

Benchmarking Equilibrium Expectation as an Estimation Method for Exponential Random Graph Models

Youhan Cheery

April 12, 2021

1 Abstract

Exponential Random Graph Models (ERGMs) are a family of network models that are useful in modelling complex network structures, such as corporate connections and social relations. ERGMs benefit from many of the properties that make their namesake, the exponential family, such a powerful device in statistical literature.

Equilibrium Expectation, a recent development in the inference of Exponential Random Graph Models proposes an efficient means to compute the parameters of large networks, a typically challenging computational problem. Through leveraging the Markov Chain properties at equilibrium it is able to make infrequent, albeit complex, updates to the parameter values to reduce computation effort required.

In contrast, Stochastic Approximation by Robbins and Monro as well as Monte Carlo methods popularised by **x**, **y**, **etc**, have been staples in the estimation of ERGM parameters for a number of years. As such, this paper will seek to benchmark the Equilibrium Expectation algorithm against Stochastic Approximation and MCMC under **comparison methods**.

2 Introduction

- Mention something about ERGMs and how they're used
- This includes the math side - and about network statistics
- Talk a little about the intractable denominator when it comes to ERGM estimation
- Computational challenges in ERGMs

3 Literature Review

In order to benchmark the Equilibrium Expectation estimation method we start by considering the two classes from which currently popular estimation techniques lie: stochastic approximation and Monte Carlo Maximum Likelihood Estimation.

3.1 Stochastic Approximation

Stochastic Approximation methods were among the first approaches to estimating the ERGM parameters by finding the actual MLE. First introduced by Herbert Robbins in 1951, the stochastic approximation method aims to iteratively find the roots of an optimisation function as represented by an expected value. For a root proposed root $\hat{\theta}$ the Herbert Robbins algorithms states that $\lim_{n \rightarrow \infty} \hat{\theta}_n = \theta$, the true root. This is done by iterating through

$$\theta_{n+1} = \theta_n - \alpha_n (N(\theta) - \alpha) \tag{1}$$

Where

In the context of ERGMs, this is ...

Implemented in 2002 by Snijders...

We specifically note that the these techniques sought the “actual” MLE, as it is worth beginning our discussion with the Maximum Pseudo Likelihood Estimation of ERGM parameters Strauss and Ikeda [2].

Due to the normalising constant, generating the parameters of the ERGM is usually intractable. Strauss and Ikeda proposed an approach where calculating the maximum likelihood function wasn't done on the direct likelihood, rather on a 'pseudo' likelihood.

For a dyad independent network model, that is, a model where every dyad (or pair of ties) is independent from the next - maximum likelihood can be calculated. While this assumption simplifies the mathematics of estimation, it limits the practicality of the network model greatly. For dyad dependent networks, Strauss and Ikeda proposed the pseudo-likelihood as the product of the probabilities of the y_{ij} , with each of the probabilities conditional on the rest of the data Strauss and Ikeda [2].

When not conditioning on data, the ERGM takes the log linear form

$$Pr(G) = \frac{1}{Z(\theta)} e^{\theta' x(G)} \quad (2)$$

where θ is a vector of parameters and $x(G)$ is a vector of graph statistics (on the observed graph).

Conditioning on the rest of the graph produces a model form that now does not depend on θ

$$\begin{aligned} Pr(y_{ij} = 1|C) &= \frac{Pr(G^-)}{Pr(G^-) + Pr(G^+)} \\ \text{logit}(Pr(y_{ij} = 1|C)) &= \theta' x(G^+) - x(G^-) \\ \text{logit}(Pr(y_{ij} = 1|C)) &= \theta' \delta x_{ij} \end{aligned} \quad (3)$$

3.2 Monte Carlo Maximum Likelihood Estimation (MCMLE)

3.3 Techniques for handling Missing Data

4 Notes from papers

4.1 Using contrastive divergence to seed Monte Carlo MLE for exponential-family random graph models (Krivitsky, 2015)

- This paper aims to address the starting point problem of MCMC methods using something known as contrastive divergence which uses a series of abridged MCMC runs instead of running one to stationarity.
- Intractable normalizing constants make life in the ERGM space challenging hence the need for MCMC methods. Sampling techniques need a starting point. Performance and feasibility depend on this value.
- Combined with importance sampling MC MLE to find the initial values
- What are curved models? What is a binary ERGM? (i.e. $h(\theta)$ proportional 1)
- The body of techniques currently used to find the MLE can be split into two broad categories: stochastic approximation (SA) and MCMLE.
- **Stochastic approximation methods**
 - The first methods were these types
 - Given a guess θ^t , these techniques simulate a sample from $ERGM(\theta^t)$ and update
 - We require an initial guess
 - A bad initial guess may result in issues of near-degeneracy, concentrating on the edge of the convex hull of the set of attainable statistics.
 - If all possible values for the parameter θ are not in the realm of Real (q dim) numbers, then MCMC sampling for θ^0 will diverge in the first place
 - Poor choice of θ^0 may induce a very dense network (after simulation), thus resulting in a large impact on computation
 - SA methods are inefficient as each new θ^t requires a burn-in and a sample to estimate $U(\theta^t)$ - **what does this mean? why do we need a new MCMC for each θ^t ?** Note here that U is the score function used to find the MLE.
 - SA will fail if the entirety of the sample lies on the edge of the convex hull.
- **Monte Carlo MLE methods**
 - Uses importance sampling integration to find $\frac{\kappa(\theta')}{\kappa(\theta)}$
 - Given a sample from the ERGM at point θ^t we then update θ^{t+1}
 - There are a number of benefits in using the MCMLE approach - namely that it uses the entire distribution of y^{θ^t} instead of just the first moment, incorporates nonlinear effects on θ to determine the next guess, and automatically determines the step length = less steps to convergence.
 - That said, it is even more susceptible to a bad θ^0 guess.
 - Can also fail for non-curved ERGMs when the convex hull of the simulated statistic does not contain $g(y^{obs})$ (i think this just means the observed graph itself?)
 - Hummel et. al (2012) proposed two ways to address the bad first guess problem. Firstly, assuming $g(Y)$ is log-normal, then $\exp[\eta(\theta') - n(\theta)^T g(Y)]$ is lognormal and we can approximate the expectation such that the maximiser depends only on the first two moments of y^{θ^t} which has a closed form solution for non-curved ERGMs. The second fix is a partial stepping technique that shifts the observed statistic to the centroid of the simulated statistic, reducing the step but preserving its direction. This approach *survives* poor starting values but does not make immune to them.

- There are two approaches to finding good values for y^{θ^0} . One is the maximum pseudo/composite likelihood estimation (MPLE/MCLE) and the second is contrastive divergence.
- **Composite likelihood methods**
 - Before simulation methods became popularised we had pseudolikelihood estimators (Strauss + Ikeda, 1974). Here the likelihood was approximated and $y_{i,j}$ is an indicator for the presence of ties. This led to a score function that is a nonlinear logistic regression.
 - MPLE can be quite different from the MLE and it's mainly used to initialise them (what is meant by "them"?)
 - MPLE/MCLE result in a multinomial model to enumerate the set of possible pairs of toggles - thus adding more burden to the modeller.
- **Contrastive divergence**
 - Hinton (2002) suggested we don't run the MCMC simulation to convergence, but instead make parallel updates starting at the observed data and calculating the gradient based on that.
- **Questions for Pavel**
 - Emphasis in this research on methods that do well for θ^0 ? "Poor choice may result in a dense network" = bad simulation? Are some methods more sensitive to initial values e.g. importance sampling vs. Robins Munro
 - Often in the CD paper we refer to curved vs. non-curved ERGMs - what is the difference and is this something that I should be separately learning about? When you have an exponential family model you have exponential form - you can transform theta before taking the dot product. Now you have model parameter and the mapping mapped to a vector it needs to be the same length of the statistics. size of theta ; g = strict definition of curved family. From the point of view of computation, whether the theta is shorter or not/transformation. Not curved ones are special cases of curved ones. Helpful for the algorithm to accommodate some transformations of theta
 - In what instances can a model only be done via MCMC simulation and not the MPLE/MCLE methods?
 - Confusing result on page 7 regarding contrastive divergence as discovered by Asuncion et al. (2010).
 - What is CD_∞ Confirm my understanding: CD_1 is the MPLE equivalent (it finds the MLE under the pseudo method), while CD_∞ is the actual MLE?
 - The three approaches to finding θ^0
- **Papers off the back of this I think I should read**
 - Hinton (2002)
 - Wang et. al. (2014)
 - Hummel (2011) perhaps?
 - Hummel et. al. (2012) (partial stepping techniques)

How well do different algorithms behave if say you start with a vector of 0s for theta, or try different configurations for the starting values

4.2 Fast Maximum Likelihood Estimation via Equilibrium Expectation for Large Network Data

- This paper proposes a fast algorithm for MLE which allows for larger network estimation. It leverages properties of the Markov chain at equilibrium.
- Existing approaches for MLE such as MCMCMLE, MoM, and Bayesian estimation use iterative algorithms that successively update θ until the expectation of the network under the statistics is equivalent to the statistic of the observed network. To do this, MCMC draws a large number of simulated networks for various values of θ , $x_s(\theta)$. The simulated network x_s is a network drawn from probability distribution $\pi(x, \theta)$.
- Each time a new simulated network is drawn, we need to satisfy the convergence criteria $\Sigma_{x,x'} \pi(x, \theta) P(x \rightarrow x', \theta) (z_A(x') - z_A(x))$. This is very expensive computationally.
- Consider the equation $E_{\pi(\theta)}(\Delta_{Z_A}(x, \theta)) = 0$. This suggests that if network x is drawn from probability distribution $\pi(x, \theta)$ then the expected value of $\Delta_{Z_A}(x, \theta)$ is 0. Only valid when at the limiting distribution.
- $\Delta_{Z_A}(x, \theta)$ is found using Monte Carlo integration. Assuming that are s_i Monte Carlo sample of networks that are i.i.d from $\pi(x, \theta)$ then the expectation can be calculated. This is obviously not a realistic assumption which will later be removed.
- I.e. $E_{\pi(\theta)}(\Delta_{Z_A}(x, \theta)) = \frac{1}{n} \Sigma_i \Delta_{Z_A}(x_{s_i}, \theta)$
- With enough sample of networks we can efficiently compute the LHS and then solve with respect to θ .

- No burn-in time using MC integration so it's much more efficient.
- Ergodicity of systems suggests that if the network x_s is very large then the true θ^* may be estimated from $f_A(x_{obs}, \theta)$, thus dropping the need to sum. This leads to the result $\Delta Z_A(x_s, \theta^{EE}) = 0$.
- If the network x_s is very large then $\theta^{EE} = \theta^*$
- **However in reality we only have ONE x_{obs}**
- In reality - this looks like Contrastive Divergence (CD) as applied to ERGM parameter estimation (instead of for finding θ^0)
- So what's done instead is that if a solution to the method of moments exists then we can draw an x_s from $\pi(x, \theta)$ such that we satisfy $z_A(x_s) = z_A(x_{obs})$. The network x_s can be drawn from an MCMC simulation. θ is updated iteratively until we have a solution for $Z_A(x_s, \theta^{EE}) = 0$. It employs Metropolis-Hastings
- Since we are not drawing many simulated networks from various θ the algorithm is considerably faster.
- To start the EE algorithm, we use CD1 estimate as the starting point. EE is more sensitive to K_A than it is to θ_0 . K_A is subtracted from the old θ to produce the next θ . This value is used to make sure that θ_A increases when dz_A is negative and decrease when dz_A is positive. We also require that the exponential family expectation is a monotonically increasing function in θ
- **Questions for Pavel**
 - In reality we only have one x_{obs} - but why can't we simulate more networks based on a starting θ and use those?

4.3 Markov Chain Monte Carlo Estimation of Exponential Random Graph Models

- This paper is about the simulation and MCMC estimation of exponential random graph models.
- A major problem with ERGMs is that for certain parameter values they can have bimodal (or multimodal) distributions for the sufficient statistics such as the number of ties
- The modality is reflected in the outcome space being divided into two or more regions whereby the usual MCMC approaches have long sojourn times, with a negligible probability of moving from one region to the next. This leads to slow convergence
- MCMC algorithms must be able to make transitions from a given graph to a very different graph. It is proposed to include transitions to the graph complement as updating steps to improve the speed of convergence to the target distribution
- Recall that a random graph, according to Strauss and Frank (1986) is a Markov graph if there are finite numbers of nodes and if edges between disjoint pairs are independent conditional on the rest of the graph
- Frank and Strauss proposed a simulation based approach to approximate the MLE of any one of the three θ_k (talking here about the standard p^* triad model), given that the other two are fixed at 0. They also proposed a logistic regression method to estimate the full θ (pseudo MLE).
- While the pseudo MLE approach is nice computationally and intuitively appealing, the properties of the resulting estimator for ERGMs are unknown.
- In addition, the pseudo MLE is not a function of the entire sufficient statistic $u(Y)$ which implies that it is not an admissible estimator for squared error loss function (lehmann, 1983)
- Many of the older MCMC estimation methods relied on MC simulations of the Markov graph at current parameter values. However, simulation algorithms for the ERGM methods can suffer from severe convergence problems not really outlined previously
- For many models (including triad), there is a large region in which the demarcation between the subset of parameters θ leads to graphs with relatively low expected densities, and the subset of parameters θ that lead to graphs with high expected densities; is quite sharp. Parameters in or near the demarcation zone can have bimodal density.
- With more nodes, the demarcation becomes more marked. (demarcation = fixing the boundary)
- One MCMC process: assuming you are at $\theta^{(n)}$, a MC simulation of the Markov graph is made which is used to estimate the moments of the distribution, which are then used to make an expansion approximating $\mu(\theta)$ for θ in the neighbourhood of $\theta^{(n)}$. This is then used to solve the moment equation to provide an estimate for $\theta^{(n+1)}$
- This paper uses a version of the Robbins-Monro algorithm, which is itself considered by some a MC variant of the Newton-Raphson algorithm.
- Robbins-Monro is a stochastic iterative algorithm intended to solve equations of the form $E\{Z_\theta\} = 0$. In the ERGM case, consider $Z_\theta = u(Y) - u_0$ where $u_0 = u(y)$ is the observed value of the sufficient statistic.

- The iteration step in the Robbins-Monro algorithm with step-size a_n is $\hat{\theta}^{(n+1)} = \hat{\theta}^{(n)} - a_n D_n^{-1} Z(n)$ where $Z(n)$ are random variables such that the conditional distribution of $Z(n)$ given $Z(1), \dots, Z(n-1)$ is the distribution of $Z(\theta)$ obtained for $\theta = \hat{\theta}^{(n)}$
- The stochastic process $\hat{\theta}^{(n)}$ generated by this MC simulation is a Markov chain, because for any n_0 the sequence $\{\hat{\theta}^{(n)}\}$ for $n > n_0$ depends on the past value via the last value, $\hat{\theta}^{(n_0)}$
- **Questions for Pavel**
 - There is a lot of mention of the triad model in these older papers. With simulation based methods, how often are they impervious to model specification?
 - Is simulating an ERGM a different problem to estimating an ERGM?

5 Meeting 24/03

- **Meeting preparation**
 - Papers read in the last fortnight: New specifications for ERGM models (Snijders et. al.), inference for curved exponential family random models (Hunter and Handcock), Modelling Social Networks from Sampled Data (Handcock and Gile)
 - A lot of material in the themes of simulation, inference and missing data. Need to start boiling down all of this raw material into a coherent thesis that I can start writing.
 - High level I am doing a bench marking of the equilibrium expectation algorithm - can we incorporate missing data effectiveness into the tests? Are there other algorithms we want to benchmark?
 - Any meta analysis papers that are worth looking into to gain some kind of inspiration for benchmarking various inference algorithms? What are the types of things people look for including speed, scalability, robustness, etc
 - Different types of error: are we only looking at MCMC error and MLE error? In the case of curved ERGMs are there any other errors we are interested in testing? What about finding the expectation and variance in curved ERGMs?
 - On the topic of curved ERGM inference - I understand the traditional Thompson and Geyer approach to MCMC. Generate initial guess, approximate the difference between current and initial by a MC integration, then approximate r and continue. But now are we talking iterative methods like Newton-Raphson?
 - Is there any relation with this work to simulation of ERGMs?
- **Meeting notes**
 - MCMC and Robbins Monro are the two basis methods for fitting these models, MCMCMLE, while Robbins Monro draws small samples and downweights them... These are competitors to equilibrium expectation
 - MCMC errors are more obvious, versus Robbins Monro not really having a good error approximator. The same goes for equilibrium expectation.
 - Okabayashi and Geyer
 - Missing data you can still calculate the likelihood - two normalising constants on top of each other. Each parameter updates requires you to run MCMC twice.
 - For Robbins Monro - there's a question of how do you run that... two MCs running at once? one for constrained and one for unconstrained etc... implementing this is a stretch goal
 - Reach goal for equilibrium expectation handling data
 - Go to `ergm.R` - check the `mainfit j- switch(...)`. This is where the equilibrium expectation algorithm will go
 - Method to prefix in `control.ergm.R` which contains the parameters
 - EE is infrequent complex updates so it should probably be written in C
 - Go to `SEXP_DISPATCH_MCMCPhase12` as an example for the implementation of the algorithm
 - Taking the outputs of MCMC and using them
 - Check out `MCMC.c.template.do_not_include_directly.h`
 - Testing a hybrid method - starting with Robbins Monro and then going into MCMCMLE

References

- [1] S. M. Herbert Robbins. A stochastic approximation method. *The Annals of Mathematical Statistics*, 1951.
- [2] D. Strauss and M. Ikeda. Pseudolikelihood estimation for social networks. *Journal of the American Statistical Association*, 1990.