

# A Simulation Study Using Random Graph Models to Fit Social Networks

Levi Lee  
Advisor: Amy Wagaman  
Amherst College

April 12, 2017

# A Highly Connected World

- Networks are everywhere
- Informational, biological, technological, and social applications
- Focus is on social networks

- Travers & Milgram (1967)
- Letter correspondence between strangers in Nebraska and Massachusetts
- Required only six people on average for the letter to be delivered
- “Six degrees of separation”

# Social Networks

- Involve a lot of people
- Form highly condensed groups
- Have relatively short distances between people

# Question

- Can we simulate this? If so, how?
- Conduct a simulation study utilizing different graph models

# Basic Terminology

- A *graph*, denoted  $G$  or  $G(V, E)$ , consists of a set of *vertices*  $i, j, k, \dots \in V$  and a set of *edges*  $\{i, j\}, \{i, k\}, \{j, k\} \dots \in E$ .
- Our focus is on *simple*, *undirected*, and *connected* graphs.

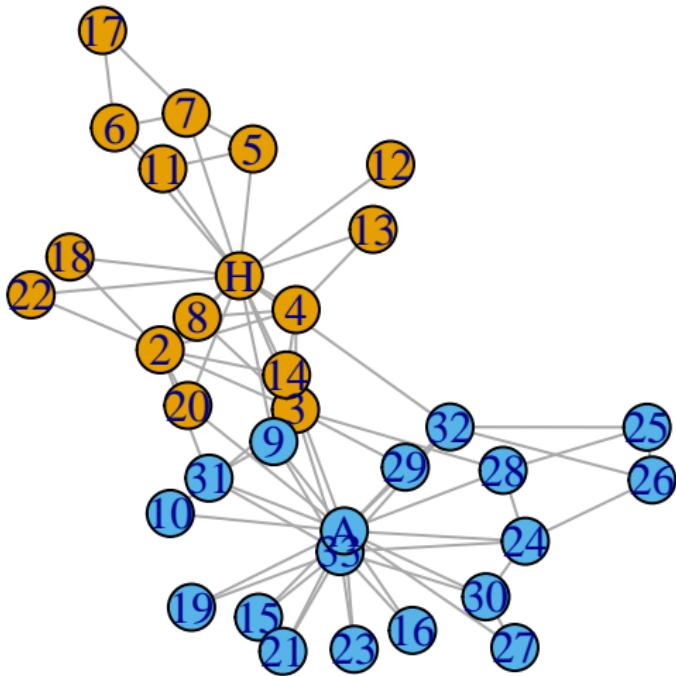
# Other Representations of Graphs

- *Adjacency matrix*: denoted  $\mathbf{A}$  is an  $N_V \times N_V$  matrix where each element denotes the existence of edges between pairs

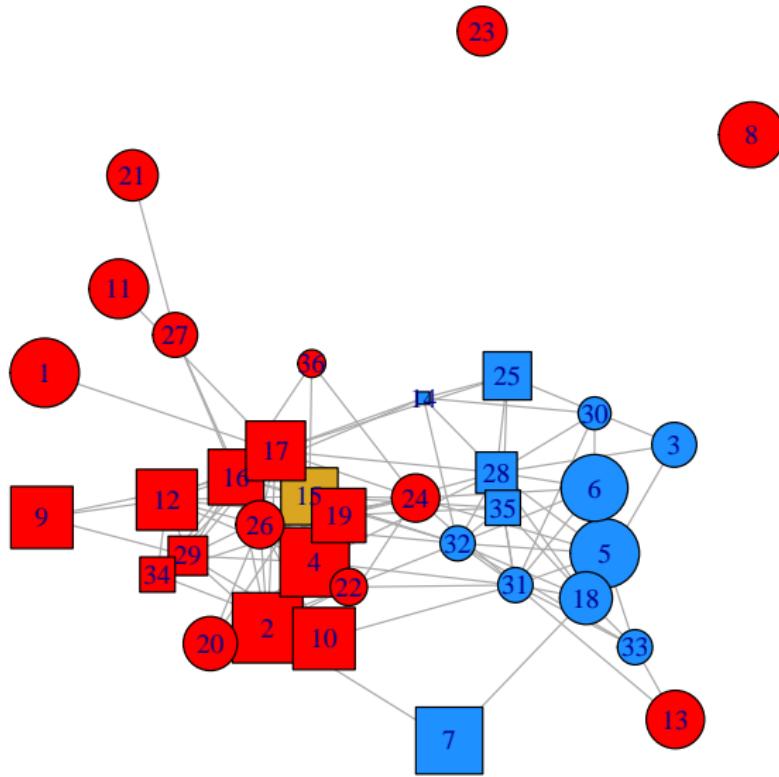
$$A_{ij} = \begin{cases} 1 & \text{if } \{i, j\} \in E \\ 0 & \text{otherwise.} \end{cases}$$

- *Edge list*: two-column list of all the edges in a graph denoted by their corresponding vertices

# Example: Karate Club of Zachary (1977)



# Example: Lazega's Law Firm (1999)



# Network Statistics

- Analogous to statistics seen in elementary statistics
- Characterize a given network with a numerical summaries

# Transitivity/Clustering Coefficient

- Ratio of triangles to connected triples
- *Triangle*: three vertices connected by three edges
- *Connected triple*: three vertices connected by two edges

$$C = \frac{(\text{number of triangles}) \times 3}{\text{number of connected triples}}$$

- Local versus global

# Notions of Distance

*Average path length:* average of the shortest paths of all distinct pairs of vertices in the network

*Diameter:* longest of all the shortest paths between distinct pairs of vertices

## Example: Karate Club and Law Firm

Network Statistic	Zachary's Karate Club	Lazega's Law Firm
Transitivity (Global)	0.256	0.389
Transitivity (Local)	0.588	0.487
Avg. Path Length	2.408	2.144
Diameter	13	5

# Centrality

- Measure of importance for each vertex in the graph
- Many different types of centralities exist

# Degree Centrality

- Based on the number of edges that are connected to a vertex
- Most common form of centrality

# Closeness Centrality

- Measures how close a vertex is to other vertices based on the inverse of the total distance of the vertex from all others

$$c_{Cl}(i) = \frac{1}{\sum_{j \in V} d(i, j)}$$

- $dist(i, j)$  is the geodesic distance between the vertices  $i, j \in V$

# Betweenness Centrality

- Measures the extent to which a vertex is located between other pairs of vertices

$$c_B(i) = \sum_{g \neq h \neq i \in V} \frac{\sigma(g, h|i)}{\sigma(g, h)}$$

- $\sigma(g, h|i)$  is the total number of shortest paths between  $g$  and  $h$  that pass through  $i$ , and  $\sigma(g, h) = \sum_{i \in V} \sigma(g, h|i)$

# Eigenvector Centrality

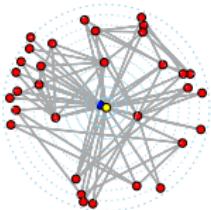
- Based on the idea of “status,” “prestige,” or “rank;” the more central the neighbors of a vertex are, the more central that vertex itself is

$$c_{Ei}(i) = \alpha \sum_{\{i,j\} \in E} c_{Ei}(u)$$

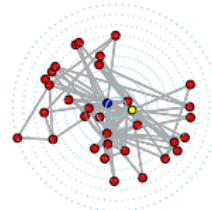
- $c_{Ei} = (c_{Ei}(1), \dots, c_{Ei}(N_V))^T$  is the solution to the eigenvalue problem  $\mathbf{Ac}_{Ei} = \alpha^{-1}\mathbf{c}_{Ei}$ , where  $\mathbf{A}$  is the adjacency matrix for network graph  $G$ .

# Example: Karate Club

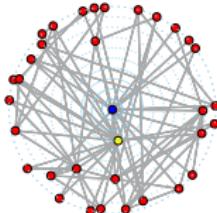
Degree



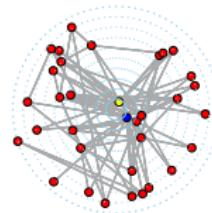
Closeness



Betweenness



Eigenvector



## Example: Karate Club and Law Firm

Network Statistic	Zachary's Karate Club	Lazega's Law Firm
Degree	4.588	6.389
Closeness Cen.	0.005	0.007
Betweenness Cen.	26.194	17.833
Eigenvector Cen.	0.377	0.458

# Graph Models

- A *graph model* takes in fixed parameters and generates a graph that varies in structure with each iteration
- Equivalently, it is a collection, or *ensemble* of graphs, denoted by

$$\{\mathbb{P}_\theta(G), G \in \mathcal{G} : \theta \in \Theta\}$$

-  $\mathcal{G}$  is a collection or ensemble of possible graphs,  $P_\theta$  is a *probability distribution* on the random graph  $G$ , and  $\theta$  is a vector of parameters that describe the graphs that  $G$  can be, ranging over possible parameters in  $\Theta$

# Erdős–Rényi Model (1959)

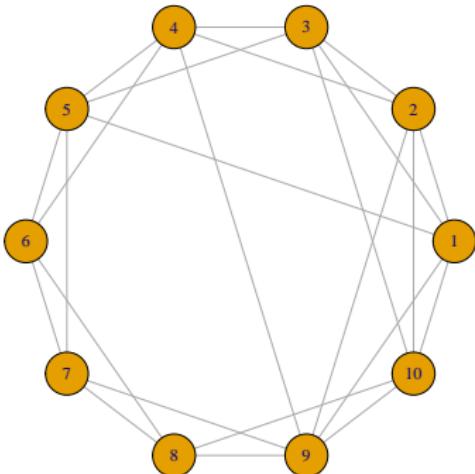
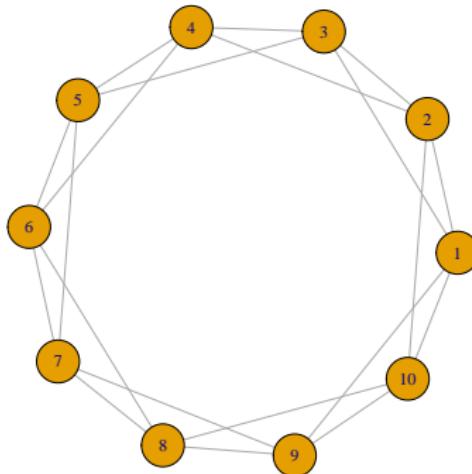
- Model with parameters:  $N_V$ , and  $N_E$  or  $p$
- Model of the form  $G(N_V, p)$  and  $G(N_V, N_E)$

# Properties of the Erdős–Rényi Model

- Short average path lengths
- Low clustering coefficient
- For the  $G(N_V, p)$  model, a particular simple graph  $g$  with exactly  $N_V$  vertices has probability  
$$P(G = g) = p^{N_E} (1 - p)^{\left(\binom{N_V}{2} - N_E\right)}$$

# Watts-Strogatz Model (1998)

- Model with parameters:  $N_V$ ,  $r$ ,  $p$



# Properties of the Watts-Strogatz Model

- High clustering coefficient
- Small average path length

# Exponential random graph models (ERGMs) I

- Exponential random graph models (ERGMs) are a class of models that can be used to generate probability distributions
- Flexible in design; we can decide our parameters
- Conduct goodness-of-fit tests for model assessment

# Exponential random graph models (ERGMs) II

- The general form for an ERGM is as follows:

$$P_{\theta, \mathcal{G}}(\mathbf{G} = \mathbf{g}) = \frac{\exp(\theta^T \mathbf{s}(\mathbf{g}))}{\kappa(\theta, \mathcal{G})}, \mathbf{g} \in \mathcal{G}$$

- **G** is the random variable representing a random graph and **g** is the particular adjacency matrix we observe. **s(g)** is the vector of model statistics for **g**,  $\theta$  is the vector of coefficients for those statistics, and  $\kappa(\theta, \mathcal{G})$  is the quantity in the numerator summed over all possible networks

# Properties of ERGMs I

- Deriving the Erdős-Rényi Model from ERGMs
- Suppose we have a particular graph  $g$  and the only statistic we have is  $L(G)$ , the number of edges in  $g$

$$P_{\theta, \mathcal{G}}(g) = \frac{\exp(\theta_L L(g))}{\sum_{g' \in \mathcal{G}} \exp(\theta_L L(g'))} = \frac{\exp(\theta_L L(g))}{\kappa(\theta, \mathcal{G})}, g \in \mathcal{G}$$

## Properties of ERGMs II

- Consider the probability distribution for a particular graph  $g$  with  $N_E$  edges again (from the Erdős-Rényi model). Using the fact that  $N_E = L(g)$ , taking the equation as a power of base  $e$ , we get the following:

$$\begin{aligned} P(g) &= p^{N_E} (1-p)^{\binom{N_V}{2} - N_E} \\ &= p^{L(g)} (1-p)^{\left(\frac{N_V(N_V-1)}{2} - L(g)\right)} \\ &= \left(\frac{p}{1-p}\right)^{L(g)} (1-p)^{\frac{N_V(N_V-1)}{2}} \\ &= \exp\left(L(g)\log\left(\frac{p}{1-p}\right) - \frac{N_V(N_V-1)}{2}\log\left(\frac{1}{1-p}\right)\right) \end{aligned}$$

# Properties of ERGMs III

$$\begin{aligned} P(g) &= \exp(\theta_L L(g) - c) \\ &= \frac{\exp(\theta_L L(g))}{\exp(c)} \end{aligned}$$

- This is the form of the ERGM seen earlier

$$P_{\theta, \mathcal{G}}(g) = \frac{\exp(\theta_L L(g))}{\kappa(\theta, \mathcal{G})}, g \in \mathcal{G}$$

- Constant  $c$  is exactly the normalizing constant we need to scale this into a probability distribution

# Data Set

- A component of the Facebook network
- 4039 vertices and 88234 edges
- Simple, connected, and undirected

# The Facebook Component

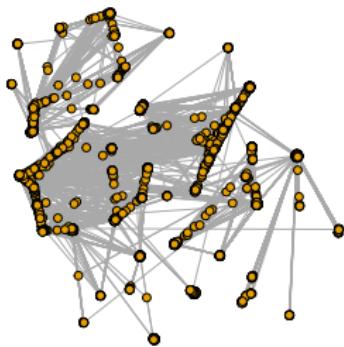
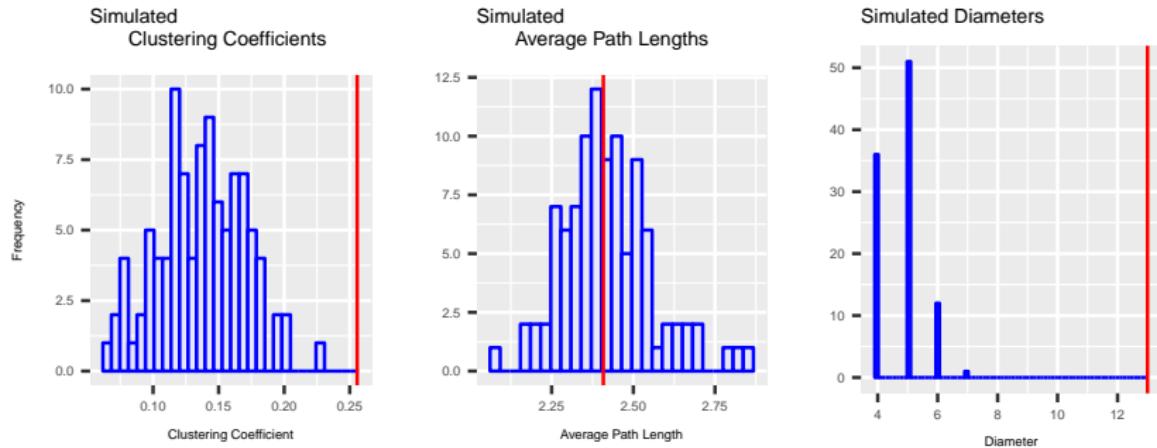


Figure 1: Overview of our component of the Facebook network

# Generating Random Graphs

- Create graphs of parable magnitude using certain parameters
- Simulate 1000 random graphs for and record network statistics for each graph
- Create distribution of these values and/or see a table of averages
- Compare network statistics of the Facebook network with that of the graphs we generated

# The Picture in Mind



# Erdős–Rényi Model

- $N_V$  and  $p$
- Number of vertices is  $N_V = 4039$
- Estimate the probability by taking the number of observed edges and dividing by the number of possible edges  
$$\hat{p} = \frac{88234}{\binom{4039}{2}} = 0.011$$
- For every possible edge among the 4039 vertices, determine if an edge will form based on the estimated probability

# Watts-Strogatz Model

- $N_V, r, p$
- Will not use  $p$
- Start with a circular model with 4039 vertices
- Assign a number of edges to the vertices equal to the smallest degree observed in our Facebook network;  $deg_{min} = 1$
- Randomly add  $88234 - 4039 = 84195$  edges until we have the same number as our observed network
- Simplify our simulated graph to eliminate multi-edges and loops

# Results for Erdős–Rényi and Watts-Strogatz Models

---

Network Statistic	Observed	Erdős-Rényi	Watts-Strogatz
Transitivity	0.617	$0.0108 \pm 0.0001$	$0.0107 \pm 9e-05$
Avg. Path Length	4.338	$2.606 \pm 0.002$	$2.609 \pm 0.0002$
Diameter	17	$3.96 \pm 0.21$	$3.95 \pm 0.22$
Degree Cen.	43.691	$43.69 \pm 0.14$	$43.45 \pm 0.01$
Betweenness Cen.	2072.642	$3242 \pm 4$	$3249.2 \pm 0.4$
Closeness Cen.	$8.881e-08$	$9.507e-05 \pm 7e-08$	$9.494e-05 \pm 7e-09$
Eigenvector Cen.	0.040	$0.620 \pm 0.022$	$0.6235 \pm 0.0227$

---

# ERGMs

- Four different ERGMs—labeled as ERGM 1a, ERGM 2a, ERGM 2b, and ERGM 3a
- ERGM 1a: one parameter: edges
- ERGM 2a: two parameters: edges and triangles
- ERGM 2b: two parameters: edges and k-stars (of size 3)
- ERGM 3a: three parameters: edges, triangles, and k-stars (of size 3)
- Estimate the coefficients ( $\theta$ 's) of the probability model
- Simulate

# Results for ERGMs I

---

Network Statistic	Observed	ERGM1a	ERGM2a
Transitivity	0.617	$0.3696 \pm 0.0012$	$0.4823 \pm 0.0020$
Avg. Path Length	4.338	$2.885 \pm 0.004$	$3.052 \pm 0.009$
Diameter	17	$5.216 \pm 0.412$	$6.098 \pm 0.035$
Degree	43.691	$43.66 \pm 0.05$	$44.54 \pm 0.03$
Betweenness Cen.	2072.642	$3805 \pm 9$	$4140 \pm 18$
Closeness Cen.	8.881e-08	$8.286e-05 \pm 8.35e-06$	$6.008e-05 \pm 1.514e-05$
Eigenvector Cen.	0.040	$0.0417 \pm 0.0008$	$0.0410 \pm 5e-05$

---

# Results for ERGMs II

---

Network Statistic	Observed	ERGM2b	ERGM3a
Transitivity	0.617	$0.3787 \pm 0.0013$	$0.4891 \pm 0.0020$
Avg. Path Length	4.338	$2.851 \pm 0.003$	$3.056 \pm 0.008$
Diameter	17	$5.095 \pm 0.293$	$6.161 \pm 0.368$
Degree Cen.	43.691	$44.06 \pm 0.06$	$44.52 \pm 0.03$
Betweenness Cen.	2072.642	$3736 \pm 6$	$4148 \pm 16$
Closeness Cen.	$8.881e-08$	$8.558e-05 \pm 6.03e-06$	$5.875e-05 \pm 1.5035e-05$
Eigenvector Cen.	0.040	$0.0392 \pm 0.0002$	$0.0410 \pm 3e-05$

---

# Conclusions and Future Work

- The models were bad, so now what?
- Choose other models
- Choose different network statistics
- Choose other data sets

# Implications

- Understand the flow of information
- Better access to jobs through networking
- Better leads to resources in research
- Improving traffic
- Understanding biological systems
- Product recommendations

# Sources |

- Baumer, B., Kaplan, D., & Horton, N. (2017). Modern Data Science With R: With Digital Download. Taylor & Francis.
- Butts, C. & others. (2008). Social network analysis with sna.
- Butts, C. & others. (2015, June). Introduction to Exponential-family Random Graph modeling with ergm.
- Hunter, D. & others. (2008). Ergm: A package to fit, simulate and diagnose exponential-family models for networks.
- Leskovec, J. & others. (2005). Realistic, Mathematically Tractable Graph Generation and Evolution, Using Kronecker Multiplication.
- Jackson, M. O. (2013). Social and economic networks: Models and analysis.
- Kolaczyk, E. D. (2009). Statistical Analysis of Network Data: Methods and Models.
- Kolaczyk, E. D. & Csárdi, G. (2014). Statistical Analysis of Network Data with R.

## Sources II

- Lazega, E. & Pattison, P. E. (1999). Multiplexity, generalized exchange and cooperation in organizations: A case study.
- Leskovec, J. & Krevl, A. (2014, June). SNAP Datasets: Stanford large network dataset collection.
- Leskovec, J. & others. (2005). Realistic, Mathematically Tractable Graph Generation and Evolution, Using Kronecker Multiplication.
- Leskovec, J. & others. (2010). Kronecker Graphs: An Approach to Modeling Networks.
- Newman, M. (2010). Networks: An Introduction.
- Travers, J., & Milgram, S. (1967). The small world problem.
- Watts, D. J., & Strogatz, S. H. (1998). Collective dynamics of “small-world” networks.
- Zachary, W. W. (1977). An information flow model for conflict and fission in small groups.