

Clustering coefficients*

Winfried Just[†] Hannah Callender[‡] M. Drew LaMar[§]

December 23, 2015

In this module we introduce several definitions of so-called clustering coefficients. A motivating example shows how these characteristics of the contact network may influence the spread of an infectious disease. In later sections we explore, both with the help of IONTW and theoretically, the behavior of clustering coefficients for various network types.

1 A motivating example

Recall¹ that the one-dimensional nearest-neighbor networks $G_{NN}^1(N, d)$ are k -regular for $k = 2d$ when $N \geq 2d + 1$. One might expect that diseases would spread on such networks in similar ways as on random regular graphs $G_{Reg}(N, k)$. Let us see whether simulations confirm this prediction.

Open IONTW, click **Defaults**, and change the following parameter settings:

infection-prob: 0.5
end-infection-prob: 1
network-type → **Random Regular**
num-nodes: 200
lambda: 4
d: 2
auto-set: **On**

Press **New** to initialize a next-generation *SIR*-model on a network $G_{Reg}(200, 4)$ with one index case in an otherwise susceptible population. Press **Metrics** to verify that $R_0 = 2$. By the result of our module *The replacement number* at this website we should have $R_t^{st} \approx 1.5 > 1$ for sufficiently small positive t . Thus we would expect to see a significant proportion of major outbreaks in addition to some minor ones. You may want to run a few exploratory simulations to check whether this is what you will see in the **World** window and the **Disease Prevalence** plot.

*©Winfried Just, Hannah Callender, M. Drew LaMar 2015

[†]Department of Mathematics, Ohio University, Athens, OH 45701 E-mail: mathjust@gmail.com

[‡]University of Portland E-mail: callende@up.edu

[§]The College of William and Mary E-mail: drew.lamar@gmail.com

¹See Exercise 4 of our module *A quick tour of IONTW* at this website <http://www.ohio.edu/people/just/IONTW>

Now let us confirm the preliminary observations with a larger number of simulations. Using the template that is provided in the instructions at this website, set up a batch processing experiment for the current parameter settings with the following specifications:

Define a **New** experiment.

Repetitions: 100

Measure runs using these reporters:

count turtles with [removed?]

Setup commands:

new-network

Exercise 1 (a) Run the experiment and analyze your data by sorting the output column from lowest to highest. If you see a distinct gap between minor and major outbreaks, report the number of minor outbreaks, as well as the mean and maximum values of the output variable for these outbreaks. Also report the minimum, mean, and maximum for the major outbreaks, as well as the overall mean.

(b) Are the results consistent with your expectations?

Now choose

network-type → **Nearest-neighbor 1**

Press **New** to initialize a next-generation *SIR*-model on a network $G_{NN}^1(200, 2)$ with one index case in an otherwise susceptible population. Press **Clear** on the bar of the **Command Center**, then **Metrics** to verify that $\langle k \rangle = 4$ and $R_0 = 2$.

You may want to run a few exploratory simulations to check whether you see similar results in the **World** window and the **Disease Prevalence** plot as for the previous network type.

Now let us confirm the preliminary observations with setting up and running a batch processing experiment with 100 runs for the current parameter settings. You may either define a **New** experiment or **Edit** the previous one by replacing

["network-type" "Random Regular"]

with

["network-type" "Nearest-neighbor 1"]

Exercise 2 (a) Run the experiment and analyze your data as in your solution of Exercise 1.

(b) Are the results similar to the one in the previous experiment? If not, does the structure of $G_{NN}^1(200, 4)$ appear to increase or decrease the severity of outbreaks relative to the corresponding random regular graph?

What is going on here? Let us take a closer look at R_t^{st} for small positive t . For the sake of argument, let us assume $t = 1$ and 3 nodes j_1, j_2, j_3 are infectious in state st . Each of these nodes will have one neighbor (the index case) who infected this node and is no longer susceptible at time $t = 1$ in state st . Let $\mathcal{N}_1(j)$ denote the set of nodes i that are adjacent

to j . In a large random 4-regular graph, with high probability it will be the case that the union of the neighborhoods $\mathcal{N}_1(j_1), \mathcal{N}_1(j_2), \mathcal{N}_1(j_3)$ contain a total of 9 susceptible nodes to whom the pathogen could be transmitted by time step 2.

Now let us see how the situation differs in graphs $G_{NN}^1(N, 2)$. For better visualization, set

num-nodes: 12

Create a **New** network with one index case in an otherwise susceptible population. In the **World** window you will see that $\mathcal{N}_1(j^*)$ contains 4 nodes.

Move the speed slider to a very slow setting; adjust for comfortable viewing as needed. Start a simulation with **Go** and stop it by pressing **Go** again. Repeat a few times if need be until you see a state with exactly 1 removed and 3 infectious nodes in the **World** window. Count the number of green nodes at the end of red edges that could become infectious at the next time step. It will be less than 9. This effect results from the special structure of the networks and explains the discrepancies that you observed.

There are various ways to quantify this effect. One measure that is popular in the literature and goes some way towards predicting the decrease in severity of outbreaks are so-called *clustering coefficients*. In your **World** window you will see at least one white edge with two red endpoints. Look at one of these edges. No effective contact between its endpoints by time $t = 2$ can be successful, and in some sense this edge decreases the number of potential nodes that can become infectious at the next time step by 2 (one for each endpoint). Clustering coefficients indicate whether we should expect many or relatively few such edges. They explain some of the decrease in the number of candidates for infection at the next step from 9 to the one you just found.

Each of the white edges with red endpoints that you see is an edge of a triangle whose third endpoint is the grey node that represents the index case. Clustering coefficients can be defined by counting the number of potential triangles; high clustering coefficients indicate that there are a lot of them; low clustering coefficients indicate few.

Let us see how this works. Change

infection-prob: 1

Create a **New** network, run a simulation in slow motion for exactly one time step. Count the number of white edges that connect two red nodes. There should be 3 of them; each one is part of a triangle whose third vertex is the index case and whose other two edges are grey.

Are 3 white edges *a lot*? To make sense of the phrases *a lot* or *few* we need to compare the observed numbers with some benchmark. In the case of clustering coefficients, the benchmark is the complete graph.

Choose

network-type → **Complete Graph**

num-nodes: 5

Create a **New** network. Run a simulation in slow motion for exactly one time step and count the number of white edges that connect red nodes. This number gives us the

benchmark; it is the number of edges in a complete graph K_n , where n is the size of $\mathcal{N}_1(j^*)$ in the previous experiment. In our case $n = 4$ and the number of edges in the complete graph is 6.

If we divide the number of white edges that we observed in the previous experiment by 6, we obtain the *node clustering coefficient* of the index case. The formal definition will be given in the next section.

2 Definitions of clustering coefficients

Several subtly different notions of *clustering coefficient* aka *transitivity* have been studied in the literature. One always needs to carefully read the definition to see what, exactly, these terms mean in the given source. We will work with four such notions in this chapter. Here we give only their definitions and briefly describe their properties. In later sections we will explore these notions at a more leisurely pace.

Consider a node i in a graph G . Recall that $\mathcal{N}_1(i)$ denotes the set of i 's neighbors, that is, nodes that are adjacent to i . Let $tr(i)$ denote the number of edges $\{j_1, j_2\} \in E(G)$ such that $j_1, j_2 \in \mathcal{N}_1(i)$. The number $tr(i)$ is exactly the number of triangles that node i forms with two of its neighbors. Let k_i denote the degree of node i .

Watts and Strogatz [4] define the *node² clustering coefficient* $C(i)$ of i by dividing $tr(i)$ by its maximum possible value $k_i(k_i - 1)/2$.

$$C(i) = \frac{2tr(i)}{k_i(k_i - 1)}. \quad (1)$$

If G represents friendships among people, the clustering coefficient $C(i)$ measures the ratio of the number of friendships between any two of i 's friends relative to a situation where all these friends would induce a complete subgraph of G . Mathematicians actually refer to such sets that induce complete subgraphs as *cliques*.

The *network clustering coefficient* C is defined as the mean of the node clustering coefficients $C(i)$:

$$C = \frac{1}{N} \sum_{i=1}^N C(i). \quad (2)$$

Unfortunately, this definition of C only makes sense if all nodes have degree $k_i \geq 2$. If $k_i < 2$, then $C(i)$ is undefined. While one could define C in this case by taking the sum in (2) only over those nodes for which $k_i \geq 2$, and replacing N in the factor $\frac{1}{N}$ by their number, here we take a different route. If $k_i \leq 1$ and (1) does not give a definition of $C(i)$, then we interpret $C(i)$ in (2) as the *edge density*, that is, the probability $\frac{2|E(G)|}{N(N-1)}$ that two randomly chosen nodes are adjacent in G .

In Section 1 we have already seen one interpretation of clustering coefficients. Here is an alternative interpretation that is often given in the literature. Consider a network G .

²Some authors refer to node clustering coefficients as *local* clustering coefficients.

Suppose we randomly pick a node i , and then we randomly pick two nodes j_1, j_2 that are adjacent to i . Does this procedure make it *more likely* or *less likely* that the pair $\{j_1, j_2\}$ forms an edge in a given graph, relative to a completely random choice of j_1, j_2 ?

To make sense of this question, let us first observe that our procedure requires that $j_1, j_2 \in \mathcal{N}_1(i)$. If $\{j_1, j_2\}$ is an edge, then the subgraph $G^{ind}(\{i, j_1, j_2\})$ of G that is induced³ by the set of nodes $\{j_1, j_2, j_3\}$ will form a triangle. If G contains relatively many triangles, as will be the case in large nearest neighbor graphs with $d > 1$, we might expect that the answer will be “more likely.” On the other hand, we should expect the answer to be “less likely” if G contains only relatively few triangles. If G does contain some edges but no triangles at all, as in trees or graphs $G_{NN}^1(N, 1)$ for $N > 3$, then $\{j_1, j_2\}$ simply cannot be an edge and the answer will definitely be “less likely.”

Intuitively, one would expect the answer “more likely” for networks of social contacts. Two randomly chosen friends of yours are more likely to be friends of each other than two randomly chosen persons. The set of all your friends is unlikely to induce a complete subgraph, or form a single clique, in the contact network, but it is rather likely that there will be cliques among them (in the mathematical sense, not in the sense of the colloquial overtones of the word). You and your friends will form a *cluster* in the friendship graph.

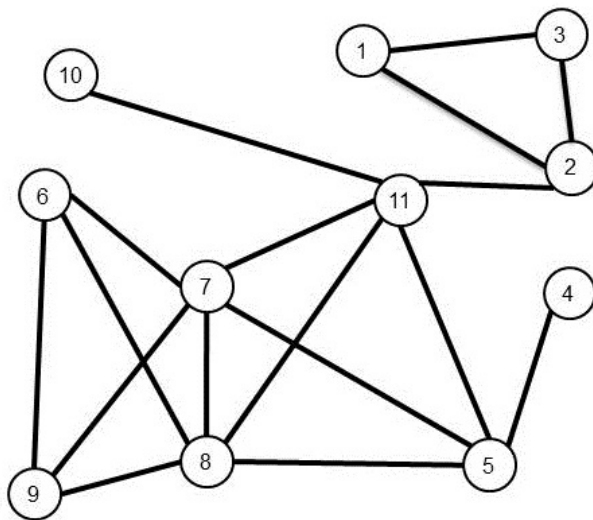


Figure 1: The graph G_1 .

Consider, for example, the graph G_1 of Figure 1. It is supposed to model a social contact network as described in Exercise 9.19 of [1]. Note that in this network, for example, $\mathcal{N}_1(9) = \{6, 7, 8\}$ is a single clique, while $\mathcal{N}_1(11)$ is the union of two cliques $\{5, 7, 8\}$ and $\{2, 10\}$.

³The subgraph $G^{ind}(V^-)$ of a graph G that is induced by a subset $V^- \subseteq V(G)$ is the graph with vertex set V^- and whose edges are exactly those pairs of vertices in V^- that are edges in G .

But what, exactly, do the phrases “relatively few” and “relatively many” triangles (or cliques of size 3) mean? As we will illustrate in Section 3, the network clustering coefficient C does not all by itself tell us whether two randomly chosen nodes are *more* likely, on average, to be adjacent in G if they share a common neighbor. To remedy this drawback of the network clustering coefficient C , let us introduce *normalized* clustering coefficients. These are obtained by dividing by the edge density, that is, by the probability $\frac{2|E(G)|}{N(N-1)}$ that two randomly chosen nodes are adjacent in G :

$$C_{norm}(i) = C(i) \frac{N(N-1)}{2|E(G)|} \quad \text{and} \tag{3}$$

$$C_{norm} = C \frac{N(N-1)}{2|E(G)|} = C \frac{N-1}{\langle k \rangle} = \frac{1}{N} \sum_{i=1}^N C_{norm}(i).$$

The second equality in the definition of C_{norm} follows from the fact that $2|E(G)| = \langle k \rangle N$. If $C(i)$ is undefined in the sense of (1) for some node i , in (3) we again interpret $C(i)$ as the edge density and obtain $C_{norm}(i) = 1$ in this case.

While $C(i), C$ are numbers between 0 and 1, the normalized clustering coefficients can take any nonnegative rational numbers as values. A value $C_{norm}(i) > 1$ indicates that the nodes in $\mathcal{N}_1(i)$ are *more* likely than average to be adjacent; a value $C_{norm}(i) < 1$ indicates that for average i the nodes in $\mathcal{N}_1(i)$ are *less* likely to be adjacent than randomly chosen nodes. If $C_{norm} > 1$ we will say the graph *exhibits clustering*; if $C_{norm} < 1$ we will say that the graph *avoids clustering*.

Exercise 3 Find the clustering coefficients $C(i), C, C_{norm}(i), C_{norm}$ for (each node i of) each of the following graphs and determine whether the graph exhibits or avoids clustering.

- (a) For the graph $G_{NN}^1(9, 2)$.
- (b) For the graph $G_{NN}^2(15, 1)$.
- (c) For the graph G_1 of Figure 1.

Exercise 4 Give an intuitive argument that for large N the normalized clustering coefficient C_{norm} in $G_{ER}(N, \lambda)$ should be very close to 1.

Many large contact networks G of interest in the study of disease transmission are *sparse*, which means that the edge density is very low. For such networks the values of the network clustering coefficient will be very close to 0 and become informative only if we compare them with a benchmark. The commonly accepted benchmark is the graph $G_{ER}(N, \lambda)$ with the same number of nodes N and the same mean degree $\langle k \rangle = \lambda$ as the network G . As you can see from Exercise 4, our normalized network clustering coefficients C_{norm} are defined in such a way that they directly give this comparison.

For many empirically studied networks the values C_{norm} are very large. This seems to be especially true if the number of nodes is large. For example, a study of the connectivity

of 6,374 servers of the internet [3] found a network with $C_{norm} = 400$, a study of the collaborations of 449,913 film actors [2] found a network with $C_{norm} = 800$, and a study of the network of 282 neurons of *C. elegans* found a network with $C_{norm} = 5.7$ [4].

This indicates that these networks exhibit some form of *strong clustering*. A mathematically precise definition of this notion poses a new mathematical challenge: How large would C_{norm} need to be so that we could confidently say that the clustering in this network is “strong”? For any given network size N , there is a theoretical upper bound on C_{norm} , but no finite upper bound exists if we allow N to be arbitrarily large. A mathematically meaningful definition of strong clustering will require us to consider a class of graphs that contains graphs of arbitrarily large size N . We can then say that this class of graphs *exhibits strong clustering* if $C_{norm} \rightarrow \infty$ a.a.s. (asymptotically almost surely), which means here that for every probability $q < 1$ and fixed C_{target} there exists $N(q, C_{target})$ such that a randomly drawn network of size $N > N(q, C_{target})$ in this class will with probability $> q$ satisfy the inequality $C_{norm} > C_{target}$.

Note that in our terminology it makes sense to say that a given network exhibits or avoids clustering. But the phrase “strong clustering” does not make sense for an individual network; it applies only to classes of networks. By Exercise 4, for any given λ the class of Erdős-Rényi networks $G_{ER}(N, \lambda)$ does not exhibit strong clustering. In the next section you will see examples of classes that do.

3 Exploring clustering coefficients of selected networks

Open IONTW and press **Defaults**. Work with the following parameter settings:

network-type → **Erdos-Renyi**
lambda: 8
num-nodes: 20, 40, 80, 160, 320

For each of the specified network sizes create one network with **New** and then press **Metrics** before creating the next network. When you are done, use the double arrow on the bar **Command Center** to enlarge this window and look at the statistics that you collected.

Exercise 5 (a) Which limit do the values of **Edge density** appear to approach as N increases?

(b) Which limit do the values of **Clustering coefficient** appear to approach as N increases?

(c) Which limit do the values of **Normalized clustering coefficient** appear to approach as N increases?

(d) Are these results consistent with what you learned in Section 2?

Press **Clear** to clean up the **Command Center** and minimize this window. Change **network-type** → **Nearest-neighbor 1**

d: 2

Repeat the steps of the data collection that you did for Erdős-Rényi networks of the sizes specified above. As you proceed, you may want to visualize the distribution of the values of $C(i)$ by choosing

plot-metric → **Normalized Coeffs**

and pressing **Update**.

Exercise 6 (a) *How would you describe the behavior of the values of Edge density as N increases?*

(b) *How would you describe the behavior of the values of Clustering coefficient as N increases?*

(c) *How would you describe the behavior of Normalized clustering coefficient as N increases?*

(d) *Does the class of networks $G_{NN}^1(N, 2)$ appear to exhibit strong clustering?*

Retain your statistics for reference and change

network-type → **Nearest-neighbor 2**

num-nodes: 25, 36, 100, 225

Repeat the steps that you did for the previous types of networks to collect data on networks $G_{NN}^2(N^2, 2)$ of the specified sizes N^2 . Inspect the data.

Exercise 7 (a) *How would you describe the behavior of the values of Edge density as N increases?*

(b) *How would you describe the behavior of the values of Clustering coefficient as N increases? How is the behavior different from the one that you observed for nearest neighbor 1 networks and how would you explain the difference?*

(c) *How would you describe the behavior of Normalized clustering coefficient as N increases?*

(d) *Does the class of networks $G_{NN}^2(N^2, 2)$ appear to exhibit strong clustering?*

Now set

d: 1

Press **New** and then **Metrics**. The command center will show you that both clustering coefficients C and C_{norm} are 0. This should be expected from the definitions, as the graph in the **World** window contains no triangles whatsoever.

Next let us explore what kind of information the different types of clustering coefficients give us about the relative likelihood that two neighbors of a randomly chosen node will form an edge compared with two randomly chosen nodes. Consider the graphs in Figures 2 and 3. We can calculate their clustering coefficients according to formulas of Subsection 2 as follows: For $i \leq 10$ we get $tr(i) = 0 = C(i)$; for $i > 10$ we get $tr(i) = 1$ and $C(i) = 1$.

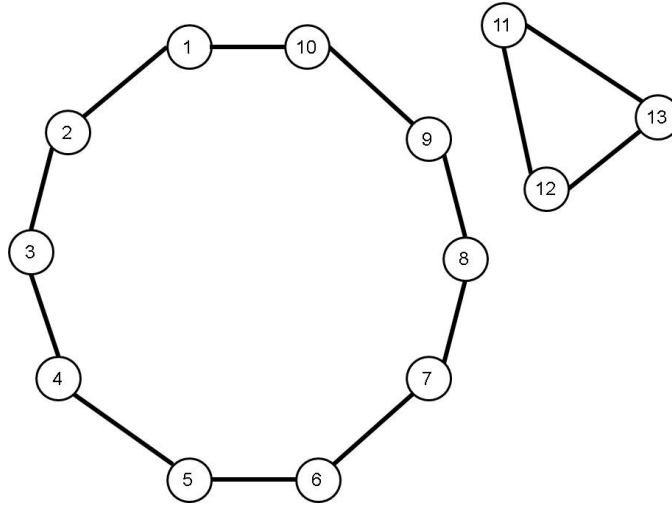


Figure 2: The graph G_2 .

By taking the mean we get the same network clustering coefficient $C = \frac{3}{13} \approx 0.23$ for both graphs.

Since the graph G_2 contains a total of 13 edges, the probability that two randomly chosen nodes are adjacent is equal to $\frac{13}{\binom{13}{2}} \approx 0.17$. Thus the average probability that two randomly chosen neighbors of a randomly chosen node are adjacent, as represented by the clustering coefficient C , is *larger* than for two nodes that are chosen completely randomly. In contrast, in G_3 the probability that two randomly chosen nodes are adjacent is equal to $\frac{33}{\binom{13}{2}} \approx 0.29$, which is larger than the clustering coefficient. It follows that in G_3 , on average, the probability that two nodes in the neighborhood $\mathcal{N}_1(i)$ of a randomly chosen node i will be *smaller* than for completely randomly chosen nodes. Thus all by itself, the network clustering coefficient C does not give us information whether two friends of one's friends are more likely to be friends than two randomly chosen people.

In contrast, the *normalized* network clustering coefficients give you this information. A value $C_{norm}(i) > 1$ indicates that the nodes in $\mathcal{N}_1(i)$ are *more* likely than average to be adjacent; a value $C_{norm} < 1$ indicates that for average i the nodes in $\mathcal{N}_1(i)$ are *less* likely to be adjacent than randomly chosen nodes. Let us check how this works out for the graphs G_2 and G_3 above.

Exercise 8 Calculate $C_{norm}(i)$ for all nodes i and C_{norm} in the graphs G_2 and G_3 of Figures 2 and 3.

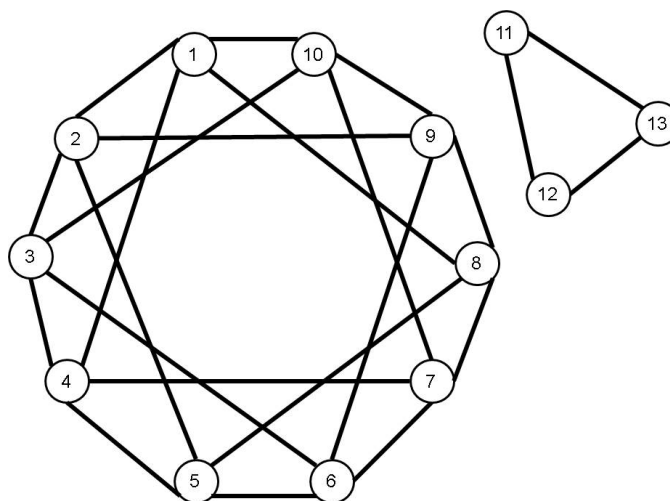


Figure 3: The graph G_3 .

4 Mathematical explorations of normalized clustering coefficients

In this section we present four theoretical results that will be of interest primarily to students with a strong mathematical background. The topics range from normalized clustering coefficients in trees, Erdős-Rényi graphs, and generic random graphs for given degree distributions to a general theorem on strong clustering. The four topics are independent of each other.

As a warm-up, we recommend the following exercise.

Exercise 9 *Suppose that G is a tree. Find a formula for $C_{norm}(G)$ in terms of the number of leaves and show that we always have $C_{norm} < 1$ if G has $N > 2$ nodes.*

Exercise 4 of Section 2 shows that the normalized clustering coefficient C_{norm} in $G_{ER}(N, \lambda)$ will be very close to 1. Clustering coefficients $C_{norm}(i)$ for individual nodes may substantially differ from 1 though. To see this, choose

network-type → **Erdos-Renyi**
num-nodes: 200
lambda: 2

Create a **New** network and press **Metrics** to look up the normalized clustering coefficient. It should be close to 1. The clustering coefficient should on average be a very small positive number. If it is 0 in your example, try again until you see a small positive number.

Now choose

plot-metric → **Normalized Coeffs**

and press **Update** to see the distribution of the normalized node clustering coefficients $C_{norm}(i)$. Most of them will be equal to 0, but a few may take very high values.

In view of our results on the degree distribution in Erdős-Rényi networks, most neighborhoods $\mathcal{N}_1(i)$ will have size on the order of $\lambda = p(N - 1)$. If λ is very small relative to N , we should expect that $C(i) = 0$ for most i . In this situation even a single edge between nodes in $\mathcal{N}_1(i)$ will result in very large values of $C_{norm}(i)$. Thus, in general, the values of $C_{norm}(i)$ will show a large range, but the effect will mostly cancel out if we compute the mean C_{norm} .

Exercise 10 Assume N is very large, but $\lambda = 2$.

(a) Derive a rough estimate of the expected maximum value of these coefficients.

(b) Check whether your estimate matches the values that IONTW displays in the plot **Network Metrics** for option **plot-metric** \rightarrow **Normalized Coeff**.

For generic random graphs, there is an interesting relation between C_{norm} and the excess $\langle k_f \rangle - \langle k \rangle$ in the friendship paradox (for definitions, see our module *The friendship paradox* at this website⁴).

Exercise 11 (a) Suppose \bar{q} is a degree distribution with $q_0 = q_1 = 0$. Show that generic random graphs $G_D(N, \bar{q})$ of large size N will satisfy

$C_{norm} > 1$ if $\langle k_f \rangle - \langle k \rangle > 1$ and

$C_{norm} < 1$ if $\langle k_f \rangle - \langle k \rangle < 1$.

(b) What can you deduce about C_{norm} for generic k -regular graphs $G_{Reg}(N, k)$ with $k \geq 2$?

The results of Exercises 6 and 7 suggest that the classes of networks $G_{NN}^1(N, 2)$ and $G_{NN}^2(N^2, 2)$ exhibit strong clustering. Let us now state and prove a general theorem that implies that this is indeed the case.

Theorem 1 Suppose we are given a class of graphs $G(N)$ that contains representatives of arbitrarily large sizes N . Moreover, assume that the mean degree $\langle k \rangle$ approaches a finite limit as N increases without bound, and $tr(i) \geq 1$ for each node i in all graphs $G(N)$. Then this class exhibits strong clustering.

First note that for all fixed $d \geq 1$ and sufficiently large N the graphs $G_{NN}^1(N, d)$ are $2d$ regular and thus satisfy the first assumption of Theorem 1. The graphs $G_{NN}^2(N^2, d)$ are not regular, but one can show that they still satisfy this first assumption for any fixed value of d . Thus in all these classes there exists some finite upper bound k_{max} on the degrees so that $k_i \leq k_{max}$ for all nodes i in all graphs $G(N)$ of the class.

Exercise 12 (a) Find an upper bound on the degree of any node in $G_{NN}^2(N^2, d)$ that does not depend on N .

(b) Show that the graphs $G_{NN}^1(N, d)$ and $G_{NN}^2(N^2, d)$ with $N > 2$ and $d > 1$ have the property that $tr(i) \geq 1$ for each node i .

⁴<http://www.ohio.edu/people/just/IONTW/>

Thus Theorem 1 implies that all classes $G_{NN}^1(N, d)$ and $G_{NN}^2(N^2, d)$ with fixed $d > 1$ exhibit strong clustering.

Exercise 13 *Prove Theorem 1 under the additional assumption that there exists some finite upper bound k_{max} on the degrees so that $k_i \leq k_{max}$ for all nodes i in all graphs $G(N)$.*

References

- [1] Winfried Just, Hannah Callender, and M Drew LaMar. Disease transmission dynamics on networks: Network structure *vs.* disease dynamics. In Raina Robeva, editor, *Algebraic and Discrete Mathematical Methods for Modern Biology*, pages 217–235. Academic Press, 2015.
- [2] Mark EJ Newman, Steven H Strogatz, and Duncan J Watts. Random graphs with arbitrary degree distributions and their applications. *Physical Review E*, 64(2):026118, 2001.
- [3] Romualdo Pastor-Satorras, Alexei Vázquez, and Alessandro Vespignani. Dynamical and correlation properties of the internet. *Physical Review Letters*, 87(25):258701, 2001.
- [4] Duncan J Watts and Steven H Strogatz. Collective dynamics of ‘small-world’ networks. *Nature*, 393(6684):440–442, 1998.