

# Random Graphs

Kent Quanrud

April 28, 2021

## 1 Introduction

Paul Erdős, inspired by Ramsey [5] before him, had a series of work analyzing *random graphs*, producing a large body of results that can mostly be grouped into two broad categories. First, he designed elaborate randomized constructions of graphs and showed that with nonzero probability, they can possess certain counterintuitive, seemingly impossible properties. This general approach is now called *Ramsey theory*. Second, he showed that for natural random graph models, these graphs, however random, tend to be extremely consistent about certain properties. Today we will study the  $G(n, p)$  random graph, sometimes called Erdős-Rényi graphs based on work by Erdős and Rényi [2, 3]. A random graph from  $G(n, p)$  is an undirected graph over  $n$  vertices, where every edge is sampled independently with probability  $p$ . By now there is a large catalog of nontrivial and useful properties that, depending on  $p$ , are almost certain to appear or not appear in such a graph (for sufficiently large  $n$ ). Most interestingly, Erdős-Rényi showed that these properties can vary dramatically with very small changes in  $p$ . Consider the following theorem.

**Theorem 1.1.** *Consider a random graph  $G \sim G(n, p)$  for  $p = c/n$ , where  $c$  is a constant.*

1. *If  $c > 1$ , then with high probability, there is exactly one connected component of  $G$  with  $\Omega(n)$  vertices, and all other components have size  $\leq O(\log n)$ .*
2. *For  $c < 1$ , then with high probability, all connected components of  $G$  has size  $< O(\log n)$ .*

The parameter  $c$  above models the average degree (in expectation). The drama lies in the fact that a tiny change in the average degree  $c$  – from .999 to 1.0001 – flips the qualitative nature of a typical random graph from one of many tiny components to essentially one giant component. This is an example of a *threshold phenomena*; alternatively, a *nonlinear dynamic*. Such phenomena is not rare: it occurs in many situations in physics, as well as in models for epidemiology and social networks. Let us briefly mention - without claiming to be very precise - that the sensitivity to  $c$  gives some motivation for controlling the “reproductive number” when analyzing and preventing the spread of infectious diseases. The reproductive number is the expected number of healthy individuals that a sick individual effects.

We note that we have a much more refined and detailed understanding then stated in Theorem 1.1. We refer the reader to [1, Chapter 7] for further details and other results in this area.

### 1.1 Overview of the proof

We will prove part 1 of Theorem 1.1 in roughly three parts.

**Part 1: the gap theorem.** Note that in theorem above, regardless of the value of  $p$ , there are simply no “medium”-size components, like a component of size  $\sqrt{n}$  of  $n/\log(n)$ . The intermediate sizes are ruled out by the following “gap theorem”. We analyze the following theorem in Section 2.

**Theorem 1.2.** *There is a universal constant  $C > 0$ , such that for  $n$  sufficiently large,  $\epsilon = \epsilon(n) \leq 1$ , and  $p = (1 + \epsilon)/n$ , we have the following. For a random graph  $G(n, p)$ , with probability of error  $\leq 1/n^2$ , no component has  $k$  vertices for any value  $k$  in the interval*

$$\frac{C \log(n)}{\epsilon^2} \leq k \leq \frac{\epsilon n}{C}.$$

**Part 2: Existence of a large component.** We prove the following theorem in Section 3

**Theorem 1.3.** *Let  $p = (1 + \epsilon)/n$  for  $\epsilon > 0$ . For any vertex  $v$ , Then for all  $3 \leq h \leq \epsilon n$ , with probability at least  $1/h$ ,  $v$  has at least  $1 + h$  vertices in its component.*

Let  $h = c \log(n)/\epsilon^2$  for a sufficiently large constant  $c$ , and let  $q = 1/h = \Omega(\epsilon^2/\log(n))$ . Call a component “small” if it has at most  $h$  vertices. We want to argue that, for  $p > 1/n$ , there is at least one component that is not small.

Any vertex  $v$  has at least a probability  $q$  of not being in a small component. Now imagine a process where we first randomly select a vertex  $v$  and inspect its component. If it is small, then we throw out  $v$  and its component, and randomly select another vertex as  $v$ , and repeat. Each vertex we inspect has a  $q$  chance of not being in a small component. We would have to fail on the order of  $n/h$  consecutive samples to conclude there is no small component - which happens with diminishingly small probability

$$(1 - q)^{\Omega(n/h)} \leq e^{-\Omega(\epsilon^4 n / \log^2 n)}.$$

Thus with very high probability, there is at least one component that is not small. In conjunction with the gap theorem, which rules out all intermediate sizes, there exists at least one giant component of size  $\Omega(\epsilon^2 n)$ .

**Part 3: Uniqueness of the giant component.** Can there be two giant components? The answer is no (with high probability) and here is a quick explanation. Instead of sampling from  $G(n, p)$  directly, we can first sample two graphs  $G_1 = (V_1, E_1)$  and  $G_2 = (V_2, E_2)$  from  $G(n/2, p)$ . In the second stage we can sample each cross-edge  $(v_1, v_2)$ , where  $v_1 \in V_1$  and  $v_2 \in V_2$ , independently with probability  $p$ . Now, by applying the theory we have already developed to  $G_1$  and  $G_2$ ,  $G_1$  and  $G_2$  will have some giant components, each of size  $\Omega(\epsilon^2 n)$ . Note that each graph can only have  $O(1/\epsilon^2)$  of them. Let  $C_1$  be a giant component in  $G_1$  and let  $C_2$  be a giant component in  $G_2$ . We sample  $|C_1||C_2| \geq \Omega(\epsilon^4 n^2)$  potential edges between  $C_1$  and  $C_2$ . Recalling that  $p$  is greater than  $1/n$ , the odds of all  $n^2$  edges failing to be sampled is vanishingly small - that is, we almost certainly connect  $C_1$  and  $C_2$ . Since there is a limited number of giant components we will almost certainly connect all of them together. Thus, for  $p > (1 + \epsilon)/n$  for  $\epsilon > 0$ , we get a *unique* giant component. This establishes Theorem 1.1 for  $c > 1$ .

**$c < 1$ .** The proof for  $c < 1$  is simpler and only requires the ideas underlying Theorem 1.2. See Section 2.1.

## 1.2 Directed graphs

One could naturally ask the same questions for directed graphs. Let  $D(n, p)$  denote the distribution over *directed graphs* where every directed edge appears independently with probability  $p$ . We might similarly ask for the maximum number of vertices reachable from any component, or the size of the maximum strongly component.

It turns out that the analysis of directed graphs can be largely reduced to undirected graphs, as shown by Karp [4] in the following delightfully simple way.

**Theorem 1.4.** *Let  $G \sim G(n, p)$  and  $D \sim G(n, p)$ , and fix a vertex  $v$ . Then the size of the connected component of  $v$  in  $G$ , and the number of vertices reachable from  $v$  in  $D$ , are identically distributed.*

*Proof.* Let us introduce a second distribution of directed random graphs. Let  $B(n, p)$  be the distribution of directed graphs where we sample each *undirected* edge  $\{u, v\}$  independently with probability  $p$ , and for each sampled edge, add both directions  $(u, v)$  and  $(v, u)$  to the graph. Clearly for a fixed vertex  $v$ , the size of the  $v$ 's (undirected) component in  $G(n, p)$  is distributed identically to the number of vertices reachable from  $v$  in  $B(n, p)$ . We claim that the number of vertices reachable from  $v$  in  $B(n, p)$  is identically distributed as in  $D(n, p)$ . At this point let us simply quote Karp [4, Lemma 1] (with minor changes in notation) whose proof is very elegant.

...To see that the last two random variables are identically distributed, note that the probability spaces  $B(n, p)$  and  $D(n, p)$  differ in only one respect: a digraph  $G$  drawn from  $B(n, p)$ , arc  $(u, v)$  is present if and only if arc  $(v, u)$  is present, while, in a digraph  $D$  drawn from  $D(n, p)$ , the event that  $(v, u)$  is present is independent of the event that  $(u, v)$  is present. Thus no experiment based on checking for the presence or absence of arcs can distinguish between the two probability spaces unless it checks both an arc and its reversal. But any standard sequential algorithm, such as breadth-first search or depth-first search, for building a search tree containing exactly the vertices reachable from vertex 1, checks for the presence of arc  $(u, v)$  only if vertex  $u$  is in the search tree and  $v$  is not; thus it never checks both an arc and its reversal, and accordingly cannot distinguish  $B(n, p)$  from  $D(n, p)$ .

To summarize the excerpt, standard search algorithms for reachability do not distinguish  $B(n, p)$  and  $D(n, p)$  anyway, so the number of reachable vertices is identically distributed. ■

## 2 A gap in component size

In this section we prove Theorem 1.2, which asserts that when  $p = (1 + \epsilon)/n$  for a constant  $\epsilon > 0$ , then with high probability, all components are either very small or very large. Our analysis follows an approach due to Karp [4]. His proof is also described in [1]. We will also reuse some of the ideas in the proof to analyze the  $p < 1/n$  in Section 2.1.

**Theorem 1.2.** *There is a universal constant  $C > 0$ , such that for  $n$  sufficiently large,  $\epsilon = \epsilon(n) \leq 1$ , and  $p = (1 + \epsilon)/n$ , we have the following. For a random graph, with probability of error  $\leq 1/n^2$ , no component has  $k$  vertices for any value  $k$  in the interval*

$$\frac{C \log(n)}{\epsilon^2} \leq k \leq \frac{\epsilon n}{C}.$$

For a vertex  $v \in V$ , let  $C(v) \subset V$  be the (randomized) component of  $v$ .

1.  $A_0 = \emptyset$ , and  $B_0 = \{v\}$ .
2. In the first iteration, set  $v_1 = v$ , set  $A_1 = \{v_1\}$ , and set  $B_1 = B_0 \cup N(v_1)$ , where  $N(v_1)$  is the (randomized) neighborhood of  $v_1$ .
3. In the  $i$ th iteration, if  $A_{i-1} \neq B_{i-1}$ , then select (any)  $v_i \in B_{i-1} \setminus A_{i-1}$ . Set  $A_i = A_{i-1} \cup \{v_i\}$  and  $B_i = B_{i-1} \cup N(v_i)$ . Otherwise we terminate with  $C(v) = A_{i-1} = B_{i-1}$ .

Above, each  $B_i$  are the set of vertices we know to be reachable from  $v$  after the  $i$ th iteration.  $A_i$  is the set of vertices whose edges have been inspected in the first  $i$  iterations.

The process terminates when  $A_k = B_k$ . But since  $A_k \subseteq B_k$  and  $|A_k| = k$ , this is precisely when  $|B_k| = k$ . As long as  $|B_i| \neq i$ ,  $B_{i+1}$  is generated by taking the union of  $B_i$  and a random sample of  $V - B_i$  where each vertex is included with probability  $p$ .

1. Initially set  $B_0 = \{v\}$ .
2. For each  $i \in \mathbb{N}$ , let  $S$  sample each vertex in  $V \setminus B_{i-1}$  independently with probability  $p$  and set  $B_i = B_{i-1} \cup S$ .
3. Let  $i$  be the first index such that  $|B_i| = i$ , and return  $C(v) = B_i$ .

Fix  $i$ . The alternative (but equivalent) process described above exposes a simple distribution for  $B_i$ . For any vertex  $x \neq v$ , we have  $x \notin B_i$  iff  $x$  failed to be added in each of the first  $i$  rounds, which occurs with probability exactly  $(1 - p)^i$ . Moreover this event is independent across vertices. Thus  $|B_i|$  is distributed exactly as the binomial distribution with  $n - 1$  coins and probability  $1 - (1 - p)^i$ .

**Lemma 2.1.** *Let  $i \leq \epsilon n/2((1 + \epsilon))$ . Then*

$$\mathbf{E}[|B_i|] \geq (1 + \epsilon)i.$$

*Proof.* We have

$$(1 - p)^i \leq e^{-ip} \leq 1 - ip + (ip)^2 \leq 1 - ip + \epsilon ip/2 = 1 - (1 - \epsilon/4)ip.$$

where (a) is because  $ip = (1 + \epsilon)i/n \leq \epsilon/4$ . Thus

$$\mathbf{E}[|B_i|] = 1 + \left(1 - (1 - p)^i\right)(n - 1) \geq (1 - \epsilon/4)ipn \geq (1 + \epsilon)i.$$

■

**Lemma 2.2.**

$$\mathbf{P}[|B_i| \leq i] \leq e^{-\epsilon^2 i/8}.$$

## 2.1 Probabilities < 1

Suppose instead that  $p = (1 - \epsilon)/n$ . Then we have  $\mathbf{E}[|B_i|] \leq (1 - \epsilon/2)i$  unless  $i$  is very close to  $n$ . In particular, for  $i = O(\log(n)/\epsilon^2)$ , the probability of  $|B_i| > i$  is  $1/\text{poly}(n)$ . Thus we see that all components will have size  $O(\log(n)/\epsilon^2)$  with high probabilities.

## 3 Galton-Watson branching processes

In the simplest case, imagine a population of size 1. Each generation, each member of the current generation flips  $k$  coins, each heads with probability  $1/k$ . For each heads, we generate another member of the next generation. The probabilities and number of coins are configured so that each member expects to have one child.

What is the probability that the population survives for  $h$  iterations, for a given parameter  $h$ ? This is answered by the following.



Figure 1: A complete binary tree of height 3, where each edge was deleted with probability  $1/2$ .

**Theorem 3.1.** *Let  $T$  be a complete  $k$ -ary tree of height  $h$ , and suppose every edge is deleted independently with probability  $(1 - 1/k)$ . Then the probability that there is a leaf connected to the root is  $\geq 1/h$  for  $h \geq 3$ , and  $\geq (1 - e^{-1})^h$  for  $h \leq 2$ .*

An example of the case  $k = 2$  is drawn in Figure 1.

*Proof.* For  $i \in \mathbb{N}$ , let  $p_i$  be the probability that a particular node at height  $i$  is connected to a subleaf. We have  $p_0 = 1$ . For a node at height  $i + 1$ , the probability that there is no path to a leaf via a particular child is

$$1 - \frac{1}{k} + \frac{1}{k}(1 - p_i) = 1 - \frac{p_i}{k}.$$

By independence, we have

$$p_{i+1} = 1 - \left(1 - \frac{p_i}{k}\right)^k.$$

Observe that the RHS is increasing in  $p_i$ ; thus to lower bound  $p_{i+1}$ , we can substitute any lower bound for  $p_i$ . We have

$$\begin{aligned} p_0 &= 1, \\ p_1 &= 1 - (1 - 1/k)^k \geq 1 - e^{-1} \geq .63, \\ p_2 &= 1 - (1 - .63/k)^k \geq 1 - e^{-.63} \geq .467, \\ p_3 &\geq 1 - (1 - .467/k)^k \geq 1 - e^{-.467} \geq .373 \geq 1/3. \end{aligned}$$

We claim by induction on  $i$  that  $p_i \geq 1/i$  for all  $i \geq 3$ . The base case  $i = 3$  was just proven. For the general case,

$$p_{i+1} \stackrel{(a)}{\geq} 1 - (1 - 1/ik)^k \geq 1 - e^{-1/i} \stackrel{(b)}{\geq} \frac{1}{i} - \frac{1}{2i^2} \geq \frac{1}{i+1}$$

Here (a) is by induction. (b) applies the inequality  $e^x \leq 1 + x + \frac{1}{2}x^2$  for  $x \leq 0$ . ■

### 3.1 Likelihood of small components

We can use the above branching process to analyze the probability that a given vertex  $v$  is in a component of size  $\leq h$ , for any  $h \leq \epsilon n / (1 + \epsilon)$ . Recall the sets  $B_0, B_1, B_2, \dots$  from Section 2. Given that  $|B_i| \leq h$ , we can think of  $B_{i+1} - B_i$  as adding (at least)  $n / (1 + \epsilon)$  children each with probability  $p = (1 + \epsilon) / n$ . Either we find new elements for all  $h$  rounds - which forces  $|B_h| \geq h$  - or we hit  $|B_i| = h$  at some point  $i < h$ . Thus the odds of  $v$  acquiring  $h$  vertices in its connected component is at least the odds produced by Theorem 3.1 for this value of  $h$ ; namely,  $1/h$ . This gives us Theorem 1.3.

### 3.2 A more general analysis

One can consider a more general model as follows. Let  $X \in \mathbb{Z}_{\geq 0}$  be a random variable taking nonnegative integer values. For each node at a generation  $i$ , we sample an independent copy of  $X$ , and generate this many children in the next generation. Here we have the following remarkably precise theorem.

**Theorem 3.2.** *Let  $n_i$  denote the number of children in the  $i$ th generation in the process described above.*

1. *If  $\mathbf{E}[X] < 1$ , then  $\lim_{i \rightarrow \infty} \mathbf{P}[n_i = 0] = 1$ .*
2. *If  $\mathbf{E}[X] > 1$ , then  $\lim_{i \rightarrow \infty} \mathbf{P}[n_i = 0] = q$ , where  $q$  is the unique solution to  $x = \sum_{i \geq 0} \mathbf{P}[X = 0]x^i$ .*

See the lecture notes by Sinclair [6] for a proof.

### References

- [1] Béla Bollobás. *Modern Graph Theory*. 1st ed. Graduate Texts in Mathematics 184. Springer-Verlag New York, 1998.
- [2] Paul Erdős and Alfred Rényi. “On Random Graphs I”. In: *Publicationes Mathematicae Debrecen* 6 (1959), p. 290.
- [3] Paul Erdős and Alfred Rényi. “On the evolution of random graphs”. In: *Publ. Math. Inst. Hungary. Acad. Sci.* 5 (1960), pp. 17–61.
- [4] Richard M. Karp. “The Transitive Closure of a Random Digraph”. In: *Random Struct. Algorithms* 1.1 (1990), pp. 73–94. DOI: 10.1002/rsa.3240010106. URL: <https://doi.org/10.1002/rsa.3240010106>.
- [5] Frank Plumpton Ramsey. “On a Problem of Formal Logic”. In: *Proceedings of the London Mathematical Society* s2-30.1 (Jan. 1930), pp. 264–286.
- [6] Alistair Sinclair. “Lecture 16”. Available at <https://people.eecs.berkeley.edu/~sinclair/cs271/n16.pdf>. 2020.