clustering and degree distribution over a wide range of parameters. Figures 2.2(right), 3.1, and 3.5 illustrate the effect of clustering on the structure of a random networks with poisson degree distributions ( $z = 3$ ) as clustering is increased from 0 to 1.00. As  $C$  is increased, nodes tend to disaggregate into smaller tightly connected clusters of nodes with similar degree. This has the overall effect of decreasing the giant component size as clustering is increased. In the limit as  $C$  goes to 1, we find that the network breaks down into many small completely connected cliques with each node in a clique sharing a common degree.

Figure 3.6 shows the effects of clustering on the size of the giant component for a poisson random network. Clustering varies from 0.05 to 0.90. The giant component seems to undergo a phase transition at a critical level of clustering around  $C = 0.60$ . In the next section we will find that the critical clustering value is actually  $C^* = 0.618$ . At this point, nodes suddenly disaggregate into much smaller, tightly inter-connected groups. Similar phase transitions have been observed throughout the networks literature, particularly concerning the targeted deletion of links and nodes in percolation phenomena [24]. This algorithm has similar disconnecting results without modifying the degree distribution of the network.

Regarding power-law networks (see figure 3.7), we note the striking tendency for moderate levels of clustering to limit the size of the giant component. Because the number of potential triads connected to a node scales as  $k^2$ , the high degree vertices account for most of the clustering. In networks with highly skewed degree distributions, the high-degree nodes must connect to one another in order to realize

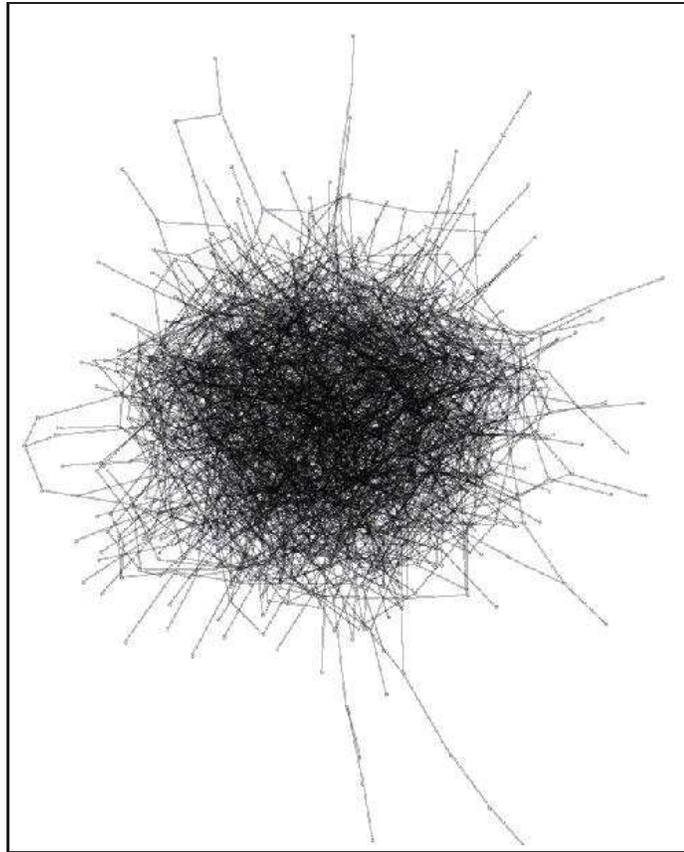


Figure 3.1: Random network on 1500 nodes, poisson degree distribution ( $z = 4$ ),  $C = 0.00$ . Compare with figures 2.2(right) and 3.5.

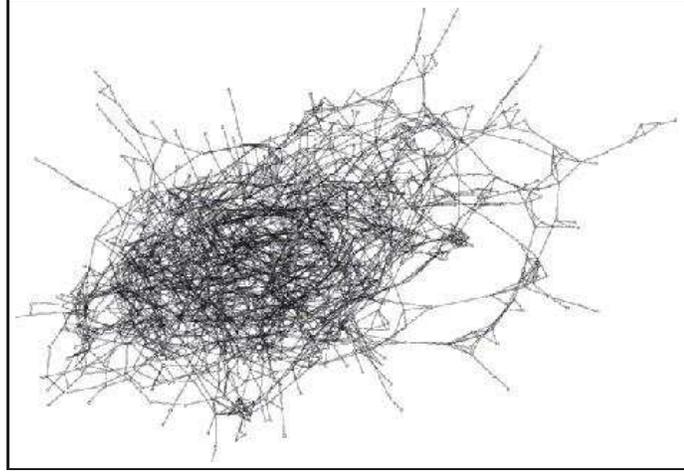


Figure 3.2: Random network on 1500 nodes, poisson degree distribution ( $z = 4$ ),  $C = 0.30$

the required number of triads. This has the effect of limiting the ability to act as hubs for low-degree vertices, and consequently the network disconnects into many small components. Large components can be preserved under much higher clustering with distributions such as the poisson.

The phase transition also undergoes major changes with the introduction of clustering, although this effect seems to depend sensitively on the degree distribution. In figure 3.8 we see that the phase transition where a giant component forms is not significantly affected by the introduction of clustering for networks with power law degree distributions. In contrast to the poisson random networks, there is no sharp phase transition between the regime with a giant component and without. This bears some resemblance to percolation phenomena, where the phase transition disappears for true power-laws and an exponent of 2. But in figure 5.1 we see that the point of the phase transition was dramatically shifted forward for the poisson random network. It is somewhat surprising to observe the phase tran-

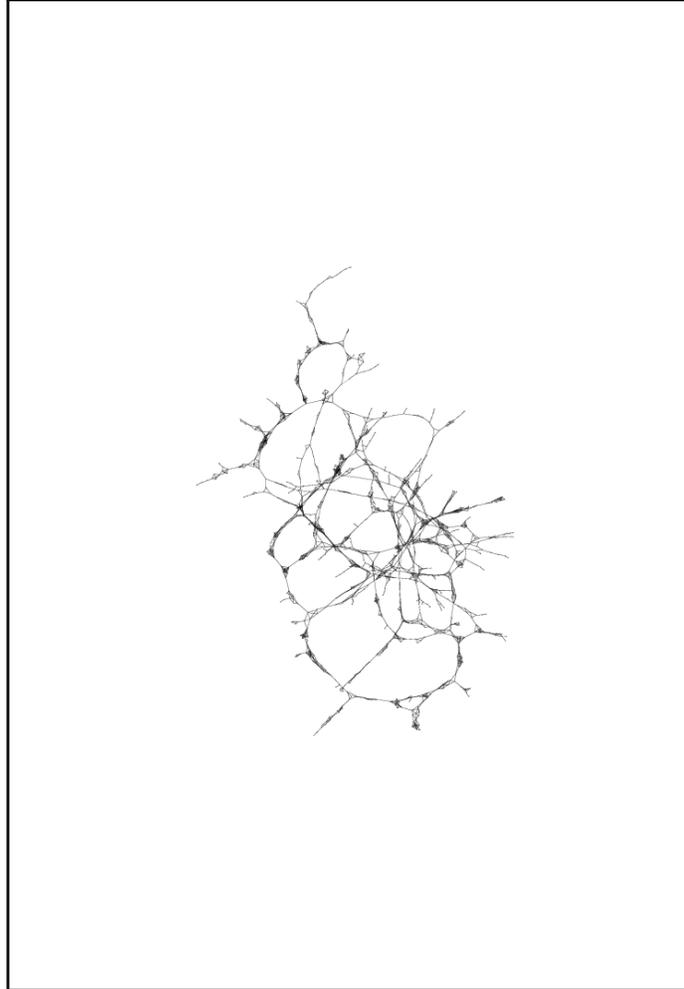


Figure 3.3: Random network on 1500 nodes, poisson degree distribution ( $z = 4$ ),  $C = 0.40$ .

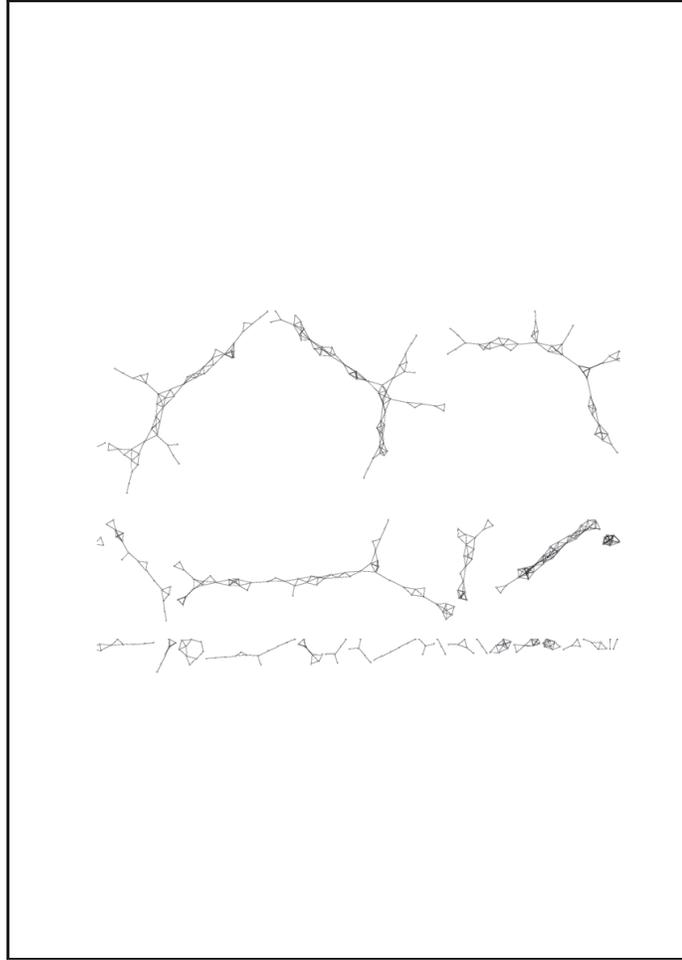


Figure 3.4: Random network on 1500 nodes, poisson degree distribution ( $z = 4$ ),  $C = 0.60$ . The image is zoomed on several of the largest components.

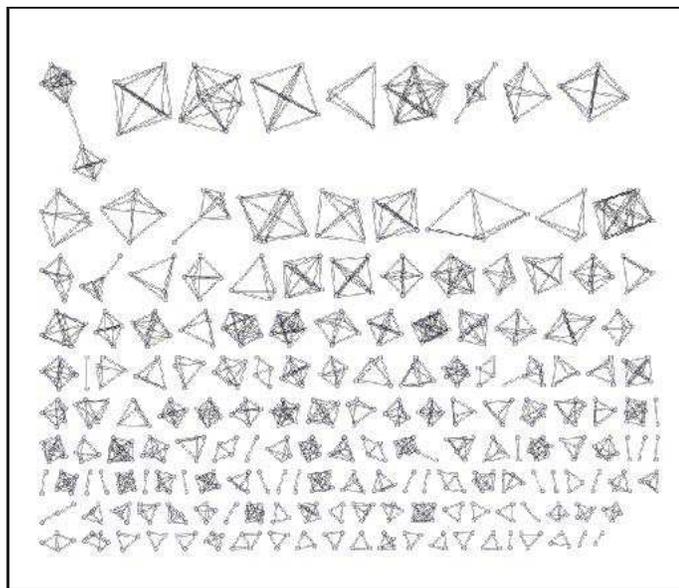


Figure 3.5: Random network on 1500 nodes, poisson degree distribution ( $z = 4$ ),  $C = 0.97$

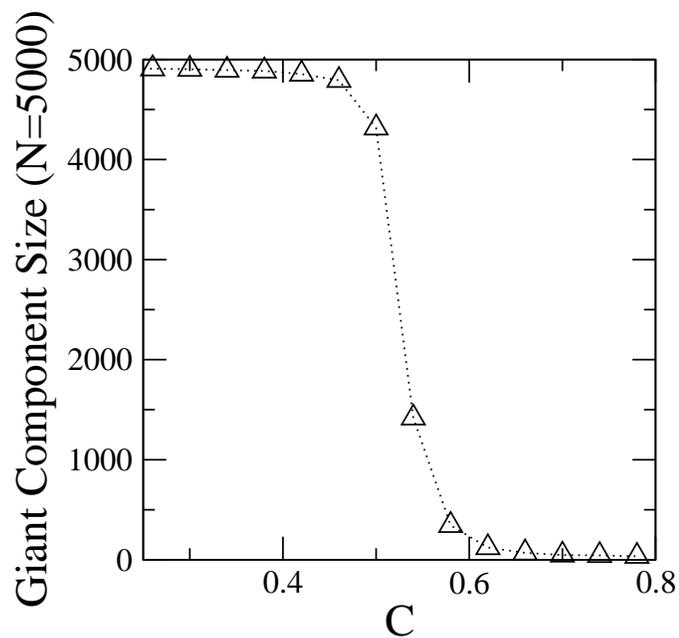


Figure 3.6: Size of the giant component versus the clustering coefficient in a poisson random network,  $z = 3$ . Each point represents the average of 40 trials.

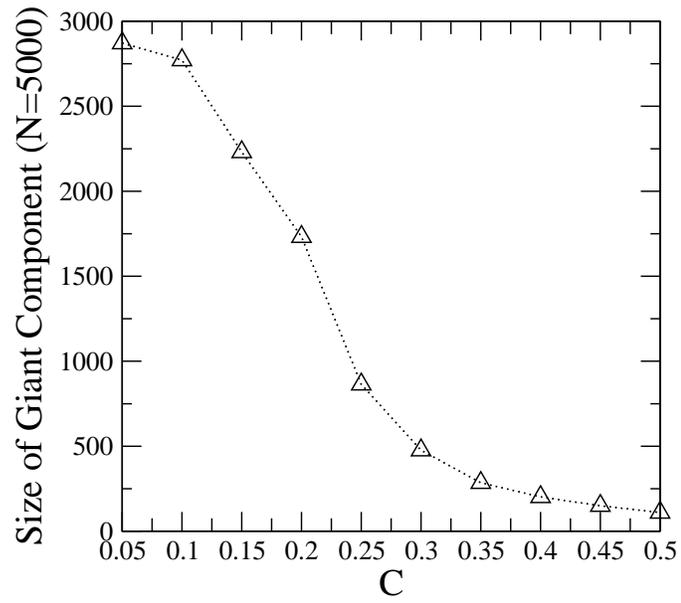


Figure 3.7:  $N=5,000$  nodes. Power law with parameters  $\kappa = 10$  and  $\gamma = 2$ . Each point represents the average of 40 trials. Compare this with 3.6. The phase transition is much less sharp than for the poisson random networks.

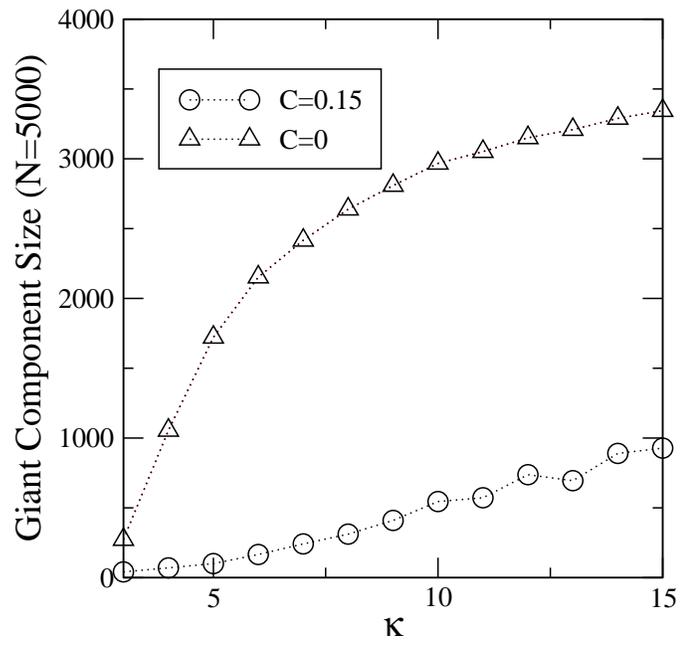


Figure 3.8: Two random networks are compared over a range of parameter values for the power law degree distribution with parameters  $\kappa$  and  $\gamma = 2$ . Each point represents the average of 40 trials.

sition being shifted *forwards* as our algorithm features the introduction of degree assortativity into the network. Previous research has shown the tendency of degree assortativity to shift the point of the phase transition backwards [25].

# Chapter 4

## Variations on the algorithm

We have proposed a very simple example of how network-growth, degree-assortativity and preferential attachment can be combined to generate networks with desirable properties. In fact, many features of this algorithm can be changed to give different and interesting results. It may be that some features of our algorithm are sub-optimal. Variations on this algorithm may be more effective at generating networks with the desired properties (e.g. a desired level of clustering, see section 7). There may be more effective ways to introduce degree assortativity, or to form a list of nodes for preferential attachment. This paper is almost certainly not the final word on this subject.

While the present algorithm was being designed, numerous similar growth algorithms were tried. This section will outline some processes similar to what we have focussed on this paper.

### 4.1 Methods for generating degree assortativity

In our initial network growth experiments, we did not introduce any degree-assortativity at all. As mentioned above, degree assortativity plays an important part in our ability to form triads to a network.

The response of the size of the giant component to the input clustering parameter  $C_{input}$  was very different, and is shown in figure 4.1. The relationship is approximately linear, and should be contrasted with the sharp decline in the size of the giant component observed above at the phase transition  $C^*$  (fig. 3.6).

Another variation on degree assortativity concerns the formulation of  $p_{ij}^c$  as

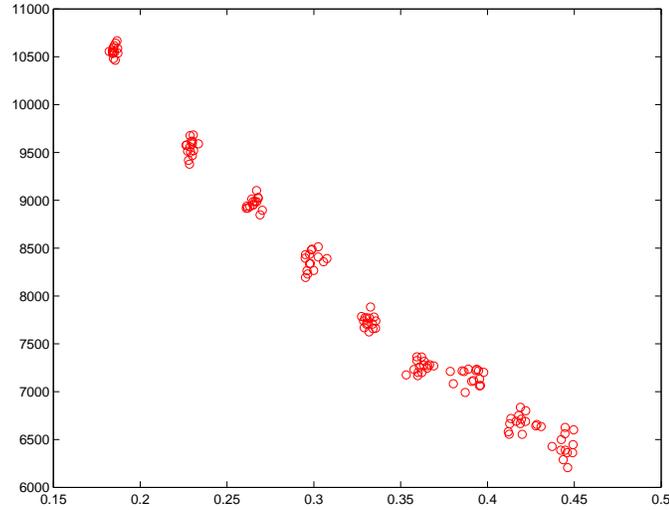


Figure 4.1: The size of the giant component is shown versus the input clustering parameter  $C_{input}$ . The network is Exponential(4),  $n = 20000$

described in the text. This is not the only "Probability of compatibility" we can devise. An alternative is clear from the way our growth algorithm works.

Let *depth* refer to the distance of a node from the initial node in the current component of a growing network. Let  $parents(i)$  denote the set of nodes at a lower *depth* than node  $i$  which are connected to node  $i$ .  $|parents(i)|$  will be the number of parents node  $i$  possesses. Let  $descendants(i)$  denote the set of nodes connected to node  $i$  which are also at a strictly greater depth than node  $i$ . In practice, a descendant of node  $i$  can never be connected to a parent of node  $i$ . This is because the parents of node  $i$  have already had their free connections "reserved" by the time a descendant of node  $i$  is designating its own connections. Hence it is not most likely (sometimes even impossible) for a descendant of node  $i$  to connect to  $C\delta(i)$  of  $i$ 's neighbors. Rather the average number of triadic connections in common with  $i$  will be  $C(\delta(i) - |parents(i)|)$ . The "probability of compatibility" between nodes

i and j then becomes:

$$p_{ij}^c = \sum_{x=0}^{\min\{d(i)-|parents(i)|, d(j)-|parents(j)|\}} p(\tau_{ij} = x)p(\tau_{ji} = x)$$

This modified degree-assortativity was not used in the experiments reported in this paper, but can be found in the clustering code released on the author’s website [41].

## 4.2 Methods for generating lists of potential triads

There are various systems of preferential attachment which can be defined for growth networks. So long as every connected triple in the network becomes a triad with probability C, the input clustering parameter will correspond to the output clustering. Therefore our preferential attachment rule should encourage the creation of triads as uniformly as possible for all connected triples. Unfortunately, a perfect way of accomplishing this has yet to be devised.

Sometimes the fate of two or more triples depends on the allocation of a single connection. This occurs whenever there are two or more paths of length two to a node which is represented in the list *PotentialTriads*. In these cases we have achieved the best results by allowing such a node to have multiple occurrences in *PotentialTriads* and therefore to form a triad with probability greater than  $C_{input}$ . This method was in fact used for the experiments reported in this paper.

Another problem concerns nodes which are two steps away, but which nevertheless have no free connections; hence a triad could never be formed with that node. We have had some success with a method which compensates for this problem. Every time such a node is encountered, a random node is chosen from the *ProspectiveTriads* list, and is re-added to the list, such that it occurs with prob-

ability greater than  $C_{input}$ . This goes some way to compensating with new triads for triads which never had a chance to exist.

# Chapter 5

## Phase transitions

It is a necessary condition for a giant component to exist that if we pick a node at random, the average number of neighbors two steps away,  $s_2$ , exceeds the number of neighbors one step away,  $s_1$  [26]. This is intuitive, since if it were not the case, the number of neighbors  $n$  steps away would decrease to zero on average, and the component would be finite in the limit of large network size. We can use this to approximate the point of the phase transition as clustering is varied in our random networks. Formally, we will solve for the point where

$$s_1 = s_2 \tag{5.1}$$

The necessary condition (5.1) will not quite be a sufficient condition in the presence of clustering as described below. Thus, our solution will only be a lower bound on the point of the phase transition, but in practice, this will serve as an excellent approximation.

For the poisson degree distribution, the average number of nodes one step away is equal to the parameter of the distribution  $z$ , so we have  $s_1 = z$ . As is well known [1], the number of edges emanating from a node if we pick an edge at random and follow it to one of its ends is also  $z$  for the poisson degree distribution. Thus, in the absence of clustering we would have simply  $s_2 = s_1 z = z^2$ , where  $s_2$  is the average number of nodes two steps away from a randomly chosen node.

In the presence of clustering, things become more complicated. Lets pick a node uniformly at random in the network and call this node  $v_0$ . A neighbor of this node,  $v_1$  will have on average  $z$  connections not in common with  $v_0$ . Furthermore, there will be on average  $Cz$  triadic connections between  $v_0$  and  $v_1$  as each of those

connections has a probability  $C$  of being a triad. We can simply deduct the triadic connections from  $s_2$ , so that we have

$$s_2 > z^2 - Cz^2 = z^2(1 - C) \quad (5.2)$$

There is not equality in equation 5.2 because there is an additional force limiting the number of second neighbors: Once two neighbors of  $v_0$ , say  $v_1$  and  $v'_1$  share a triadic connection, it becomes more likely that a node two steps away from  $v_0$ , say  $v_2$ , is a common neighbor of both  $v_1$  and  $v'_1$ . In fact, such connections exist with probability  $C$ . Then, the number of connections we should deduct from every neighbor at distance two due to common connections of nodes at distance one is equal to  $C$  times the average number of triadic connections at distance one, or in other words  $z^2C^2$ . Thus, we have

$$s_2 = z^2 - Cz^2 - C^2z^2 = z^2(1 - C - C^2)$$

We can use this to solve for the critical  $z_C^*$  where a giant component forms given a level of clustering  $C$ :

$$z = z^2(1 - C - C^2) \quad (5.3)$$

The non-zero root of this equation is given by

$$z_C^* = \frac{1}{1 - C - C^2} \quad (5.4)$$

Note that when  $C=0$ , we retrieve the well known result that a giant component forms when  $z = 1$  in the absence of clustering. Unfortunately, we can only say that this is a lower bound for the phase transition due to that the nodes at distance two are not identical to  $v_0$ . The number of outgoing connections from such nodes (to nodes not already counted) is less than  $z - C^2z$  on average.

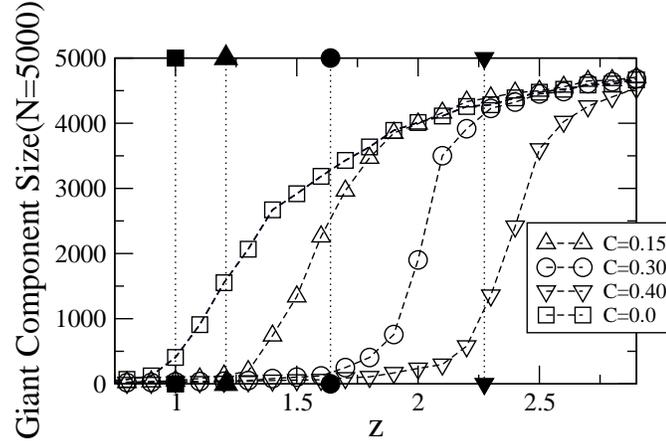


Figure 5.1: The size of the giant component is shown vs.  $z$ , the parameter of the poisson degree distribution, for four levels of clustering ( $C = 0.0, C = 0.15, C = 0.30, C = 0.40$ ). The vertical lines indicate the point of the phase transition for each level of clustering predicted by equation 5.4

In figure 5.1 we have plotted the size of the giant component versus the parameter  $z$  for several levels of clustering. The vertical lines correspond to the phase transitions  $z_C^*$  as given by (5.4). We find good agreement between theory and simulation.

There is a singularity in (5.4) where  $1 - C - C^2 = 0$ . At this point,  $C^* = 0.618$ , the giant component disappears regardless of the average degree  $z$  of the degree distribution.  $C^*$  represents the critical level of clustering that can coexist in a network with a giant component.

# Chapter 6

## Finite size effects

During the execution of the algorithm, it occasionally happens that a node cannot find a suitable neighbor due to the absence of a node left in the network which has free stubs and the correct degree to satisfy the degree assortativity requirements. This imperfection is due to the finite size of the network. In the limit of large size, it would always be possible to find a scale such that every node can find just the right profile of neighbors with the right degree. There is no perfect way to deal with such discrepancies. For the simulations used in this article, we have simply truncated the degree of that node so that it does not have to seek a new neighbor. Even with networks of only 5000 nodes, the number of corrections made is quite small.

Figures 6.1 and 6.2 show the effects of network size and clustering on the amount of degree-corrections made by the algorithm. Figure 6.1 shows the effects of clustering on the number of corrections made for two networks. Note that the total number of "stubs" in the network is equal to the average degree of the nodes times the population size. The corrections made is shown as the proportional reduction in the number of "stubs". Even at 90% clustering, the poisson random network only undergoes less than 5% reduction in its "stubs".

Figure 6.2 shows the effects of network size on the number of corrections made. As expected, the number of corrections drops with the number of nodes in the network. For 7000 nodes and 80% clustering, a poisson random network undergoes less than a 0.4% reduction in its "stubs".

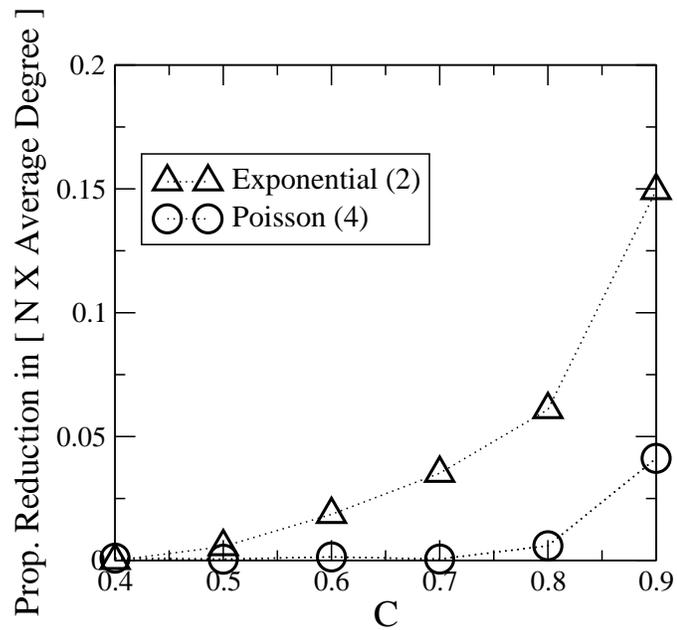


Figure 6.1: The percentage reduction in the number of "stubs" is shown versus the Clustering Coefficient for two networks: (i) Poisson degree distribution with parameter = 4, (ii) Exponential degree distribution with parameter = 2.  $N=5000$  for both networks. Each point is based on the average of 20 trials.

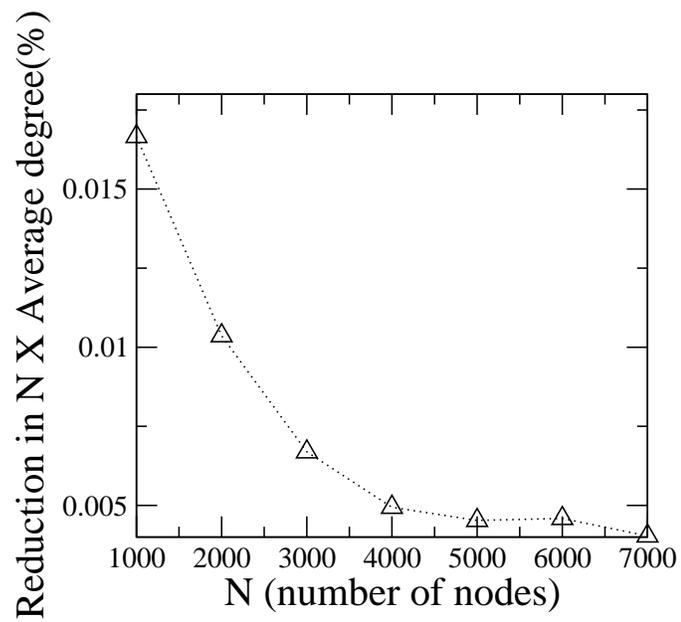


Figure 6.2: The percentage reduction in the number of "stubs" is shown versus the network size. The network has a Poisson degree distribution with parameter  $\lambda = 4$ ,  $C = 0.80$ . Each point is based on the average of 20 trial networks.

## Chapter 7

# Dependence of the clustering coefficient on input parameter $C_{input}$

We have demonstrated a random network model which can generate any desired level of clustering for any degree distribution. Getting a desired level of clustering  $C$  is not always as simple as setting the parameter  $C_{input} = C$ . In general the "input" clustering will be very close to the "output" clustering, though there are sometimes systematic differences. Figure 7.1 shows the value of the clustering coefficient achieved over a broad range of values of  $C_{input}$  for a Poisson random network. Although the  $C$  values do not always fall on the diagonal, they nevertheless cover the full spectrum of  $C = 0$  to  $C = 1.00$  making it possible to achieve any desired level of clustering.

It would be desirable for the input clustering to correspond exactly to the output clustering. The causes of the discrepancy are not fully understood as of the writing of this manuscript, but are probably related to inaccurate degree-assortativity and improperly allocated "prospective triad" lists. Improving the algorithm so that  $C_{input}$  more closely corresponds to  $C$  would be worthy subject for future research.

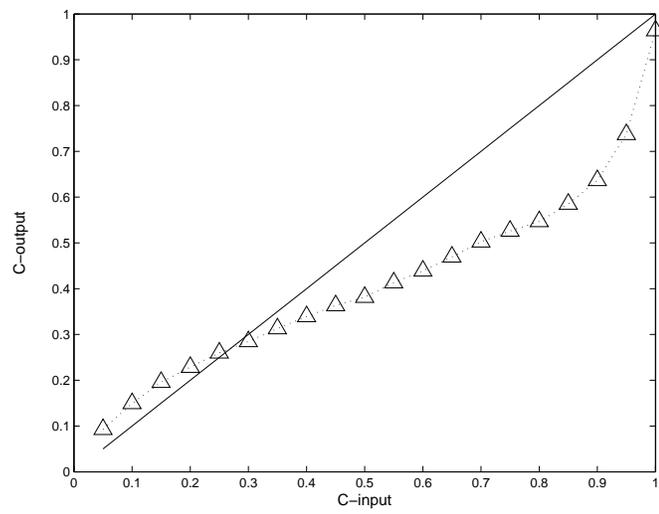


Figure 7.1: The clustering realized versus the input clustering parameter  $C_{input}$ .

The random network has a poisson degree distribution with  $z = 8$ .  $N = 2500$ .

## Chapter 8

### Implications for sociology

The statistical properties of large social networks have been neglected by most social networks researchers in favor of the study of small networks which feature complete information about nodes and ties. This has begun to change in recent years as researchers from other disciplines have made great strides in the mathematics of large random networks—discoveries with direct applications to social networks. Indeed these advances were largely stimulated by a sociological question, the small-world problem, which was expertly investigated by Duncan Watts, an applied mathematician-turned sociologist. Now the methods developed by mathematicians and physicists are returning home to sociology where they may find new applications and facilitate our understanding of a broad range of large social networks, everything from markets and supply chains to internet-dating communities [39].

The present work aims to be a part of this quickly growing literature on large, complex social networks. From the very beginning of this literature—Duncan Watt’s investigation of the small-world problem—transitivity of network connections has been a primary feature of interest. Duncan Watts explained how high transitivity can co-exist with short average path length. This was accomplished with a simple network model which featured random connections and transitivity which was built into a specified lattice topology and a constant degree distribution. Watts did not, however, have a network model which allowed him to smoothly interpolate between various levels of clustering for any degree distribution. One significant aspect of this research is that it allows sociologists to explore broad

ranges of clustering with realistic degree distributions. The degree distribution can even be taken directly from empirical data.

Another aim of this paper is to bring recognition to the multitude of mechanisms for injecting desired topologies into large random networks. Indeed, social networks researchers have been developing network models which feature transitivity for more than a decade [38]. In more recent years, *exponential random network* models have gained a strong foothold in the discipline. Network growth models have received less attention, and perhaps should receive more. Growth models are very flexible in the range of topologies they can produce. They are also suggestive of the mechanisms which produce the topologies we observe. For example, we have demonstrated that network growth and degree-assortativity coupled with preferential attachment to neighbors-of-neighbors is *alone* capable of generating large amounts of clustering.

Finally, a major contribution of this research to sociology is to clarify the relationship between transitivity and the connectivity of social networks. We have shown how increasing transitivity decreases the size of the giant component. Furthermore, there is an upper bound to transitivity, beyond which a giant component will not exist in a random network. It is unlikely that transitivity reaches such extremes in large social networks, as connectivity is an important feature to most of its constituents.

## Chapter 9

### Discussion

We have presented a method for generating random networks which unite two frequently modeled topological features— clustering and the degree distribution.

Random network models can serve several important purposes. First, they can serve as a null hypothesis about the structure of a real-world network. Significant deviations in the structure of the real-world network from a corresponding random graph indicate that there are more forces at work shaping the network than are being accounted for in the random graph model. These deviations can then motivate further inquiry into the forces shaping real-world networks [1].

Secondly, real-world networks are very often of a scale that it is impossible to map them entirely. Various network sampling techniques have been devised to estimate features of the network topology in the absence of data on the entire network [27, 28, 29]. Given reliable estimates about network topology, a random network can then be generated which reproduces this topology. The random network may be used as a stand-in for modeling various dynamic models on networks.

Lastly, the family of random networks we have presented here enables the exploration of a huge parameter space for models on networks. There are a growing number of models which describe dynamic processes on networks. Examples are models of diffusion processes, such as models of epidemics [30, 31, 32], models of fads [33, 34], the spread of rumors [35, 36], and the migration of species among connected habitats [37]. Other models explore interactions among nodes embedded in a network. Examples include spin-glasses, kuramoto oscillators, and disordered neural networks [12]. There are many applications for exploring the effects of clustering and degree distributions on these and other models.

## BIBLIOGRAPHY

- [1] M.E.J Newman, S.H. Strogatz, and D.J. Watts, *Phys. Rev. E* 64(2) 2001.
- [2] M. E. J. Newman, D. J. Watts, and S. H. Strogatz, *Proc. Natl. Acad. Sci. USA* 99, 2566-2572 (2002).
- [3] R. Albert, A.L. Barabasi, *Statistical Mechanics of Complex Networks*, *Rev. Mod. Phys.* 74, 47 (2002)
- [4] R. Albert, H. Jeong, and A.L. Barabasi, *Nature (London)*, 406, 6794 (2000); 406 378 (2000)
- [5] L. Barabasi, *Linked*, Perseus, Cambridge 2002
- [6] E. Ravasz and A.-L. Barabási, Hierarchical organization in complex networks. Preprint cond-mat/0206130 (2002).
- [7] Caldarelli, G., Capocci, A., De Los Rios, P., and Muñoz, M. A., Scale-free networks from varying node intrinsic fitness, *Phys. Rev. Lett.* **89**, 258702 (2002).
- [8] D. J. Watts, S.H.Strogatz, *Nature* 393, 440-442 (1998)
- [9] D. J. Watts, *Small worlds: The dynamics of Networks between Order and Randomness*, Princeton University
- [10] M. E. J. Newman, *SIAM Review* 45, 167-256 (2003).
- [11] L. A. N. Amaral, A. Scala, M. Barthélémy, and H. E. Stanley, Classes of small-world networks. *Proc. Natl. Acad. Sci. USA* **97**, 11149–11152 (2000).
- [12] B. J. Kim, Performance of networks of artificial neurons: The role of clustering. *Phys. Rev. E* **69**, 045101(R) 2004
- [13] S.N. Dorogovtsev and J.F.F. Mendes, *The evolution of networks: from biological nets to the Internet and WWW*, Oxford : Oxford University Press, 2003.
- [14] Krapivsky, P. L. and Redner, S., Organization of growing random networks, *Phys. Rev. E* **63**, 066123 (2001).
- [15] J. Davidsen, H. Ebel, S. Bornholdt, *Phys. Rev. Lett.* 88 (2002) 128701, cond-mat/0108302
- [16] E. M. Jin, Michelle Girvan, and M. E. J. Newman, *Phys. Rev. E* 64, 046132 (2001). Press, Princeton (1998)
- [17] S.N.Dorogovtsev, J.F.F.Mendes, A.N.Samukhin, Structure of growing networks with preferential linking, *Phys. Rev. Lett.* **85**, 4633 (2000)

- [18] P. Holme and B. J. Kim, Growing scale-free networks with tunable clustering. *Phys. Rev. E* **65**, 026107 (2002).
- [19] Guillaume J.-L. and Latapy M., A realistic model for complex networks, (2003) cond-mat/0307095
- [20] R. Milo, N. Kashtan, S. Itzkovitz, M. E. J. Newman, and U. Alon, Uniform generation of random graphs with arbitrary degree sequences, cond-mat/0312028
- [21] S.N.Dorogovtsev, J.F.F.Mendes,A.N.Samukhin, How to generate a random growing network, cond-mat/0206132 (2002)
- [22] M.E.J. Newman, Properties of highly clustered networks, *Phys. Rev. E* **68**, 026121 (2003)
- [23] Ravasz, E., Somera, A. L., Mongru, D. A., Oltvai, Z., and Barabási, A.-L., Hierarchical organization of modularity in metabolic networks, *Science* **297**, 1551–1555 (2002).
- [24] D. Stauffer and A. Aharony, *Introduction to percolation theory*, Taylor & Francis, London , 1992
- [25] M.E.J. Newman, Mixing patterns in networks, *Phys. Rev. E* **67**, 026126 (2003).
- [26] Molloy and Reed, *Random Struct. and Algorithms* **6**, 161 (1995)
- [27] D. Heckathorn, Respondent-Driven Sampling II: Deriving valid population estimates from chain-referral samples of hidden populations. *Social Problems*, 2002.
- [28] M. Salganik, D. Heckathorn, Sampling and estimation in hidden populations using respondent-driven sampling, forthcoming in *Sociological Methodology* (2004)
- [29] S.P. Blythe, C. Castillo-Chavez, G. Casella, Empirical methods for the estimation of the mixing probabilities for socially structured populations from a single survey sample, *Math. Pop. Studies* **3** 199-225 (1992)
- [30] L. Ancel-Meyers, M.E.J.Newman, Applying network theory to epidemics: control measures for mycoplasma pneumoniae outbreaks, *Emerging Inf. Dis.* **9**(2) 2003.
- [31] M. Kretzschmar, M. Morris, Measures of concurrency in networks and the spread of infectious diseases, *Math. Biosciences* **133**:165-195 (1996)
- [32] R. Cohen, K. Erez, D. ben-Avraham, S. Havlin, Breakdown of the internet under intentional attack, *Phys. Rev. E.*, v. **86** (16) pg. 3682 (2001)

- [33] D.J.Watts, A simple model of global cascades on random networks, PNAS **9**(9) 2002.
- [34] D.Centola, M. Macy, and V. Eguiluz, Maybe its not such a small world afterall. unpublished manuscript.
- [35] X. Guardiola, A. Diaz-Guilera, C.J. Perez, A. Arenas, M. Llas, arxiv.org:cond-mat/0204141 (2002)
- [36] D. Zanette, Dynamics of rumor-propagation on small-world networks, Phys. Rev. E **65** 041908
- [37] S. Ellner, Effects of successional dynamics on metapopulation persistence, Ecology **84**(4) 882-889 2003
- [38] J. Weesie, A transitive random network model, Social Networks 11 (1989) 363-386
- [39] P. Holme, C. Edling and F. Liljeros, Structure and time evolution of an internet dating community, Social Networks 26 (2004) 155-174
- [40] All networks were rendered with yEd © <http://www.yworks.com>, free for academic use.
- [41] Code for generating networks like those described in this article can be found at: <http://www.people.cornell.edu/pages/emv7/clustering>