

Basic Properties of Erdős Renyi Model

Ravi Niure

Seminar “Theoretical Topics in Data Science”

Supervisor: León Bohn

RWTH Aachen University

Email: *ravi.niure@mail.rwth-aachen.de*

February 5, 2021

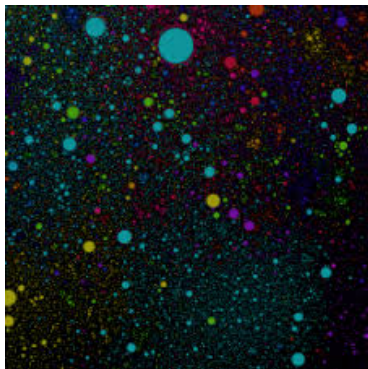
Overview

- 1 Motivation
- 2 Definition
- 3 Additional Concepts
- 4 Degree Distribution
- 5 Phase Transitions
- 6 Existence of Triangles
- 7 Threshold for Diameter of two
- 8 Conclusion

Motivation

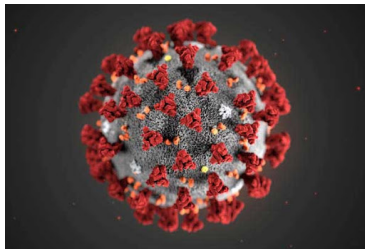
Networks are everywhere!

- COVID-19, World Wide Web, The Internet, Transportation networks, Gene interactions, Social networks, scientific collaborations, etc.



(a) Map of Internet

Source: <https://internet-map.net>



(b) SARS-CoV-2 virus

- **Mathematical representation:** Graphs with several *nodes* or *vertices* and *edges* or *links*
- Real networks are large and complex.
- Complete description of such networks and graphs is almost impossible.
- Statistical approach: link between microscopic properties and macroscopic phenomena!

This leads to **Random Graphs!**

Random Graphs

- A **random graph** is a result of **random process**.
- **Random process:** Connection between two elements is a random event determined by the sum of a very large number of unpredictable events. [2]
- A simple random graph model is **Erdős-Rényi Model**.

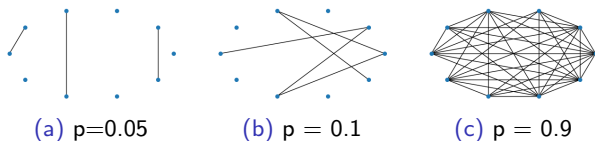


Figure: Graphs with $N = 10$ with different edge probability

Definition

Erdős-Rényi Model

Definition

The $G(n, p)$ model is a graph-valued random variable with two parameters, $n \in \mathbb{N}$, the total number of vertices and $p \in [0, 1]$, the probability of edge formation between any pair of vertices (v, w) . [3]

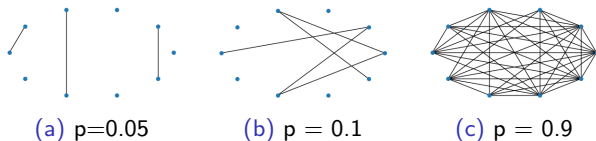


Figure: Graphs with $N = 10$ with different edge probability

- **Potential neighbours**

$$n - 1 \cong n \text{ for } n \rightarrow \infty$$

- **Edge Probability** as a function of n

$$p = \frac{d}{n}, \text{ where } d \text{ is some constant in many cases.}$$

- **Expected Degree**

$$E[\text{degree}] = (n - 1) \frac{d}{n} \approx n \frac{d}{n} = d$$

Additional Concepts

Linearity of Expectation

Theorem (Linearity of expectation)

For any random variables X_1, X_2, \dots, X_n and constants $\alpha_1, \alpha_2, \dots, \alpha_n$,

$$\mathbb{E} \left(\sum_{i=1}^n \alpha_i X_i \right) = \sum_{i=1}^n \mathbb{E} (\alpha_i X_i).$$

Markov's Inequality

Theorem (Markov's inequality)

Let x be a non-negative random variable. Then for $\alpha > 0$,

$$\text{Prob}(x \geq \alpha) \leq \frac{E(x)}{\alpha}.$$

If $\alpha = 1$, then,

$$\text{Prob}(x \geq 1) \leq E(x).$$

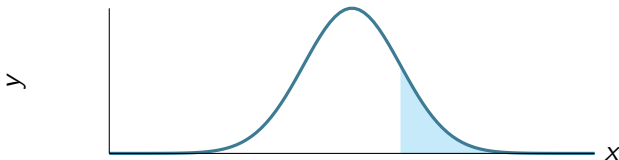


Figure: Markov's Inequality

Chebychev's Inequality

Theorem (Chebychev's inequality)

Let x be a random variable. Then for $c > 0$,

$$\text{Prob}(|x - E(x)| \geq c) \leq \frac{\text{Var}(x)}{c^2}.$$

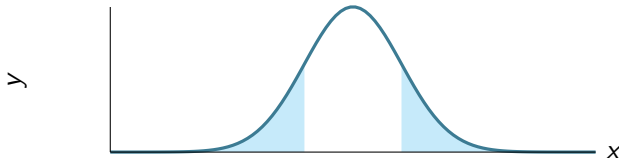


Figure: Chebychev's Inequality

Degree Distribution

Degree Distribution

- The degree distribution of $G(n, p)$ is **binomial**.

$$\begin{aligned}\text{Prob}(\text{vertex has degree } k) &= \binom{n-1}{k} p^k (1-p)^{n-1-k} \\ &\approx \binom{n}{k} p^k (1-p)^{n-k}\end{aligned}$$

- As $n \rightarrow \infty$, the **binomial** distribution approaches **poisson** distribution

$$\lim_{n \rightarrow \infty} \binom{n}{k} \left(\frac{d}{n}\right)^k \left(1 - \frac{d}{n}\right)^{n-k} = \frac{n^k}{k!} \frac{d^k}{n^k} e^{-d} = \frac{d^k}{k!} e^{-d}$$

for $p = \frac{d}{n}$ and constant d .

Degree Distribution in Real Networks

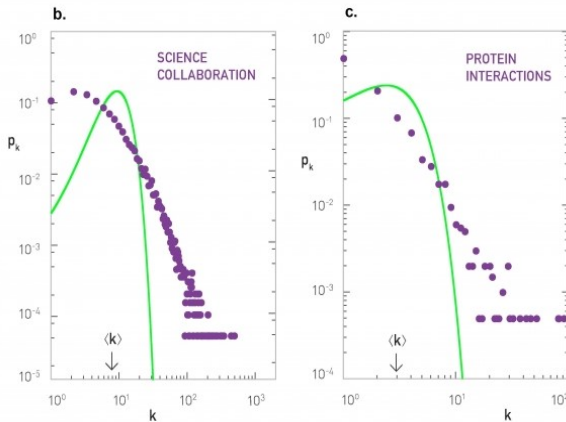


Figure: Degree Distribution in various networks and Poisson approximation [1]

Limitations of ER Model

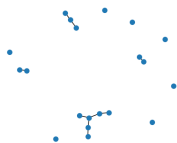
- Degree distribution is concentrated around expected degree.
- Independent edge formation doesn't reflect real networks.

Phase Transitions

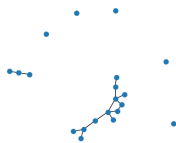
Phase Transitions



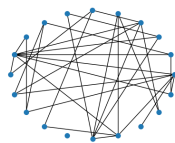
(a) $p=0.01$



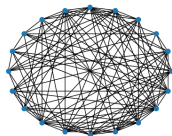
(b) $p = 0.05$



(c) $p = 0.06$



(d) $p = 0.15$



(e) $p = 0.55$

Figure: Graphs with $N = 20$

Threshold

Definition

For any monotonically increasing property, if there exists a function $p(n)$ such that:

- For all $p_1(n) \in o(p(n))$, the graph $G(n, p_1(n))$ almost surely does not have the property.
- For any $p_2(n) \in \omega(p(n))$, the graph $G(n, p_2(n))$ almost surely has the property.

Then we say that a phase transition occurs and $p(n)$ is the *threshold*.

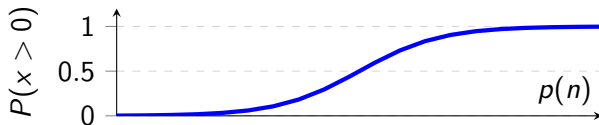


Figure: Threshold

Sharp Threshold

Definition

For threshold of $p(n)$, if there exists a constant c such that:

- for $c < 1$, the graph $G(n, cp(n))$ almost surely does not have the property.
- for $c > 1$, the graph $G(n, cp(n))$ almost surely has the property.

Then $p(n)$ is a *sharp threshold*.

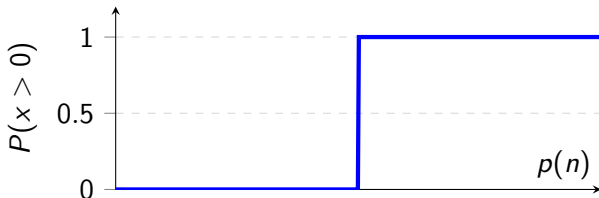


Figure: Sharp Threshold

Theorem (First Moment Method)

Let $x \in \mathbb{Z}_{\geq 0}$ be a random variable representing the number of occurrence of a certain property. Then, as $\mathbb{E}(x)$ approaches 0, with n approaching infinity, the probability of the existence of the property tends to 0.

Theorem (Second Moment Method)

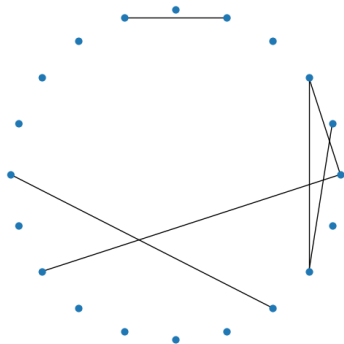
Let $x \in \mathbb{Z}_{\geq 0}$ be a random variable representing the number of occurrence of a certain property, and $E(x) > 0$. If $\text{Var}(x) = o(E^2(x))$, then x is almost surely greater than zero.

Corollary

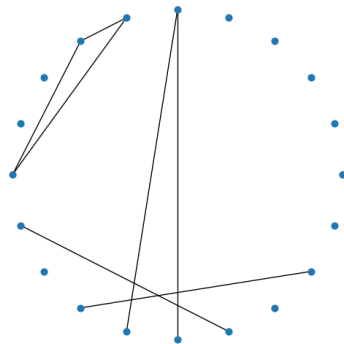
If $E(x^2) \leq E^2(x)(1 + o(1))$, then x is almost surely greater than zero.

Existence of Triangles

Triangles



(a) Graph without a triangle



(b) Graph with a triangle

Lemma

The expected number of triangles in a random graph $G(n, \frac{d}{n})$ is $\frac{d^3}{6}$.

Proof:

Let x be the number of triangles. Then x is given by,

$x = \sum_{ijk} \Delta_{ijk}$, where Δ_{ijk} is an indicator variable.

$$\begin{aligned} E(x) &= E\left(\sum_{ijk} \Delta_{ijk}\right) = \sum_{ijk} E(\Delta_{ijk}) = \binom{n}{3} \left(\frac{d}{n}\right)^3 \\ &= \frac{n!}{3!(n-3)!} \left(\frac{d}{n}\right)^3 \approx \frac{n^3}{3!} \left(\frac{d}{n}\right)^3 = \frac{d^3}{6} \end{aligned}$$

Introvert vs Extrovert Rooms

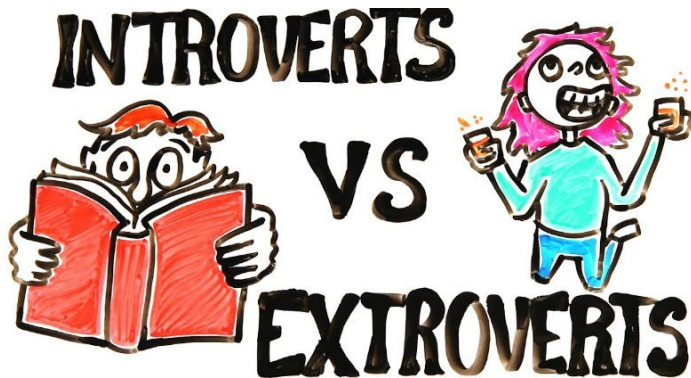


Figure: Introverts and Extroverts

Source: <https://www.asapscience.com/>

Lemma

Let x be the number of triangles in $G(n, \frac{d}{n})$, then $\text{Var}(x) \leq E(x) + o(1)$.

Proof:

Lets start by calculating $E(x^2)$. We can write x as

$$x = \left(\sum_{ijk} \Delta_{ijk} \right)$$

Expanding the squares of sum,

$$E(x^2) = E \left(\sum_{i,j,k} \Delta_{ijk} \right)^2 = E \left(\sum_{\substack{i,j,k \\ i',j',k'}} \Delta_{ijk} \Delta_{i',j',k'} \right)$$

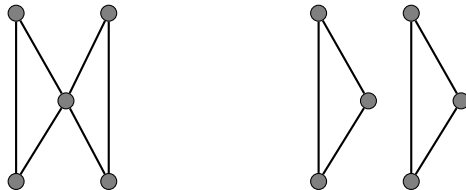


Figure: Triangles share 1 or less vertex

$$\mathbb{E} \left(\sum \Delta_{ijk} \Delta_{i'j'k'} \right) = \sum \mathbb{E}(\Delta_{ijk}) \mathbb{E}(\Delta_{i'j'k'}) \leq \mathbb{E}^2(x)$$

Note: The independence of the triangles and linearity of expectation!

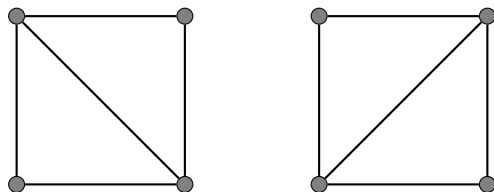


Figure: Triangles share 2 vertices

Note: We can't make independence assumption with the triangles. However, we can reformulate the problem as "rectangle with a diagonal" as opposed to two triangles.

let y be the number of rectangles with diagonal.

$$\text{Prob}(\square_{abcd}) = \left(\frac{d}{n}\right)^4 \left(\frac{d}{n} + \frac{d}{n}\right) - \left(\frac{d}{n}\right)^6 \approx 2 \left(\frac{d}{n}\right)^5$$

The expected number of such rectangles with either diagonal becomes,

$$E(y) = E\left(\sum_{abcd} \square_{abcd}\right) = \sum_{abcd} (E \square_{abcd}) = 2 \binom{n}{4} \left(\frac{d}{n}\right)^5 \approx \frac{d^5}{6n}$$

As n approaches ∞ , this expectation will approach zero, i.e $o(1)$.

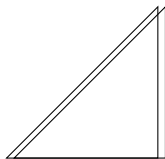


Figure: Two triangles sharing 3 vertices

For part 3, i, j, k and i', j', k' are the same sets. The contribution of this part of the summation to $E(x^2)$ is $E(x)$.

Thus, putting all the parts together,

$$E(x^2) = E \left(\sum_{\substack{i,j,k \\ i',j',k'}} \Delta_{ijk} \Delta_{i',j',k'} \right) \leq E^2(x) + E(x) + o(1)$$

which implies,

$$\text{Var}(x) = E(x^2) - E^2(x) \leq E(x) + o(1)$$

Threshold for Existence of Triangles

Corollary (Threshold for triangles)

The threshold for the existence of triangles in $G(n, \frac{d}{n})$ is $p(n) = \frac{1}{n}$.

Proof. For x to be zero, x would have to differ from its expected value by at least the expected value. Thus,

$$\begin{aligned}\text{Prob}(x = 0) &\leq \text{Prob}(|x - E(x)| \geq E(x)) \leq \frac{\text{Var}(x)}{E^2(x)} \leq \frac{E(x) + o(1)}{E^2(x)} \\ &\leq \frac{6}{d^3} + o(1)\end{aligned}$$

For $d > \sqrt[3]{6} \cong 1.8$, $\text{Prob}(x = 0) < 1$

For $d < 1.8$, $E(x) = \frac{d^3}{6} < 1$

Therefore, $p = \frac{d}{n} \approx \frac{1}{n}$ is the threshold.

Simulation with networkx

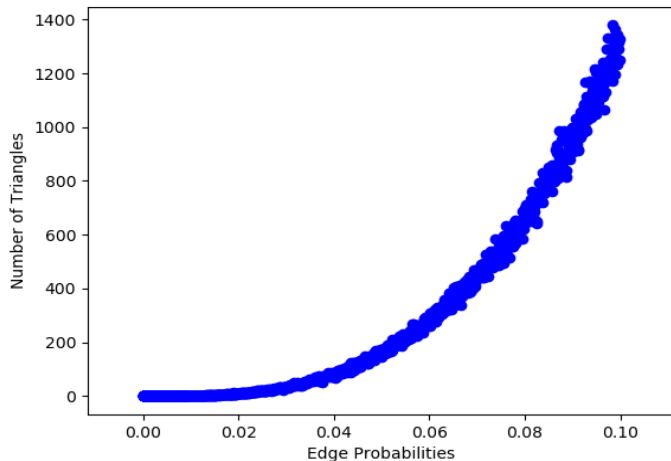


Figure: Number of Triangles when $N = 200$

Threshold for existence of triangles when $N = 200$

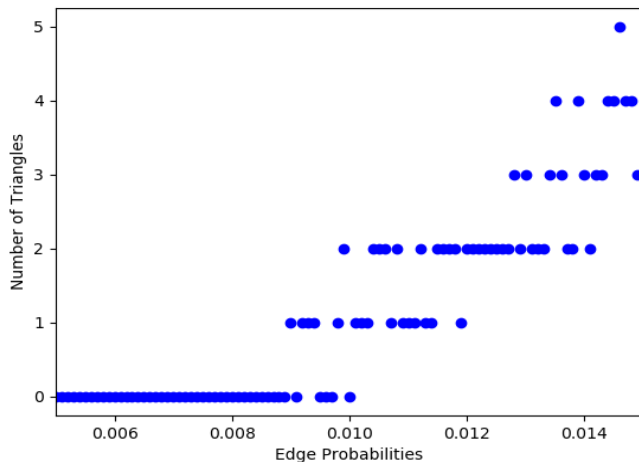
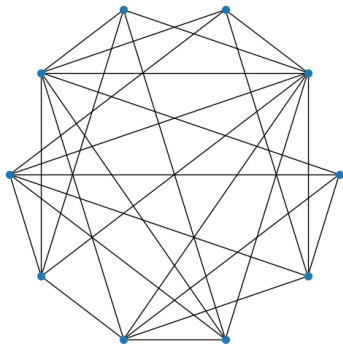


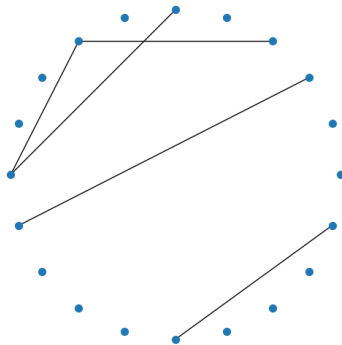
Figure: Number of Triangles when $N = 200$

Threshold for Diameter of two

Different diameter graphs



(a) Graph with diameter 2



(b) Graph with undefined diameter

Definition of diameter two

Definition

A graph has diameter two if and only if for each pair of vertices i and j , either there is an edge between them or there is another vertex k to which both i and j have an edge.

Good Pair Bad Pair!

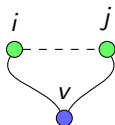
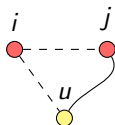
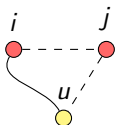
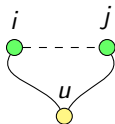
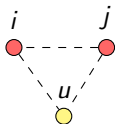


Figure: Good Pair Bad Pair

Threshold for diameter two

Theorem (Threshold for diameter two)

The property that $G(n, p)$ has diameter two has a sharp threshold at

$$p = \sqrt{2} \sqrt{\frac{\ln n}{n}}.$$

- $\frac{1}{\sqrt{n}}$: threshold for finding a common neighbor
- $\sqrt{\ln n}$: ensures every pair has a common neighbor

Proof - First moment argument

Number of bad pair of vertices $(x) = \sum_{i < j} l_{ij}$

$$E(x) = \binom{n}{2} (1-p) (1-p^2)^{n-2}$$

Let $p = \sqrt{c} \sqrt{\frac{\ln n}{n}}$ and for large n ,

$$E(x) \cong \frac{n^2}{2} \left(1 - c \sqrt{\frac{\ln n}{n}}\right) \left(1 - c^2 \frac{\ln n}{n}\right)^n \cong \frac{n^2}{2} e^{-c^2 \ln n} \cong \frac{1}{2} (n)^{2-c^2}$$

Expectation

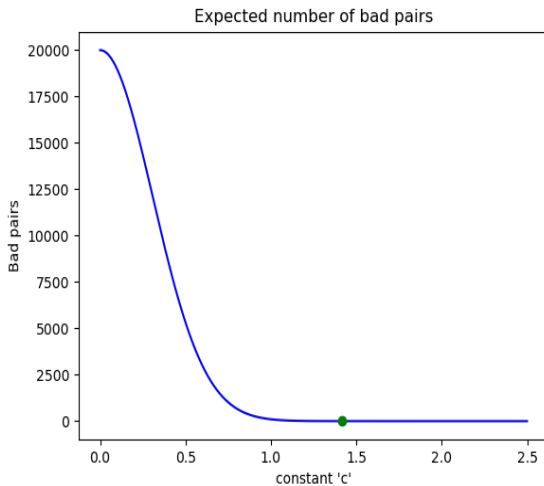


Figure: Expected bad pairs

For $c < \sqrt{2}$, Second moment argument

$$\begin{aligned} E(x^2) &= E \left(\sum_{i < j} l_{ij} \right)^2 = E \left(\sum_{i < j} l_{ij} \sum_{k < l} l_{kl} \right) = E \left(\sum_{\substack{i < j \\ k < l}} l_{ij} l_{kl} \right) \\ &= \sum_{\substack{i < j \\ k < l}} E(l_{ij} l_{kl}). \end{aligned}$$

The above sum can be partitioned into:

$$E(x^2) = \underbrace{\sum_{\substack{i < j \\ k < l}} E(l_{ij} l_{kl})}_{a=4} + \underbrace{\sum_{\substack{\{i,j,k\} \\ i < j}} E(l_{ij} l_{ik})}_{a=3} + \underbrace{\sum_{i < j} E(l_{ij}^2)}_{a=2}$$

Second moment argument - cont'd

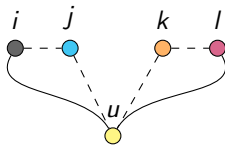


Figure: Four distinct vertices

For the first part,

$$\mathbb{E}(I_{ij}I_{kl}) \leq (1 - p^2)^{2(n-4)} \leq n^{-2c^2} (1 + o(1))$$

Across all $\{i, j, k, l\}$,

$$\sum_{\substack{i < j \\ k < l}} \mathbb{E}(I_{ij}I_{kl}) \leq \frac{1}{4} n^{4-2c^2} (1 + o(1)) = \mathbb{E}^2(x) (1 + o(1))$$

Second moment argument - cont'd

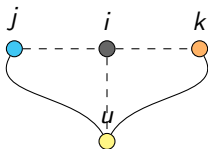


Figure: Three distinct vertices

For the second part of the sum,

$$\sum_{\substack{\{i,j,k\} \\ i < j}} \mathbb{E}(l_{ij}l_{ik}) \leq n^{3-2c^2}$$

Second moment argument - cont'd

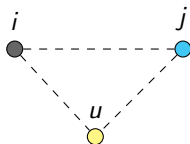


Figure: Three distinct vertices

For the third summation,

$$\sum_{ij} E(l_{ij}^2) = E(x).$$

Combining all three parts,

$$E(x^2) \leq \frac{1}{4}n^{4-2c^2} + n^{3-2c^2} + n^{2-c^2}$$

Second moment argument - cont'd

For $c < \sqrt{2}$,

$$E(x^2) \leq E^2(x) (1 + o(1))$$

This gives,

$$E(x^2) - E^2(x) \leq E^2(x) o(1)$$

Thus,

$$\text{Var}(x) \leq E^2(x) o(1) \implies \text{Var}(x) = o(E^2(x))$$

As such, from second moment argument, a bad pair almost surely exists and therefore there is no graph with diameter ≤ 2 .

Simulation with networkx

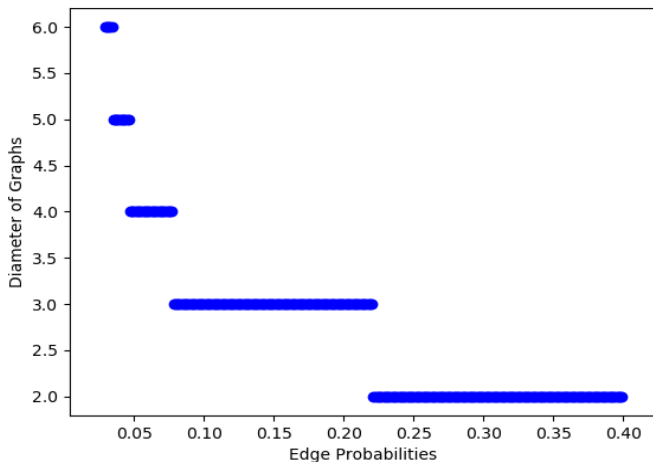


Figure: Diameter of Graphs with different probability when $N = 200$

Conclusion

- Random Graphs provide essential frameworks for studying large and complex networks.
- Erdős-Rényi model is the earliest and simple random graph model.
- **Limitations:** Thin tail, independent edge formation, etc
- Concentration inequalities provide expectation and variance relationship; used to argue for/against existence of properties.
- Triangles are essential to compute clustering coefficients.
 - Expected number of triangles, independent of nodes, is $\frac{d^3}{6}$.
 - The sharp threshold for existence is $\frac{1}{n}$.
- The threshold for two degrees of separation is $\sqrt{2\frac{\ln n}{n}}$.



Albert Barabási. *Random Networks Network Science*. URL: <http://networksciencebook.com/chapter/3> (visited on 02/05/2021).



Alain Barrat, Marc Barthélemy, and Alessandro Vespignani. *Dynamical Processes on Complex Networks*. 1st. New York, USA: Cambridge University Press, 2008. ISBN: 9780521879507.



Avrim Blum, John Hopcroft, and Ravindran Kannan. *Foundations of Data Science*. pre-publication. Cambridge, UK: Cambridge University Press, 2020. ISBN: 9781108485067.