

Exact Statistics and Semi-Parametric Tests for Small Network Data

George G. Vega Yon, MS^{*} Andrew Slaughter, PhD[†] Kayla de la Haye, PhD^{*}

^{*}University of Southern California [†]U.S. Army Research Institute for the
Department of Preventive Medicine Behavioral and Social Sciences

IC²S² 2019, Amsterdam
June 20, 2019

Acknowledgements



This material is based upon work support by, or in part by, the U.S. Army Research Laboratory and the U.S. Army Research Office under grant number W911NF-15-1-0577

Computation for the work described in this paper was supported by the University of Southern California's Center for High-Performance Computing (hpc.usc.edu).



The views expressed in this presentation are those of the authors, and do not represent the official policy or positions of the Department of the Army, the DoD, or the U.S. Government.



Context: Social abilities and team performance

Two research questions

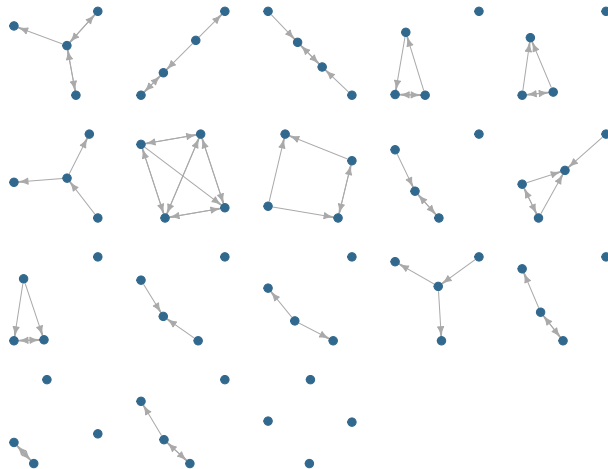
How do **social abilities** impact **network structure**?

How does **collective intelligence** collective intelligence affect team (network)
performance performance?

To answer this question, we have the following experimental data:

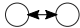
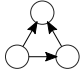

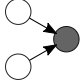
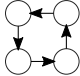
- ▶ 42 mixed-gender teams,
- ▶ Which completed 1 hour of group tasks (MIT's IQ test for teams)
- ▶ Individual survey capturing information regarding socio-demographics **and**:
 - ▶ **Social Intelligence**: Social Perception (measured by RME), Social Accommodation, Social Gregariousness, and Social Awareness
 - ▶ **Social Networks**: Advice Seeking, Leadership, Influence (among others).

Context (cont'd)



We can do a lot of simple statistics: density, % of *[blank]*, etc. but... **how can we go beyond that?**

Exponential random graph models

Representation	Description
	Mutual Ties (Reciprocity) $\sum_{i \neq j} y_{ij} y_{ji}$
	Transitive Triad (Balance) $\sum_{i \neq j \neq k} y_{ij} y_{jk} y_{ik}$
	Homophily $\sum_{i \neq j} y_{ij} \mathbf{1}(x_i = x_j)$
	Covariate Effect for Incoming Ties $\sum_{i \neq j} y_{ij} x_j$
	Four Cycle $\sum_{i \neq j \neq k \neq l} y_{ij} y_{jk} y_{kl} y_{li}$

ERGMs can do the job.

Exponential random graph models (a crash course)

A vector of
model parameters

A vector of
sufficient statistics

$$\Pr(\mathbf{Y} = \mathbf{y} \mid \theta, \mathbf{X}) = \frac{\exp\{\theta^t s(\mathbf{y}, \mathbf{X})\}}{\sum_{\mathbf{y}' \in \mathcal{Y}} \exp\{\theta^t s(\mathbf{y}', \mathbf{X})\}}, \quad \forall \mathbf{y} \in \mathcal{Y}$$

Observed data

The normalizing
constant

All possible
networks

There is one problem with this model ...

$$\Pr(\mathbf{Y} = \mathbf{y} \mid \theta, \mathbf{X}) = \frac{\exp\{\theta^t \mathbf{s}(\mathbf{y}, \mathbf{X})\}}{\sum_{\mathbf{y}' \in \mathcal{Y}} \exp\{\theta^t \mathbf{s}(\mathbf{y}', \mathbf{X})\}}, \quad \forall \mathbf{y} \in \mathcal{Y}$$

A vector of
model parameters A vector of
sufficient statistics

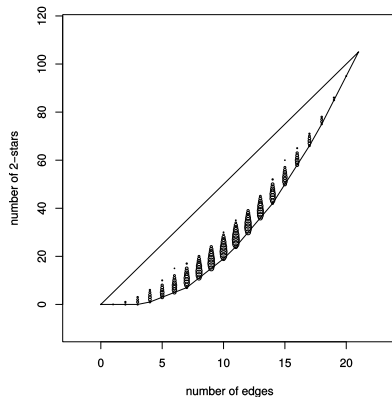
Observed data The normalizing
constant

All possible
networks

because of \mathcal{Y} , the **normalizing constant** is
a summation of $2^{n(n-1)}$ terms 🤯!

Revising model degeneracy and existence of MLE

Following Handcock (2003), the key question is: Where do the sufficient statistics live?



- ▶ In the interior: **Good**, we (possibly) get nice estimates in both MC-MLE and MLE
- ▶ Not in the interior: **We are in trouble**, MLE may not exist

ERGMs for small networks

- Calculating the likelihood function for a directed graph means (at some point) enumerating $2^{n(n-1)}$ **terms**.

$$\Pr(\mathbf{G} = \mathbf{g} \mid \theta, \mathbf{X}) = \frac{\exp\{\theta^t s(\mathbf{g}, \mathbf{X})\}}{\sum_{\mathbf{g}' \in \mathcal{G}} \exp\{\theta^t s(\mathbf{g}', \mathbf{X})\}}$$

- So, if $n = 6$, then we have approx 1,000,000,000 terms 🤯.
- This has lead the field to aim for (very neat) simulation based methods
- But, if our small networks have (at most) 6 nodes...

We can go back to the good-old-fashion MLE

Keeping $n \leq 6$ we can

- ▶ Compute the likelihood function exactly, and hence use ``simple'' optimization to get MLEs.
- ▶ Obtain more **accurate** estimates **faster** (in most cases).
- ▶ Since (usually) small networks come in many...obtain pooled estimates. Which helps with **power and degeneracy**)
- ▶ And more:
 - ▶ All MLE goodies, e.g., LRT
 - ▶ Enhanced simulation methods: resampling, cross-validation
 - ▶ Trivially extend ERGM: mixed-effects models, dependency structures across net
 - ▶ etc.

This and more has been implemented in the `ergmito` (`lifecycle` `experimental`) R package (available at <https://github.com/muriteams/ergmito>)

(built on top of Statnet's amazing `ergm` Hunter et al. (2008); Handcock et al. (2018) R package)

Sidetrack...

ito, ita: From the latin *-ītus*. suffix in Spanish used to denote small or affection.
e.g.:

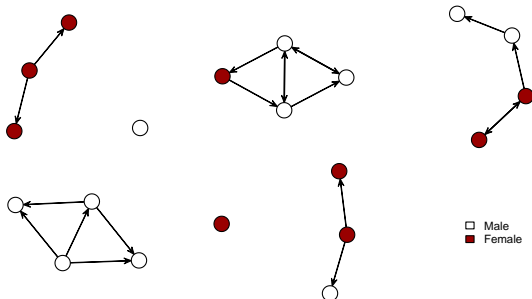
¡Qué lindo ese perrito! / What a beautiful little dog!

¿Me darías una tacita de azúcar? / Would you give me a small cup of sugar?

Special thanks to George Barnett who proposed the name during the 2018 NASN!

Quick example

Suppose that we have 5 networks (as in the R package network)



And we would like to fit a model using the edgecount and number of gender-homophilic ties.

How can we do it?

ergmito example (cont'd)

The same as you would do with the `ergm` package

```
model1 <- ergmito(fivenets ~ edges + nodematch("female"))
```

```
summary(model1) #
```

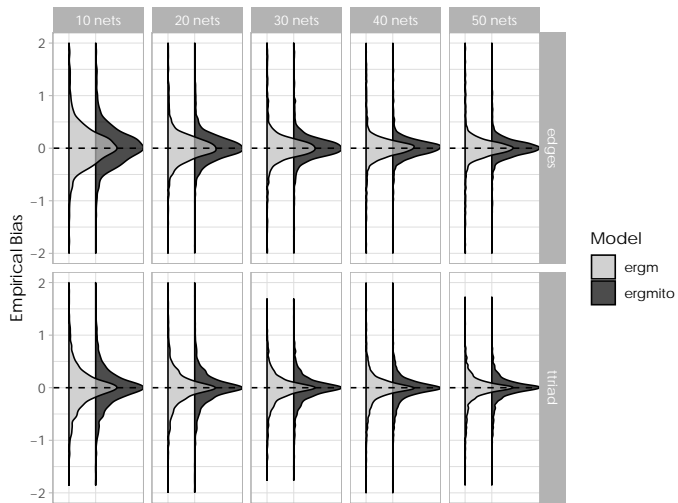
```
##
## ERGMito estimates
##
## formula:  fivenets ~ edges + nodematch("female")
##
##              Estimate Std. Error z value Pr(>|z|)
## edges          -1.70475     0.54356 -3.1363 0.001711 **
## nodematch.female  1.58697     0.64305  2.4679 0.013592 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## AIC: 73.34109    BIC: 77.52978    (Smaller is better.)
```

Go to <https://github.com/muriteams/ergmito> for more on this R package.

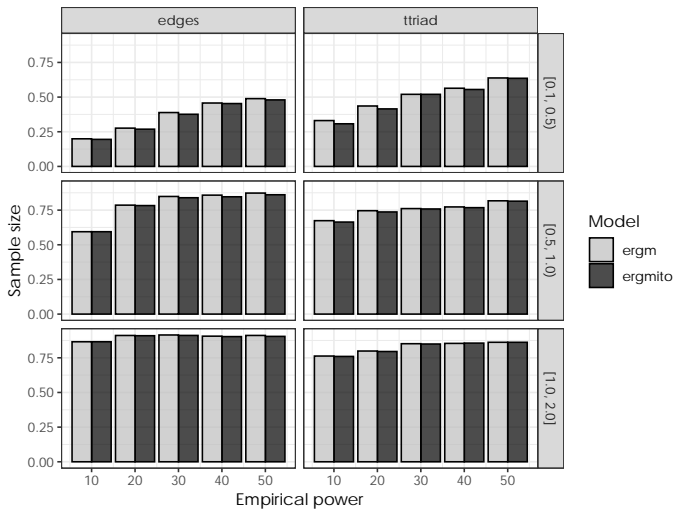
How many networks?

- ▶ Thinking about power and unbiasedness, we did a simulation study
- ▶ Simulated 20,000 samples of networks using the following steps:
 1. Draw parameters for the model based on the terms edges and ttriad (transitive triples) from a $\text{uniform}(-2, 2)$.
 2. Draw group sizes by randomly selecting the number of networks of size 4 and size 5. Each sample has any of $\{10, 20, \dots, 50\}$ networks.
 3. Using 1. and 2., simulate networks using an ERGM model
- ▶ We fitted the models using both MC-MLE (block-diagonal ergm) and MLE (ergmito).
- ▶ We looked (are looking) at empirical bias (sanity check), power and elapsed time.

How many networks? Bias

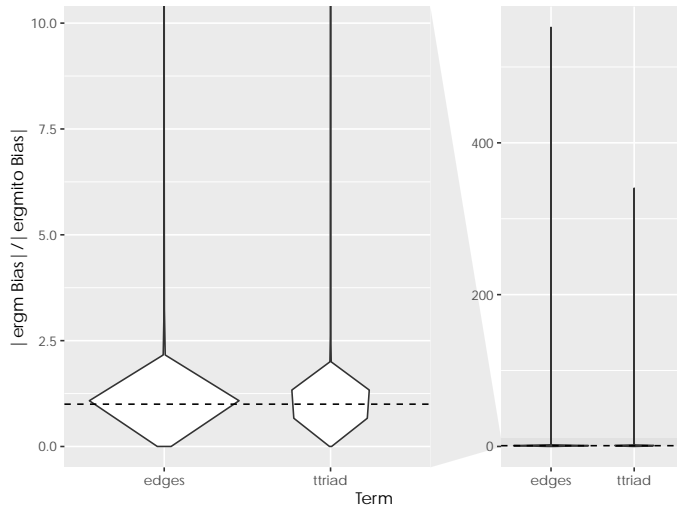


How many networks? Power

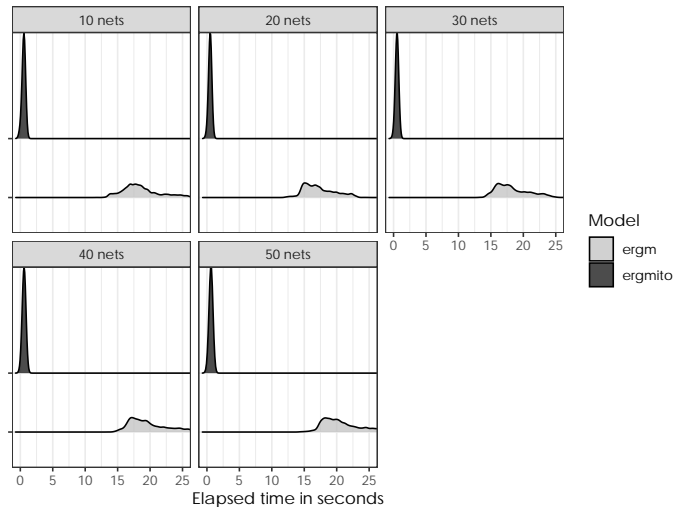


What about a real data set?

How many networks? improvements?



How many networks? improvements? (contd')



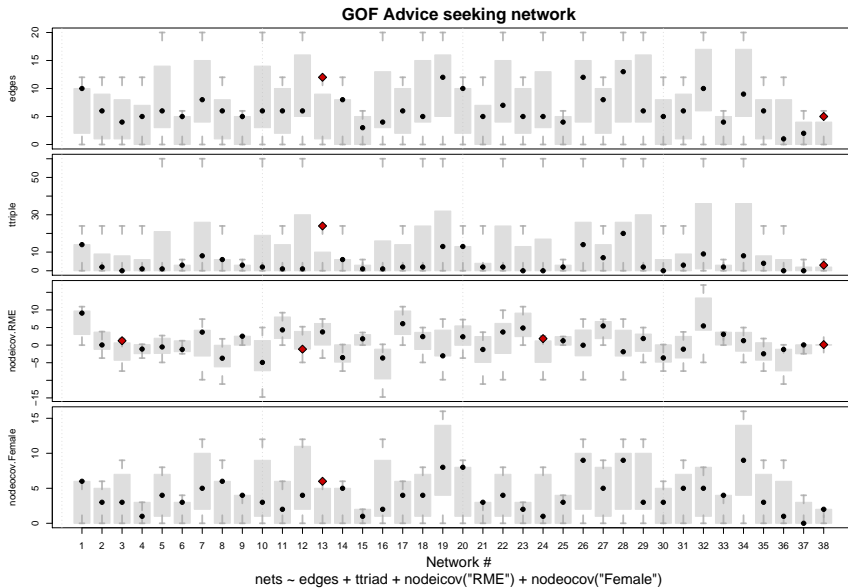
Preliminary results

From our sample of 42 small networks:

	Advice	Dislike	Influence	Leader	Trust
edges	-0.85*** (0.17)	-2.30*** (0.20)	-0.77*** (0.13)	-0.53*** (0.14)	-0.47*** (0.14)
ttriple	0.24*** (0.06)		0.21** (0.08)		0.20*** (0.06)
nodeicov.RME	0.40*** (0.09)		0.21* (0.09)	0.42*** (0.11)	0.25** (0.09)
nodeocov.Female	0.53** (0.18)				
nodematch.Female		0.56* (0.27)			
nodeicov.SI3Fac1		-0.35* (0.15)			
nodeicov.Female				-0.52** (0.20)	
nodeocov.RME				-0.32** (0.11)	
nodeocov.SI3Fac1					0.31*** (0.09)
AIC	695.07	381.72	756.84	637.01	776.82
BIC	712.13	394.52	769.92	654.07	794.25
Log Likelihood	-343.54	-187.86	-375.42	-314.50	-384.41
Num. networks	38	38	41	38	41
Convergence	0	0	0	0	0

*** $p < 0.001$, ** $p < 0.01$, * $p < 0.05$

Table 1: Selected models for each one of the studied networks. Results presented here correspond to a forward selection process.



Context: Social abilities and team performance

Two research questions

~~How do **social abilities** impact **network structure**?~~

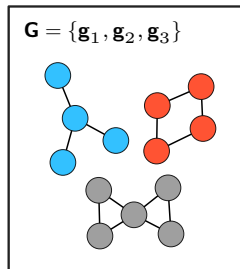
How does **collective intelligence** affect team (network) **performance**?

Networks and team performance

Suppose we have the following:

- ▶ Data on structure, nodes, and an outcome: $(\mathbf{g}, \mathbf{x}, y)$
- ▶ In general, we are interested on assessing the following: $\mathbf{g} \perp y$
- ▶ Three scenarios:
 - a. Direct association
 - b. Known (hypothesized) mediated association (a known confounder)
 - c. Unknown (hypothesized) mediated association (an unknown confounder)

Step 1:
Fit the ERGMito



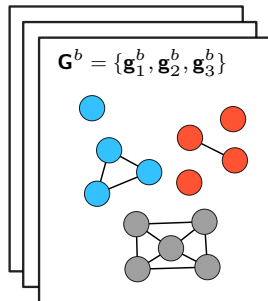
Fit the ERGMito,
This will give us $\mathcal{D}(\hat{\theta}, X_j)$

Step 2:
Calculate $t_0 =$

$$t \left(\begin{bmatrix} \text{blue path} \\ \text{red path} \\ \text{gray cycle} \end{bmatrix}, \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} \right)$$

Throughout the simulations
the only part that changes is
the networks, not Y

Step 3:
For $b \in 1, \dots, B$ do



- 3.1) For $j \in \{1, 2, 3\}$ draw a new network from \mathcal{D}
- 3.2) Use the new sample to calculate $t_b = t(\mathbf{G}^b, Y)$

We are still working (thinking) about this...

Discussion

- ▶ ERGMItos... This is not new. What's new is the set of tools to apply it
- ▶ Taking this approach we can improve our estimates (power) and help with degeneracy
- ▶ The tool is working (according to the simulation study...)
- ▶ Need to conduct more simulations using nodal attributes and compare with ERGM block diagonal models.
- ▶ What about goodness-of-fit? Still need to better think about it

Discussion (contd')

- ▶ The simplicity of the estimation procedure allows us to think of:
 - ▶ Separable Temporal ERGMitos, a.k.a. TERGMitos
 - ▶ Mixture models and Bayesian inference (if you are into that kind of stuff)
 - ▶ More flexible formulas (e.g. interactions between terms and graph-level attributes)
 - ▶ Better odds ratios (not simply exponentiating the coefficients)
 - ▶ Simulation based methods (small size \implies sampling from in-memory data, and exact tests)
 - ▶ Cross-validation/model selection in ERGMs (thank you, Nolan 🙏!)
- ▶ Still thinking about how to test for association between network structure and group outcome

Thanks!

We thank members of our MURI research team, USC's Center for Applied Network Analysis, Garry Robins, Carter Butts, Johan Koskinen, Noshir Contractor, and attendees of the NASN 2018 conference for their comments.



George G. Vega Yon, MS

Andrew Slaughter, PhD

Kayla de la Haye, PhD

vegayon@usc.edu

<https://ggvy.cl>

 gvegayon  gvegayon

References

Handcock, Mark S. 2003. ``Assessing Degeneracy in Statistical Models of Social Networks." Working Paper No. 39 76 (39): 33--50. <https://doi.org/10.1.1.81.5086>.

Handcock, Mark S., David R. Hunter, Carter T. Butts, Steven M. Goodreau, Pavel N. Krivitsky, and Martina Morris. 2018. Ergm: Fit, Simulate and Diagnose Exponential-Family Models for Networks. The Statnet Project (<http://www.statnet.org>). <https://CRAN.R-project.org/package=ergm>.

Hunter, David R., Mark S. Handcock, Carter T. Butts, Steven M. Goodreau, and Martina Morris. 2008. ``Ergm: A Package to Fit, Simulate and Diagnose Exponential-Family Models for Networks." Journal of Statistical Software 24 (3): 1--29.