

GRAPH DATA ANALYSIS

MARIA CAMERON

CONTENTS

1. Introduction	1
2. Basic algorithms for graph exploration	1
2.1. Breadth-first search	2
2.2. Depth-first search	5
3. Random graph models	8
3.1. Poisson random graphs	8
3.2. Generating functions	9
3.3. Poisson random graphs: mean size of non-giant component	10
3.4. Random graphs with arbitrary degree distributions	12
References	15

1. INTRODUCTION

See my slides of Nov. 12, 2020, and read: (i) chapters 1 and 2 in [M. Newman's review paper \[1\]](#), and (ii) a [3-page paper by Watts and Strogatz in Nature 1998 \[2\]](#).

2. BASIC ALGORITHMS FOR GRAPH EXPLORATION

Suppose we have a graph $G(V, E)$ where V is the set of vertices and E is the set of edges. Often, the graph is given by the list of its vertices $1, \dots, n$, and a the list of edges. The list of edges, sorted according to the tail vertices of the edges is easily recasted into the adjacency lists. Another way to define a graph is to provide its $n \times n$ *adjacency matrix* $A = (A_{ij})$ where $A_{ij} = 1$ if there is an edge from i to j and $A_{ij} = 0$ otherwise. These ways of defining a graph is suitable for both, directed and undirected graph, and easily can be extended to a weighted graph by adding weights to the list of edges.

Given a graph, the first natural question is whether it is connected or not. If, not, what are its *connected components*?

Definition 1. *A connected component of a graph $G(V, E)$ is a subset of its vertices $V_1 \subset V$ such that*

- (1) *for any pair of vertices $i, j \in V_1$ there is a path in $G(V, E)$ from i to j , and*
- (2) *if there is a path from $i \in V_1$ to $j \in V$, then $j \in V_1$, i.e., the connected component is a maximal subset.*

If the graph $G(V, E)$ is directed, and we restrict ourselves to accounting only for directed paths, then such a component is called strongly connected.

The two straightforward algorithms that compute all connected components in time $O(|V| + |E|)$ are the *breadth-first search* (BFS) and *depth-first search* (DFS) – see the so-called CLRS textbook [3] (Chapter 22).

2.1. Breadth-first search. The breadth-first search (BFS) is one of the simplest algorithms for searching a directed or undirected unweighted graph and is a building block for many important graph algorithms. It is suitable for both directed and undirected graphs. It computes shortest paths from the *source* (any designated vertex) to all other vertices reachable from it.

The BFS uses the data structure called the *queue*. The key feature of a queue is the principle that *the first element entering it is the first one who exits it*. This structure is denoted by the symbol Q in the pseudocode below (Algorithm 1). The color attribute of a vertex indicates if the vertex

- has never been in the queue (**WHITE**),
- is either in the queue now, or is just out of it and its adjacency list being explored now (**GRAY**),
- is out of the queue and adjacency list has been already explored (**BLACK**).

The other attributes of a vertex v are $v.d$, the shortest distance to the source, and $v.parent$, the predecessor of v in a shortest path from the source to it. The *predecessor subgraph* for a given graph $G(V, E)$ and a source vertex $s \in V$ is defined as $G_\pi = (V_\pi, E_\pi)$ where

$$V_\pi = \{v \in V \mid v.parent \neq \text{NIL}\} \cup \{s\}, \quad E_\pi = \{(v.parent, v) \mid v \in V_\pi \setminus \{s\}\}.$$

If not all vertices are discovered in a run of BFS, then we can start a new run with an undiscovered vertex chosen as a source.

Remark A *tree* is a special kind of undirected graph that contains no cycles. A tree with n vertices contains $n - 1$ edges. Often one of the vertices is specified as a *root*. All vertices that are not a root and have only one adjacent edge are called *leaves*. A union of trees is called a *forest*.

Let us prove the correctness of the BFS [3] (Section 22.2).

Theorem 1. *Let $G(V, E)$ be a directed or undirected graph, and suppose BFS is run on G from a given source $s \in V$.*

- (1) *BFS discovers every vertex $v \in V$ that is reachable from the source s ;*
- (2) *upon termination, for all $v \in V$, $v.d = \delta(s, v)$, the shortest path length from s to v .*
- (3) *For any vertex $v \neq s$ that is reachable from s , one of the shortest paths from s to v is a shortest path from s to $v.parent$ followed by the edge $(v.parent, v)$.*

First we will prove some useful inequalities.

Algorithm 1: Breadth-first search.

Input: An unweighted graph $G(V, E)$, and a root vertex s .

```

for each vertex  $u \in V$  do
    |  $u.color = WHITE$ ;
    |  $u.d = \infty$ ;
    |  $u.parent = NIL$ ;
end
 $s.color = GRAY$ ;
 $s.d = 0$ ;
 $Q = \emptyset$ ;
Add  $s$  to  $Q$ ;
while  $Q$  is not empty do
    | Take the first vertex  $u$  from  $Q$ ;
    | for each  $v$  adjacent to  $u$  do
    | | if  $v.color == WHITE$  then
    | | |  $v.color = GRAY$ ;
    | | |  $v.d = u.d + 1$ ;
    | | |  $v.parent = u$ ;
    | | | Add  $v$  to the end of  $Q$ ;
    | | end
    | end
    |  $u.color = BLACK$ ;
end

```

Output: The connected component containing s , distances and paths from s to all vertices in each connected component, and the *breadth-first tree* (the predecessor subgraph).

Lemma 1. *In the settings of Theorem 1:*

- (1) $\delta(s, v) \leq \delta(s, u) + 1 \quad \forall (u, v) \in E.$
- (2) $v.d \geq \delta(s, v) \quad \forall v \in V.$

Proof. Prove inequality (1). If u and v are both not reachable from s , then $\infty \leq \infty + 1$ which is true. If both u and v are reachable from s , the shortest path from s to v cannot be longer than the shortest path from s to u followed by the edge (u, v) .

Inequality (2) will be proven by induction on the number of additions to the queue. The basis of the induction holds as $s.d = 0 \geq \delta(s, s) = 0$, and $v.d = \infty \geq \delta(s, v)$ for all $v \in V$. Assume that (2) holds right after the first k additions to the queue. Suppose a white vertex v is discovered during the search from u just dequeued. By induction hypothesis,

$$u.d \geq \delta(s, u).$$

Hence, using (1), we obtain (2):

$$v.d \geq u.d + 1 \geq \delta(s, u) + 1 \geq \delta(s, v)$$

The vertex v is then enqueued and made gray. Therefore, $v.d$ no longer can change. Hence the induction hypothesis is maintained. \square

Lemma 2. *Suppose that vertices v and w are enqueued during the execution of BFS, and that v is enqueued before w . Then $v.d \leq w.d$ at the time that w is enqueued.*

We will prove that if the queue contains vertices (v_1, \dots, v_r) where v_1 is the head and v_r is the tail, then

$$(3) \quad v_r.d \leq v_1.d + 1 \quad \text{and} \quad v_i.d \leq v_{i+1}.d \quad \forall i = 1, \dots, r-1.$$

Once we prove (3), the statement of Lemma 2 will follow immediately.

Proof. The proof is conducted by induction on the number of queue operations. Initially, when the queue contains only s , (3) holds automatically. Suppose that (3) holds at some step of the execution of BFS. To prove the induction step, we want to show that then it also holds after (i) dequeuing v_1 and (ii) enqueueing a new vertex v_{r+1} . Note that v_{r+1} may happen before dequeuing v_1 .

As v_1 is dequeued, $v_1.d \leq v_2.d$ by induction assumption. Hence $v_r.d \leq v_1.d + 1 \leq v_2.d + 1$. Hence (3) is maintained.

Let (v_1, \dots, v_r) be the current queue, and we are enqueueing a new vertex v_{r+1} discovered from a dequeued vertex u . Since u used to be in the queue right before v_1 ,

$$v_1.d \geq u.d, \quad \text{and} \quad v_r.d \leq u.d + 1.$$

Hence

$$v_{r+1}.d = u.d + 1 \leq v_1.d + 1, \quad \text{and} \quad v_r.d \leq u.d + 1 = v_{r+1}.d.$$

Therefore, (3) is maintained as a result of any queue operation. \square

Now we are ready to prove theorem 1.

Proof. We will prove statement (2) from converse. Let v be the vertex with the minimum $\delta(s, v)$ (length of the shortest path from s) that receives an incorrect d value. Clearly, $v \neq s$. Moreover, v must be reachable from s , since otherwise $\delta(s, v) = \infty \geq v.d \geq \delta(s, v)$ implying that $v.d = \delta(s, v)$.

Thus, $v.d > \delta(s, v)$. Let u be the vertex immediately preceding v in a shortest path from s to v . Then

$$\delta(s, v) = \delta(s, u) + 1.$$

Since $\delta(s, u) < \delta(s, v)$ and because of our choice of v , we have $u.d = \delta(s, u)$ which implies that

$$(4) \quad v.d > \delta(s, v) = \delta(s, u) + 1 = u.d + 1.$$

Now consider the moment of dequeuing u . At this time, v is either white, or gray, or black.

- If v is white, then $v.d = u.d + 1$ which contradicts (4).

- If v black, it means that it was removed from the queue before u , hence by Lemma 2, $v.d \leq u.d$ which contradicts (4).
- If v is gray, then it was discovered from a vertex t and painted gray, and t was removed from the queue before u . Hence, by Lemma 2,

$$t.d \leq u.d \quad \text{and} \quad v.d = t.d + 1 \leq u.d + 1$$

which contradicts (4).

Therefore, in all scenarios, we come to a contradiction. Hence $v.d = \delta(s, v)$ for all $v \in V$.

All vertices reachable from s must be discovered, since otherwise we would have

$$\infty = v.d > \delta(s, v).$$

Finally, we observe that if $v.parent = u$ then $v.d = u.d + 1$. Thus, we obtain a shortest path from s to v by adding $(v.parent, v)$ to a shortest path from s to $v.parent$. □

2.2. Depth-first search. The depth-first search (DFS), like BFS is suitable for unweighted graphs, directed or undirected. DFS creates a *predecessor subgraph* for a given graph $G(V, E)$ defined a bit differently from the one for BFS: $G_\pi = (V, E_\pi)$ where

$$E_\pi = \{(v.parent, v) \mid v \in V, v.parent \neq \text{NIL}\}.$$

This subgraph is called a *depth-first forest* comprising several *depth-first trees*. In an undirected graph, each depth-first tree spans a connected component of the graph. In general, each depth-first tree contains all vertices reachable from its root. In a directed graph, we say that v is reachable from s if and only if there is a directed path from s to v .

In DFS, the vertices are colored the same way as in BFS. This technique guarantees that each vertex ends up in exactly one depth-first tree, so that the trees are disjoint.

Besides the depth-first forest, DFS also *timestamps* each vertex. Each vertex has two timestamps.

- The first timestamp $v.d$ records when it is first *discovered*; $v.color$ changes from **WHITE** to **GRAY**.
- The second timestamp $v.f$ records when DFS *finishes* examining all v 's adjacency list, i.e., discovers all vertices hanging down from v ; $v.color$ changes from **GRAY** to **BLACK**.

Note that all timestamps are between 1 and $2|V|$ as there is exactly one discovery event and one finishing event for each vertex. For every vertex u ,

$$u.d < u.f.$$

Algorithm 2–3 constitutes a pseudocode for DFS. Algorithm 3 is a *recursive function* as it calls itself. I'd like to remark that `time` can be made into a global variable. In C language, I would define it as a local variable of pointer-to-integer type. In the pseudocodes, I also define it as a local variable but in a manner suitable for Matlab implementation.

Note that the output of the DFS depends on the order in which the vertices are discovered. The cost of DFS is $O(|V| + |E|)$.

Let us discuss some properties of DFS.

Algorithm 2: DFS(G)

Input: An unweighted graph $G(V,E)$.
for *each vertex* $u \in V$ **do**
 | $u.\text{color} = \text{WHITE};$
 | $u.\text{parent} = \text{NIL};$
end
 $\text{time} = 0;$
for *each vertex* $u \in V$ **do**
 | **if** $u.\text{color} == \text{WHITE}$ **then**
 | $\text{time} = \text{DFS-VISIT}(G,u,\text{time});$
 | **end**
end

Algorithm 3: $\text{time} = \text{DFS-VISIT}(G,u,\text{time})$

Input: An unweighted graph $G(V,E)$, a selected vertex u , and the time variable.
 $\text{time} = \text{time} + 1;$
 $u.d = \text{time};$
 $u.\text{color} = \text{GRAY};$
for *each* v *adjacent to* u **do**
 | **if** $v.\text{color} == \text{WHITE}$ **then**
 | $v.\text{parent} = u;$
 | $\text{time} = \text{DFS-VISIT}(G,u,\text{time});$
 | **end**
end
 $u.\text{color} = \text{BLACK};$
 $\text{time} = \text{time} + 1;$
 $u.f = \text{time};$

Theorem 2. (Parenthesis theorem) *Let DFS be applied to any directed or undirected graph $G(V,E)$. Then for any two vertices u and v , exactly one of the following three conditions holds:*

- $[u.d, u.f] \cap [v.d, v.f] = \emptyset$, and neither u nor v is a descendant of the other in the depth-first forest,
- $[u.d, u.f]$ is entirely inside $[v.d, v.f]$, and u is a descendant of v in a depth-first tree, or
- $[v.d, v.f]$ is entirely inside $[u.d, u.f]$, and v is a descendant of u in a depth-first tree.

Furthermore, v is a descendant of u in the depth-first forest if and only if

$$u.d < v.d < v.f < u.f.$$

The proof directly follows from the structure of Algorithm 2-3.

Theorem 3. (White-path theorem) *In a depth-first forest of a directed or undirected graph $G(V, E)$, vertex v is a descendant of vertex u if and only if at time $u.d$ (when the search discovers u), there is a path from u to v consisting entirely of white vertices.*

Proof. \implies : If $v = u$ then the path consists of a single vertex u which is still white at the moment when we set $u.d$. If v is a proper descendant of u then $u.d < v.d$, so v is white at time $u.d$. Since v is any descendant of u , all vertices on the unique simple path from u to v in the depth-first forest are white at time $u.d$.

\impliedby : Suppose that there is a path of white vertices from u to v at time $u.d$, but v does not become a descendant of u in the depth-first tree. Without the loss of generality, we assume that every vertex other than v along the path becomes a descendant of u . Let w be the predecessor of v in the path, so that w is a descendant of u . Then $w.f < u.f$ by Theorem 2. Since v must be discovered after u is discovered but before w is finished, by Theorem 2 we have

$$u.d < v.d < w.f \leq u.f.$$

Hence, by Theorem 2, $[v.d, v.f]$ is entirely inside $[u.d, u.f]$ and hence v is a descendant of u . \square

Another property of DFS is that it can be used to classify edges of $G(V, E)$. We define the following types of edges in terms of the depth-first forest G_π produced by DFS on G :

- **Tree edges** are those in G_π . Edge (u, v) is a tree edge iff v is discovered by exploring edge (u, v) .
- **Back edges** are those edges (u, v) connecting a vertex u to an ancestor v in a depth-first tree. We consider self-loops, which occur in directed graphs, to be back edges.
- **Forward edges** are those nontree edges (u, v) that connect a vertex u to its descendant in a depth-first tree.
- **Cross edges** are all other edges. They can go between vertices in the same depth-first tree, as long as one vertex is not an ancestor of the other, or they can go between vertices in different depth-first trees.

Suppose DFS is exploring an edge (u, v) . Note that then u is gray. There are the following possibilities for v :

- $v.color = \text{WHITE}$, hence (u, v) is a tree edge;
- $v.color = \text{GRAY}$, hence (u, v) is a back edge;
- $v.color = \text{BLACK}$, hence (u, v) is a forward or cross edge.

Theorem 4. *In DFS of an undirected graph G , every edge of G is either a tree edge or a back edge.*

Proof. Let (u, v) be an arbitrary edge of G , and suppose without the loss of generality that $u.d < v.d$. Then, since v is in the adjacency list of u ,

$$u.d < v.d < v.f < u.f.$$

Hence v is white as it is discovered from u and hence (u, v) is a tree edge.

If $v.d < u.d$, then (u, v) is a back edge since v is still gray at the time this edge is explored. \square

Theorem 5. *A directed graph G is acyclic if and only if DFS of G yields no back edges.*

Proof. \implies : Suppose that a DFS produces a back edge (u, v) . Then vertex v is an ancestor of vertex u in the depth-first forest. Thus, G contains a path from v to u , and the back edge (u, v) completes a cycle.

\impliedby : Suppose that G contains a cycle c . We show that DFS of G yields a back edge. Let v be the first vertex to be discovered in c , and let (u, v) be the preceding edge in c . At time $v.d$, the vertices of c form a path of white vertices from v to u . By Theorem 3, vertex u becomes a descendant of v in the depth-first forest. Therefore, (u, v) is a back edge. \square

3. RANDOM GRAPH MODELS

Network models help us to understand the meaning of such statistical network properties as average shortest path length, clustering coefficient, degree distribution – how they relate to the way the network has been formed and how they relate to each other.

3.1. Poisson random graphs. The study of network models dates back to 1950s. Solomonoff and Rapoport (1951) [4] and, independently, Erdős and Renyi (1959) [5] proposed random graph models with similar properties:

- [4]: the random graph $G(n, p)$ which is an ensemble of random graphs with n vertices and in which each of $(1/2)n(n-1)$ possible edges exists with probability p ;
- [5]: the random graph $G(n, N)$ which is an ensemble of random graphs with n vertices and N edges randomly selected out of $(1/2)n(n-1)$ possible edges.

In both works, the limit $n \rightarrow \infty$ was considered, and a number of asymptotic estimates were obtained. The model $G(n, p)$ makes analysis a bit easier, so, we will focus on it. We will also let n tend to infinity in a manner where *the mean degree*

$$(5) \quad z = 2 \frac{pn(n-1)}{2n} = p(n-1)$$

remains constant. In this case, the model has a Poisson degree distribution. Indeed, the probability that a vertex has degree k is given by:

$$(6) \quad p_k = \binom{n}{k} p^k (1-p)^{n-k}.$$

Letting $n \rightarrow \infty$, fixing k , and setting $p = z/(n-1)$, we obtain the Poisson distribution:

$$(7) \quad p_k = \lim_{n \rightarrow \infty} \frac{n(n-1) \dots (n-k+1)}{k!} \frac{z^k}{(n-1)^k} \left(1 - \frac{z}{n-1}\right)^{n-1} \left(1 - \frac{z}{n-1}\right)^{1-k} = \frac{z^k e^{-z}}{k!}.$$

Both Solomonoff and Rapoport and Erdős and Renyi discovered the most interesting and important property of the Poisson random graph: it possesses a phase transition. If the mean degree $z < 1$, all connected components of the graph are small (with high probability). In contrast, if $z > 1$, there appears the so-called *giant component* containing

$O(n)$ vertices. Let u be the probability that a random vertex does not belong to the giant component. It is equal to the probability that none of its neighbors belongs to the giant component. This argument allows us to derive an equation for u :

$$(8) \quad u = \sum_{k=0}^{\infty} p_k u^k = e^{-z} \sum_{k=0}^{\infty} \frac{(uz)^k e^{-z}}{k!} = e^{z(u-1)}.$$

Observing that $S = 1 - u$ is the probability for a random vertex to belong to the giant component, we obtain an equation for S :

$$(9) \quad S = 1 - e^{-zS}.$$

The graphs of $1 - S - \exp(-zS)$ for $z = 1/2, 1$, and 2 are shown in Fig. 1.

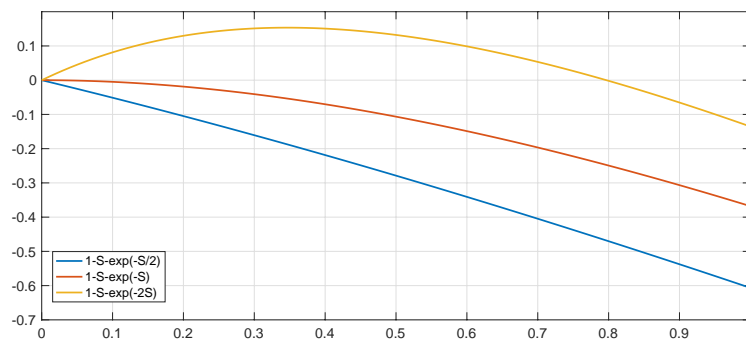


FIGURE 1. Graphs of $f(S) = 1 - S - \exp(-zS)$ where S is the fraction of vertices in the giant component for various mean degrees z . Using elementary calculus, one can show that $S = 0$ is the only root for $z \leq 1$, and there appears an additional solution $S > 0$ for $z > 1$.

In order to obtain more theoretical results for Poisson random graphs, we need to enhance our arsenal of analytical tools. A powerful tool is the method of generating functions widely used by M. Newman and collaborators. A beautiful book “[Generatingfunctionology](#)” by H. Wilf is available online.

3.2. Generating functions. This brief review follows the one in [6]. Let p_k be a discrete probability distribution. The index k runs from 0 to ∞ . If the number of outcomes is finite, then the tail of the distribution is merely 0, and all results obtained in this section will still apply. A generating function for this probability distribution is defined as a power series:

$$(10) \quad G(x) := \sum_{k=0}^{\infty} p_k x^k.$$

Since p_k is a probability distribution, we have

$$(11) \quad G(1) = \sum_{k=0}^{\infty} p_k = 1.$$

The generating function allows us to restore the distribution according to the formula:

$$(12) \quad p_k = \frac{1}{k!} \left. \frac{d^k G(x)}{dx^k} \right|_{x=0}.$$

Therefore, the function $G(x)$ encapsulates all information available about the distribution. In particular, it “generates” the distribution.

The generating function allows us to calculate all moments for the distribution. In particular, the mean is given by

$$(13) \quad \langle k \rangle = \sum_{k=0}^{\infty} k p_k = \sum_{k=1}^{\infty} k p_k 1^{k-1} = G'(1).$$

The m th moment is given by

$$(14) \quad \langle k^m \rangle = \sum_{k=0}^{\infty} k^m p_k = \left[\left(x \frac{d}{dx} \right)^m G(x) \right]_{x=1}.$$

If p_k is the distribution for a property of an object (i.e., the degree distribution for the vertices of a random graph) and $G(x)$ is its generating function, then the distribution of the sum of this property over m independent realizations is generated by $G(x)^m$. For example, let $m = 2$:

$$\begin{aligned} [G(x)]^2 &= \left[\sum_{k=0}^{\infty} p_k x^k \right]^2 = \sum_{j,k} p_j p_k x^{j+k} \\ &= p_0 p_0 x^0 + (p_0 p_1 + p_1 p_0) x + (p_0 p_2 + 2 p_1 p_1 + p_2 p_0) x^2 + \dots \end{aligned}$$

3.3. Poisson random graphs: mean size of non-giant component. We will follow the discussion in [6]. For the Poisson random graph with mean degree z , the generating function for the degree distribution is given by

$$(15) \quad G_0(x) = \sum_{k=0}^{\infty} e^{-z} \frac{z^k}{k!} x^k = e^{z(x-1)}.$$

As it should be, $G_0(1) = e^0 = 1$. Let us check that the mean degree is z . Indeed,

$$G'_0(1) = z e^{z(x-1)} \Big|_{x=1} = z.$$

Now we will calculate the average size of non-giant component, i.e., the expected size of the component where a randomly picked vertex that does not lie in the giant component belongs to. Let us pick a vertex v and an edge emanating from it. Suppose that this edge leads to another vertex w . We want to obtain the distribution for the so-called *excess degree* of w , i.e, the probability q_k that k edges other than (v, w) are incident to w . Since

the probability to arrive to w from v is proportional to the degree of w which is $k + 1$, we get:

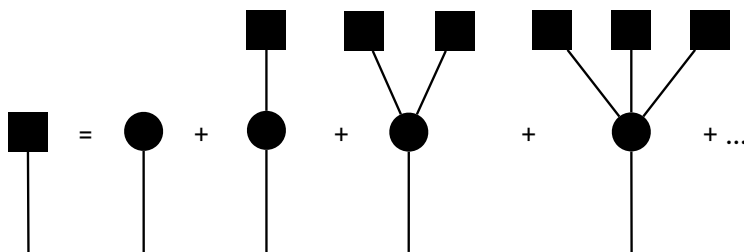
$$(16) \quad q_k = \frac{(k + 1)p_{k+1}}{\sum_{k=1}^{\infty} kp_k} = \frac{(k + 1)p_{k+1}}{z} = e^{-z} \frac{(k + 1)z^{k+1}}{z(k + 1)!} = e^{-z} \frac{z^k}{k!} = p_k.$$

Therefore, for the Poisson random graph, the distributions p_k and q_k coincide! Hence, for the Poisson model, the generating function $G_1(x)$ for the excess degree distribution coincides with $G_0(x) = e^{z(x-1)}$.

Let $H_1(x)$ be the generating function for the total number of vertices reachable by choosing a random edge, not belonging to the giant component if any, and following to one of its ends. It must satisfy the following self-consistency equation obtained under the assumption that all small components are tree-like, i.e., contain no cycles:

$$(17) \quad H_1(x) = xq_0 + xq_1H_1(x) + xq_2[H_1(x)]^2 + xq_3[H_1(x)]^3 + \dots = xG_1(H_1(x)).$$

Let us clarify why we think about each non-giant component as a tree. Remember that we assume that n is very large and the probability of an edge between any pair of vertices scales as $O(1/n)$. The number of vertices in a non-giant component is a small fraction of n . Hence the probability that there is an extra edge connecting two vertices of the same small component also scales as $O(1/n)$ and hence tends to zero as $n \rightarrow \infty$. In [6], (17) is illustrated with the following diagram:



Assume for now that there is no giant component. We randomly select a vertex v and look at the probability distribution for the size of the connected component containing it. It is generated by the function

$$(18) \quad H_0(x) = xp_0 + xp_1H_1(x) + xp_2[H_1(x)]^2 + xp_3[H_1(x)]^3 + \dots = xG_0(H_1(x)).$$

The mean size of a non-giant component is given by

$$(19) \quad \langle s \rangle = H'_0(1) = G_0(H_1(1)) + G'_0(H_1(1))H'_1(1) = 1 + G'_0(1)H'_1(1).$$

The derivative $H'(1)$ can be found from (17):

$$(20) \quad H'_1(1) = G_1(H_1(1)) + G'_1(H_1(1))H'_1(1) = 1 + zH'_1(1).$$

Therefore,

$$(21) \quad H'_1(1) = \frac{1}{1 - z}.$$

Hence, we find that

$$(22) \quad \langle s \rangle = 1 + \frac{z}{1-z} = \frac{1}{1-z}$$

Now suppose that there is a giant component occupying the fraction S of vertices. The probability u that a vertex is not in the giant component satisfies (8):

$$u = e^{z(u-1)} = G_0(u).$$

We also observe that $H_0(1) = H_1(1) = u$, the fraction of vertices not in the giant component. Therefore, we need to renormalize all probabilities for component sizes in the series $H_0(x)$:

$$h_k \mapsto \frac{h_k}{H_0(1)},$$

so that the renormalized probabilities sum up to 1 and hence are suitable for computing the expected size for non-giant components. Hence, (19) becomes

$$(23) \quad \langle s \rangle = \frac{H'_0(1)}{H_0(1)} = \frac{G_0(H_1(1)) + G'_0(H_1(1))H'_1(1)}{H_0(1)}.$$

In the presence of the giant component,

$$(24) \quad H'_1(1) = \frac{G_1(u)}{1 - G'_1(u)} = \frac{u}{1 - ze^{z(u-1)}} = \frac{u}{1 - zu}.$$

Hence, we obtain:

$$(25) \quad \langle s \rangle = \frac{1}{u} \left[G_0(u) + \frac{G'_0(u)zu}{1 - zu} \right] = 1 + \frac{zu}{1 - zu} = \frac{1}{1 - zu} = \frac{1}{1 - z + zS}.$$

3.4. Random graphs with arbitrary degree distributions. While Poisson random graphs offer an analytically solvable model in the limit $n \rightarrow \infty$, they do not resemble real-life networks in a number of aspects. One such an aspect is the degree distribution. In many real-world networks, a power degree distribution is observed – see Fig. 3.2 in [1]: the internet, the World-Wide Web, protein interactions, collaborations in mathematics, citations networks all exhibit degree distribution of the form

$$(26) \quad p_k = \frac{k^{-\alpha}}{\zeta(\alpha)}, \quad \text{where} \quad \zeta(\alpha) = \sum_{k=1}^{\infty} k^{-\alpha} \text{ is the Riemann zeta function.}$$

Note that the Riemann zeta function is finite for all α such that $\text{Re}(\alpha) > 1$. Another degree distribution that is observed in real-world networks is exponential, e.g., in the power-grid network (Fig. 3.2 in [1]).

Motivated by this fact, random graphs with a specified degree distribution were introduced – see [6] and references therein. These graphs can be sampled by generating vertices with numbers of “stubs” distributed according to the given distribution p_k and then randomly matching the stubs. The only restriction of this approach is that the total number of stubs must be even.

The method of generating functions used in the previous section is straightforwardly transferable to a model with specified degree distribution. As before, we obtain the excess degree distribution

$$q_k = \frac{(k+1)p_{k+1}}{\sum_{j=1}^{\infty} j p_j} = \frac{(k+1)p_{k+1}}{z} \quad \text{where} \quad z = \sum_{j=1}^{\infty} j p_j \quad \text{is the mean degree.}$$

The distributions p_k and q_k are generated by $G_0(x)$ and $G_1(x)$ respectively which no longer need to coincide. Note that

$$(27) \quad G_1(x) = \frac{G_0'(x)}{z}.$$

Random graphs with a given degree may or may not have a giant component, and, what is most interesting for me, a criterion for the existence of giant component can be derived [6].

3.4.1. *Phase transition.* Let us pick a vertex v and consider its first neighbors, second neighbors, ..., m th neighbors, and so on as in M. Newman, "Random graphs as models for networks" (2002). The excess degree distribution for the first neighbors is q_k . The mean excess degree is given by

$$(28) \quad \sum_{k=0}^{\infty} k q_k = \frac{1}{z} \sum_{k=0}^{\infty} k(k+1)p_{k+1} = \frac{1}{z} \sum_{k=0}^{\infty} (k-1)k p_k = \frac{\langle k^2 \rangle - \langle k \rangle}{\langle k \rangle},$$

where $z \equiv z_1 \equiv \langle k \rangle$ is the mean degree. The expected number of second neighbors is equal to the mean excess degree of the first neighbors times the expected number of first neighbors:

$$(29) \quad z_2 = \langle k \rangle \frac{\langle k^2 \rangle - \langle k \rangle}{\langle k \rangle} = \langle k^2 \rangle - \langle k \rangle$$

Then we take an arbitrary second neighbor of v and apply the same argument to calculate the expected number of third neighbors. The average excess degree for second neighbors is still given by (28). The probability that any of these excess edges reconnect to the first neighbor of v or to v itself tends to zero as $n \rightarrow \infty$. Therefore, the expected number of third neighbors is the expected number of second neighbors times the mean excess degree:

$$z_3 = z_2 \frac{\langle k^2 \rangle - \langle k \rangle}{\langle k \rangle} = \frac{z_2}{z_1} z_2 = \left[\frac{z_2}{z_1} \right]^2 z_1.$$

By a similar argument, we find the expected number of m th neighbors:

$$(30) \quad z_m = \frac{z_2}{z_1} z_{m-1} = \left[\frac{z_2}{z_1} \right]^{m-1} z_1.$$

Equation (30) allows us find a criterion for existence of the giant component. Summing (30) over all m and adding the vertex v , we find the expected size of a connected component

containing v :

$$(31) \quad \langle s \rangle = 1 + z_1 \sum_{m=1}^{\infty} \left[\frac{z_2}{z_1} \right]^{m-1} = \begin{cases} 1 + \frac{z_1^2}{z_1 - z_2}, & z_1 > z_2 \\ \infty, & \text{otherwise.} \end{cases}$$

If this expression is finite, i.e., if $z_2 < z_1$, there is no giant component. Otherwise, there is. Recalling (29) for z_2 , we write the condition for the phase transition: $z_1 = z_2$, i.e.,

$$(32) \quad \boxed{0 = z_2 - z_1 = \langle k^2 \rangle - 2\langle k \rangle = \sum_{k=0}^{\infty} k(k-2)p_k.}$$

We remark that (32) for the phase transition was originally derived by [Molloy and Reed \(1995\)](#) [7] using a different approach.

3.4.2. Average shortest path length. Equation (30) allows us to estimate the average shortest path length in the giant component of a random graph with a specified degree distribution provided that there is one. Let us assume that the giant component embraces almost all vertices in the graph. We assume that $z_2 \gg z_1$ and set the expected size of l th neighborhood to be equal to the whole graph n :

$$z_l = \left[\frac{z_2}{z_1} \right]^{l-1} z_1 \simeq n.$$

For here, we obtain the following estimate for the average shortest path length:

$$(33) \quad l \simeq 1 + \frac{\log(n/z_1)}{\log(z_2/z_1)}.$$

For the special case of Erdős-Renyi random graph we have $z_2 = z_1^2 = z^2$ (as $G_1 = G_0$), hence

$$l \simeq 1 + \frac{\log n - \log z}{\log z} = \frac{\log n}{\log z}.$$

3.4.3. Clustering coefficient. Another important quantity of interest is the *clustering coefficient*. There are different non-equivalent definitions of clustering coefficient. Here, we will define it as the ratio of the tripled number of triangles in the network to the number of connected triples:

$$(34) \quad C = \frac{3 \times \text{number of triangles in the network}}{\text{number of connected triples of vertices}}.$$

Let v be an arbitrary vertex, and let k_i and k_j be excess degrees of two of its neighbors. The distributions for the excess degrees are q_k . The probability that these two neighbors are connected to each other is

$$(35) \quad \frac{k_i k_j}{nz}.$$

Indeed, we have k_i stubs emanating from i . For each of them, the probability to connect to one out of k_j stubs of j is k_j/nz as nz is the total number of stubs. Averaging this probability

over the joint distribution for k_i and k_j which is merely the square of the distribution for the excess degree due to the independence of k_i and k_j and recalling (28), we get:

$$(36) \quad C = \frac{\langle k_i k_j \rangle}{nz} = \frac{1}{nz} \sum_{k,l=0}^{\infty} klq_k q_l = \frac{1}{nz} \left[\sum_{k=0}^{\infty} kq_k \right]^2 = \frac{1}{nz} \left[\frac{\langle k^2 \rangle - \langle k \rangle}{\langle k \rangle} \right]^2.$$

This expression can be rewritten in terms of the coefficient of variation of the degree distribution c_v defined by

$$(37) \quad c_v^2 := \frac{\langle k^2 \rangle - \langle k \rangle^2}{\langle k \rangle^2},$$

i.e., c_v is the ratio of the standard deviation to the mean. Using it, we obtain:

$$(38) \quad C = \frac{z}{n} \left[\frac{\langle k^2 \rangle - \langle k \rangle^2 + \langle k \rangle^2 - \langle k \rangle}{\langle k \rangle^2} \right]^2 = \frac{z}{n} \left[c_v^2 + \frac{z-1}{z} \right]^2.$$

Thus, we see that C scales with n as $O(z/n)$, however, the factor by which z/n is multiplied can be large.

3.4.4. *Average size of non-giant components.* In a manner similar to the one used for the Poisson random graphs, one can calculate the average size of a non-giant component [6]:

$$(39) \quad \langle s \rangle = \frac{H'_0(1)}{H_0(1)} = 1 + \frac{zu^2}{[1-S][1-G'_1(u)]},$$

where $u \equiv H_1(1)$ is the smallest non-negative solution to $u = G_1(u)$, and S is the fraction of graph occupied by the giant component: $S = 1 - G_0(u)$.

Exercise Give a detailed derivation of (39).

REFERENCES

- [1] M. E. J. Newman, "The structure and function of complex networks," *SIAM Review*, vol. 45, no. 2, pp. 167–256, 2003.
- [2] D. J. Watts and S. H. Strogatz, "Collective dynamics of "small-world" networks," *Nature*, vol. 393, pp. 440–442, 1998.
- [3] T. H. Cormen, E. Leiserson, Charles, R. L. Rivest, and C. Stein, *Introduction to algorithms, Third edition*. The MIT press, 2009.
- [4] R. Solomonoff and A. Rapoport, "Connectivity of random nets," *Bulletin of Mathematical Biophysics*, vol. 13, pp. 107–117, 1951.
- [5] P. Erdős and A. Renyi, "On random graphs," *Publ. Math. Debrecen*, vol. 6, pp. 290–297, 1959.
- [6] M. E. J. Newman, H. Strogatz, Steven, and D. J. Watts, "Random graphs with arbitrary degree distributions and their applications," *Physical Review E*, vol. 64, p. 026118, 2001.
- [7] M. Molloy and B. Reed, "A critical point for random graphs with a given degree sequence," *Random Structures and Algorithms*, vol. 6, pp. 161–180, 1995.