

RANDOM GRAPHS PROJECT

SHANE LUBOLD

Description: In this project, we will analyze a common random graph model, the Erdős–Rényi model. We will use basic probability theory and simulations to understand the properties and behavior of this random graph model, with an emphasis on these properties as the size of the graph grows. To begin, we will first review basic graph theory and probability theory. Then, we will examine the degree distribution of nodes in this model, the expected number of cliques, and study an important property in random graph theory: the “threshold property”. In addition, we will examine the probabilistic method, an interesting proof technique that is frequently used to prove statements about random graph models.

Contact: Shane Lubold, sl223@uw.edu, Padelford A-318.

Overview of Project

- (1) Week 1: Review graph theory fundamentals and basic probability.
- (2) Week 2: Introduce Erdős–Rényi graph model and basic simulations
- (3) Week 3: Degree distribution of model
- (4) Week 4: Number of triangles/cliques in ER.

1. GRAPH THEORY AND PROBABILITY REVIEW (WEEK 1)

We first review basic graph theory concepts.

Definition 1.1 (Graph). A graph $G = (V, E)$ is a set of vertices (or nodes) $V = \{1, \dots, n\}$ and a set of edges $E \subseteq V \times V$.

A simple way to represent the edges in a graph is with an adjacency matrix.

Definition 1.2 (Adjacency matrix). For a graph $G = (V, E)$ on n nodes, the adjacency matrix A of G is an $n \times n$ matrix such that $A_{ij} = 1$ if there is an edge between nodes i and j and zero otherwise.

Exercise 1.1. Let $V = \{1, 2, 3\}$ be the vertex set of a graph G . How many graphs are there on this vertex set? Draw all possible graphs on this vertex set. Determine the total number of graphs on $V = \{1, 2, 3, 4\}$ and try to extend this formula to the general case when $V = \{1, 2, \dots, n-1, n\}$ (This is more challenging. To do this, think about how many $n \times n$ symmetric matrices there are with zero on the diagonal and only 0 or 1 on the off-diagonal. Then note that for each adjacency matrix, there is a corresponding graph).

Definition 1.3 (Path). A path between nodes i and j is a sequence of distinct edges (e_1, \dots, e_m) connecting i and j such that all nodes along the sequence of edges are distinct.

Informally, a path is a sequence of edges that allows one to “walk” from node i to node j along edges in G .

Definition 1.4 (Connected Nodes and connected Graph). Two nodes i and j are connected if there exists a path between i and j . A graph is called connected if all $\binom{n}{2}$ pairs of nodes are connected.

Exercise 1.2. Draw an example of a graph on 5 nodes $V = \{1, 2, 3, 4, 5\}$ such that three of the nodes are connected and the other two are not connected.

Definition 1.5 (Isolated Node). A node is called isolated if it does not have an edge with any other node in the graph.

Definition 1.6 (Clique). A clique $C \subseteq V$ of a graph $G = (V, E)$ is a set of nodes such that for every pair (i, j) in C , i and j are connected. That is, the subgraph induced by C is complete.

Exercise 1.3. Draw an example of a graph on 5 nodes $V = \{1, 2, 3, 4, 5\}$ with a 3 clique and a graph on 5 nodes with a 5 clique. (Note that a graph on n nodes containing an n clique is called a complete graph on n nodes, denoted by K_n).

We now introduce the degree of a node.

Definition 1.7. Let $i \in V$ be a node of a graph $G = (V, E)$ with adjacency matrix A . Then, the degree of node i , denoted by $d(i)$, is the number of edges connecting i to other nodes in the graph; that is,

$$d(i) = \sum_{j \in V} A_{i,j}.$$

We will be using R throughout this project, so let's get familiar with using R to make graphs. A good reference for this is Section 2 of <https://kateto.net/netscix2016.html>. To get started, run the following commands in R.

```
library(igraph)
g1 <- graph( edges=c(1,2, 2,3, 3, 1), n=3, directed=F )
plot(g1)
```

The edges specified are (1,2), (2, 3), (3,1) on $n = 3$ nodes. The notation “directed = F” means that we are creating an undirected graph. By changing these values, we can create graphs with different edge sets. For example,

```
library(igraph)
g1 <- graph( edges=c(1,2, 2,3, 3, 1, 1, 5, 2, 4, 4, 5), n=5, directed=F)
plot(g1)
```

We now review some basic probability theory that will be used throughout this project.

Definition 1.8 (Expected value). *The expected value of random variable $X \in \{0, 1, 2, 3, \dots\}$ with probability mass function p is defined as*

$$E(X) = \sum_{x=0}^{\infty} x \cdot p(x)$$

For a sequence $\{X_n\}$ of random variables, it always holds that $E(\sum_{n=1}^N X_n) = \sum_{n=1}^N E(X_n)$, even if the $\{X_n\}$ are dependent.

Proposition 1.1. *Let X be a random variable and let I be an indicator random variable that is 1 when $X \in A$ and 0 when $X \notin A$. Then,*

$$E(I) = P(X \in A) .$$

Exercise 1.4. *Prove this result.*

As an example, let $X \sim N(0,1)$ and let $A = [0, \infty)$. Then, if I is an indicator random variable that is 1 when $X \in A$ and 0 otherwise, then $P(X \in A) = 1/2$ (check this!), so that $E(I) = 1/2$ too.

Definition 1.9. *Let n be an integer and $p \in [0,1]$. A discrete random variable X has a binomial distribution $\text{Binomial}(n,p)$ when the probability mass function of X is*

$$\mathbb{P}(X = k) = \binom{n}{k} p^k (1-p)^{n-k}, \quad k \in \{0, 1, \dots, n-1, n\} .$$

Exercise 1.5. *Show that if X has a $\text{Binomial}(n,p)$ distribution then $E(X) = np$.*

Definition 1.10. *Let $\lambda > 0$. A discrete random variable X has a $\text{Poisson}(\lambda)$ distribution when the probability mass function of X is*

$$\mathbb{P}(X = k) = \frac{\lambda^k \exp(-\lambda)}{k!} .$$

Exercise 1.6. *Prove that if $X \sim \text{Poisson}(\lambda)$, then $E(X) = \lambda$.*

Next week, we will show that the Poisson and Binomial distributions are closely related as $n \rightarrow \infty$.

Definition 1.11. *Let A and B be events. A and B are said to be independent when*

$$P(A \text{ and } B) = P(A)P(B)$$

2. ERDŐS–RÉNYI MODEL (WEEK 2)

We now introduce the Erdős–Rényi (ER) model. To do this, we fix an integer n and a value $p \in [0, 1]$, which controls the probability nodes connect to each other. To simulate an ER model, we do the following: For each pair of nodes $(i, j) \in V^2$ with $i < j$, add an edge between nodes i and j independently with probability p . There are a total of $\binom{n}{2} = \frac{n(n-1)}{2}$ possible edges in the graph.

In the previous section, we described the degree of a node i for a deterministic graph. Now, when discussing a random graph, the degree of a node for each simulation is random, and so we therefore want to talk about the *degree distribution* of a node.

Exercise 2.1. *Simulate 1000 ER graphs with $n = 100$, $p \in \{1/10, 1/2, 9/10\}$. For each graph and each value of p , compute the degree of the 1st node and plot a histogram of these values. Based on your previous probability courses, can you guess what the distribution of the degree of a node is?*

We have the following result:

Proposition 2.1. *The degree distribution of any vertex i is Binomial($n - 1, p$); that is,*

$$\mathbb{P}(d(i) = k) = \binom{n-1}{k} p^k (1-p)^{n-1-k}.$$

Exercise 2.2. *Prove this result.*

From this result, we know a lot about the expected degree of a node and its variance. In particular, we know that $E(d(i)) = (n-1)p$ (which you proved last week) and $\text{Var}(d_i) = p(n-1)(1-p)$. We can also compute the probability that a node is *isolated*; that is, that it connects to no other node. By setting $k = 0$, it is given by $\mathbb{P}(i \text{ is isolated}) = (1-p)^{n-1}$. Similarly, by setting $k = n-1$, we see that the probability that node i is connected to all other nodes in the graph is given by p^{n-1} .

Exercise 2.3. *Simulate 1000 ER graphs with $n = 10$, $p \in \{1/10, 1/2, 9/10\}$. For each graph and each value of p , compute the probability that the first node is isolated and that it is connected to all other nodes in the graph, and compare to the values given in the paragraph above.*

As is the case for most of the results in this project, we are interested in the behavior of the ER model as the graph gets larger (i.e., as $n \rightarrow \infty$). To explore this, we assume that $p = p(n)$ is a function of n and we simulate $G \sim G(n, p(n))$ as $n \rightarrow \infty$. We have the following result about the degree distribution of such a graph:

Proposition 2.2. *Assume that $np(n) \rightarrow c$ as $n \rightarrow \infty$. Then,*

$$\mathbb{P}(d(i) = k) \rightarrow \frac{c^k e^{-c}}{k!}.$$

That is, the degree distribution of any node approaches a Poisson distribution with parameter c .

Exercise 2.4. *Prove this result. (Hint: You will need to use the fact that $(1 - \frac{x}{n})^n \rightarrow e^{-x}$ as $n \rightarrow \infty$.)*

Exercise 2.5. Assume that we take $p(n) = \frac{c}{n}$ for some real number c . Prove that $\mathbb{P}(\text{node 1 is isolated}) \rightarrow e^{-c}$ as $n \rightarrow \infty$. Now, set $c = 1$ and simulate 1000 graphs with $n \in \{100, 1000, 10000\}$. For each value of n , count the number of times node 1 is isolated and compare this to e^{-1} . Does it seem to be getting closer to e^{-1} ?

Exercise 2.6. Pick a sequence $p_1(n)$ such that $np_1(n) \rightarrow 1$ and a sequence $p_2(n)$ such that $np_2(n) \rightarrow 2$. For each sequence, simulate 1000 ER graphs with $n \in \{100, 1000, 10000\}$. For each value of n , plot a histogram of the degree of node 1 and compare it to distribution from the proposition above.

Finally, we will prove two basic probability theory tool that will be used often in the following sections.

Proposition 2.3 (Markov's Inequality). If X is a non-negative random variable and $a > 0$, then $P(X \geq a) \leq \frac{E(X)}{a}$.

Exercise 2.7. In this exercise we will prove Markov's inequality for discrete random variables (i.e., random variables taking values in $\{0, 1, \dots\}$.) To do this, first express $E(X)$ as

$$E(X) = \sum_{x=0}^{\infty} xp(x) = \sum_{x=0}^{a-1} xp(x) + \sum_{x=a}^{\infty} xp(x).$$

Then, drop one of the two terms above and bound the remaining term to obtain the term $aP(X \geq a)$. Then, divide by a to finish the proof.

Exercise 2.8 (Boole's Inequality). The following inequality holds for any events A_1, \dots, A_n :

$$\mathbb{P}\left(\bigcup_{i=1}^n A_i\right) \leq \sum_{i=1}^n \mathbb{P}(A_i).$$

To prove that this inequality holds, note that if the $\{A_i\}$ are mutually disjoint (meaning that for each pair A_i, A_j it holds that $A_i \cap A_j = \emptyset$), then the above inequality is in fact an equality. Now, by representing events as circles in a Venn diagram, argue why the above inequality always holds. See the Wikipedia page on "Inclusion-Exclusion principle" to see what it meant by Venn diagram and how it is helpful for this problem.