

Directed Random Geometric Graphs

Jesse Michel,¹ Sushruth Reddy,¹ Rikhav Shah,¹ Sandeep Silwal,¹ and Ramis Movassagh^{2,*}

¹Massachusetts Institute of Technology, Cambridge MA, 02139

²IBM Research, MIT-IBM AI Lab, Cambridge MA, 02142

(Dated: August 8, 2018)

Many real-world networks are intrinsically directed. Such networks include activation of genes, hyperlinks on the internet, and the network of followers on Twitter among many others [1–3]. The challenge, however, is to create a network model that has many of the properties of real-world networks such as powerlaw degree distributions and the small-world property [4, 5]. To meet these challenges, we introduce the *Directed* Random Geometric Graph (DRGG) model, which is an extension of the random geometric graph model. We prove that it is scale-free with respect to the indegree distribution, has binomial outdegree distribution, has a high clustering coefficient, has few edges and is likely small-world. These are some of the main features of aforementioned real world networks. We empirically observe that word association networks have many of the theoretical properties of the DRGG model.

Contents

I. Introduction	2
II. Previous work	2
III. Directed Random Geometric Graph (DRGG) Model	3
IV. Mathematical results	4
A. Degree distributions	4
B. Clustering Coefficient of $G(n, \alpha, d)$	6
C. Undirected edges and graph Limits	9
D. Diameter	10
V. An Application to real world networks	11
A. Degree Distribution	12
B. Other Graph Statistics	12
C. Analysis of Hubs	13
VI. Conclusions and Future Work	13
VII. Acknowledgements	14
VIII. Appendix	14
A. Proof of Theorem 2	14
1. Upper Bound Calculation	15
2. Lower Bound Calculation	15
B. Proof of Theorem 3	16
1. Computing the expectation	16
2. Computing the integral expression for \bar{C}_{in} for odd d	17
C. Analysis of Hubs	18
References	18

I. INTRODUCTION

The widespread availability of rich data on complex networks has spurred mathematical research into the properties of such networks. Researchers frequently use random graphs to model real networks [1, 3, 5–9]. There are many properties of real networks that appear in a diverse set of fields. The most prevalent of these properties are the following:

Scale-free: the distribution of the degrees of nodes is given by a power law, i.e. $P[\text{degree of } v = k] \sim k^{-\gamma}$ [7]. In general, power-laws frequently occur in systems where ‘rich’ items become richer, though there are many potential explanations [2, 7]. For networks, this means that nodes with high degree are more likely to attract connections from additional nodes added to the network. This frequently occurs when the degree is a rough measure of popularity, for example, websites that are linked to by a lot of other websites are seen as credible, so are linked to by more websites [10]. In directed graphs, the indegree and outdegree distributions can be examined separately. Networks can exhibit a power law indegree distribution but a skinny tail (e.g. Gaussian) outdegree distribution [11].

Few-edges: Networks are usually quite sparse [12]. Every node has a small expected degree that doesn’t increase much as the number of nodes scales. For example, if Twitter doubles its number of users, it’s unlikely that the average user would increase the number of users they follow by any constant factor. Real networks tend to have $\tilde{O}(n)$ edges¹.

Small-world: Most nodes are connected via a relatively short number of hops in the graph [13, 14]. Specifically, the expected length of the shortest path between two randomly selected points is $O(\log n)$. This property is related to the diameter, which is the longest shortest path in the graph. Clearly, the diameter is an upper bound on the expected length of a shortest path.

Clustering: The *clustering coefficient* is the average likelihood that given a node, a randomly selected pair of its neighbors will be connected via a single edge. Intuitively, agents tend to make connections with those they are already ‘close’ to. We can write an expression for the clustering coefficient,

$$\frac{1}{n} \sum_u P[v \rightsquigarrow w | u \rightsquigarrow v, u \rightsquigarrow w],$$

where $v \rightsquigarrow w$ is the event that there is a directed edge from v to w . This paper finds a class of random graphs which satisfies all four of these key properties.

II. PREVIOUS WORK

We briefly describe classes of random graphs that have attracted much mathematical attention. Perhaps the most prominent are Erdős-Renyi random graphs, denoted $G(n, p)$. It is natural to extend this model to the directed case by including the directed edges (v, u) and (u, v) independently with probability p . Erdős-Renyi graphs are tractable and small-world but are not scale-free.

Another well-known network with heavy-tailed power law node degree distributions is the preferential attachment model. While this model exhibits the desired power law distribution for degrees and is mathematically tractable, this model lacks certain key properties of some real world networks such as a realistic clustering coefficient, which measures the commonality of small communities [15, 16].

Random geometric graphs (RGG) are another class of random graphs that have become more popular in recent years due to their simple formulation [17]. In the simple RGG model, one uniformly distributes n points over some space (often taken to be the unit d -dimensional cube or torus), and takes them to be the vertices of a random graph. One connects any pair of vertices with distance less than a fixed R . Intuitively nodes which are closer are more connected, which is a desired feature of some real world networks. Much is known about this class of graphs, such as the sizes of connected components and the minimum value of R which results in a connected graph [18]. These types of random graphs are commonly used to model radio broadcasting towers and there has been much work on distributed algorithms on these models (see [8, 9]). Some variations on the basic RGG model have also been studied. One variant, termed the k -nearest neighbor model, samples points as in the basic RGG model but only connects each

¹ $\tilde{O}(n) = \mathcal{O}(n \text{ poly}(\log n))$

point to its k nearest neighbors [19]. Another variant generalizes the RGG model by sampling points as before, but connecting nodes with a probability dependent on their distance [20].

III. DIRECTED RANDOM GEOMETRIC GRAPH (DRGG) MODEL

Consider n points (nodes/vertices) uniformly distributed on a d -dimensional cube $[0, L]^d$. To make the model fully translation-invariant, we impose periodic boundary conditions, that is, for any i ,

$$(x_1, \dots, x_{i-1}, x_i + L, x_{i+1}, \dots, x_d) = (x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_d).$$

This is equivalent considering points on a d -dimensional torus. For simplicity, below we take the unit cube (i.e., set $L = 1$), but our results can be extended to a general L . Each point v is assigned a radius r_v distributed according to a Pareto distribution [21],

$$f(r) = \begin{cases} \eta/r^\alpha & r_0 \leq r \leq 1/2 \\ 0 & \text{otherwise} \end{cases},$$

where $\alpha > d + 1$ is a parameter of the graph and r_0 is the minimum allowed radius, chosen based on n so that the resulting graph is almost surely connected [18],

$$r_0 = \left(\frac{\log n}{V_d n} \right)^{1/d}.$$

where V_d be the volume of the unit ball in d -dimensions (i.e. $v_1 = 2, v_2 = \pi, v_3 = 4\pi/3$). $r \leq 1/2$ ensures that balls of selected radii do not intersect themselves. η is the normalizing factor that ensures $\int_{r_0}^{1/2} f(r) dr = 1$:

$$\eta = \frac{(\alpha - 1)r_0^{\alpha-1}}{1 - (2r_0)^{\alpha-1}}.$$

We refer to the ball of radius r_v centered at v as “ v ’s circle” and denote the distance on the torus between v and u as $d(u, v) = d(v, u)$. The directed edges are added as follows. Add a directed edge from u to v (denoted by $u \rightsquigarrow v$) if u is in v ’s circle (i.e. $d(v, u) \leq r_v$). Intuitively, r_v is a measure of the popularity of v . If v is unpopular, then only nodes close by will connect; conversely very popular nodes attract the attention of nodes that are far away. Lastly, since we want to model real-world networks, we can think of d as being a small integer, such as $1 \leq d \leq 5$. See Figures (1) and (2) for reference.

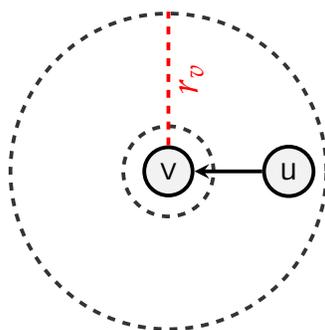


Figure 1: Illustration of $u \rightsquigarrow v$: There is a directed edge from vertex u to vertex v since u lies inside v ’s circle of radius r_v . There is no directed edge from v to u since v lies outside of u ’s circle.

We refer to this as the Directed Random Geometric Graph (DRGG) model, denoted by $G(n, \alpha, d)$. In summary the notation is:

- $G(n, \alpha, d)$ is a DRGG where n is the number of nodes and d is the dimension of the space.
- $f(r) = \eta/r^\alpha$ probability density function for the selection of radii.

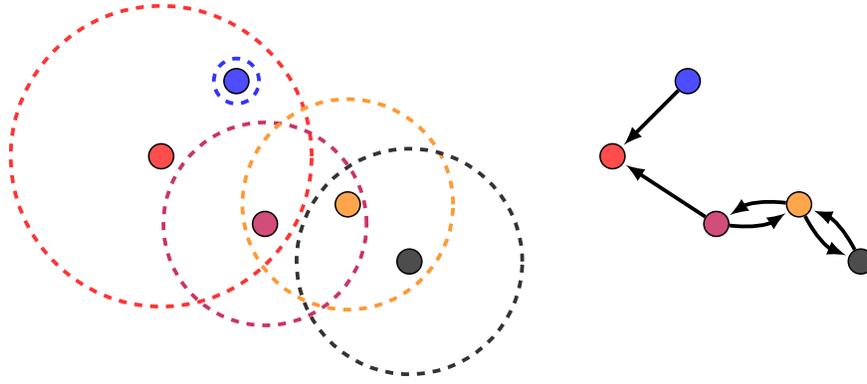


Figure 2: Left: Vertices with their respective circles plotted in the same color. Right: The corresponding DRGG.

- $[r_0, 1/2]$ is the support of $f(r)$.

IV. MATHEMATICAL RESULTS

For any class of random graphs, there are a variety of properties of theoretical and practical interests as outlined in the introduction. In this section, we prove that the indegree distribution of DRGGs is given by a power law with parameter $\frac{\alpha-1}{d} + 1$ while the outdegree distribution is a binomial with parameter $\Theta(\frac{\log n}{n})$ where n is the number of vertices. We also show that the total number of edges in DRGGs is $\Theta(n \log n)$ and show that the clustering coefficient of DRGGs approaches a constant as $n \rightarrow \infty$. Finally, we show how the standard geometric random graphs can be thought of as a graph limit of the DRGG model. We end with a conjecture about the diameter of DRGGs.

A. Degree distributions

As explained in Section I, motivated by real-world networks, we want a graph model whose indegree distribution follows a power law, and whose outdegree distribution decays exponentially. Here we prove that indeed in DRGG, the indegrees follow a power law distribution, and outdegrees a binomial distribution.

Theorem 1. *Let $\delta_{in}(k)$ and $\delta_{out}(k)$ be the density function for the indegree and outdegree of the random graph $G(n, \alpha, d)$ respectively. Then, asymptotically as $n \rightarrow \infty$,*

$$\delta_{in}(k) \propto k^{-\frac{\alpha-1}{d}-1} \quad (1)$$

for fixed d , $\alpha > d + 1$, and $\frac{\alpha-1}{d} \ll k \ll n$ and

$$\delta_{out}(k) = \binom{n-1}{k} z^k (1-z)^{n-1-k} \quad (2)$$

where $z = \Theta(\log n/n)$.

Proof. We first prove the indegree case. Consider the probability that a vertex v has indegree k . For this to happen, there must be exactly k other vertices in v 's circle (of radius r_v). The probability of this is

$$P[\delta_{in}(v) = k] = \int_{r_0}^{1/2} P[r_v = r] P[\delta_{in}(v) = k \mid r_v = r] dr. \quad (3)$$

Since the points are distributed uniformly on the plane, $P[\delta_{in}(v) = k \mid r_v = r]$ is the binomial random variable that describes the probability of exactly k points lying within the volume of the d -dimensional sphere. This gives

$$P[\delta_{in}(v) = k] = \eta \binom{n}{k} \int_{r_0}^{1/2} \frac{1}{r^\alpha} (V_d r^d)^k (1 - V_d r^d)^{n-k} dr. \quad (4)$$

Making the substitution $u = V_d r^d$ gives us

$$P[\delta_{in}(v) = k] = \frac{\eta}{d} \binom{n}{k} V_d^{\frac{\alpha-1}{d}} \int_{V_d r_0^d}^{V_d/2^d} u^{k-1-\frac{\alpha-1}{d}} (1-u)^{n-k} du. \quad (5)$$

We let $\beta = \frac{\alpha-1}{d}$ for convenience and express the integrand in terms of the exponential function obtaining

$$P[\delta_{in}(v) = k] = \frac{\eta}{d} \binom{n}{k} V_d^\beta \int_{V_d r_0^d}^{V_d/2^d} \exp \left[n \left(\frac{k-\beta-1}{n} \ln u + \left(1 - \frac{k}{n}\right) \ln(1-u) \right) \right] du.$$

Now let $f(u) = \frac{k-\beta-1}{n} \ln u + \left(1 - \frac{k}{n}\right) \ln(1-u)$ so that the integrand is of the form $e^{nf(u)}$. Treating n as the large part and setting $f'(u) = 0$, the saddle point is $u_{sp} = \frac{k-\beta-1}{n-\beta-1}$ and

$$f''(u_{sp}) = -\frac{(n-\beta-1)^3}{n(k-\beta-1)(n-k)}.$$

Note that for large n , $f''(u_{sp})$ approaches a negative constant. The steepest descent method gives (see [22])

$$P[\delta_{in}(v) = k] \approx \frac{\eta}{d} \binom{n}{k} V_d^\beta e^{nf(u_{sp})} \sqrt{\frac{2\pi}{n |f''(u_{sp})|}}, \quad (6)$$

which after inserting $f(u_{sp})$ and $f''(u_{sp})$ reads

$$P[\delta_{in}(v) = k] \approx \frac{\eta}{d} \binom{n}{k} V_d^\beta \frac{\sqrt{2\pi}(k-\beta-1)^{k-\beta-1/2} (n-k)^{n-k+1/2}}{(n-\beta-1)^{n-\beta+1/2}}. \quad (7)$$

In fact, as $n \rightarrow \infty$, the two sides of Eq. (6) become asymptotically equivalent [23]. As show in figure 3, the steepest descent result gives an excellent approximation. One can ignore the term $\frac{\eta}{d}$ since it do not depend on k so they only contribute to the normalizing constant for δ_{in} . We now show that for large n , Eq. (7) gives rise to a power law in k . Using Stirling's approximation one obtains

$$\frac{1}{\sqrt{2\pi}(n-\beta-1)} \left(\frac{n}{k-\beta-1} \right)^{\beta+1} \leq \binom{n}{k} \frac{(k-\beta-1)^{k-\beta-1/2} (n-k)^{n-k+1/2}}{(n-\beta-1)^{n-\beta+1/2}} \leq \frac{1}{\sqrt{2\pi}(n-\beta-1)} \left(\frac{n}{k} \right)^{\beta+1}.$$

Ignoring factors that do not depend on k , we see that $P[\delta_{in}(v) = k] \propto k^{-\beta-1}$ as desired.

We now analyze the outdegree case. Consider a vertex v . v has outdegree k if it lies in exactly k circles of other vertices. Now consider any other vertex $u \neq v$. The probability that v lies inside u 's circle is given by

$$z = \int_{r_0}^{1/2} P[v \text{ in } u\text{'s ball} | r_u = r] P[r_u = r] dr,$$

where $P[v \text{ in } u\text{'s ball} | r_u = r] = V_d r^d$ and $P[r_u = r] = \eta/r^\alpha$. Note that z does not depend on v or u . Since the radii of all circles are picked independently, it follows that the probability that the outdegree of vertex v is k is $\binom{n-1}{k} z^k (1-z)^{n-k}$. One can calculate z exactly to be

$$z = \int_{r_0}^{1/2} \frac{\eta}{r^\alpha} V_d r^d dr = \frac{\eta V_d}{d-\alpha+1} \left(\frac{1}{2^{d-\alpha+1}} - r_0^{d-\alpha+1} \right).$$

Defining $\alpha = \beta d + 1$ for $\beta > 1$ the above expression reads

$$z = \frac{\beta V_d}{\beta-1} \left(\frac{r_0^{\beta d} 2^{-d} - 2^{-\beta d} r_0^d}{r_0^{\beta d} - 2^{-\beta d}} \right).$$

Using Taylor expansion, and recalling that $\beta = \frac{\alpha-1}{d}$ and $1 \gg z$, the equation above reduces to

$$z = \left(\frac{\alpha - 1}{\alpha - 1 - d} \right) \frac{\log n}{n} + \mathcal{O} \left(\frac{\log n}{n} \right)^{\beta-1}$$

and we are done. □

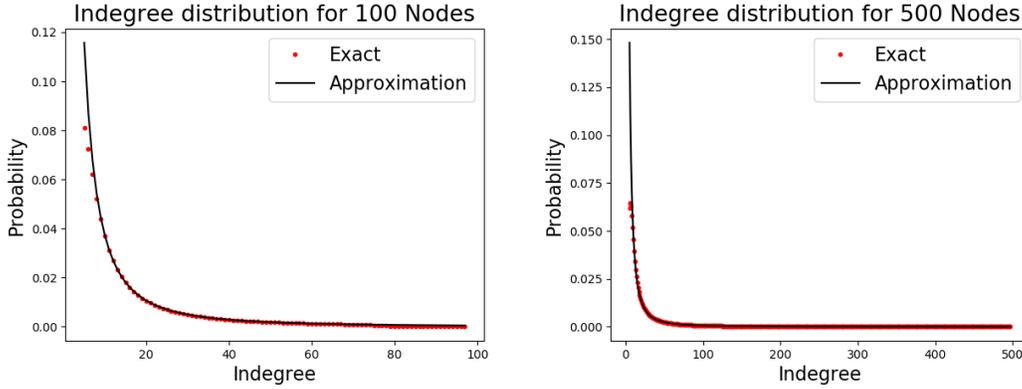


Figure 3: Illustration of the approximation of Eq. (7) for 100, and 500 vertices for $G(n, \alpha, d)$. The approximation is very accurate for both cases. For 100 nodes we see that the fit is tight except at the head and tail, which aligns with the fact that steepest descent makes a Gaussian approximation. We see a tighter fit for $n = 500$, which aligns with the asymptotic nature of our approximation in 1. A heavy tail distribution can also be seen which numerically confirms the results proven in 1.

Corollary 1. *The expected number of edges in $G(n, \alpha, d)$ is $\Theta(n \log n)$.*

Proof. To find the total number of edges, it suffices to count all the outdegrees edges. Since there are $n - 1$ nodes that v can connect to with independent probabilities, Theorem 1 implies that the expected outdegree of any vertex is $\Theta(\log n)$. Hence, the total number of edges is $\Theta(n \log n)$. □

B. Clustering Coefficient of $G(n, \alpha, d)$

The clustering coefficient of $G(n, \alpha, d)$ is discussed in this section. This metric is important because it measures the propensity of nodes with common neighbors to themselves be connected [4, 24, 25]. Hence, real world networks often exhibit large clustering coefficients. We begin by showing that $G(n, \alpha, d)$ is rich in triangles, namely, $G(n, \alpha, d)$ has $\Theta(n \log^2(n))$ triangles. Triangles are important because they are a measure of connectivity and contribute to the clustering coefficient that will be discussed later in this section. In addition, counting the total number of triangles in a network is itself a well-studied problem [26]. In a directed graph, there are fundamentally two types of triangles: **Type 1** and **Type 2** triangles. These triangles are shown in Figure (4). Due to the geometric nature of $G(n, \alpha, d)$, it

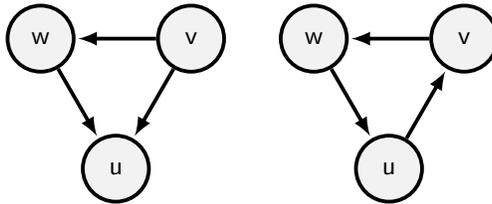


Figure 4: A type 1 triangle on the left and a type 2 triangle on the right.

turns out that we do not need to account for both types of triangles. Say that vertices u, v, w form a type 1 or 2 triangle if the edge relationships between u, v and w match Figure (4). Note that we say $v \rightsquigarrow w$ if there is a directed edge from v to w .

Lemma 1. *If vertices u, v and w form a type 2 triangle in $G(n, \alpha, d)$, then they also form a type 1 triangle.*

Proof. We note that if u, v and w form a type 2 triangle then $d(w, v) \leq r_w, d(v, u) \leq r_v, d(u, w) \leq r_u$. Without loss of generality, suppose that r_v is the minimum radius: $r_v = \min(r_u, r_v, r_w)$. Then, we have $r_v \leq r_u$ so it follows that

$$d(v, u) \leq r_v \leq r_u.$$

Hence, there must also be a directed edge from v to u . That is, u, v and w also form a type 1 triangle, as desired. \square

Note that type 1 triangles measure the following phenomenon in real networks: If two people follow the same person, what is the probability that one follows the other? Indeed, in real world networks, the number of triangles is ‘large’, which indicates that models such as preferential attachment do not accurately model real world networks [4]. However, the total number of triangles of type 1 in $G(n, \alpha, d)$ is quite ‘large.’

Theorem 2. *The expected number of directed, acyclic (type 1) triangles in a graph, $|S|$, for $S = \{u, v, w : v \rightsquigarrow u, w \rightsquigarrow u, v \rightsquigarrow w\}$, is $\mathbb{E}[|S|] = \Theta(n \log^2(n))$ as $n \rightarrow \infty$ for fixed α, d under the constraint that $\alpha > 2d$.*

Proof. We give an abridged version of the proof. Please see the appendix for full details. Throughout the proof, $P[\cdot]$ denotes the probability of an event while $p(\cdot)$ denotes the probability density of a random variable. By linearity of expectation,

$$\mathbb{E}[|S|] = n(n-1)(n-2) P(v \rightsquigarrow u, w \rightsquigarrow u, v \rightsquigarrow w)$$

for randomly selected vertices u, v, w . Without loss of generality, we may take u to be located at the origin. Then condition on the locations of v and w , denoted by \mathbf{x}_v and \mathbf{x}_w , as well as the radii of circles centered at u and w , denoted by r_u and r_w .

It is clear that the radii of circles corresponding to u and w , along with the locations of v and w are independent. Additionally, edges are conditionally independent given these radii and their locations. The probability density of the points x_v and x_w is 1 since v, w are chosen uniformly over the unit torus: i.e. $p(\mathbf{x}_v) = p(\mathbf{x}_w) = 1$. Furthermore, $P[v \rightsquigarrow u | \mathbf{x}_v, r_u] = \mathbb{1}[d(\mathbf{x}_v, 0) < r_u]$, as we draw an edge from v to u if and only if \mathbf{x}_v is inside the circle of radius r_u centered at u (the origin). Similarly, $P[w \rightsquigarrow u | \mathbf{x}_w, r_u] = \mathbb{1}[d(\mathbf{x}_w, 0) < r_u]$. Finally, we have that $P[v \rightsquigarrow w | \mathbf{x}_w, \mathbf{x}_v, r_w] = \mathbb{1}[d(\mathbf{x}_v, \mathbf{x}_w) < r_w]$, as we draw an edge from v to w if and only if the distance between their coordinates is less than the radius of w 's circle. Note that $\mathbb{1}[d(\mathbf{x}_v, \mathbf{x}_w) < r_w] \leq 1$ is always true. Using these ideas, we arrive at the following inequality

$$P[v \rightsquigarrow u, w \rightsquigarrow u, v \rightsquigarrow w] \leq V_d^2 \eta \int_{r_0}^{\frac{1}{2}} r_u^{2d-\alpha} = V_d^2 \eta \left[\frac{r_u^{2d-\alpha+1}}{2d-\alpha+1} \right]_{r_0}^{\frac{1}{2}}.$$

As n becomes large, for any fixed d we have that r_0 goes to 0. This then implies (for $\alpha > 2d$) that the dominant term in the above (and a valid, asymptotically tight upper bound) is:

$$V_d^2 \eta \frac{r_0^{2d-\alpha+1}}{\alpha-2d-1} = C_1 \frac{\log^2(n)}{n^2}$$

where C_1 is some constant independent of n . Using similar reasoning as above, we can get a lower bound on $P[v \rightsquigarrow u, w \rightsquigarrow u, v \rightsquigarrow w]$. Specifically, in the appendix, it is shown that

$$\begin{aligned} P[v \rightsquigarrow u, w \rightsquigarrow u, v \rightsquigarrow w] &\geq V_d^2 \eta^2 \int_{r_0}^{\frac{1}{4}} dr_u r_u^{2d-\alpha} \left(\frac{(2r_u)^{1-\alpha} - (1/2)^{1-\alpha}}{\alpha-1} \right) \\ &= \frac{V_d^2 \eta^2}{\alpha-1} \left[\frac{2^{-\alpha} r_0^{2d-2\alpha-2}}{\alpha-d-1} \left(1 - (4r_0)^{2(\alpha+1)-2d} \right) - \frac{2^{\alpha-1} r_0^{2d-\alpha+1}}{\alpha-2d-1} \left(1 - (4r_0)^{\alpha-1-2d} \right) \right]. \end{aligned}$$

For fixed d , as n goes to infinity, r_0 goes to 0, so that the leading order term here is the one containing the factor $r_0^{2d-2\alpha-2}$ (as $\alpha > 2d > d-1$). Thus, we have the lower bound

$$P[v \rightsquigarrow u, w \rightsquigarrow u, v \rightsquigarrow w] \geq (C_2 - \epsilon) \frac{\log^2(n)}{n^2}$$

for any $\epsilon > 0$, for large enough n . This proves that, for fixed α, d :

$$C_1 \frac{\log^2(n)}{n^2} \geq P[v \rightsquigarrow u, w \rightsquigarrow u, v \rightsquigarrow w] \geq (C_2 - \epsilon) \frac{\log^2(n)}{n^2}.$$

Multiplying by $n(n-1)(n-2)$, we get:

$$C_1 n \log^2(n) \geq \mathbb{E}[|S|] \geq (C_2 - \epsilon) n \log^2(n)$$

for any ϵ , given that n is sufficiently large. This proves that $\mathbb{E}[|S|] \in \Theta(n \log^2(n))$, as desired. \square

We now turn to the clustering coefficient. For a directed graph, there exist many different possible definitions of the clustering coefficient [24, 25]. One natural choice is \overline{C}_{in} , which is defined as:

$$\overline{C}_{\text{in}} = \frac{1}{|V|} \sum_{v \in V} c_v$$

where c_u is the local clustering coefficient for vertex u and is defined as:

$$c_u = \frac{|\{v, w : v \rightsquigarrow u, w \rightsquigarrow u, v \rightsquigarrow w\}|}{d_u(d_u - 1)}$$

where d_u is the degree of vertex u . In this section, we find the asymptotic expectation value of this clustering coefficient as $n \rightarrow \infty$, and show it approaches a constant (dependent on d and α). Note that the expectation is over all possible generated graphs. This is in contrast to other random graph models such as Erdős-Renyi random graphs and preferential attachment graphs where the clustering coefficient are very small or approach 0 as $n \rightarrow \infty$ [16, 27]. We then examine the case $\alpha \rightarrow \infty$, and show that it matches results from the random geometric graph model with radius r_0 . This is in line with the correspondence outlined in the previous section, as in this limit, all edges become undirected and the radius concentrates at r_0 , as in the case of random geometric graphs. Our main theorem is as follows:

Theorem 3. *As $n \rightarrow \infty$ for fixed d , α satisfying $\alpha > 2d + 1$, the expectation of the clustering coefficient approaches a constant which is **independent** of n (the number of vertices). For **odd** d , the constant can be calculated exactly:*

$$(\alpha - 1)^2 \frac{d!!}{(d-1)!!} \sum_{k=0}^{\frac{d-1}{2}} \frac{(-1)^k}{2k+1} \binom{\frac{d-1}{2}}{k} \left[\frac{1}{(\alpha-1)^2 - d^2} - \frac{d}{2^{2k+1}(2k+d+1)} \left(\frac{1}{(\alpha-1)^2 - (2k+d+1)^2} \right) \right].$$

Proof. We give an abridged proof. Please see the appendix for full details. By linearity of expectation (where the expectation is over all possible DRGG), $\mathbb{E}[\overline{C}_{\text{in}}] = \mathbb{E}[c_u]$ for a randomly selected vertex u . Moreover,

$$\mathbb{E}[c_u] = \mathbb{E} \left[\frac{|\{v, w : v \rightsquigarrow u, w \rightsquigarrow u, v \rightsquigarrow w\}|}{d_u(d_u - 1)} \right] = \sum_k P(d_u = k) \frac{\mathbb{E}[|\{v, w : v \rightsquigarrow u, w \rightsquigarrow u, v \rightsquigarrow w\}| | d_u = k]}{k(k-1)}$$

from the law of iterated expectation. Furthermore,

$$\mathbb{E}[|\{v, w : v \rightsquigarrow u, w \rightsquigarrow u, v \rightsquigarrow w\}| | d_u = k] = k(k-1)P(v \rightsquigarrow w | v \rightsquigarrow u, w \rightsquigarrow u, d_u = k)$$

by linearity of expectation. Then, we can substitute this expression into the sum given for $\mathbb{E}[c_u]$:

$$\mathbb{E}[\overline{C}_{\text{in}}] = \sum_k P[d_u = k] P[v \rightsquigarrow w | v \rightsquigarrow u, w \rightsquigarrow u, d_u = k].$$

After some computation, detailed in the Appendix, one obtains:

$$\mathbb{E}[\overline{C}_{\text{in}}] = \int dr_u p(r_u) \int dr_w p(r_w) \int_{\|\mathbf{x}_v\| < r_u} d^d \mathbf{x}_v \int_{\|\mathbf{x}_w\| < r_u} d^d \mathbf{x}_w \frac{\mathbb{1}[\|\mathbf{x}_v - \mathbf{x}_w\| < r_w]}{V_d^2 r_u^{2d}}.$$

For odd d , we can compute this explicitly (work detailed in the Appendix):

$$\mathbb{E}[\overline{C}_{in}] = \frac{dd!!}{(d-1)!!} \sum_{k=0}^{\frac{d-1}{2}} \frac{(-1)^k}{2k+1} \binom{\frac{d-1}{2}}{k} \left(\frac{1}{d} \frac{(\alpha-1)^2}{(\alpha-1)^2 - d^2} - \frac{1}{2^{2k+1}(2k+d+1)} \frac{(\alpha-1)^2}{(\alpha-1)^2 - (2k+d+1)^2} \right). \quad (8)$$

For even d , we cannot get nice closed forms since the density function of two randomly chosen points does not have a nice form like the one that exists for odd d . However, we can perform asymptotic analysis and show that the clustering coefficient also approaches a constant as $n \rightarrow \infty$ for even d as well. \square

Remark 1. Note that from the expression in Eq. (8), it is clear that our value of $\mathbb{E}[\overline{C}_{in}]$ is a constant independent of n . This is different than other standard models such as the Erdős-Renyi random graphs which are known to have low clustering coefficient [27] and the preferential attachment model which have the property that the clustering coefficient approaches 0 as $n \rightarrow \infty$ [16].

It is interesting to take the limit as $\alpha \rightarrow \infty$ in Eq.(8) which gives:

$$\lim_{\alpha \rightarrow \infty} \mathbb{E}[\overline{C}_{in}] = \frac{dd!!}{(d-1)!!} \sum_{k=0}^{\frac{d-1}{2}} \frac{(-1)^k}{2k+1} \binom{\frac{d-1}{2}}{k} \left(\frac{1}{d} - \frac{1}{2^{2k+1}(2k+d+1)} \right).$$

This limit **matches the known clustering coefficient derived by Dall and Christensen** of the standard RGG model [17]. This is not surprising since taking $\alpha \rightarrow \infty$ in $G(n, \alpha, d)$ results in the standard RGG model, as explored in the next section. The clustering coefficient for various odd dimensions along with their values as $\alpha \rightarrow \infty$ is shown in Table III in the appendix.

C. Undirected edges and graph Limits

In this section, we prove a lemma which shows that, given the existence of an edge, there is an asymptotically constant probability of an edge in the opposite direction. This shows that we have a positive fraction of what we have termed ‘undirected edges.’

Lemma 2. For fixed α, d , and for any randomly selected vertices u, v , as $n \rightarrow \infty$, we have that $P[u \rightsquigarrow v | v \rightsquigarrow u] \rightarrow \frac{2\beta-2}{2\beta-1}$ where $\beta = \frac{\alpha-1}{d}$.

Proof. The conditional probability can be written as $\frac{P[u \rightsquigarrow v, v \rightsquigarrow u]}{P[v \rightsquigarrow u]}$. We now compute the numerator and denominator separately. As before, we can without loss of generality situate u at the origin. Then,

$$P[u \rightsquigarrow v, v \rightsquigarrow u] = \int_{r_0}^{\frac{1}{2}} dr_u \int_{[0,1]^d} d^d \mathbf{x}_v \int_{r_0}^{\frac{1}{2}} dr_v P[u \rightsquigarrow v, v \rightsquigarrow u | r_u, \mathbf{x}_v, r_v] \cdot p(r_u, \mathbf{x}_v, r_v).$$

Since u, v share an undirected edge if and only if their separation is less than the radii of both of their circles, we have

$$P[u \rightsquigarrow v, v \rightsquigarrow u | r_u, \mathbf{x}_v, r_v] = \mathbb{1}[d(\mathbf{x}_v, 0) < r_u] \mathbb{1}[d(\mathbf{x}_v, 0) < r_v]$$

Moreover, $p(r_u, \mathbf{x}_v, r_v) = p(r_u)p(\mathbf{x}_v)p(r_v)$ as each of these 3 quantities is chosen independently. Then, the integral is rewritten:

$$\int_{\|\mathbf{x}_v\| < \frac{1}{2}} d^d \mathbf{x}_v \int_{\|\mathbf{x}_v\|}^{\frac{1}{2}} dr_u p(r_u) \int_{\|\mathbf{x}_v\|}^{\frac{1}{2}} dr_v p(r_v) = \int_{\|\mathbf{x}_v\| < \frac{1}{2}} d^d \mathbf{x}_v P[r_u > \|\mathbf{x}_v\|] P[r_v > \|\mathbf{x}_v\|].$$

The integrand only depends on the norm of \mathbf{x}_v ; we use spherical coordinates to write it as:

$$S_{d-1} \int_0^{\frac{1}{2}} r^{d-1} P[r_u > r] P[r_v > r] dr$$

where $S_{d-1} = \frac{2\pi^{d/2}}{\Gamma(\frac{d}{2})}$. We can split this integral into portions that go from 0 to r_0 , where $P[r_u > r]$ and $P[r_v > r]$ each

evaluate to 1, and a portion going from r_0 to $\frac{1}{2}$. By direct computation and for $b \geq r_0$, one obtains

$$\int_b^{\frac{1}{2}} dr p(r) = \frac{\eta}{\alpha - 1} \left(\frac{1}{b^{\alpha-1}} - 2^{\alpha-1} \right),$$

and we have

$$P[u \rightsquigarrow v, v \rightsquigarrow u] = S_{d-1} \left[\frac{r_0^d}{d} + \frac{\eta^2}{(\alpha - 1)^2} \int_{r_0}^{\frac{1}{2}} r^{d-1} \left(\frac{1}{r^{\alpha-1}} - 2^{\alpha-1} \right)^2 dr \right].$$

We are interested in the behavior of this integral as $n \rightarrow \infty$, or equivalently $r_0 \rightarrow 0$. Therefore, for sufficiently large n ,

$$\int_{r_0}^{1/2} r^{d-1} \left(\frac{1}{r^{\alpha-1}} - 2^{\alpha-1} \right) dr = \frac{r_0^{d-2\alpha+2}}{2\alpha - d - 2} + \mathcal{O}(r_0^{d-\alpha+1}) + C$$

for some constant C . Using the Taylor expansion of η around 0 and noting that $\alpha > d + 1$ (see III), we get

$$P[u \rightsquigarrow v, v \rightsquigarrow u] = S_{d-1} \left(\frac{r_0^d}{d} + \frac{r_0^d}{2\alpha - d - 2} \right) + \mathcal{O}(r_0^{d+\alpha-1}).$$

For the denominator of the conditional probability, assuming that u is located at the origin, we compute

$$P[v \rightsquigarrow u] = \int_{r_0}^{\frac{1}{2}} dr_u p(r_u) P[v \rightsquigarrow u | r_u].$$

Note that $P[v \rightsquigarrow u | r_u] = P[||\mathbf{x}_v|| < r_u] = V_d r_u^d$. Hence,

$$P[v \rightsquigarrow u] = V_d \eta \int_{r_0}^{\frac{1}{2}} r^{-\alpha+d} dr = V_d r_0^d \left(\frac{\alpha - 1}{\alpha - d - 1} \right) + \mathcal{O}(r_0^{\alpha-1}).$$

Dividing $P[u \rightsquigarrow v, v \rightsquigarrow u]$ by $P[v \rightsquigarrow u]$ and using the fact that $r_0 \rightarrow 0$, we get that

$$\lim_{n \rightarrow \infty} P[u \rightsquigarrow v | v \rightsquigarrow u] = \frac{S_{d-1}}{V_d} \left(\frac{1}{d} + \frac{1}{2\alpha - d - 2} \right) \frac{\alpha - d - 1}{\alpha - 1} = \frac{2\beta - 2}{2\beta - 1},$$

since $\frac{S_{d-1}}{V_d} = d$ where $\beta = \frac{\alpha-1}{d}$. □

Remark 2. We can use Lemma 2 to understand the limiting behavior of $G(n, \alpha, d)$ for fixed n and d and as $\alpha \rightarrow \infty$. In this case, the probability density function for the radii of the vertices converges to a delta distribution at the minimum radius r_0 . However, we can actually say something stronger. In the $\alpha \rightarrow \infty$ case, $G(n, \alpha, d)$ actually converges to a **standard undirected random geometric graph** with fixed radius r_0 . Note that r_0 is the sharp connectivity threshold for undirected random geometric graphs with fixed radius. To show this, note that the asymptotics we arrived at $P[u \rightsquigarrow v | v \rightsquigarrow u]$ are also valid for fixed n and d and $\alpha \rightarrow \infty$ since the r_0 term still dominates. Thus,

$$\lim_{\alpha \rightarrow \infty} P[u \rightsquigarrow v | v \rightsquigarrow u] = \frac{S_{d-1}}{dV_d} = 1.$$

This proves the following corollary.

Corollary 2. For fixed n and d , $\lim_{\alpha \rightarrow \infty} G(n, \alpha, d)$ converges to the standard random geometric graph model (RGG).

D. Diameter

The diameter of a graph is defined as the longest path among the set of shortest paths over all pairs of vertices. In our case, we are only concerned with directed paths. Similar to the clustering coefficient, the diameter of a graph is

a measure of connectivity. It is common for real-world networks to have a small diameter, as can be seen from the popular ‘six-degrees of separation’ phenomenon. Based on numerical results, we conjecture that the diameter of the DRGG model is $\mathcal{O}(\log n)$. It appears that the ‘fat tail’ property of the radii distribution contributes to the significant reduction in the diameter of the DRGG model. However, even though we are not able to prove our conjecture for the diameter, it is still possible to prove a related result which hints that the diameter is indeed $\mathcal{O}(\log n)$.

Lemma 3. *Let α and d be fixed in $G(n, \alpha, d)$. Pick two vertices u and v uniformly at random in $G(n, \alpha, d)$. Let a_k denote the number of directed paths of length k from u to v . If $k \geq \frac{\log n}{\log \log n}$, then $\mathbb{E}[a_k] \geq 1$.*

Proof. First consider three vertices w_1, w_2 and w_3 chosen uniformly at random. If there is a directed edge from w_1 to w_2 , then this does not affect the probability of an edge from w_2 to w_3 . This is because once the locations of w_1 and w_2 are fixed, the location of w_3 is still uniform, and having an edge from w_2 to w_3 only depends on the location of w_3 . Thus, given a length k directed path from u to v , the edges of this path are independent of one another. Hence, if a_k denotes the expected number of directed paths of length k from u to v , we have

$$\mathbb{E}[a_k] = \binom{n-2}{k-1} z^k (k-1)! \tag{9}$$

where z is the probability of an edge in $G(n, \alpha, d)$ which is given by $\frac{C \log n}{n}$ (Theorem 1) where $C = \frac{\alpha-1}{\alpha-1-d}$. Taking n large, shifting $k \rightarrow k+1$, and using $\binom{n}{k} \geq \frac{n^k}{k^k}$, we have

$$\mathbb{E}[a_k] \geq \sqrt{2\pi k} \left(\frac{C \log n}{e} \right)^k \frac{C \log n}{n}.$$

Hence, it suffices to find a k such that $\left(\frac{\log n}{e}\right)^k \geq n$. This is equivalent to finding a k such that $k \log \log n - k \geq \log n$. Rearranging, we see that $k \geq \frac{\log n}{\log \log n}$ works for large n , as desired. Therefore, the expected number of paths of length $k = \frac{\log n}{\log \log n}$ is at least

$$\mathbb{E}[a_k] \geq C \sqrt{2\pi k} \log k \geq 1.$$

□

Remark 3. In fact we can take the \log of the expression in Eq. (9), it can be shown that $k = \Omega\left(\frac{\log n}{\log \log n}\right)$ is the threshold for the expected number of paths being asymptotically greater than 1.

The above lemma tells us that we can expect to find a short path between any two vertices. However, this result is still far from establishing bounds on the diameter or even the length of the shortest path between two vertices chosen uniformly at random. We end this section with the conjecture.

Conjecture 1. *Let α and d be fixed in $G(n, \alpha, d)$. The length of the diameter of $G(n, \alpha, d)$ is $\mathcal{O}(\log n)$.*

V. AN APPLICATION TO REAL WORLD NETWORKS

We tested our model on a variety of real world networks. Our code is available at <https://github.com/martinjm97/DRGG>. Interestingly, we empirically observed that networks created through word association resulted in networks that had binomial outdegree distribution and power law indegree distribution. An example of this is the University of South Florida Word Association Network. To create this network, researchers asked participants to write the first word that came to mind that was meaningfully related or strongly associated to words that were presented to them. Then a directed edge was drawn between the word said by the participant and the word that was presented to them. This network has approximately 10^4 vertices and 7.2×10^4 edges. For more information about this network, see [28]. In this section we investigate this network and see how its properties compare to that of the DRGG model.

A. Degree Distribution

We begin by exploring how the indegree and outdegree distributions for the word association networks compare to the predictions of DRGG. The outdegree and indegree distributions of the network along with the best fits according to DRGG are shown in Figure 5. Note that for degree distributions, DRGG essentially has **one free parameter**, namely $\beta = \frac{\alpha-1}{d}$. Therefore, we fit both the outdegree and indegree distributions using β . As shown in Figure 5, DRGG is a close fit, especially considering that there was only one free parameter to tune. We discovered that the value of $\beta = 7/3 \approx 2.33$ ($\alpha = 8, d = 3$) resulted in the best fit.

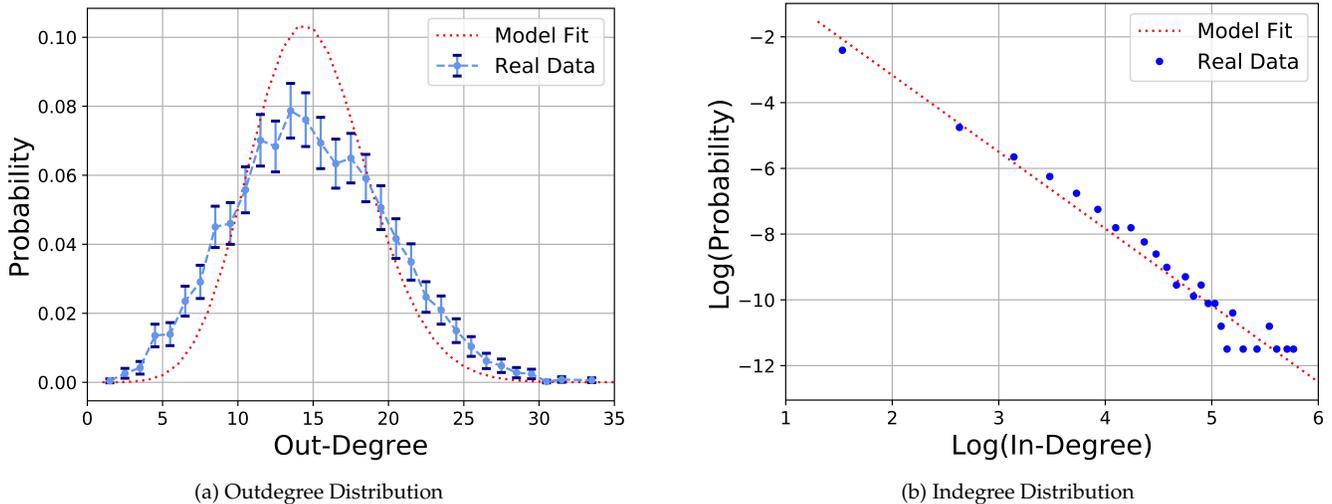


Figure 5: Fit for outdegree and indegree distributions for the word association network. β for both fits was $\beta = 7/3 \approx 2.33$. The outdegree distribution on the left also has error bars of approximately two standards of deviation (we estimate one standard deviation of a bin by the square root of the number of items that fall in the bin).

B. Other Graph Statistics

It is interesting to investigate how other statistics quantities compare to the values predicted by DRGG. In this section we use the model parameters of $\alpha = 8, d = 3$ to fit the degree distributions.

We compared the average clustering coefficient (the clustering coefficient averaged over all nodes), the diameter, and the average path length. The average clustering coefficient gives information about how tightly-knit small communities are in the graph [14]. Likewise, diameter and average path length provides insight into how closely the simulation models the real-world data. The comparison between our model predictions and actual real data values are shown in Table I. Note that the average clustering coefficient is noticeably lower in our model predictions than in the actual data. This suggests that the word association data set has a stronger clique behavior than our model. Overall, the empirical degree distributions and the empirical graph statistics do not quite match the predictions given by DRGG. This may be explained by the fact our results are asymptotic in n and the network we studied only has 10^4 vertices.

	Avg. Clustering Coefficient	Diameter	Avg. Path Length
Simulation	0.512 ± 10^{-2}	9.67 ± 0.94	5.149 ± 0.65
Real Data	0.119	7	4

Table I: The results we averaged over 100 simulations of DRGG and were computed on the undirected version of the graph i.e., the directed edges were made by bidirectional (undirected). This was done for computational simplicity. Since the DRGG model is strongly connected, this should not have a large impact on the results. The same procedure was applied to the real data set. Finally, we only used the giant component in the real data set to calculate these statistics. The mean values for the three graph statistics along with two standard deviations of error are displayed.

C. Analysis of Hubs

We analyze the words in the ‘hubs’ of the word association network [28]. By hubs, we specifically mean words that have a high indegree and represent the tail-end of the power law indegree distribution. The 20 words with the largest indegrees are shown in Table II. A lot of these hub words are **emotional** words such as Love, Good, Bad, Pain, and Happy. In addition, the hubs also include words that are ubiquitous in everyday life, such as money, water, car, work, and people.

Further analysis was performed on the largest 50 hubs to understand their significance in relation to the remainder of the graph. One way of understanding the relationships between nodes in the graph is to look at their semantic similarity, or closeness in meaning. Several metrics have been proposed to quantify semantic similarity in words. We used two, based on the WordNet database and Word2Vec model. The WordNet database contains a hand-catalogued tree-like hierarchy of words. Given a word at node n , hypernyms (words with broader meanings) are located higher in the tree relative to n , while hyponyms (words with more specific meanings) are located lower in the tree. Two words can be judged to be similar in meaning if they are close together in the graph induced by these word relationships [29]. One particular such measure of similarity is Wu-Palmer similarity:

$$\text{Similarity score}(w_1, w_2) = \frac{2 \cdot \text{depth}(\text{least common ancestor}(w_1, w_2))}{\text{depth}(w_1) + \text{depth}(w_2)}.$$

The Wu-Palmer similarity always lies between 0 and 1, and is higher if the words are closer in semantic similarity [30]. We found that the average similarity between a hub and its neighbors is 0.359, the average similarity between an arbitrary node in the graph and its neighbors is 0.354, and that the average similarity between two arbitrarily selected nodes in the graph is 0.241. The average similarity between hubs is 0.263. Thus, according to this similarity metric, hubs are on average slightly closer to their neighbors than other nodes. As expected, randomly chosen nodes have lower similarity as they may not be related, while nodes that are connected would be expected to have higher similarity. Hubs are also not very closely related: they are thus all common yet distinct words that serve as distinct “sinks” in the word graph.

To confirm this trend, we compared these results to those obtained from a different similarity metric. The Word2Vec model maps different words to continuous vector representations of a desired dimension. For more information, see the appendix. Similarity between vectors is then supposed to capture semantic similarity: indeed, relationships such as King – Man + Woman = Queen seem to approximately hold between the vectors representing these respective words. We found that the average dot product similarity between a hub and its neighbors is 0.452, the average similarity between an arbitrary node and its neighbors is 0.422, and the average similarity between randomly selected nodes is 0.338. The average similarity between different hubs is 0.415. Thus, as before, hubs are on average slightly more similar to neighbors than arbitrarily chosen nodes, and much more similar than random words are to each other. Interestingly, this metric describes hubs as fairly similar to each other- this might be an artifact of this metric, as all hubs describe fairly common words that might appear together often in many text corpora.

We also used a library named TextBlob to perform sentiment analysis to determine whether hubs expressed significantly different emotions than arbitrary nodes in the graph. For more information about TextBlob, please see the appendix.

Phrases can be constructed by concatenating words in the association network with their neighbors. TextBlob was used to analyze these phrases. Phrases that contain hub words have an average polarity of 0.0970 and subjectivity of 0.553, with variances of 0.0731 and 0.0752 respectively. Phrases made from arbitrary graph nodes have an average polarity of 0.00511 and subjectivity of 0.532, with variances of 0.0982 and 0.0611 respectively. Thus, hubs tend to be slightly more “positive” than “negative”, and slightly more “subjective” than “objective” compared to the average node in the graph, but these differences are on the order of the variances in these numbers. It thus seems that hubs are generally fairly neutral.

In summary, it seems that hubs are distinguished by the fact that they are short words that are easily memorable. This explains why they have high indegree: people have an easy time remembering them, regardless of their semantic content.

VI. CONCLUSIONS AND FUTURE WORK

We have introduced a new model of random graphs, the Directed Random Geometric Graph (DRGG) model that has the property of being scale free in its indegree distribution, has few edges, and has a high clustering coefficient.

Food: 324	Money: 302	Water: 276	Car: 259	Good: 255
Bad: 229	Work: 195	House: 185	School: 183	Love: 181
Man: 171	Paper: 163	Pain: 158	Animal: 156	People: 154
Fun: 151	Book: 149	Clothes: 147	Happy: 145	Hard: 144

Table II: Words with the largest indegrees in the word association network. The respective indegrees are shown next to the words.

Furthermore, we have displayed that this model can be applied to real world networks such as word association networks. Future work includes further theoretical investigation of the DRGG model, such as proving Conjecture 1 which states that the diameter of $G(n, \alpha, d)$ is $\mathcal{O}(\log n)$.

A potential future application of DRGGs is to model power grid networks. Power grids are comprised of different types of nodes such as generators and transformers. Edges are directed because some nodes produce energy, others transfer energy, and yet others use up energy. Therefore, relationship among the nodes is asymmetrical. This matches the asymmetric indegrees and outdegrees of our model. Furthermore, similar to our model, connections are highly correlated with distance, since nodes that are far apart will be impractical to connect. However, a challenge for this analysis is the lack of directed data that is available, although there has been some recent progress to construct such a network by pulling data from multiple sources [31]. We envision the analysis of directed power grid networks as a direction for future research with potential for broad applicability. In addition, we hope to see further applications of the DRGG model to other real world networks.

VII. ACKNOWLEDGEMENTS

Jesse Michel, Sushruth Reddy, Rikhav Shah, and Sandeep Silwal would like to thank IBM Research for the opportunity to do an internship at the MIT-IBM AI lab in Cambridge MA in the winter of 2018, which was awarded to them for winning the 2017 HackMIT competition. They would also like to thank Ramis Movassagh for mentoring them during this period.

VIII. APPENDIX

A. Proof of Theorem 2

Proof. Throughout the proof, $P[\cdot]$ denotes a probability while $p(\cdot)$ denotes a probability density. By linearity of expectation, we can write $\mathbb{E}[|S|] = n(n-1)(n-2)P(v \rightsquigarrow u, w \rightsquigarrow u, v \rightsquigarrow w)$ for randomly selected vertices u, v, w . Without loss of generality, as we are working on a torus, we may take u to be located at $(0, 0, \dots, 0)$. We can then condition on the locations of v and w , which we denote \mathbf{x}_v and \mathbf{x}_w , as well as the radii of circles centered at u and w , which we denote r_u and r_w , respectively, to obtain:

$$P[v \rightsquigarrow u, w \rightsquigarrow u, v \rightsquigarrow w] = \int_{r_0}^{\frac{1}{2}} dr_u \int_{r_0}^{\frac{1}{2}} dr_w \int_{[0,1]^d} d^d \mathbf{x}_v \int_{[0,1]^d} d^d \mathbf{x}_w P[v \rightsquigarrow u, w \rightsquigarrow u, v \rightsquigarrow w | r_u, r_w, \mathbf{x}_v, \mathbf{x}_w] p(r_u, r_w, \mathbf{x}_v, \mathbf{x}_w).$$

We now note that the radii of circles corresponding to u and w , along with the locations of v and w are independent. Additionally, edges are conditionally independent given these radii and locations. Then, we may rewrite the above probability as:

$$\int_{r_0}^{\frac{1}{2}} dr_u \int_{r_0}^{\frac{1}{2}} dr_w \int_{[0,1]^d} d^d \mathbf{x}_v \int_{[0,1]^d} d^d \mathbf{x}_w P[v \rightsquigarrow u | \mathbf{x}_v, r_u] P[w \rightsquigarrow u | \mathbf{x}_w, r_u] P[v \rightsquigarrow w | \mathbf{x}_w, \mathbf{x}_v, r_w] p(r_u) p(r_w) p(\mathbf{x}_v) p(\mathbf{x}_w).$$

Note firstly that the probability density of the points x_v and x_w is 1 since v, w are chosen uniformly from the unit torus which means $p(\mathbf{x}_v) = p(\mathbf{x}_w) = 1$. Furthermore, $P[v \rightsquigarrow u | \mathbf{x}_v, r_u] = \mathbb{1}[d(\mathbf{x}_v, 0) < r_u]$, as we draw an edge from v to u if and only if \mathbf{x}_v is inside the circle of radius r_u centered at u (the origin). Similarly, $P[w \rightsquigarrow u | \mathbf{x}_w, r_u] = \mathbb{1}[d(\mathbf{x}_w, 0) < r_u]$. Finally, we have that $P[v \rightsquigarrow w | \mathbf{x}_w, \mathbf{x}_v, r_w] = \mathbb{1}[d(\mathbf{x}_v, \mathbf{x}_w) < r_w]$, as we draw an edge from v to w if and only if the

distance between their coordinates is less than the radius of w 's circle. Substituting, we obtain:

$$P[v \rightsquigarrow u, w \rightsquigarrow u, v \rightsquigarrow w] = \int_{r_0}^{\frac{1}{2}} dr_u p(r_u) \int_{r_0}^{\frac{1}{2}} dr_w p(r_w) \int_{d(\mathbf{x}_v, 0) < r_u} d^d \mathbf{x}_v \int_{d(\mathbf{x}_w, 0) < r_u} d^d \mathbf{x}_w \mathbb{1}[d(\mathbf{x}_v, \mathbf{x}_w) < r_w].$$

We now upper and lower bound this expression in order to show that it is of order $\Theta(n \log^2(n))$. We first prove an upper bound.

1. Upper Bound Calculation

Note that $\mathbb{1}[d(\mathbf{x}_v, \mathbf{x}_w) < r_w] \leq 1$ is always true, and so our probability satisfies:

$$P[v \rightsquigarrow u, w \rightsquigarrow u, v \rightsquigarrow w] \leq \int_{r_0}^{\frac{1}{2}} dr_u p(r_u) \int_{r_0}^{\frac{1}{2}} dr_w p(r_w) \int_{d(\mathbf{x}_v, 0) < r_u} d^d \mathbf{x}_v \int_{d(\mathbf{x}_w, 0) < r_u} d^d \mathbf{x}_w 1.$$

Now, the last two integrals individually evaluate to the volume of the d -dimensional ball with radius r_u . Furthermore, the integral over r_w evaluates to 1, as $p(r_w)$ is a normalized probability density function. The remaining integral can be evaluated as follows:

$$P[v \rightsquigarrow u, w \rightsquigarrow u, v \rightsquigarrow w] \leq V_d^2 \eta \int_{r_0}^{\frac{1}{2}} r_u^{2d-\alpha} = V_d^2 \eta \left[\frac{r_u^{2d-\alpha+1}}{2d-\alpha+1} \right]_{r_0}^{\frac{1}{2}}.$$

As n becomes large, for any fixed d we have that r_0 goes to 0. This then implies (for $\alpha > 2d$) that the dominant term in the above (and a valid upper bound) is:

$$V_d^2 \eta \frac{r_0^{2d-\alpha+1}}{\alpha-2d-1} = C_1 \frac{\log^2(n)}{n^2}$$

where C_1 is some constant independent of n .

2. Lower Bound Calculation

We now show a lower bound on the integral expression. We claim that:

$$P[v \rightsquigarrow u, w \rightsquigarrow u, v \rightsquigarrow w] \geq \int_{r_0}^{\frac{1}{2}} dr_u p(r_u) \int_{r_0}^{\frac{1}{2}} dr_w p(r_w) \int_{d(\mathbf{x}_v, 0) < r_u} d^d \mathbf{x}_v \int_{d(\mathbf{x}_w, 0) < r_u} d^d \mathbf{x}_w \mathbb{1}[r_w > 2r_u]$$

This is because for choices of r_w, r_u such that $r_w > 2r_u$, we have $d(\mathbf{x}_v, \mathbf{x}_w) \leq d(\mathbf{x}_v, 0) + d(\mathbf{x}_w, 0) \leq 2r_u \leq r_w$ by the triangle inequality. Then, this integrand is always nonnegative and is identical to the original integrand when it is nonzero. This integral is then a lower bound on the original, as desired. Now, the integrals over $\mathbf{x}_v, \mathbf{x}_w$ can be done as before to give:

$$P[v \rightsquigarrow u, w \rightsquigarrow u, v \rightsquigarrow w] \geq V_d^2 \int_{r_0}^{\frac{1}{2}} dr_u r_u^{2d} p(r_u) \int_{r_0}^{\frac{1}{2}} dr_w p(r_w) \mathbb{1}[r_w > 2r_u]$$

We can eliminate the indicator function by rewriting the bounds as:

$$P[v \rightsquigarrow u, w \rightsquigarrow u, v \rightsquigarrow w] \geq V_d^2 \int_{r_0}^{\frac{1}{4}} dr_u r_u^{2d} p(r_u) \int_{2r_u}^{\frac{1}{2}} dr_w p(r_w)$$

Doing the integral over r_w , this becomes:

$$\begin{aligned} P[v \rightsquigarrow u, w \rightsquigarrow u, v \rightsquigarrow w] &\geq V_d^2 \eta^2 \int_{r_0}^{\frac{1}{4}} dr_u r_u^{2d-\alpha} \left(\frac{(2r_u)^{1-\alpha} - (1/2)^{1-\alpha}}{\alpha - 1} \right) \\ &= \frac{V_d^2 \eta^2}{\alpha - 1} \left[\frac{2^{-\alpha} r_0^{2d-2\alpha-2}}{\alpha - d - 1} \left(1 - (4r_0)^{2(\alpha+1)-2d} \right) - \frac{2^{\alpha-1} r_0^{2d-\alpha+1}}{\alpha - 2d - 1} \left(1 - (4r_0)^{\alpha-1-2d} \right) \right]. \end{aligned}$$

Note that for fixed d , as n goes to infinity, r_0 goes to 0, so that the leading order term here is the one containing the factor $r_0^{2d-2\alpha-2}$ (as $\alpha > 2d > d - 1$). Thus, we have the lower bound

$$P[v \rightsquigarrow u, w \rightsquigarrow u, v \rightsquigarrow w] \geq (C_2 - \epsilon) \frac{\log^2(n)}{n^2}$$

for any $\epsilon > 0$, for large enough n . This then proves that, for fixed α, d , that:

$$C_1 \frac{\log^2(n)}{n^2} \geq P[v \rightsquigarrow u, w \rightsquigarrow u, v \rightsquigarrow w] \geq (C_2 - \epsilon) \frac{\log^2(n)}{n^2}$$

Multiplying by $n(n-1)(n-2)$, we get:

$$C_1 n \log^2(n) \geq \mathbb{E}[|S|] \geq (C_2 - \epsilon) n \log^2(n)$$

for any ϵ , given that n is sufficiently large. This then proves that $\mathbb{E}[|S|] \in \Theta(n \log^2(n))$, as desired. \square

B. Proof of Theorem 3

1. Computing the expectation

Proof. Now, note that

$$P[v \rightsquigarrow w | v \rightsquigarrow u, w \rightsquigarrow u, d_u = k, r_u, r_w, \mathbf{x}_v, \mathbf{x}_w] = P[v \rightsquigarrow w | \mathbf{x}_v, \mathbf{x}_w, r_w] = \mathbb{1}[\|\mathbf{x}_v - \mathbf{x}_w\| < r_w]$$

(i.e. v 's lying in w 's circle is independent of all variables but for the positions of v, w and w 's radius). Furthermore,

$$p(r_u, r_w, \mathbf{x}_v, \mathbf{x}_w | v \rightsquigarrow u, w \rightsquigarrow u, d_u = k) = p(r_w) p(r_u, \mathbf{x}_v, \mathbf{x}_w | v \rightsquigarrow u, w \rightsquigarrow u, d_u = k),$$

from the conditional independence of r_w , as none of the other variables being considered involve edges pointing to w . Furthermore, this equals $p(r_w) p(r_u | v \rightsquigarrow u, w \rightsquigarrow u, d_u = k) p(\mathbf{x}_v, \mathbf{x}_w | r_u, v \rightsquigarrow u, w \rightsquigarrow u)$, by the chain rule of probability. Note that $p(r_u | v \rightsquigarrow u, w \rightsquigarrow u, d_u = k) = p(r_u | d_u = k)$ as r_u is independent of the fact that two things lie within u 's circle given the indegree of u . Furthermore, note that

$$p(\mathbf{x}_v, \mathbf{x}_w | r_u, v \rightsquigarrow u, w \rightsquigarrow u) = \frac{P[v \rightsquigarrow u, w \rightsquigarrow u, \mathbf{x}_v, \mathbf{x}_w | r_u] p(\mathbf{x}_v, \mathbf{x}_w | r_u)}{P[v \rightsquigarrow u, w \rightsquigarrow u | r_u]}$$

by Bayes' Rule. Note that $p(\mathbf{x}_v, \mathbf{x}_w | r_u) = p(\mathbf{x}_v) p(\mathbf{x}_w) = 1$, that $P[v \rightsquigarrow u, w \rightsquigarrow u | r_u] = (V_d r_u^d)^2$, and that

$$P[v \rightsquigarrow u, w \rightsquigarrow u | r_u, \mathbf{x}_v, \mathbf{x}_w] = \mathbb{1}[\|\mathbf{x}_w\| < r_u] \mathbb{1}[\|\mathbf{x}_v\| < r_u].$$

Now, we decompose the conditional probability above as:

$$\begin{aligned} &P[v \rightsquigarrow w | v \rightsquigarrow u, w \rightsquigarrow u, d_u = k] \\ &= \int dr_u \int dr_w \int d^d \mathbf{x}_v \int d^d \mathbf{x}_w P[v \rightsquigarrow w | v \rightsquigarrow u, w \rightsquigarrow u, d_u = k, r_u, r_w, \mathbf{x}_v, \mathbf{x}_w] p(r_u, r_w, \mathbf{x}_v, \mathbf{x}_w | v \rightsquigarrow u, w \rightsquigarrow u, d_u = k). \end{aligned}$$

Then, our integral becomes:

$$P[v \rightsquigarrow w | v \rightsquigarrow u, w \rightsquigarrow u, d_u = k] = \int dr_u p(r_u | d_u = k) \int dr_w p(r_w) \int_{\|\mathbf{x}_v\| < r_u} d^d \mathbf{x}_v \int_{\|\mathbf{x}_w\| < r_u} d^d \mathbf{x}_w \frac{\mathbb{1}[\|\mathbf{x}_v - \mathbf{x}_w\| < r_w]}{V_d^2 r_u^{2d}}.$$

Substituting these results back into the original expression, we get the desired result in the main body of the paper. \square

2. Computing the integral expression for \bar{C}_{in} for odd d

Proof. Note that the inner 2 integrals give the probability that 2 points randomly chosen in a sphere of radius r_u are less than r_w apart. The answer to this question can be derived from a result that can be found in [32]. Namely, for **odd** d , the probability distribution function for the distance between two points being exactly r apart in a ball of radius R is:

$$P(r) = \frac{dr^{d-1}}{R^d} \frac{d!!}{(d-1)!!} \sum_{k=0}^{\frac{d-1}{2}} \frac{(-1)^k}{2k+1} \binom{\frac{d-1}{2}}{k} \left(1 - \left(\frac{r}{2R}\right)^{2k+1}\right).$$

(The even case is much harder to work with and does not result in a nice closed form. Thus, we will only work with the odd d case). The cumulative distribution function is then (for odd d):

$$D(r) = \frac{dd!!}{(d-1)!!} \sum_{k=0}^{\frac{d-1}{2}} \frac{(-1)^k}{2k+1} \binom{\frac{d-1}{2}}{k} \left(\frac{1}{d} \left(\frac{r}{R}\right)^d - \frac{1}{2^{2k+1}(2k+d+1)} \left(\frac{r}{R}\right)^{2k+d+1}\right).$$

Then, when we substitute back into the integral, we wish to compute integrals of the form:

$$\int dr_u p(r_u) \int dr_w p(r_w) \left(\frac{r_w}{r_u}\right)^m = \frac{\eta^2 r_0^{2-2\alpha}}{(1-\alpha)^2 - m^2} \left[\left(1 - \left(\frac{1}{2r_0}\right)^{m-\alpha+1}\right) \left(1 - \left(\frac{1}{2r_0}\right)^{-m-\alpha+1}\right) \right]$$

where m is an arbitrary integer. Note that since $m \leq (d-1)/2$, and $\alpha > 2d+1$, we have that $\pm m - \alpha + 1 < 0$. Then, as $n \rightarrow \infty$ (and, thus, as $r_0 \rightarrow 0$), the terms in parentheses both go to 1, so that the dominant term is proportional to $r_0^{2-2\alpha}$. We then get:

$$\int dr_u p(r_u) \int dr_w p(r_w) \left(\frac{r_w}{r_u}\right)^m \approx \frac{\eta^2}{(1-\alpha)^2 - m^2} r_0^{2-2\alpha} \rightarrow \frac{(\alpha-1)^2}{(\alpha-1)^2 - m^2}.$$

Putting this all back together again, we get the final expression in the main body of the paper.

d	$\mathbb{E}[\bar{C}_{in}]$	$\lim_{\alpha \rightarrow \infty} \mathbb{E}[\bar{C}_{in}]$
1	$\frac{(\alpha-1)^2}{4} \left(\frac{4}{\alpha^2-2\alpha} + \frac{1}{-\alpha^2-2\alpha+1} \right)$	$\frac{3}{4}$
3	$\frac{3(\alpha-1)^2(5\alpha^4-20\alpha^3+9\alpha^2+22\alpha-72)}{32(\alpha^2-2\alpha-8)(\alpha^2-2\alpha-5)(\alpha^2-2\alpha-3)}$	$\frac{15}{32}$
5	$\frac{(\alpha-1)^2(159\alpha^6-954\alpha^5+5364\alpha^4-15096\alpha^3-73679\alpha^2+175006\alpha+392040)}{512(\alpha^2-2\alpha-24)(\alpha^2-2\alpha-9)(\alpha^2-2\alpha-7)(\alpha^2-2\alpha-5)}$	$\frac{159}{512}$

Table III: Values of $\mathbb{E}[\bar{C}_{in}]$ for various odd dimensions. The limiting value as $\alpha \rightarrow \infty$ is also shown. We only have values for odd d since we could not calculate a closed form expression for the probability of the distance between two randomly chosen points being exactly r apart in a ball of radius R .

C. Analysis of Hubs

The mapping utilized by Word2Vec is obtained by training a neural net. Specifically, the vectors are the solutions of an optimization problem which roughly attempts to maximize the dot products of vectors corresponding to words that are located close to each other in some text corpus. Ideally, words that are located close to each other often have similar meaning, and thus higher dot products, giving some indication of semantic similarity [33]. For the purposes of these experiments, we used a set of word vectors that were pre-trained on Google News articles. TextBlob uses a probabilistic model, tending to classify words as positive if they occur in many positive movie reviews, and negative if they occur in low-rated reviews. The library returns polarity and subjectivity values, which measure how negative/positive (on a $[-1, 1]$ scale) and objective/subjective (on a $[0, 1]$ scale) a given phrase is, respectively.

* Electronic address: R. Movassagh: ramis@us.ibm.com, Rest: mithack@mit.edu

- [1] Bin Zhang and Steve Horvath. A general framework for weighted gene coexpression network analysis. In *STATISTICAL APPLICATIONS IN GENETICS AND MOLECULAR BIOLOGY 4: ARTICLE 17*, 2005.
- [2] A. Capocci, V. D. P. Servedio, F. Colaiori, L. S. Buriol, D. Donato, S. Leonardi, and G. Caldarelli. Preferential attachment in the growth of social networks: The internet encyclopedia wikipedia. *Phys. Rev. E*, 74:036116, Sep 2006. doi: 10.1103/PhysRevE.74.036116. URL <https://link.aps.org/doi/10.1103/PhysRevE.74.036116>.
- [3] D. Ediger, K. Jiang, J. Riedy, D. A. Bader, and C. Corley. Massive social network analysis: Mining twitter for social good. In *2010 39th International Conference on Parallel Processing*, pages 583–593, Sept 2010. doi: 10.1109/ICPP.2010.66.
- [4] M. Kaiser and C. C. Hilgetag. Spatial growth of real-world networks. *Phys. Rev. E*, 69(3):036103, March 2004. doi: 10.1103/PhysRevE.69.036103.
- [5] S. Boccaletti, V. Latora, Y. Moreno, M. Chavez, and D.-U. Hwang. Complex networks: Structure and dynamics. *Physics Reports*, 424(4):175 – 308, 2006. ISSN 0370-1573. doi: <https://doi.org/10.1016/j.physrep.2005.10.009>. URL <http://www.sciencedirect.com/science/article/pii/S037015730500462X>.
- [6] S. Tanimoto. Power laws of the in-degree and out-degree distributions of complex networks. *ArXiv e-prints*, December 2009.
- [7] Albert-Laszlo Barabasi and Reka Albert. Emergence of scaling in random networks. *Science*, 286(5439):509–512, 1999. ISSN 0036-8075. doi: 10.1126/science.286.5439.509. URL <http://science.sciencemag.org/content/286/5439/509>.
- [8] Robert Elsasser, Leszek Gasieniec, and Thomas Sauerwald. On radio broadcasting in random geometric graphs. *International Symposium on Distributed Computing*, 5218:212–226, 01 1970.
- [9] Xingde Jia. Wireless networks and random geometric graphs. In *7th International Symposium on Parallel Architectures, Algorithms and Networks, 2004. Proceedings.*, pages 575–579, May 2004. doi: 10.1109/ISPAN.2004.1300540.
- [10] Albert-Laszlo Barabasi, Reka Albert, and Hawoong Jeong. Scale-free characteristics of random networks: The topology of the world-wide web, 2000.
- [11] Stefan Kambiz Behfar, Ekaterina Turkina, Patrick Cohendet, and Thierry Burger-Helmchen. Directed networks different link formation mechanisms causing degree distribution distinction. *Physica A: Statistical Mechanics and its Applications*, 462(C): 479–491, 2016. URL <https://EconPapers.repec.org/RePEc:eee:phsmap:v:462:y:2016:i:c:p:479-491>.
- [12] B. A. Huberman, D. M. Romero, and F. Wu. Social networks that matter: Twitter under the microscope. *ArXiv e-prints*, December 2008.
- [13] Mark Newman. *Networks: An Introduction*. Oxford University Press, Inc., New York, NY, USA, 2010. ISBN 0199206651, 9780199206650.
- [14] Duncan Watts and Steven H. Strogatz. Collective dynamics of small world networks. *Nature*, 393:440–2, 07 1998.
- [15] N. Durak, A. Pinar, T. G. Kolda, and C. Seshadhri. Degree Relations of Triangles in Real-world Networks and Models. *ArXiv e-prints*, July 2012.
- [16] B. Bollobas. Mathematical results on scale-free random graphs. In *Handbook of Graphs and Networks*, pages 1–37. Wiley, 2003.
- [17] Jesper Dall and Michael Christensen. Random geometric graphs. *Phys. Rev. E*, 66:016121, Jul 2002. doi: 10.1103/PhysRevE.66.016121. URL <https://link.aps.org/doi/10.1103/PhysRevE.66.016121>.
- [18] M. Bradonjic and W. Perkins. On Sharp Thresholds in Random Geometric Graphs. *ArXiv e-prints*, August 2013.
- [19] Paul Balister, Amites Sarkar, and Bela Bollobas. *Percolation, Connectivity, Coverage and Colouring of Random Geometric Graphs*, pages 117–142. Springer Berlin Heidelberg, Berlin, Heidelberg, 2008. ISBN 978-3-540-69395-6. doi: 10.1007/978-3-540-69395-6_2. URL https://doi.org/10.1007/978-3-540-69395-6_2.
- [20] Carl P. Dettmann and Orestis Georgiou. Random geometric graphs with general connection functions. *Phys. Rev. E*, 93: 032313, Mar 2016. doi: 10.1103/PhysRevE.93.032313. URL <https://link.aps.org/doi/10.1103/PhysRevE.93.032313>.
- [21] Eric W. Weisstein. "pareto distribution." from mathworld—a wolfram web resource. <http://mathworld.wolfram.com/ParetoDistribution.html>, 2018.
- [22] Richard Chapling. Asymptotic methods, April 2016.

- [23] A. Erdelyi. *Asymptotic Expansions*. Dover Books on Mathematics. Dover Publications, 1956. ISBN 9780486603186. URL <https://books.google.com/books?id=aedk-OHdmNYC>.
- [24] G. Fagiolo. Clustering in complex directed networks. *Phys. Rev. E*, 76(2):026107, August 2007. doi: 10.1103/PhysRevE.76.026107.
- [25] Benjamin Tabak, Marcelo Takami, J. Rocha, and Daniel Cajueiro. Directed clustering coefficient as a measure of systemic risk in complex banking networks. Working Papers Series 249, Central Bank of Brazil, Research Department, 2011. URL <https://EconPapers.repec.org/RePEc:bcb:wpaper:249>.
- [26] Mohammad Al Hasan and Vachik S. Dave. Triangle counting in large networks: a review. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, pages e1226–n/a, 2017. ISSN 1942-4795. doi: 10.1002/widm.1226. URL <http://dx.doi.org/10.1002/widm.1226>. e1226.
- [27] N. Meghanathan. A random network model with high clustering coefficient and variation in node degree. In *2015 8th International Conference on Control and Automation (CA)*, pages 54–57, Nov 2015. doi: 10.1109/CA.2015.20.
- [28] Schreiber Nelson, McEvoy. The university of south florida word association, rhyme, and word fragment norms. <http://www.usf.edu/FreeAssociation/>, 1998.
- [29] Christiane Fellbaum. *WordNet*. Springer, 2010.
- [30] Zhibiao Wu and Martha Palmer. Verb semantics and lexical selection. *CoRR*, abs/cmp-lg/9406033, 1994. URL <http://arxiv.org/abs/cmp-lg/9406033>.
- [31] Saleh Soltan, Alexander Loh, and Gil Zussman. A learning-based method for generating synthetic power grids. 2017.
- [32] S.-J. Tu and E. Fischbach. Random distance distribution for spherical objects: general theory and applications to physics. *Journal of Physics A Mathematical General*, 35:6557–6570, August 2002. doi: 10.1088/0305-4470/35/31/303.
- [33] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. Distributed representations of words and phrases and their compositionality. *CoRR*, abs/1310.4546, 2013. URL <http://arxiv.org/abs/1310.4546>.