

# Advanced Network Analysis

## Introduction to Exponential Random Graph Models

Shahryar Minhas [s7minhas.com]

# Reading

- Hunter DR, Handcock MS, Butts CT, Goodreau SM, Morris M. ergm: A package to fit, simulate and diagnose exponential-family models for networks. *Journal of Statistical Software*. 2008;24(3):1-29.
- Hunter, DR, Goodreau, SM, Handcock, MS. Goodness of fit of social network models. *Journal of the American Statistical Association*. 2008;103(481):248-258.
- Cranmer, S. J., & Desmarais, B. A. (2011). Inferential network analysis with exponential random graph models. *Political Analysis*, 19(1), 66-86.

# Why ERGMs?

- Test hypotheses about processes that give rise to particular network structures.
- ERGMs provide a way to:
  - estimate and evaluate the effect of exogenous covariates
  - and also the effect of endogeneous network structures, such as triangles, reciprocity, etc.

# High level logic

- Conceptualize the observed network data as just one realization of a set of possible networks with similar important characteristics produced by some unknown stochastic process
- A statistical model for a network on a given set of actors assigns a probability to all possible networks on those actors
- The range of possible networks and their probability of occurrence under the model is represented by a probability distribution on the set of all possible graphs
- Estimate model parameters using observed network as guide

# High Level Logic (TL;DR)

- **One Realization:** Our observed network is one possible outcome among many.
- **Stochastic Process:** These outcomes are generated by an unknown random process.
- **Probability Model:** Assigns likelihoods to all possible networks for a set of actors.
- **Parameter Estimation:** Use the observed network to estimate model parameters.

# ERGM Framework

- Let  $\mathcal{Y}$  be the sample space of  $Y$ , e.g.  $\{0, 1\}^N$
- Any model-class for the multivariate distribution of  $Y$  can be parameterized in the form:

$$\Pr(Y = y) = \frac{\exp(\theta^T g(y))}{k(\theta, \mathcal{Y})}, y \in \mathcal{Y}$$

- The above gives us the probability of a single graph
- $g(y)$ : vector of network statistics
- $\theta$ : vector of model parameters
- $k(\theta, \mathcal{Y})$ : normalizing constant which is summed over all possible graphs.  
The denominator represents the quantity from the numerator summed over all possible networks with  $n$  nodes, constraining the probabilities to sum to 1.

$$k(\theta, \mathcal{Y}) = \sum_{y \in \mathcal{Y}} \exp(\theta^T g(y))$$

# Vector of network statistics

- The  $g(y)$  term includes parameters to estimate the effect of "network statistics"
- These are counts of network configurations:
  - Density: # ties
  - Reciprocity: # of reciprocal ties
  - Triangles: # of triangles
  - and many more such terms
- We can also model the effect of exogenous dyadic and nodal covariates in this framework
- Many ERGM terms have been developed ([Morris et al. 2008](#) provide a comprehensive review)

# Estimation: MCMC-MLE

- ERGM computations are too difficult to perform directly, lets use an iterative method for simulating draws from a given distribution (e.g., see [Snijders 2002](#) and [Handcock 2003](#))
- ERGMs have no "closed form" or analytical solution to estimate the parameters,  $\theta$
- Basically, this procedure works by:
  - Starting with some initial values for  $\theta$
  - Simulate networks from those values
  - Compare the mean statistics to the observed
  - Repeat until difference is lower than some stopping condition



# MCMC-MLE ... more details

- Generate distributions of different  $g$ s that emerges for any given specification of  $\beta$ s
- Search over set of  $\beta$ s to find one that leads to the highest likelihood of getting a network that looks similar to the observed  $g$
- How to generate a distribution of different  $g$ 's for any given specification of  $\beta$ ?
  - Fix starting network at  $g^0$
  - Randomly pick a link,  $ij$ , to change
  - Then, based on MPLE, randomly put the link  $ij$  in or out with the appropriate probability given the profile of parameters  $\beta$  and given  $g_{-ij}^0$
  - This leads to a new network  $g^1$
  - Iterate through different links in the network, which results in a Markov chain over the resulting networks, and over time the probability that we visit any given network approaches that of its steady state distribution.

# Still some problems ...

Basic problem with MCMC-MLE is that even with an MCMC method, we are still only sampling relatively few networks relative to the huge number possible ...

- **Bhamidi et al. (2011)**: The time an MCMC has a chance to sample enough networks to gain a representative sample is generally exponential in the number of links ... unless the edges are independent, which if they are should then make you wonder why you did an ERGM in the first place.
- **Chatterjee and Diaconis (2013)**: Provide further evidence that calculating the normalizing constant for moderately sized networks is prohibitively time-consuming.
- **Shalizi & Rinaldo (2013)**: When running an ERGM inferences are only well calibrated, if you are able to view the whole network.

# Time to estimate some ERGMs



# Statnet package

statnet is a package in R that includes a range of "sub-packages" that we will be using. The `ergm` package is one such package and we will be using it to build our first inferential network model.

Project website:

```
library(statnet)
```

## statnet

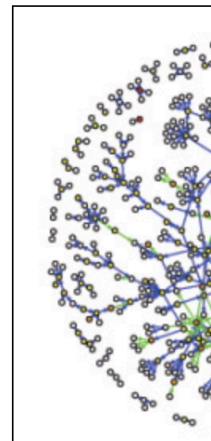
*Software tools for the analysis, simulation and visualization of network data.*

### Welcome to statnet!

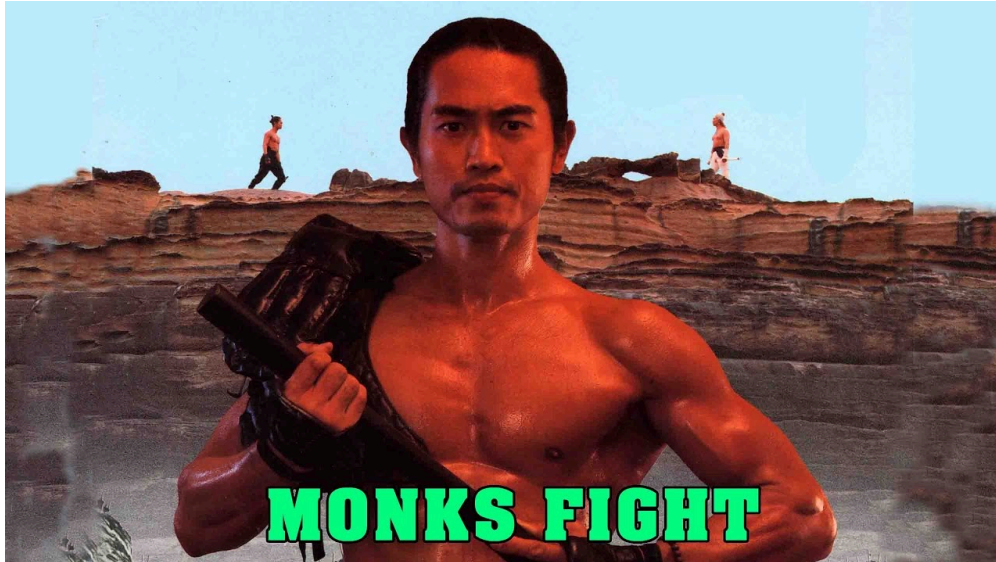
Visit the [statnet Wiki](#) for information on, background material for and access to the **statnet** suite of packages for network analysis. You can find [installation instructions](#), [tutorials](#), and [developer resources](#) at the wiki.

### What is statnet?

**statnet** is a suite of software packages for network analysis that implement recent advances in the statistical modeling of networks. The analytic framework is based on Exponential family Random Graph Models (ergm). **statnet** provides a comprehensive framework for ergm-based network modeling, including tools for model estimation, model evaluation, model-based network simulation, and network visualization. This broad functionality is powered by a central Markov chain Monte Carlo (MCMC) algorithm.



# Sampson Monastery Data



- Sampson (1969) recorded the interactions among a group of monks while he was a resident at their monastery.
- We'll be working with a cross-sectional directed dataset that flows from his work there.
- A directed edge from monk A to monk B exists if A indicated that he has positive relations with B.

# Sampson Monastery Data

```
data('sampson')  
samplike
```

```
## Network attributes:  
##   vertices = 18  
##   directed = TRUE  
##   hyper = FALSE  
##   loops = FALSE  
##   multiple = FALSE  
##   total edges= 88  
##     missing edges= 0  
##     non-missing edges= 88  
##  
## Vertex attribute names:  
##   cloisterville group vertex.names  
##  
## Edge attribute names:  
##   nominations
```

# network objects

network objects enable you to provide greater structure to adjacency matrices. To go from an adjacency matrix to a network object one just needs to run:

```
adjMat = as.matrix.network(samplike)
```

And vice versa:

```
as.network.matrix(adjMat,  
  matrix.type='adjacency',  
  directed=TRUE  
)
```

```
## Network attributes:  
##   vertices = 18  
##   directed = TRUE  
##   hyper = FALSE  
##   loops = FALSE  
##   multiple = FALSE  
##   bipartite = FALSE  
##   total edges= 88  
##   missing edges= 0
```

# Get information about attributes

To view the data stored in the attributes we can use get methods, lets see what's in the group attribute:

```
network::get.vertex.attribute(samplike, 'group')  
samplike %v% 'group' # returns the same result as above
```

```
## [1] "Turks"      "Turks"      "Outcasts"   "Loyal"      "Loyal"      "Loyal"  
## [7] "Turks"      "Loyal"      "Loyal"      "Loyal"      "Loyal"      "Turks"  
## [13] "Outcasts"   "Turks"      "Turks"      "Turks"      "Outcasts"   "Outcasts"
```

Sampson made these groups based on his observations:

- Loyal Opposition consists of the novices who entered the monastery first.
- The Young Turks arrived later, in a period of change. They questioned practices in the monastery, which the members of the Loyal Opposition defended.
- Some novices did not take sides in this debate, so they are labeled 'interstitial'.
- The Outcasts are novices who were not accepted in the group.



# Sampson data visualization

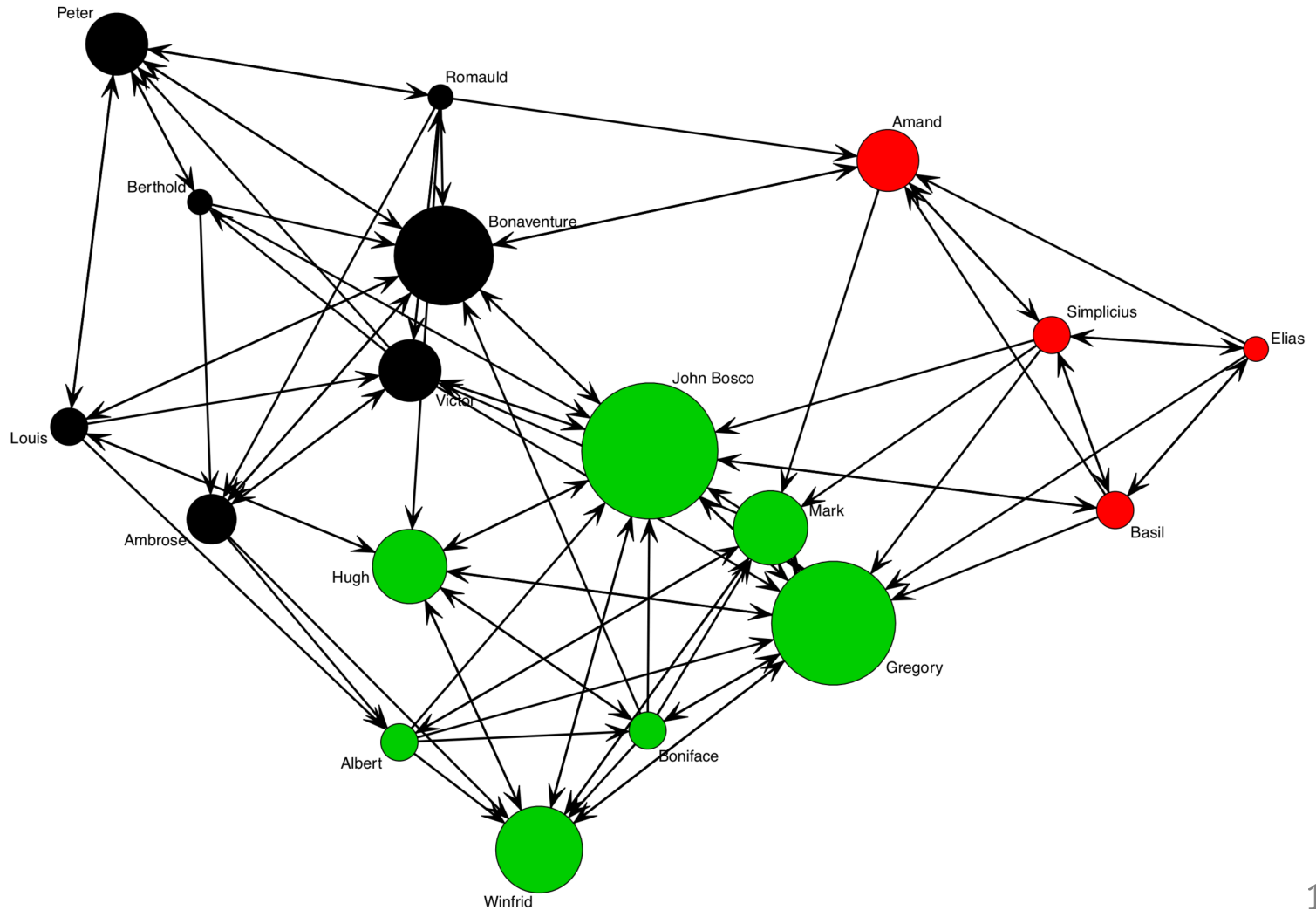
The network package comes with visualization functionality for networks as well. `igraph` is more flexible for visualization purposes, but most basic operations can still be accomplished with the network package.

```
set.seed(6886)

vertexSize = degree(samplike, cmode = 'indegree')/2

plot(samplike,
     displaylabels = TRUE,
     # size of nodes based on vector vertexSize
     vertex.cex = vertexSize,
     # color of nodes based on vertex attribute: group
     vertex.col = 'group'
)
```

# Sampson data visualization



# Exploring the Sampson data

- Lets run a simple model with no exogenous or endogenous covariates -- this is just equivalent to a GLM with an intercept term
  - (also referred to as an Erdos-Renyi model)
- The function to run an ERGM is simply `ergm`. The `statnet` package imports it from the `ergm` package.

```
m1 = ergm(samplike ~ edges)
```

# Summarizing the result from an `ergm` object

```
summary(m1)
```

```
## Call:
## ergm(formula = samplike ~ edges)
##
## Maximum Likelihood Results:
##
##           Estimate Std. Error MCMC % z value Pr(>|z|)
## edges   -0.9072      0.1263      0  -7.183   <1e-04 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##           Null Deviance: 424.2  on 306  degrees of freedom
##           Residual Deviance: 367.2  on 305  degrees of freedom
##
## AIC: 369.2  BIC: 372.9  (Smaller is better. MC Std. Err. = 0)
```

# Exploring the sampson data

- Interpretation here is straightforward, baseline probability of a tie in the network is:

```
plogis(coef(m1)[['edges']])
```

```
## [1] 0.2875817
```

# Lets add an exogenous covariate

- The network visualization showed significant clustering by group, so adding in a covariate based on that variable seems reasonable
- In this case, we do so through the nodematch function
  - nodematch creates an edge level covariate that is one between i and j when they are in the same group and zero otherwise
- This is an example of a homophilous effect

```
m2 = ergm(samplike ~ edges + nodematch('group'))
```

# Lets add an exogenous covariate

```
summary(m2)
```

```
## Call:
## ergm(formula = samplike ~ edges + nodematch("group"))
##
## Maximum Likelihood Results:
##
##              Estimate Std. Error MCMC % z value Pr(>|z|)
## edges           -2.0015      0.2131      0  -9.393  <1e-04 ***
## nodematch.group    2.6481      0.3026      0   8.751  <1e-04 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##      Null Deviance: 424.2  on 306  degrees of freedom
## Residual Deviance: 276.9  on 304  degrees of freedom
##
## AIC: 280.9  BIC: 288.3  (Smaller is better. MC Std. Err. = 0)
```

# Many other ways to add exogenous covariates

**Morris et al. (2008)** detail a variety of ways to create terms for use in `ergm`. Here are some prominent ones:

- `nodecov()`: main effect of a covariate
- `nodeocov()`: main effect of a nodal, sender covariate
- `nodeicov()`: main effect of a nodal, receiver covariate
- `absdiff()`: absolute difference between covariate value for *i* and *j*
- `edgecov()`: main effect of a dyadic covariate



# Add reciprocity term

- Now lets add an endogenous parameter, specifically, reciprocity
- This can be done by using the `mutual` term, which:
  - "adds one network statistic to the model, equaling the number of pairs of actors  $i$  and  $j$  for which  $(i \rightarrow j)$  and  $(j \rightarrow i)$  both exist" (Morris et al. (2008))
- Estimating the effect of network statistics such as these is only possible via pseudolikelihood or MCMC-MLE approaches

```
m3 = ergm(samplike ~ edges + nodematch('group') +  
  mutual  
  )
```

# Add reciprocity term

```
summary(m3)
```

```
## Call:
## ergm(formula = samplike ~ edges + nodematch("group") + mutual)
##
## Monte Carlo Maximum Likelihood Results:
##
##              Estimate Std. Error MCMC % z value Pr(>|z|)
## edges          -2.2733      0.2332      0  -9.749  < 1e-04 ***
## nodematch.group   2.0100      0.3168      0   6.345  < 1e-04 ***
## mutual           1.4404      0.4635      0   3.107  0.00189 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##      Null Deviance: 424.2  on 306  degrees of freedom
## Residual Deviance: 268.0  on 303  degrees of freedom
##
## AIC: 274  BIC: 285.1  (Smaller is better. MC Std. Err. = 0.2697)
```

# Interpretation

The baseline probability of a tie now is:

```
plogis(coef(m3)[['edges']]) #plogis is equivalent to  $\exp(xb)/(1+\exp(xb))$ 
```

```
## [1] 0.09335707
```

But if the reciprocal tie is present even if the two actors are not in the same group, then the log odds of the tie is 3.23x greater:

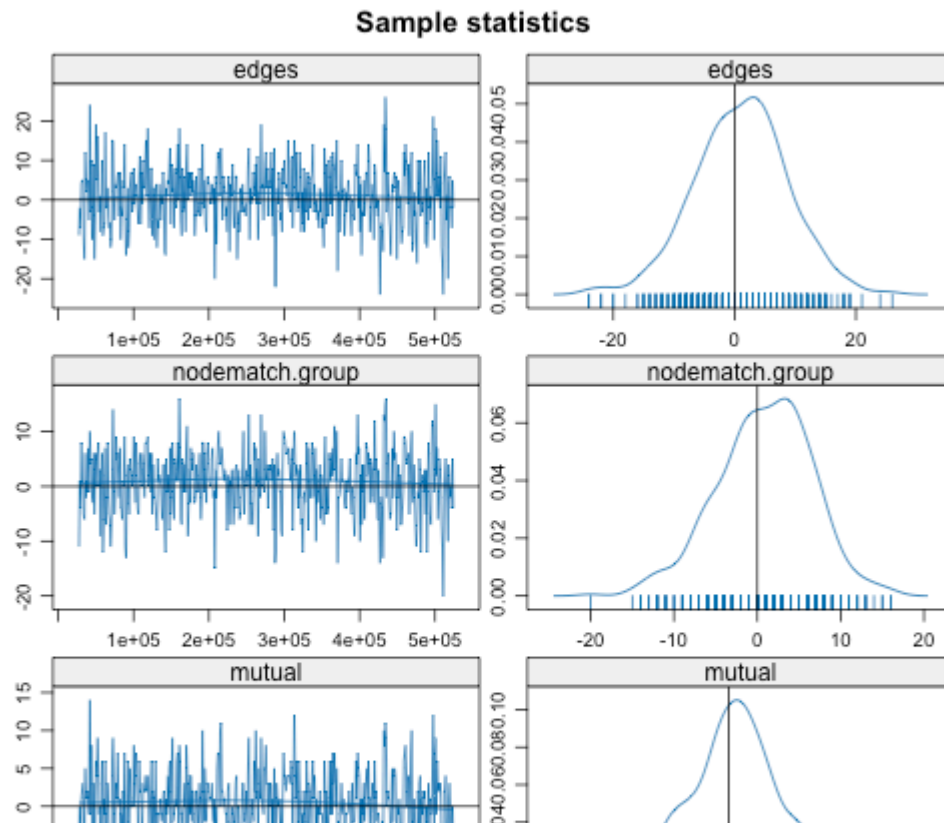
```
plogis(coef(m3)[['edges']] + coef(m3)[['mutual']])
```

```
## [1] 0.3030263
```

# Checking convergence

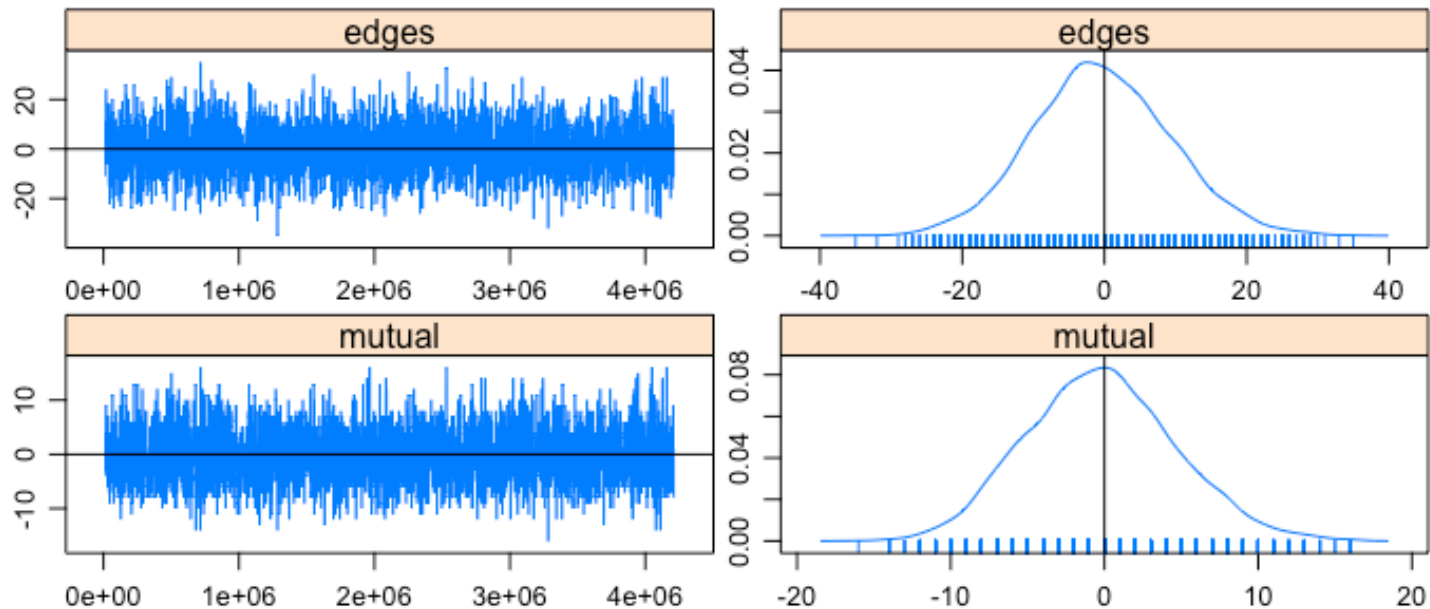
Given our model with a mutual term is estimated via MCMC-, we should always check convergence:

```
mcmc.diagnostics(m3)
```



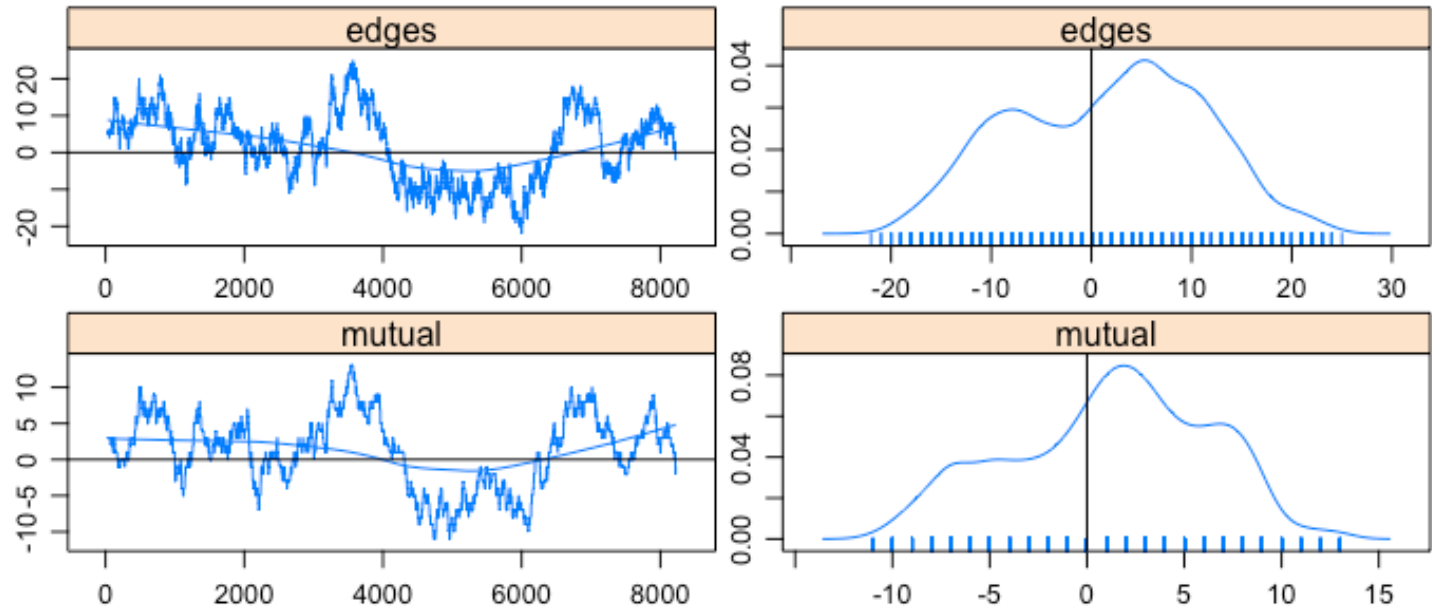
# Checking convergence

Sample statistics



# Example of bad chain

Sample statistics



# What to do?

Increase the number of iterations.

```
?control.ergm  
  
m = ergm(formula, data=data,  
  control=control.ergm(  
    seed=6886,  
    MCMC.samplesize=10000  
  )  
)
```

# Using simulation to gauge fit

- Since ERGMs are generative, given a set of coefficient values, we can simulate networks that are near the maximum likelihood realization of sufficient statistic
- This can be useful for examining fit, among other things, and is easy using `simulate`
- In addition to checking model fit, you can change parameter values, constrain the network in various ways, etc. See `?simulate.ergm` for details.



# Simulation and fit

Ideally, the results closely approximate the visualization of the Sampson network that we presented at the beginning of the application section.

```
set.seed(6886)
simNets = simulate(m3, nsim = 5)

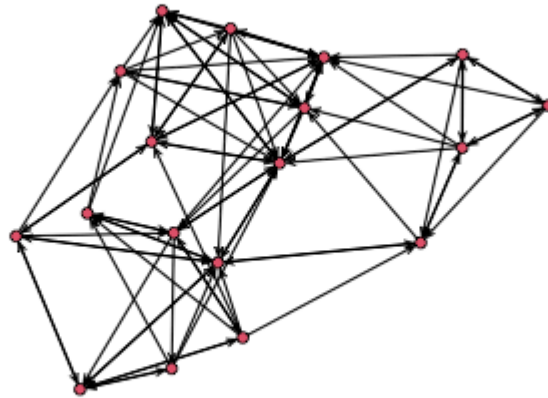
# get layout for actors
p = plot(samplike)

# Define a plotting function:
plotSimNet = function(net, label){
  set.seed(6886)
  plot(net, displaylabels = FALSE,
        vertex.cex = degree(net, cmode = 'indegree')/2, edge.col = "k",
        vertex.col = 'group', coord=p )
  title(label) }

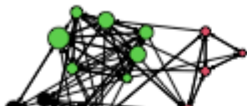
par(mfrow = c(2, 3))

# add actual network to list of sim nets
# for comparison
simNets[[6]] = samplike
labels = c(paste0("sim",1:5), 'actual')
lapply(1:length(simNets), function(i){
```

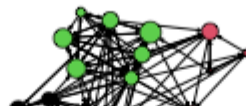
# Simulation and fit



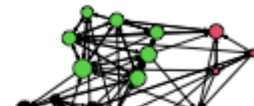
sim1



sim2



sim3



# Running a bunch of simulations

- Running a few simulations and looking at pictures is somewhat useful, but as we generate more and more we can get a more accurate sense of how well our model fits the observed network
- To do this, we can use the `gof` (goodness-of-fit) function
- `gof` simulates networks from the ERGM estimates and, for some set of network statistics, compares the distribution in the simulated networks to the observed values

# Examining model fit

- After running a bunch (exact number can be controlled by the `control.gof.ergm` function) of simulations, we want to use some criteria to compare our simulated models with the observed network
- A standard set of statistics network scholars use to compare how well their model is capturing network dependencies are:
  - `in degree`: Proportion of nodes with the same value of the attribute as the receiving node
  - `out degree`: Proportion of nodes with the same value of the attribute as the sending node
  - `edge-wise shared partners`: Similar to above except this counts the number of dyads with the same number of edges
  - `minimum geodesic distance`: The proportion of pairs of nodes whose shortest connecting path is of length  $k$ , for  $k = 1, 2, \dots$  -- also, pairs of nodes that are not connected are classified as  $k = \infty$

# Examining model fit

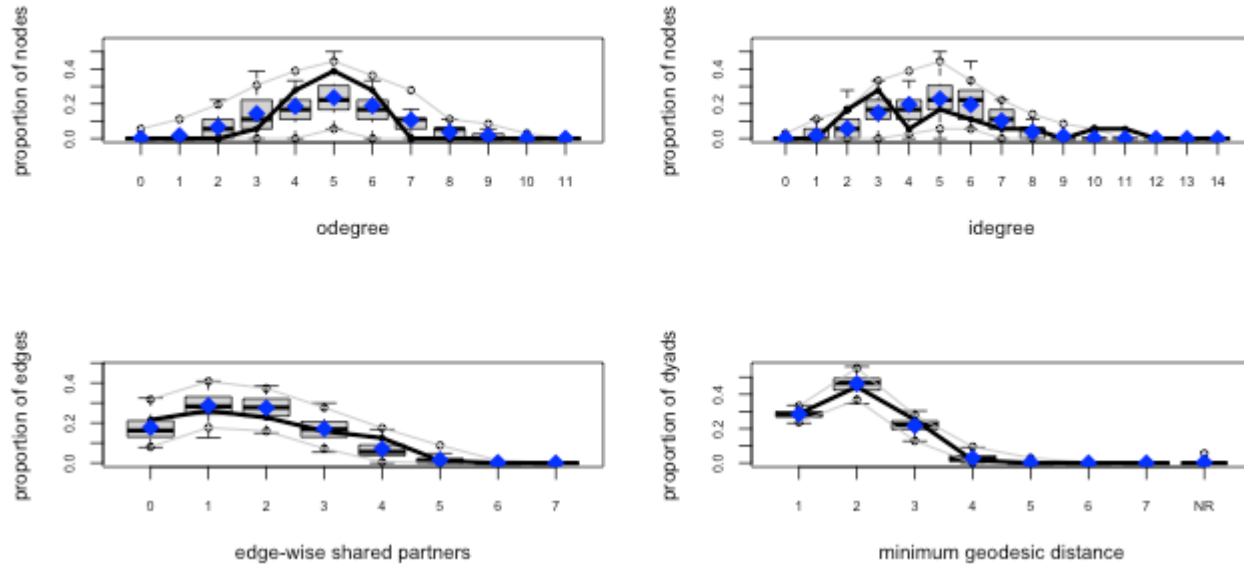
We can generate a comparison along those dimension using the following code:

```
set.seed(6886)
gofM3 = gof(
  m3,
  # specify stats to compare against (- indicates remove)
  GOF=~idegree + odegree + espartners + distance-model
)

# we'll compare against four plots, so set up plotting window
par(mfrow = c(2, 2))
plot(gofM3)
```

# Examining model fit

## Goodness-of-fit diagnostics



# Accounting for the popular monks

- Often in social networks, we find that there can be nodes that play "central" in networks
- To do this we include the `idegree1.5` network statistics
  - this "equals the sum over the actors of each actor's indegree taken to the 3/2 power (or, equivalently, multiplied by its square root)" (`ergm-terms`)

```
m4 = ergm(samplike ~  
  edges + nodematch('group') +  
  mutual + idegree1.5  
)
```

# Accounting for the popular kids

```
summary(m4)
```

```
## Call:
## ergm(formula = samplike ~ edges + nodematch("group") + mutual +
##       idegree1.5)
##
## Monte Carlo Maximum Likelihood Results:
##
##              Estimate Std. Error MCMC % z value Pr(>|z|)
## edges           -4.4593     0.6904      0  -6.459  < 1e-04 ***
## nodematch.group    2.1822     0.3547      0   6.152  < 1e-04 ***
## mutual            1.4129     0.4814      0   2.935 0.003338 **
## idegree1.5         0.6406     0.1824      0   3.512 0.000445 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##      Null Deviance: 424.2  on 306  degrees of freedom
## Residual Deviance: 259.3  on 302  degrees of freedom
##
## AIC: 267.3  BIC: 282.2  (Smaller is better. MC Std. Err. = 0.1951)
```



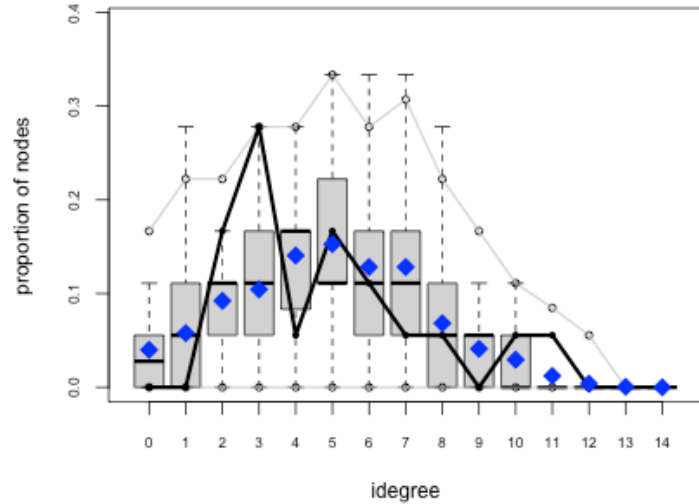
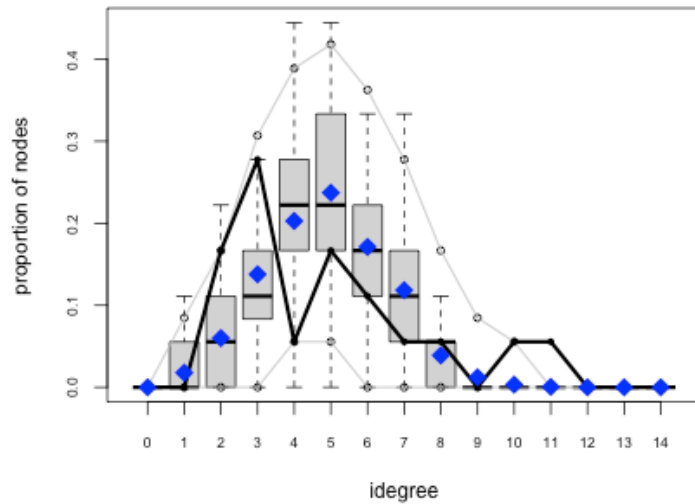
# Did we improve in degree fit?

```
gofM3_indegree = gof(m3, GOF=~indegree-model)
gofM4_indegree = gof(m4, GOF=~indegree-model)

# plot
par(mfrow=c(1,2))
plot(gofM3_indegree)
plot(gofM4_indegree)
```

# Did we improve in degree fit?

Goodness-of-fit diagnostics



# Other model selection methods

We can also use Akaike and Bayesian Information Criteria (AIC and BIC) as criteria for model selection. Accessing these statistics from `ergm` objects is simple:

```
round(sapply(list(m1, m2, m3, m4), AIC), 0)
```

```
## [1] 369 281 274 267
```

```
round(sapply(list(m1, m2, m3, m4), BIC), 0)
```

```
## [1] 373 288 285 282
```

# Triangles

- Triangle terms are a strong motivation for why ERGMs are utilized in social networks
- They are used to answer the question:
  - How does the likelihood of an interaction change if two actors already have an interaction in common?
- Modeling triangles within `ergm` can be done using the `triangles` term. This "adds one statistic to the model equal to the number of triangles in the network" (Morris et al. (2008)):
  - For an undirected network, a triangle is defined to be any set  $\{(i,j), (j,k), (k,i)\}$  of three edges
  - For a directed network, a triangle is defined as any set of three edges  $(i,j)$  and  $(j,k)$  and either  $(k,i)$  or  $(i,k)$

```
m5 = ergm(samplike ~  
  edges + nodematch('group') +  
  mutual + idegree1.5 +  
  triangles  
)
```

# Triangles can be dangerous

- Model degeneracy can often occur with ERGMs, much of the ERGM literature notes that degeneracy is a sign of model misspecification
- Degeneracy here means that the model places a large amount of probability on a small subset of networks that fall in the set of obtainable networks but share little resemblance with the observed network
- Terms like `triangles` can at times produce degenerate graphs
- [Hunter & Handcock 2006](#) introduces a set of geometrically weighted terms to help deal with this issue

# Geometrically Weighted Terms

- The implication of a triangles term is that the likelihood of tie changes proportionately to the number of shared friends two people have
  - Specifically, if having one shared friend makes a tie 25% more likely, having six shared friends makes a tie 150% more likely
- Idea behind geometric terms is to discount each additional tie
- To capture the same effect as `triangles`, we can do this via the geometrically-weighted edgewise shared partners (`gwesp`) term
  - `gwesp` takes a parameter, `decay` that controls how much to discount 2nd, 3rd, etc. shared partners
  - `ergm` will estimate a value for `decay` by default, but most applied scholars fix the `decay` parameter
  - The closer `decay` is to zero, the more dramatic the discounting applied to subsequent shared partners

# Geometrically Weighted Terms

```
m6 = ergm(samplike ~  
  edges + nodematch('group') +  
  mutual + idegree1.5 +  
  gwesp(decay = .5, fixed = TRUE)  
)
```

# Geometrically Weighted Terms

```
summary(m6)
```

```
## Call:
## ergm(formula = samplike ~ edges + nodematch("group") + mutual +
##       idegree1.5 + gwesp(decay = 0.5, fixed = TRUE))
##
## Monte Carlo Maximum Likelihood Results:
##
##              Estimate Std. Error MCMC % z value Pr(>|z|)
## edges              -5.0430     0.7584      0 -6.650 < 1e-04 ***
## nodematch.group       2.9391     0.5204      0  5.648 < 1e-04 ***
## mutual               1.4130     0.5084      0  2.779 0.00545 **
## idegree1.5           1.0307     0.2302      0  4.477 < 1e-04 ***
## gwesp.OTP.fixed.0.5  -0.5490     0.1869      0 -2.938 0.00330 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##      Null Deviance: 424.2  on 306  degrees of freedom
## Residual Deviance: 253.1  on 301  degrees of freedom
##
## AIC: 263.1  BIC: 281.7  (Smaller is better. MC Std. Err. = 0.307)
```