



UNIVERSITY
OF WOLLONGONG
AUSTRALIA

University of Wollongong
Research Online

National Institute for Applied Statistics Research
Australia Working Paper Series

Faculty of Engineering and Information Sciences

2015

Using contrastive divergence to seed Monte Carlo MLE for exponential-family random graph models

Pavel Krivitsky
University of Wollongong

Recommended Citation

Krivitsky, Pavel, Using contrastive divergence to seed Monte Carlo MLE for exponential-family random graph models, National Institute for Applied Statistics Research Australia, University of Wollongong, Working Paper 11-15, 2015, 26.
<http://ro.uow.edu.au/niasrawp/28>

Research Online is the open access institutional repository for the University of Wollongong. For further information contact the UOW Library:
research-pubs@uow.edu.au

Using contrastive divergence to seed Monte Carlo MLE for exponential-family random graph models

Abstract

Exponential-family models for dependent data have applications in a wide variety of areas, but the dependence often results in an intractable likelihood, requiring either analytic approximation or MCMC-based techniques to fit, the latter requiring an initial parameter configuration to seed their simulations. A poor value can lead to slow convergence or outright failure. The approximate techniques that could be used to seed them tend not to be as general as the simulation-based, and require implementation separate from that of the MLE-finding algorithm.

Contrastive divergence is a more recent simulation-based approximation technique that uses a series of abridged MCMC runs instead of running them to stationarity. We combine it with the importance sampling Monte Carlo MLE for a general method to obtain adequate initial values the MLE-finding techniques, describe and extend it to a wide variety of modeling scenarios, and address practical issues such as stopping criteria and selection of tuning parameters.

Our approach reuses the aspects of an MLE implementation that are model-specific, so little to no additional implementer effort is required to obtain adequate initial parameters. We demonstrate this on a series of network datasets and models drawn from ERGM computation literature.

NIASRA

NATIONAL INSTITUTE FOR APPLIED
STATISTICS RESEARCH AUSTRALIA



***National Institute for Applied Statistics Research
Australia***

University of Wollongong

Working Paper

11-15

**Using Contrastive Divergence to Seed Monte Carlo MLE for
Exponential-Family Random Graph Models**

Pavel N. Krivitsky

*Copyright © 2015 by the National Institute for Applied Statistics Research Australia, UOW.
Work in progress, no part of this paper may be reproduced without permission from the Institute.*

National Institute for Applied Statistics Research Australia, University of Wollongong,
Wollongong NSW 2522. Phone +61 2 4221 5435, Fax +61 2 4221 4845. Email:
anica@uow.edu.au

Using Contrastive Divergence to Seed Monte Carlo MLE for Exponential-Family Random Graph Models

Pavel N. Krivitsky¹

School of Mathematics and Applied Statistics and
National Institute for Applied Statistics Research Australia,
University of Wollongong

¹Calculations were performed on the National Institute for Applied Statistics Research Australia (NIASRA) High Performance Computing cluster.

Abstract

Exponential-family models for dependent data have applications in a wide variety of areas, but the dependence often results in an intractable likelihood, requiring either analytic approximation or MCMC-based techniques to fit, the latter requiring an initial parameter configuration to seed their simulations. A poor value can lead to slow convergence or outright failure. The approximate techniques that could be used to seed them tend not to be as general as the simulation-based, and require implementation separate from that of the MLE-finding algorithm.

Contrastive divergence is a more recent simulation-based approximation technique that uses a series of abridged MCMC runs instead of running them to stationarity. We combine it with the importance sampling Monte Carlo MLE for a general method to obtain adequate initial values the MLE-finding techniques, describe and extend it to a wide variety of modeling scenarios, and address practical issues such as stopping criteria and selection of tuning parameters.

Our approach reuses the aspects of an MLE implementation that are model-specific, so little to no additional implementer effort is required to obtain adequate initial parameters. We demonstrate this on a series of network datasets and models drawn from ERGM computation literature.

Keywords: curved exponential family; ERGM; network data; partial stepping

1 Introduction

Exponential family models for dependent data have found applications in point processes, social networks, statistical physics, and image analysis alike, but this dependence often produces likelihoods with intractable normalizing constants. A variety of techniques—frequentist and Bayesian—have been proposed for their estimation. Although some approximations are available, the exact techniques invariably require a starting parameter configuration $\boldsymbol{\theta}^0$, their performance and even feasibility depending on this value.

In this work, we focus on the problem of a general way of obtaining a good $\boldsymbol{\theta}^0$ with minimal additional implementer effort, particularly for the application of these models to modeling of social networks—the exponential-family random graph models (ERGMs) (Wasserman and Pattison, 1996), as extended to curved families by Snijders et al. (2006) and Hunter and Handcock (2006) and to networks with valued ties by Robins et al. (1999) and Krivitsky (2012), we consider the broad class of models defined as follows. Given a set $N = \{1, 2, \dots, n\}$ of actors of interest, let $\mathbb{Y} \subseteq N \times N$ be the set of potential relationships among them (usually a proper subset, excluding self-loops or if only ties among specific subsets of actors are of interest). Then, with \mathbb{S} being the set of relationship values (which could be simply $\{0, 1\}$ for binary networks), we define the sample space of mappings $\mathcal{Y} \subseteq \mathbb{S}^{\mathbb{Y}}$ (again, sometimes a proper subset if, say, we wish to constrain the network to have a specific number of ties or a specific degree distribution).

Then, $\mathbf{Y} \sim \text{ERGM}_{\mathcal{Y}, h, \boldsymbol{\eta}, \mathbf{g}}(\boldsymbol{\theta})$ if

$$\Pr_{\mathcal{Y}, h, \boldsymbol{\eta}, \mathbf{g}}(\mathbf{Y} = \mathbf{y}; \boldsymbol{\theta}) \stackrel{\text{def}}{=} \frac{h(\mathbf{y}) \exp\{\boldsymbol{\eta}(\boldsymbol{\theta})^\top \mathbf{g}(\mathbf{y})\}}{\kappa_{\mathcal{Y}, h, \boldsymbol{\eta}, \mathbf{g}}(\boldsymbol{\theta})}, \quad \mathbf{y} \in \mathcal{Y} :$$

an exponential family over a sample space \mathcal{Y} of networks (potentially with valued ties), parametrized by a q -vector $\boldsymbol{\theta}$, and specified by a reference measure $h(\mathbf{y})$ (with $h(\mathbf{y}) \propto 1$ being typical for binary ERGMs), a mapping $\boldsymbol{\eta}$ from $\boldsymbol{\theta}$ to the p -vector of canonical parameters (and in non-curved ERGMs, $\boldsymbol{\eta}(\boldsymbol{\theta}) \equiv \boldsymbol{\theta}$ with $p \equiv q$), and a sufficient statistic p -vector \mathbf{g} . The normalizing constant $\kappa_{\mathcal{Y}, h, \boldsymbol{\eta}, \mathbf{g}}(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \sum_{\mathbf{y}' \in \mathcal{Y}} h(\mathbf{y}') \exp\{\boldsymbol{\eta}(\boldsymbol{\theta})^\top \mathbf{g}(\mathbf{y}')\}$, is often intractable for models that seek to reproduce more complex social effects, such as triadic closure. It also identifies the *natural parameter space* of the model, $\Theta_N \stackrel{\text{def}}{=} \{\boldsymbol{\theta} : \kappa_{\mathcal{Y}, h, \boldsymbol{\eta}, \mathbf{g}}(\boldsymbol{\theta}) < \infty\}$, which equals \mathbb{R}^q for binary ERGMs, but which may be far more complex for valued ERGMs, such as if geometric or Conway–Maxwell–Poisson (CMP) distribution (Shmueli et al., 2005) is used for social

interaction counts (Krivitsky, 2012). Unless it is relevant to the discussion, we will, generally, omit “ $\mathcal{Y}, h, \boldsymbol{\eta}, \mathbf{g}$ ” from the subscript.

Given an observed network, \mathbf{y}^{obs} , it is desired to find the MLE, $\hat{\boldsymbol{\theta}} \stackrel{\text{def}}{=} \arg \max_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta}} \log \Pr(\mathbf{Y} = \mathbf{y}^{\text{obs}}; \boldsymbol{\theta})$, or, equivalently, to solve the score estimating equation,

$$\mathbf{U}(\hat{\boldsymbol{\theta}}) \stackrel{\text{def}}{=} \nabla_{\boldsymbol{\theta}} \ell(\hat{\boldsymbol{\theta}}) = \boldsymbol{\eta}'(\hat{\boldsymbol{\theta}})^\top [\mathbf{g}(\mathbf{y}^{\text{obs}}) - \mathbb{E}\{\mathbf{g}(\mathbf{Y}); \hat{\boldsymbol{\theta}}\}] = -\boldsymbol{\eta}'(\hat{\boldsymbol{\theta}})^\top \mathbb{E}\{\mathbf{z}(\mathbf{Y}); \hat{\boldsymbol{\theta}}\} = \mathbf{0}, \quad (1)$$

(Hunter and Handcock, 2006, eq. 3.1), where $\boldsymbol{\eta}'(\cdot) \stackrel{\text{def}}{=} \nabla_{\boldsymbol{\theta}} \boldsymbol{\eta}(\cdot)$, $\mathbb{E}(\cdot; \cdot)$ denotes the expectation under the model and parameter configuration in question, and $\mathbf{z}(\mathbf{y}) \stackrel{\text{def}}{=} \mathbf{g}(\mathbf{y}) - \mathbf{g}(\mathbf{y}^{\text{obs}})$.

We use $\vec{\mathbf{y}}$ as shorthand for a sample or series of networks $\mathbf{y}^1, \dots, \mathbf{y}^S$, with $\vec{\mathbf{y}}^\theta$ in particular being a sample from $\text{ERGM}(\boldsymbol{\theta})$, and we use $\mathbf{g}(\vec{\mathbf{y}})$ for a $p \times S$ matrix with sth column containing $\mathbf{g}(\mathbf{y}^s)$, with $\bar{\mathbf{g}}(\vec{\mathbf{y}}) \stackrel{\text{def}}{=} \mathbf{g}(\vec{\mathbf{y}}) \mathbf{1}_S / S$, and, analogously $\mathbf{z}(\vec{\mathbf{y}})$ and $\bar{\mathbf{z}}(\vec{\mathbf{y}})$; and we define $\mathbf{U}_{\vec{\mathbf{y}}^\theta}(\boldsymbol{\theta}) \stackrel{\text{def}}{=} -\boldsymbol{\eta}'(\boldsymbol{\theta})^\top \mathbf{z}(\vec{\mathbf{y}}^\theta)$, a $q \times S$ matrix whose sth column is the contribution to (1) from \mathbf{y}^s , so that $\bar{\mathbf{U}}_{\vec{\mathbf{y}}^\theta}(\boldsymbol{\theta})$ is the sample estimate of $\nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta})$. We also use the sample variance of a statistic $\mathbf{t}(\vec{\mathbf{y}})$,

$$\widetilde{\text{Var}}\{\mathbf{t}(\vec{\mathbf{y}})\} \stackrel{\text{def}}{=} \frac{1}{S-1} \sum_{s=1}^S \{\mathbf{t}(\mathbf{y}^s) - \bar{\mathbf{t}}(\vec{\mathbf{y}})\} \{\mathbf{t}(\mathbf{y}^s) - \bar{\mathbf{t}}(\vec{\mathbf{y}})\}^\top.$$

A body of literature exists on computational methods for finding $\hat{\boldsymbol{\theta}}$ given a starting configuration $\boldsymbol{\theta}^0$; and on approximate techniques suitable for finding such a configuration.

1.1 Techniques for finding the MLE

The currently popular MLE techniques can be broadly classified into two categories: stochastic approximation (SA) and Monte Carlo Maximum Likelihood Estimation (MCMLE). We review them in turn.

1.1.1 Stochastic Approximation Methods

Stochastic approximation methods represented the first attempts to find the actual MLE for ERGMs, starting with Snijders (2002) application of Robbins and Monro (1951) and similar algorithms, and, later, refinements such as

those of Okabayashi and Geyer (2012). Given a guess $\boldsymbol{\theta}^t$, these techniques simulate a sample $\bar{\mathbf{y}}^{\boldsymbol{\theta}^t} = (\mathbf{y}^{\boldsymbol{\theta}^t, 1}, \dots, \mathbf{y}^{\boldsymbol{\theta}^t, S})$ from $\text{ERGM}(\boldsymbol{\theta}^t)$ and updates the guess to

$$\boldsymbol{\theta}^{t+1} \stackrel{\text{def}}{=} \boldsymbol{\theta}^t - \alpha_t \bar{\mathbf{U}}_{\bar{\mathbf{y}}^{\boldsymbol{\theta}^t}}(\boldsymbol{\theta}^t),$$

for α_t a scalar or a $q \times q$ matrix that is decreasing in t .¹ (Robbins–Monro implementation as used by Snijders (2002) and the PNet software suite for ERGM inference (Wang et al., 2014) uses a scalar multiple of the inverse of the diagonal of $\widehat{\text{Var}}(\bar{\mathbf{y}}^{\boldsymbol{\theta}^t})$ in particular.)

Methods of this type require an initial guess, $\boldsymbol{\theta}^0$. In the context of network models in particular, a poor initial guess may induce a near-degenerate distribution concentrated on the edge of the convex hull of the set of attainable statistics $\text{Conv}(\{\mathbf{g}(\mathbf{y}') : \mathbf{y}' \in \mathcal{Y}\})$ (often an empty network or a complete graph). (Rinaldo et al., 2009; Hunter et al., 2012, and others) While $\mathbf{U}(\boldsymbol{\theta}^0)$ itself may not be on the edge of this convex hull, its sample value $\mathbf{U}_{\bar{\mathbf{y}}^{\boldsymbol{\theta}^0}}(\boldsymbol{\theta}^0)$ could very well be, leaving the gradient-based methods without an unambiguous direction of ascent. And, if $\Theta_{\text{N}} \neq \mathbb{R}^q$, MCMC sampling for $\boldsymbol{\theta}^0 \notin \Theta_{\text{N}}$ will diverge in the first place, and locating a $\boldsymbol{\theta}^0 \in \Theta_{\text{N}}$ may itself be a challenge. (Krivitsky, 2012)

Choice of $\boldsymbol{\theta}^0$ can affect estimation in other ways as well: while one can represent a network \mathbf{y} as an $n \times n$ matrix of relationship values, most large networks studied tend to be sparse, and sparse matrix representations are used as a result. Then, storing and processing a network with more ties is more costly in both memory and time, and if a poor choice of $\boldsymbol{\theta}^0$ induces very dense networks, computation can be slowed down severely or fail.

SA algorithms also tend to be relatively computationally inefficient: every new guess $\boldsymbol{\theta}^t$ requires a burn-in period and a sample to estimate $\mathbf{U}(\boldsymbol{\theta}^t)$, and optimal length of each step is unknown, so relatively many such steps are typically required.

1.1.2 Monte Carlo Maximum Likelihood Estimation

Introduced by Geyer and Thompson (1992), and applied to curved ERGMs by Hunter and Handcock (2006), MCMLE draws on importance sampling

¹The gradient methods cited are all specified for non-curved ERGMs, but this is a direct extension.

integration, observing that

$$\begin{aligned} \frac{\kappa(\boldsymbol{\theta}')}{\kappa(\boldsymbol{\theta})} &= \sum_{\mathbf{y} \in \mathcal{Y}} \exp[\{\boldsymbol{\eta}(\boldsymbol{\theta}') - \boldsymbol{\eta}(\boldsymbol{\theta})\}^\top \mathbf{g}(\mathbf{y})] \frac{h(\mathbf{y}) \exp\{\boldsymbol{\eta}(\boldsymbol{\theta})^\top \mathbf{g}(\mathbf{y})\}}{\kappa(\boldsymbol{\theta})} \\ &= \text{E} \left(\exp[\{\boldsymbol{\eta}(\boldsymbol{\theta}') - \boldsymbol{\eta}(\boldsymbol{\theta})\}^\top \mathbf{g}(\mathbf{Y})]; \boldsymbol{\theta} \right), \end{aligned}$$

and proposes to estimate this expectation for values of $\boldsymbol{\theta}'$ near $\boldsymbol{\theta}$ based on a sample from the model with configuration $\boldsymbol{\theta}$: given a sample $\vec{\mathbf{y}}^{\boldsymbol{\theta}^t}$ from $\text{ERGM}(\boldsymbol{\theta}^t)$, update the guess

$$\begin{aligned} \boldsymbol{\theta}^{t+1} &= \arg \max_{\boldsymbol{\theta}'} \left(\{\boldsymbol{\eta}(\boldsymbol{\theta}') - \boldsymbol{\eta}(\boldsymbol{\theta}^t)\}^\top \mathbf{g}(\mathbf{y}^{\text{obs}}) - \log \frac{1}{S} \sum_{s=1}^S \exp[\{\boldsymbol{\eta}(\boldsymbol{\theta}') - \boldsymbol{\eta}(\boldsymbol{\theta}^t)\}^\top \mathbf{g}(\mathbf{y}^{\boldsymbol{\theta}^t, s})] \right) \\ &= \arg \max_{\boldsymbol{\theta}'} \log \frac{1}{S} \sum_{s=1}^S \exp[-\{\boldsymbol{\eta}(\boldsymbol{\theta}') - \boldsymbol{\eta}(\boldsymbol{\theta}^t)\}^\top \mathbf{z}(\mathbf{y}^{\boldsymbol{\theta}^t, s})]. \end{aligned} \quad (2)$$

This is equivalent to solving

$$\hat{\mathbf{U}}_{\vec{\mathbf{y}}^{\boldsymbol{\theta}^t}}(\boldsymbol{\theta}^{t+1}) \stackrel{\text{def}}{=} -\boldsymbol{\eta}'(\boldsymbol{\theta}^{t+1})^\top \hat{\mathbf{E}}_{\vec{\mathbf{y}}^{\boldsymbol{\theta}^t}}\{\mathbf{z}(\mathbf{Y}); \boldsymbol{\theta}^{t+1}\} = \mathbf{0},$$

the MCMLE approximation of the score equation, where, for a statistic $\mathbf{t}(\cdot)$,

$$\hat{\mathbf{E}}_{\vec{\mathbf{y}}^{\boldsymbol{\theta}^t}}\{\mathbf{t}(\mathbf{Y}); \boldsymbol{\theta}'\} \stackrel{\text{def}}{=} \frac{\sum_{s=1}^S \exp[\{\boldsymbol{\eta}(\boldsymbol{\theta}') - \boldsymbol{\eta}(\boldsymbol{\theta}^t)\}^\top \mathbf{g}(\mathbf{y}^{\boldsymbol{\theta}^t, s})] \mathbf{t}(\mathbf{y}^{\boldsymbol{\theta}^t, s})}{\sum_{s=1}^S \exp[\{\boldsymbol{\eta}(\boldsymbol{\theta}') - \boldsymbol{\eta}(\boldsymbol{\theta}^t)\}^\top \mathbf{g}(\mathbf{y}^{\boldsymbol{\theta}^t, s})]},$$

the MCMLE importance sampling approximation of $\text{E}\{\mathbf{t}(\mathbf{Y}); \boldsymbol{\theta}'\}$.

The MCMLE approach has the benefit of making very efficient use of the simulated sample, compared to the SA methods (Geyer and Thompson, 1992, Sec. 1.3): it uses the entire distribution of $\vec{\mathbf{y}}^{\boldsymbol{\theta}^t}$, rather than just its first moment, incorporates nonlinear effects of $\boldsymbol{\theta}$ on $\text{E}\{\mathbf{g}(\mathbf{Y}); \boldsymbol{\theta}\}$ in determining the next guess, and automatically determines the optimal (or close) step length, requiring much fewer sampling runs before convergence.

This efficiency comes at a cost: MCMLE is highly sensitive to a poor initial guess $\boldsymbol{\theta}^0$. Whereas SA methods only fail if the sample lies entirely on the edge of the convex hull (or $\boldsymbol{\theta}^0 \notin \boldsymbol{\Theta}_N$), MCMLE for non-curved ERGMs will also fail whenever the convex hull of the simulated statistics, $\text{Conv}\{\mathbf{g}(\vec{\mathbf{y}}^{\boldsymbol{\theta}^0})\}$, does not contain $\mathbf{g}(\mathbf{y}^{\text{obs}})$. (Equivalently, $\mathbf{0} \notin \text{Conv}\{\mathbf{U}_{\vec{\mathbf{y}}^{\boldsymbol{\theta}^0}}(\boldsymbol{\theta}^0)\}$.) Then, $\boldsymbol{\theta}^{t+1}$ does not exist. (Hummel et al., 2012, p. 926)

Hummel et al. (2012) proposed two major modifications to the MCMLE algorithm that ameliorate this. The first is the lognormal approximation: if $\mathbf{g}(\mathbf{Y})$ is approximately normal, $\exp[\{\boldsymbol{\eta}(\boldsymbol{\theta}') - \boldsymbol{\eta}(\boldsymbol{\theta})\}^\top \mathbf{g}(\mathbf{Y})]$ is lognormal, and its expectation gives an approximation

$$\ell(\boldsymbol{\theta}') - \ell(\boldsymbol{\theta}) \approx \{\boldsymbol{\eta}(\boldsymbol{\theta}') - \boldsymbol{\eta}(\boldsymbol{\theta})\}^\top \{-\bar{\mathbf{z}}(\bar{\mathbf{y}}^{\boldsymbol{\theta}'})\} - \{\boldsymbol{\eta}(\boldsymbol{\theta}') - \boldsymbol{\eta}(\boldsymbol{\theta})\}^\top \widetilde{\text{Var}}\{\mathbf{z}(\bar{\mathbf{y}})\}\{\boldsymbol{\eta}(\boldsymbol{\theta}') - \boldsymbol{\eta}(\boldsymbol{\theta})\}/2, \quad (3)$$

whose maximizer in $\boldsymbol{\theta}'$ depends only on the first two moments of $\bar{\mathbf{y}}^{\boldsymbol{\theta}'}$ and has a closed form for non-curved ERGMs—the version derived by Hummel et al. (2012)—extending directly to curved models, though the maximizer no longer has a closed form (as implemented in the R (R Core Team, 2015) package `ergm` (Hunter et al., 2008; Handcock et al., 2015)).

Also introduced was the Partial Stepping technique, where a step length $0 < \gamma \leq 1$ is selected, and $\mathbf{g}(\mathbf{y}^{\text{obs}})$ is replaced with $\gamma \mathbf{g}(\mathbf{y}^{\text{obs}}) + (1 - \gamma) \bar{\mathbf{g}}(\bar{\mathbf{y}}^{\boldsymbol{\theta}'})$ in the calculation of $\hat{\mathbf{U}}_{\bar{\mathbf{y}}^{\boldsymbol{\theta}'}}(\cdot)$. In other words, the vector of observed statistics is shifted towards the centroid of the simulated statistics, reducing the length of the step while preserving its general direction. Hummel et al. choose γ adaptively, selecting a safety margin (1.05) and finding the highest $\gamma \leq 1$ such that

$$1.05\gamma \mathbf{g}(\mathbf{y}^{\text{obs}}) + (1 - 1.05\gamma) \bar{\mathbf{g}}(\bar{\mathbf{y}}^{\boldsymbol{\theta}'}) \in \text{Conv}\{\mathbf{g}(\bar{\mathbf{y}}^{\boldsymbol{\theta}'})\}. \quad (4)$$

While this approach survives poor starting values (provided $\boldsymbol{\theta}^0 \in \Theta_{\text{N}}$), it is not immune to them, in that a poor starting value is likely to result in a tiny γ and a very long optimization. And so, we turn to the question of obtaining good values for $\boldsymbol{\theta}^0$.

1.2 Techniques for Finding Starting Values

Although there have been some recent developments on asymptotic approximations for ERGMs (He and Zheng, 2015), they have only been derived for a very specific set of models, and may or may not generalize. The two major techniques for obtaining $\boldsymbol{\theta}^0$ are the *maximum pseudo-/composite likelihood estimation* (MPLE/MCLE) and the more recently proposed *contrastive divergence* (CD). (It is also possible to instead fit a simpler submodel, and initialize the remaining elements of $\boldsymbol{\theta}$ to 0, as is done by PNet (Wang et al., 2014).)

1.2.1 Composite Likelihood

Before simulation-based methods were proposed, the only practical way to fit ERGMs with intractable normalizing constants was using pseudolikelihood (Besag, 1974; Strauss and Ikeda, 1990), approximating

$$L(\boldsymbol{\theta}) \approx \tilde{L}(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \prod_{(i,j) \in \mathbb{Y}} \Pr(Y_{i,j} = y_{i,j}^{\text{obs}} | \mathbf{Y}_{-(i,j)} = \mathbf{y}_{-(i,j)}^{\text{obs}}; \boldsymbol{\theta}), \quad (5)$$

where $y_{i,j}$ is the indicator of the presence of a tie from actor i to actor j and $\mathbf{y}_{-(i,j)}$ is the set of all ties in \mathbf{y} excluding (i, j) . The pseudolikelihood is then maximized to produce the maximum pseudolikelihood estimator (MPLE) $\tilde{\boldsymbol{\theta}}$.

For binary ERGMs, this gives an estimating equation

$$\tilde{\mathbf{U}}(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \boldsymbol{\eta}'(\tilde{\boldsymbol{\theta}})^\top \sum_{(i,j) \in \mathbb{Y}} [y_{i,j}^{\text{obs}} - \text{logit}^{-1}\{\boldsymbol{\eta}(\tilde{\boldsymbol{\theta}})^\top \boldsymbol{\Delta}_{i,j} \mathbf{g}(\mathbf{y}^{\text{obs}})\}] \boldsymbol{\Delta}_{i,j} \mathbf{g}(\mathbf{y}^{\text{obs}}) = \mathbf{0},$$

a (nonlinear) logistic regression, with “covariates” $\boldsymbol{\Delta}_{i,j} \mathbf{g}(\mathbf{y}) \stackrel{\text{def}}{=} \mathbf{g}(\mathbf{y} \cup \{(i, j)\}) - \mathbf{g}(\mathbf{y} \setminus \{(i, j)\})$, the effect of adding the tie (i, j) to the network \mathbf{y} on $\mathbf{g}(\mathbf{y})$, all other ties being equal.

MPLE can be quite different from the MLE, however, (van Duijn et al., 2009) so, with growing computing power making methods of Section 1.1 feasible, today it is mainly used to initialize them. Even in that capacity, it has practical limitations. For example, consider a network drawn from a process for which the total number of ties that can be observed is fixed at c . That is, $\mathcal{Y} = \{\mathbf{y} \in 2^{\mathbb{Y}} : |\mathbf{y}| = c\}$, used in the application of Hunter and Handcock (2006). One Metropolis–Hastings algorithm for exploring such a sample space selects one tie and one non-tie in \mathbf{y}^s at random and proposes to toggle both of them, thus preserving the total number of ties. Using this algorithm to sample $\tilde{\mathbf{y}}^{\tilde{\boldsymbol{\theta}^t}$ for either MCML or SA would result in the MLE on the constrained sample space.

In contrast, MPLE, and its generalization, maximum composite likelihood estimate (MCLE) (Lindsay, 1988), would require an algorithm to enumerate, rather than explore, the set of possible pairs of toggles, and the resulting pseudolikelihood would no longer be a binary logistic regression, but rather a multinomial model. In practice, this creates an additional burden on the implementer. Other constraints—such as conditioning on the degree sequence of a graph—require as many as 4 or 6 toggles in the proposal. (Rao et al.,

1996) The resulting combinatorial explosion can be addressed by sampling, but the problem of requiring a reimplementaion of MPLE remains.

In valued ERGMs, $\Pr(Y_{i,j} = \mathbf{y}_{i,j}^{\text{obs}} | \mathbf{Y}_{\neg(i,j)} = \mathbf{y}_{\neg(i,j)}^{\text{obs}}; \boldsymbol{\theta})$ might, itself, be intractable, such as when CMP (Shmueli et al., 2005) is used, whereas MCMC-based methods require no additional implementational or computational effort. (Krivitsky, 2012)

1.2.2 Contrastive Divergence

In a model whose log-likelihood gradient could only be obtained by an MCMC simulation, Hinton (2002) proposed not to run the MCMC simulation to convergence, but rather to make a series of parallel MCMC updates, each starting at the observed data, and calculate the gradient based on that. As applied to ERGMs by Asuncion et al. (2010), given an MCMC sampling algorithm for $\text{ERGM}(\boldsymbol{\theta})$, let $\text{ERGM}_{\text{CD}_k}(\boldsymbol{\theta})$ be the distribution of random graphs produced after k MCMC transitions starting with \mathbf{y}^{obs} . Call its expectation $\text{E}_{\text{CD}_k}(\cdot; \boldsymbol{\theta})$. Then, a CD_k estimate $\tilde{\boldsymbol{\theta}}^k$ solves

$$\mathbf{U}_{\text{CD}_k}(\tilde{\boldsymbol{\theta}}^k) \stackrel{\text{def}}{=} \boldsymbol{\eta}'(\tilde{\boldsymbol{\theta}}^k)^\top [-\text{E}_{\text{CD}_k}\{\mathbf{z}(\mathbf{Y}); \tilde{\boldsymbol{\theta}}^k\}] = \mathbf{0}, \quad (6)$$

shown by Hyvriinen (2006) to be equivalent to the MPLE if only one variable (i.e. edge) is updated and the updates are full-conditional Gibbs. Asuncion et al. (2010) noted that CD_1 (the MPLE) and CD_∞ (the MLE) were endpoints of a continuum of increasingly close approximations to the latter and showed that if k variables are block-updated in each MCMC step (*blocked contrastive divergence* (BCD)), CD_1 estimate is equivalent to maximizing the composite likelihood with block size of k . Asuncion et al. then applied CD_k to a number of exponential families, including ERGMs, using SA (with $\boldsymbol{\alpha}_t$ a scalar) to find the MCLE. Carreira-Perpiñ and Hinton (2005) proposed using the CD_k estimates to seed MCMLE.

No burn-in phase is required for CD_k estimates, which means that some of the inefficiency of the SA algorithms is not as problematic, but the issues of step length remain: Asuncion et al. (2010) used very short steps, for example. Also, the sampling algorithm required is distinct from the one that one might use for MCMLE, so using BCD as initial values for MCMLE may require additional effort on the part of the implementer.

Note, however, that CD_k sampling alleviates the sensitivity issues of MCMLE: if $\boldsymbol{\eta}(\boldsymbol{\theta}^0) = \mathbf{0}$, then for sampling $\bar{\mathbf{y}}^{\boldsymbol{\theta}^0, k}$ from $\text{ERGM}_{\text{CD}_k}(\boldsymbol{\theta}^0)$ is very

unlikely to produce realizations such that $\mathbf{g}(\mathbf{y}^{\text{obs}}) \notin \text{Conv}\{\mathbf{g}(\bar{\mathbf{y}}^{\theta^0, k})\}$, and it is also immune to the problem of $\theta^0 \notin \Theta_N$. We therefore propose to combine the two approaches.

Fellows (2014) described a framework for contrastive divergence as a variational approximation, provided some guidelines on what proposal kernels are likely to perform well, and advocated using a more efficient Newton-like update of the form

$$\boldsymbol{\theta}^{t+1} = \boldsymbol{\theta}^t - [\widetilde{\text{Var}}\{\mathbf{z}(\bar{\mathbf{y}}^{\theta^t, k})\}]^{-1}[\bar{\mathbf{z}}(\bar{\mathbf{y}}^{\theta^t, k})], \quad (7)$$

for the special case of non-curved ERGMs. This approach is equivalent to lognormal approximation of Hummel et al. (2012) (with step length γ fixed at 1) and Robbins–Monro without an α_t that does not decrease in t . The author has also recently become aware of a thesis by Hummel (2011) that also discussed ERGM CD inference. Hummel focused on exploring different MCMC kernels, but some computational considerations were also discussed, and we note the overlap where it occurs.

Outline

We begin by extending the MCMLE Partial Stepping technique of Hummel et al. (2012) to curved ERGMs in Section 2. In Section 3, we motivate and describe an algorithm for using an MCMLE-like technique to efficiently obtain CD estimates, and discuss practical considerations in applying this approach to a variety of network data and model types. Finally, in Section 4, we illustrate the technique’s versatility and gain some intuition for the tuning parameters it requires through a series of applications to the network data and models previously considered in ERGM computation literature.

2 Partial Stepping for Curved ERGMs

Hummel et al. (2012) derive Partial Stepping and the adaptive selection of the step length γ for non-curved ERGMs. Using their approach with curved models is likely to result in unnecessarily conservative step lengths, however. To see why, consider a popular Geometrically Weighted Degrees (GWD) (Hunter and Handcock, 2006, eq. 4.8) ERGM term. In our notation, this term has two free parameters, θ_1 (the strength of the effect) and θ_2 (decay

rate), which map to $(n - 1)$ -subvectors of $\boldsymbol{\eta}(\cdot)$ and $\boldsymbol{g}(\cdot)$ having elements

$$\begin{aligned}\eta_i(\boldsymbol{\theta}) &= \theta_1 \exp(2\theta_2) [\{1 - \exp(-\theta_2)\}^i - 1 + i \exp(-\theta_2)] \\ g_i(\mathbf{y}) &= \sum_{j=1}^n 1_{|\mathbf{y}_j|=i},\end{aligned}$$

for $i = 1, \dots, n - 1$, with $|\mathbf{y}_j|$ being the degree of actor j . That is, for every degree value i , $\boldsymbol{\eta}(\boldsymbol{\theta})$ has an element with a coefficient proportional to θ_1 and decaying in i at a rate controlled by θ_2 , and $\boldsymbol{g}(\mathbf{y})$ has an element with the count of actors with degree exactly i .

The sufficient statistic therefore includes the full degree distribution of the network. A necessary, though not sufficient, requirement for (4) to hold for a given γ is that

$$\min_s g_i(\mathbf{y}^{\theta^t, s}) < 1.05\gamma g_i(\mathbf{y}^{\text{obs}}) + (1 - 1.05\gamma)\bar{g}_i(\bar{\mathbf{y}}^{\theta^t}) < \max_s g_i(\mathbf{y}^{\theta^t, s})$$

hold for every degree value i , and applying Partial Stepping to $\boldsymbol{g}(\cdot)$ itself would select γ accordingly, as if every element of $\boldsymbol{\eta}$ were a free parameter, even though the actual dimension of $\boldsymbol{\theta}$ is much smaller.

To address this problem, we observe that (1) can be expressed as

$$\mathbf{U}(\boldsymbol{\theta}) = \boldsymbol{\eta}'(\boldsymbol{\theta})^\top \boldsymbol{g}(\mathbf{y}^{\text{obs}}) - \boldsymbol{\eta}'(\boldsymbol{\theta})^\top \text{E}\{\boldsymbol{g}(\mathbf{Y}); \boldsymbol{\theta}\}$$

which suggests that for curved ERGMs, we might use γ such that

$$1.05\gamma \boldsymbol{\eta}'(\boldsymbol{\theta}^t)^\top \boldsymbol{g}(\mathbf{y}^{\text{obs}}) + (1 - 1.05\gamma) \boldsymbol{\eta}'(\boldsymbol{\theta}^t)^\top \bar{\boldsymbol{g}}(\bar{\mathbf{y}}^{\theta^t}) \in \text{Conv}\{\boldsymbol{\eta}'(\boldsymbol{\theta}^t)^\top \boldsymbol{g}(\bar{\mathbf{y}}^{\theta^t})\}.$$

Our generalization does not provide the same guarantees as using the raw $\boldsymbol{g}(\bar{\mathbf{y}}^{\theta^t})$, since $\boldsymbol{\eta}'(\boldsymbol{\theta}^t)$ is not constant in $\boldsymbol{\theta}^t$, but it gives each element of $\boldsymbol{g}(\cdot)$ its due weight.

3 Contrastive Divergence via Monte Carlo MLE

3.1 Motivation

Just as the algorithm in Section 1.1.2 solves the score equations (1), we can apply the importance sampling paradigm to solving (6). For the special case of CD_1 , on a Metropolis–Hastings sampler with proposal density $q(\cdot|\cdot)$,

$$\text{E}_{\text{CD}_1}\{\mathbf{z}(\mathbf{Y}); \boldsymbol{\theta}\} = \sum_{\mathbf{y}' \in \mathcal{Y} \setminus \{\mathbf{y}^{\text{obs}}\}} q(\mathbf{y}'|\mathbf{y}^{\text{obs}}) \min \left[1, \frac{q(\mathbf{y}^{\text{obs}}|\mathbf{y}')}{q(\mathbf{y}'|\mathbf{y}^{\text{obs}})} \exp\{\boldsymbol{\eta}(\boldsymbol{\theta})^\top \mathbf{z}(\mathbf{y}')\} \right] \mathbf{z}(\mathbf{y}')$$

(because for rejections, $\mathbf{y}' \equiv \mathbf{y}^{\text{obs}}$, so $\mathbf{z}(\mathbf{y}') = \mathbf{0}$), and since

$$\text{E}_{\text{CD}_1} \{ \mathbf{z}(\mathbf{Y}); \boldsymbol{\theta}' \} = \text{E}_{\text{CD}_1} \left(\frac{\min \left[\frac{q(\mathbf{Y}|\mathbf{y}^{\text{obs}})}{q(\mathbf{y}^{\text{obs}}|\mathbf{Y})}, \exp\{\boldsymbol{\eta}(\boldsymbol{\theta}')^\top \mathbf{z}(\mathbf{Y})\} \right]}{\min \left[\frac{q(\mathbf{Y}|\mathbf{y}^{\text{obs}})}{q(\mathbf{y}^{\text{obs}}|\mathbf{Y})}, \exp\{\boldsymbol{\eta}(\boldsymbol{\theta})^\top \mathbf{z}(\mathbf{Y})\} \right]} \mathbf{z}(\mathbf{Y}); \boldsymbol{\theta} \right)$$

the importance sampling estimator for it for $\boldsymbol{\theta}'$ based on a sample $\mathbf{y}^{\boldsymbol{\theta},1} = (\mathbf{y}^{\boldsymbol{\theta}^t,1,1}, \dots, \mathbf{y}^{\boldsymbol{\theta}^t,1,S})$ drawn from $\text{ERGM}_{\text{CD}_1}(\boldsymbol{\theta})$ is

$$\hat{\text{E}}_{\mathbf{y}^{\boldsymbol{\theta},1}} \{ \mathbf{z}(\mathbf{Y}); \boldsymbol{\theta}' \} = \frac{1}{S} \sum_{s=1}^S \frac{\min \left[\frac{q(\mathbf{y}^{\boldsymbol{\theta},1,s}|\mathbf{y}^{\text{obs}})}{q(\mathbf{y}^{\text{obs}}|\mathbf{y}^{\boldsymbol{\theta},1,s})}, \exp\{\boldsymbol{\eta}(\boldsymbol{\theta}')^\top \mathbf{z}(\mathbf{y}^{\boldsymbol{\theta},1,s})\} \right]}{\min \left[\frac{q(\mathbf{y}^{\boldsymbol{\theta},1,s}|\mathbf{y}^{\text{obs}})}{q(\mathbf{y}^{\text{obs}}|\mathbf{y}^{\boldsymbol{\theta},1,s})}, \exp\{\boldsymbol{\eta}(\boldsymbol{\theta})^\top \mathbf{z}(\mathbf{y}^{\boldsymbol{\theta},1,s})\} \right]} \mathbf{z}(\mathbf{y}^{\boldsymbol{\theta},1,s}). \quad (8)$$

If the ratios of $q(\cdot|\cdot)$ are recorded during the sampling, this could be implemented directly; and similarly—although with complications—for $k > 1$. In practice, MCMLE weights ($\exp[\{\boldsymbol{\eta}(\boldsymbol{\theta}') - \boldsymbol{\eta}(\boldsymbol{\theta})\}^\top \mathbf{z}(\mathbf{y}^{\boldsymbol{\theta},1,s})]$) can be used instead: the importance weight in (8) for a given $\mathbf{y}^{\boldsymbol{\theta},1,s}$ is monotonically increasing in $\boldsymbol{\eta}(\boldsymbol{\theta}')^\top \mathbf{z}(\mathbf{y}^{\boldsymbol{\theta},1,s})$, with the weights being equal (to 1) if $\boldsymbol{\theta}' = \boldsymbol{\theta}$, so using the MCMLE weights will, at worst, make the approximation somewhat worse when $\boldsymbol{\theta}'$ is far away from $\boldsymbol{\theta}$, but if (8) evaluated at $\boldsymbol{\theta}' = \boldsymbol{\theta}^{t+1}$ is close to $\mathbf{0}$ and $\boldsymbol{\theta}^{t+1} = \boldsymbol{\theta}^t$, we can be confident that the optimization has converged. For higher k , the distribution $\text{ERGM}_{\text{CD}_k}(\boldsymbol{\theta})$ of the sample will be closer to $\text{ERGM}(\boldsymbol{\theta})$, so this approximation will only improve.

3.2 Algorithm

This leads to a CD update of the form

$$\boldsymbol{\theta}^{t+1} = \arg \max_{\boldsymbol{\theta}'} \log \frac{1}{S} \sum_{s=1}^S \exp[-\{\boldsymbol{\eta}(\boldsymbol{\theta}') - \boldsymbol{\eta}(\boldsymbol{\theta}^t)\}^\top \mathbf{z}(\mathbf{y}^{\boldsymbol{\theta}^t,k,s})].^2$$

It has a number of appealing properties. From the implementation point of view, the only change required to turn MCMLE into CD is modifying the MCMC sampler to revert the chain to \mathbf{y}^{obs} every k steps, and any improvements to the sampling algorithm also improve the estimator.

²Hummel (2011, eq. 4.3) used a similar update in the context of CD for non-curved ERGMs, but did not motivate the use of MCMLE importance sampling weights explicitly.

From the computational cost point of view, in MLE methods, every new guess $\boldsymbol{\theta}^t$ requires a long burning-in period, a fixed cost that cannot be reduced by parallel processing, and $\mathbf{y}^{\boldsymbol{\theta}^t, s}$ tend to be autocorrelated, which encourages using a large S and fewer iterations. But, as $\boldsymbol{\theta}^t$ in (2) moves farther away from $\boldsymbol{\theta}^t$, the accuracy of the estimate decreases. $\vec{\mathbf{y}}^{\boldsymbol{\theta}^t, k}$, on the other hand, is a random sample, requiring a total of Sk MCMC steps per iteration, and the sampling is an embarrassingly parallel problem. This means that a series of relatively short, inexpensive CD steps can be used to obtain an initial value.

To ameliorate potential problems with using MCMLE weights rather than true weights, we propose to use the Hummel et al. (2012) Partial Stepping technique with a more conservative γ safety margin than the Hummel et al. (2012) default of 1.05.³ Whether their lognormal approximation should be used is less clear. Its Newton-like update (7) is optimal if $\mathbf{g}(\vec{\mathbf{Y}}^{\boldsymbol{\theta}^t, k, s})$ is well approximated by the multivariate normal distribution and the relationship between $\boldsymbol{\theta}$ and $\mathbf{U}_{\text{CD}_k}(\boldsymbol{\theta})$ is well approximated by linear over the magnitude of the update, but, for modest k , this is highly unlikely to be the case: for example, if $g(\mathbf{y}) = |\mathbf{y}|$, the number of edges in the network, for any MCMC step that toggles one potential tie at a time $\mathbf{g}(\mathbf{y}^{\boldsymbol{\theta}^t, 1, s})$ can be one of only three values, $\mathbf{g}(\mathbf{y}^{\text{obs}}) - 1$, $\mathbf{g}(\mathbf{y}^{\text{obs}})$, or $\mathbf{g}(\mathbf{y}^{\text{obs}}) + 1$.

At the same time, although every MCMC step reduces the Kullback–Leibler divergence between $\text{ERGM}(\boldsymbol{\theta})$ and $\text{ERGM}_{\text{CD}_k}(\boldsymbol{\theta})$ (Cover and Thomas, 1991, Thm. 15.1.10, for example), a full-conditional Gibbs sampler is likely to do so faster than a Metropolis–Hastings sampler with the same block size. MPLE is equivalent to CD with full-conditional Gibbs sampling (Hyvrinen, 2006), while Metropolis–Hastings is more practical for ERGMs (Hunter et al., 2008), so it is likely that MPLE will outperform CD_1 , and Fellows (2014), in particular, focuses on full-conditional Gibbs.

3.3 Artificial multiplicity

Fellows (2014) also shows that increasing k alone may not be sufficiently effective at improving the estimators, and suggests that CD kernels should instead be designed to “focus” on the dependencies in the model: if blocked contrastive divergence (Asuncion et al., 2010) is used for, say, a network model with triadic closure, the “blocks” should include triads.

³Hummel (2011, p. 77) CD implementation also uses 1.05. We explore its effects in Section 4.

Unfortunately, specialized proposals defeat the advantage of CD as a source of initial values: it is no longer a drop-in replacement for MCMC. Therefore, we propose an ad hoc remedy by modifying the Metropolis–Hastings algorithm to create artificial blocks of proposals. Recall that, given a proposal distribution $q(\cdot|\cdot)$, the acceptance probability

$$\alpha(\mathbf{y}^*|\mathbf{y}) = \min\left(1, \frac{q(\mathbf{y}|\mathbf{y}^*)}{q(\mathbf{y}^*|\mathbf{y})} \exp[\boldsymbol{\eta}(\boldsymbol{\theta})^\top \{\mathbf{g}(\mathbf{y}^*) - \mathbf{g}(\mathbf{y})\}]\right).$$

For MCMC, a simple proposal that toggles only one dyad, or the minimal number of dyads needed to preserve a constraint, generally suffices. A more complex proposal can be emulated by chaining m simple proposals, i.e., $\mathbf{y}^{*1} \sim q(\mathbf{y}^{*1}|\mathbf{y})$, $\mathbf{y}^{*2} \sim q(\mathbf{y}^{*2}|\mathbf{y}^{*1})$, \dots , $\mathbf{y}^{*m} \sim q(\mathbf{y}^{*m}|\mathbf{y}^{*m-1})$, then accepting \mathbf{y}^{*m} with probability

$$\alpha(\mathbf{y}^{*m}|\mathbf{y}) = \min\left(1, \frac{q(\mathbf{y}|\mathbf{y}^{*1})}{q(\mathbf{y}^{*1}|\mathbf{y})} \dots \frac{q(\mathbf{y}^{*m-1}|\mathbf{y}^{*m})}{q(\mathbf{y}^{*m}|\mathbf{y}^{*m-1})} \exp[\boldsymbol{\eta}(\boldsymbol{\theta})^\top \{\mathbf{g}(\mathbf{y}^{*m}) - \mathbf{g}(\mathbf{y})\}]\right),$$

remaining at \mathbf{y} otherwise. This is not the correct acceptance probability (a correct one would consider all possible ways to propose \mathbf{y}^{*m} from \mathbf{y}), so m is a trade-off between the correctness of the stationary distribution and incorporation of the dependence in the model.

But, an approximation is what we require. We will use $\tilde{\boldsymbol{\theta}}^{(m,k)}$ to refer to a $\text{CD}_{(m,k)}$ estimate, taking k steps with artificial multiplicity m .

3.4 Stopping Criterion

We briefly turn to the question of when to consider the optimization to be concluded. The stopping criterion of Hummel et al. (2012) is not well-suited to CD, because for small $m \times k$ in particular, it may not be possible for $\bar{\mathbf{y}}^{\boldsymbol{\theta}^t, (m,k)}$ to draw sufficiently far away from \mathbf{y}^{obs} for Hummel et al. for $\text{Conv}\{\mathbf{g}(\bar{\mathbf{y}}^{\boldsymbol{\theta}^t, (m,k)})\}$ to not contain $\mathbf{g}(\mathbf{y}^{\text{obs}})$, no matter how bad $\boldsymbol{\theta}^t$ is.

The forms of the estimating equations (1) and (6) suggest another straightforward method to determine whether a particular $\boldsymbol{\theta}^t$ is sufficiently close to $\tilde{\boldsymbol{\theta}}^{(m,k)}$ to stop. For each guess $\boldsymbol{\theta}^t$, CD draws a simple random sample $\bar{\mathbf{y}}^{\boldsymbol{\theta}^t, (m,k)}$ from $\text{ERGM}_{\text{CD}_{(m,k)}}(\boldsymbol{\theta}^t)$, $\bar{\mathbf{g}}(\bar{\mathbf{y}}^{\boldsymbol{\theta}^t, (m,k)})$ is an unbiased estimator of $\text{E}_{\text{CD}_{(m,k)}}\{\mathbf{g}(\mathbf{Y}); \boldsymbol{\theta}^t\}$, and premultiplication by $\boldsymbol{\eta}'(\boldsymbol{\theta}^t)$ is a linear transformation, so $\bar{\mathbf{U}}_{\bar{\mathbf{y}}^{\boldsymbol{\theta}^t, (m,k)}}(\boldsymbol{\theta}^t)$ is unbiased for $\mathbf{U}_{\text{CD}_{(m,k)}}(\boldsymbol{\theta}^t)$.

Therefore, we can use a Hotelling's T^2 -Test (Hotelling, 1931) to test $H_0 : \mathbf{U}_{\text{CD}_{(m,k)}}(\boldsymbol{\theta}^t) = \mathbf{0}$ based on a sample $\mathbf{U}_{\bar{\mathbf{y}}^{\boldsymbol{\theta}^t, (m,k)}}(\boldsymbol{\theta}^t)$, stopping upon a failure

to reject. The decision to terminate entails accepting a null hypothesis, but this can be ameliorated in practice by setting a very high α , because the cost of a Type I error is small: setting $\alpha = 0.5$ only entails running on average $1/\alpha = 2$ more iterations than necessary.

3.5 Choice of k and m

The choice of k is a trade-off: higher k leads to $\tilde{\boldsymbol{\theta}}^{(m,k)}$ being closer to $\hat{\boldsymbol{\theta}}$, but the computing cost increases in proportion to it, and sensitivity to poor $\boldsymbol{\theta}^0$ does as well, and a similar trade-off (up to a point) applies for m .

Our goal is to maximize the utility of the CD estimate as the starting value of MCMLE, and a simple one-dimensional metric of this utility is available: the Hummel et al. (2012) adaptive step length for the first MCMLE iteration (4). This is, essentially, a measurement of how deep in the convex hull of $\boldsymbol{\eta}'(\tilde{\boldsymbol{\theta}}^{(m,k)})^\top \mathbf{g}(\tilde{\mathbf{y}}^{\tilde{\boldsymbol{\theta}}^{(m,k)}})$ is $\boldsymbol{\eta}'(\tilde{\boldsymbol{\theta}}^{(m,k)})^\top \mathbf{g}(\mathbf{y}^{\text{obs}})$. An estimated step length of 1 or close implies that only a few steps of full MCMLE will be required, while a step length close to 0 implies that the starting value is practically useless.

We therefore propose to evaluate $\tilde{\boldsymbol{\theta}}^{(m,k)}$ for a series of (m, k) configurations, then, for each estimate $\tilde{\boldsymbol{\theta}}^{(m,k)}$, draw an MCMC sample, evaluate the adaptive step length, and initialize the MCMLE with the one giving the highest γ such that (4) holds. Because MCMLE step requires a long burn-in, this is likely to be computationally expensive, but we can, instead, use a proxy in the form of a short MCMC run that would nonetheless have a burn-in period much longer than the highest value of $m \times k$ used.

4 Examples

We illustrate the proposed technique by replicating examples found in the ERGM computational methods literature. We list them here, identifying the computational challenge of each, provide the details about the data and the models in the Appendix.

Lazega, a collaboration network of lawyers, was used by Hunter and Handcock (2006) to demonstrate inference for curved ERGMs, fitting a curved ERGM conditional having a specific number of ties—a complex constraint. (We also replicate the curved fit without the constraint, modeling edge count.)

E. coli, a transcriptional regulation network, was selected by Hummel et al. (2012) for being particularly difficult to fit.

Kapferer, a network of workers in a tailor shop in Zambia, was also used by Hummel et al. (2012).

Zachary, a valued network of counts of contexts of interactions among members of a university karate club, which we use to demonstrate immediate applicability to models for valued networks, fitting a Binomial- and a Poisson-reference ERGM. (For the latter, we include the CMP (Shmueli et al., 2005) term, deliberately initializing CD with a starting value outside of Θ_N to test the algorithm’s robustness.)

4.1 Procedure

We have implemented the proposed techniques in the R (R Core Team, 2015) package `ergm` (Hunter et al., 2008; Handcock et al., 2015) and released them on an experimental basis.

We refrain from tuning the algorithms to each specific dataset, and unless otherwise stated, we use default settings of the `ergm` package. For CD, we use $S = 1024$, start the estimation at $\mathbf{0}_q$, unless otherwise noted, and allow 60 iterations. For each example, we evaluate the MPLE (if available—no implementation of MPLE for curved ERGMs is known to the author) and CD for each combination of the following factors (five times):

(\mathbf{m}, \mathbf{k}): every combination of $k = 1, 2, 4, 16, 128$ and $m = 1, 2, 4, 8$ such that $k \times m \leq 256$;

Update type: “IS MCMLE”, the importance-sampling-based update (2) and “Lognormal”, using (3) for the likelihood; and

γ margin: 1.05 of Hummel et al. (2012), 1.5 and 2 (more conservative).

Having found the $\tilde{\boldsymbol{\theta}}$ according to each method and parameters, we measure its utility as a starting value for MCMLE by obtaining an MCMC sample from $\text{ERGM}(\tilde{\boldsymbol{\theta}})$ starting at \mathbf{y}^{obs} for each of the following settings, then evaluating adaptive step length γ as proposed by Hummel et al. (2012) or our extension in Section 2:

$\gamma_{\mathbf{F}}$: based on `ergm` (Handcock et al., 2015) package defaults (burn-in: 16384, sample size: 1024, interval: 1024), to evaluate the suitability of $\tilde{\boldsymbol{\theta}}$ as a starting value; and

$\gamma_{\mathbf{S}}$: based on a shortened MCMC run (burn-in: 8192, sample size: 1024, interval: 8), as a proxy for $\gamma_{\mathbf{F}}$ to test the suggestion of Section 3.5.

Table 1: Aggregate effects of the update type and the γ margin on quality, speed, and reliability. Means are taken after standardizing each value by its example’s overall mean and standard deviation.

Settings		γ_F		Cost (mean)		Failures	
Update	γ mar.	mean	< 0.02	Iter.	$\frac{\text{sec.}}{m \times k}$	Error	Unconv.
IS MCMLE	0.05	-0.24	17%	-0.24	-0.21	0%	13%
IS MCMLE	0.50	-0.01	8%	-0.23	-0.15	0%	5%
IS MCMLE	1.00	0.03	8%	-0.03	-0.14	1%	6%
Lognormal	0.05	0.09	5%	0.03	0.12	0%	5%
Lognormal	0.50	0.08	5%	0.14	0.18	1%	5%
Lognormal	1.00	0.05	5%	0.34	0.21	0%	6%

Error: failed with an error (typically, was trapped)

Unconv.: failed to meet the convergence criterion in 60 iterations

4.2 Results

Table 1 gives the effects of the update type and the γ margin. Importance sampling MCMLE updates as opposed to the Newton-like lognormal updates appear to be a trade-off between speed and stability, with MCMLE making more efficient steps, at a greater risk of making a poor step. A more conservative γ margin, alleviates this, while retaining the efficiency improvement.

The effects of m and k are visualized in Figure 1. The general pattern appears to be that MPLE, where available, outperforms CD with small k and m , but is, eventually, outperformed by CD, except in hard-to-sample models such as the *E. coli*. At the same time, there are diminishing returns as $m \times k$ increases, and, in valued ERGMs, they actually perform worse.

For the hard-to-sample *E. coli*, higher artificial multiplicities seem to outperform lower for the same $m \times k$, but the results are less consistent for other ERGMs, and, in particular, for the valued ERGMs and the fixed-edges model, whose proposal is already multiplicitous; this may be because there are many more possible ways to a given \mathbf{y}^{*m} from \mathbf{y} in those cases, which $\alpha(\mathbf{y}^{*m}|\mathbf{y})$ ignores.

In the CMP model, CD was able to locate an adequate $\theta^0 \in \Theta_N$ in 95% of the trials.

Lastly the relationship between γ_S and γ_F in Figure 2 is positive but very noisy. But, for our purpose of selecting an adequate CD estimate to seed the MCMLE, it is adequate: the configuration with the best (rightmost) γ_S is,

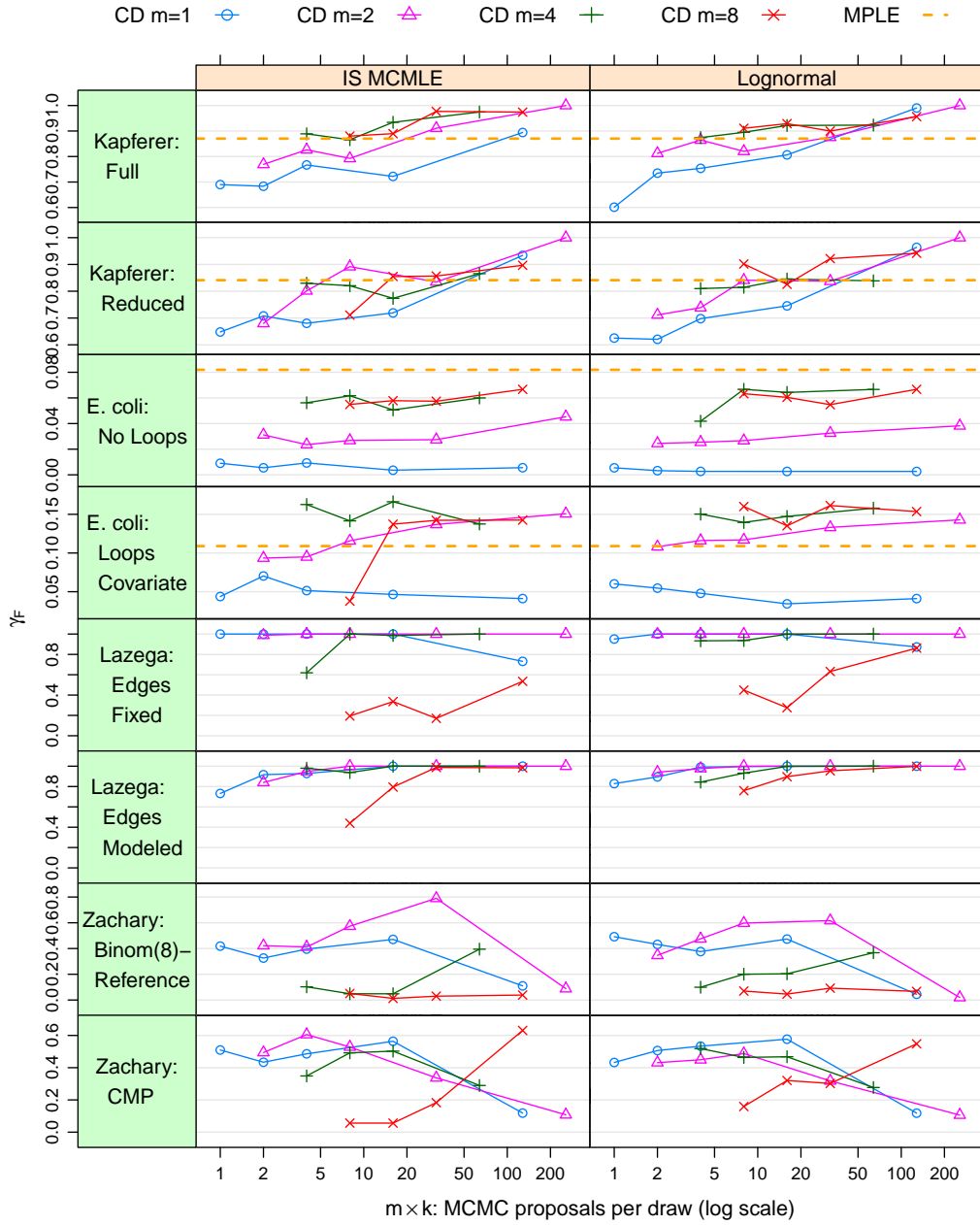


Figure 1: Effect of (m, k) and update type on the quality of the starting value as measured by γ_F (using margin 1.05). Values displayed are the means of the five replications, using CD with γ margin of 1.5. Dashed line gives γ_F from the MPLE fit.

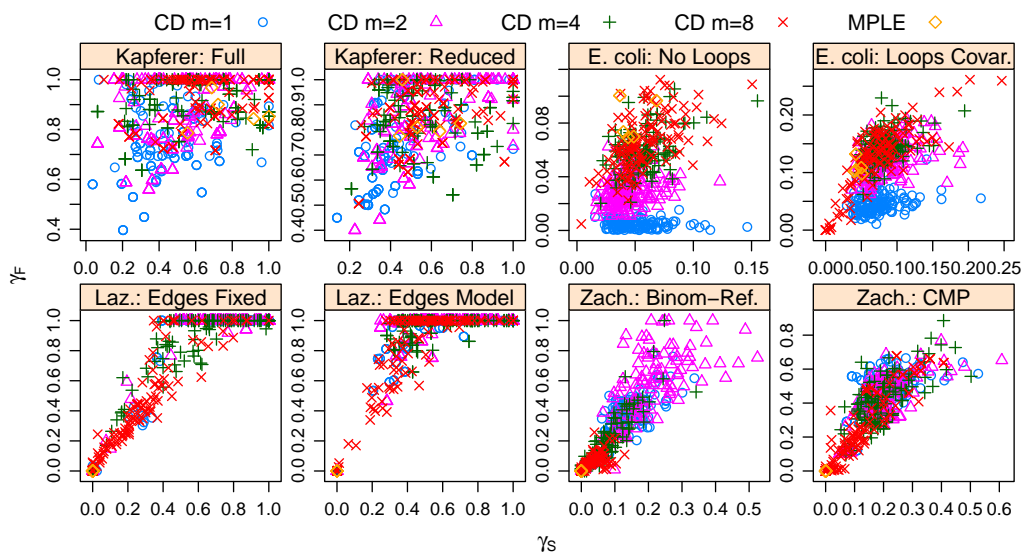


Figure 2: Using Hummel step length from a short run to predict step length for full MCMC: prediction is likely to be noisy, but valid for selection of an adequate starting value.

while not necessarily the best (topmost) γ_F , is usually among the best, and never among those that give $\gamma_F \approx 0$. At the same time, evidence from the *E. coli* fits suggests that for hard-to-sample models too short pilot runs may result in selecting a poor start.

5 Conclusion

We have reviewed the available techniques for obtaining initial values for the simulation-based MLE methods for exponential family models with intractable normalizing constants, and, combining the approaches of Monte Carlo MLE and contrastive divergence, we propose a fairly universal algorithm for obtaining these values; and we have extended the existing techniques for improving MCMLE to the curved ERGMs.

Our examples demonstrate the viability and versatility of our approach: adequate starting values are produced for a wide variety of datasets and models—some designed to be difficult—with an algorithm agnostic to the specifics of the model. In practice, this means that any implementation of MCMLE for a new valued or constrained ERGM class (e.g., rank or signed networks) acquires a source of starting values without additional effort.

Different problems call for different (m, k) , and we have shown that short pilot MCMC runs can be used to select an adequate starting value for the MCMC out of several candidates, which are themselves inexpensive to fit. Thus, this tuning can be automated.

An alternative approach to selecting (m, k) may be to use an increasing sequence of k s, initializing each at the previous one’s solution as its stopping criterion is met. This approach should be used with caution, however, because $\tilde{\boldsymbol{\theta}}$ based on a small k can be worse than $\tilde{\boldsymbol{\theta}} = \mathbf{0}$. This is subject for future research.

We focused on the case where the networks were fully observed. Handcock and Gile (2010) formulated a framework for modeling of partially observed networks—networks that have missing ties—and expressed the log-likelihood as $\ell(\boldsymbol{\theta}) = \log \Pr(\mathbf{Y} \in \mathcal{Y}(\mathbf{y}^{\text{obs}}); \boldsymbol{\theta}) = \log \sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{y}^{\text{obs}})} \Pr(\mathbf{Y} = \mathbf{y}'; \boldsymbol{\theta})$, where $\mathcal{Y}(\mathbf{y}^{\text{obs}})$ is defined as the set of networks whose partial observation could have produced \mathbf{y}^{obs} : essentially, all of the ways to impute the missing ties in \mathbf{y}^{obs} . They then proposed to maximize this likelihood by taking advantage of the fact that, if $\kappa_{\mathcal{Y}'}(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \sum_{\mathbf{y}' \in \mathcal{Y}'} h(\mathbf{y}') \exp\{\boldsymbol{\eta}(\boldsymbol{\theta})^\top \mathbf{g}(\mathbf{y}')\}$, log-likelihood can

be expressed as $\ell(\boldsymbol{\theta}) = \log \kappa_{\mathcal{Y}(\mathbf{y}^{\text{obs}})}(\boldsymbol{\theta}) - \log \kappa_{\mathcal{Y}}(\boldsymbol{\theta})$, resulting in

$$U(\hat{\boldsymbol{\theta}}) = \nabla_{\boldsymbol{\theta}} \ell(\hat{\boldsymbol{\theta}}) = \boldsymbol{\eta}'(\hat{\boldsymbol{\theta}})^\top [\mathbb{E}_{\mathcal{Y}(\mathbf{y}^{\text{obs}})}\{\mathbf{g}(\mathbf{Y}); \hat{\boldsymbol{\theta}}\} - \mathbb{E}_{\mathcal{Y}}\{\mathbf{g}(\mathbf{Y}); \hat{\boldsymbol{\theta}}\}] = \mathbf{0},$$

with MCMLE approximation also possible for the first term by sampling $\bar{\mathbf{y}}^{\theta^t} | \mathbf{y}^{\text{obs}}$ from $\text{ERGM}_{\mathcal{Y}(\mathbf{y}^{\text{obs}})}(\boldsymbol{\theta}^t)$. Partial Stepping can be extended to this case as well, by translating $\mathbf{g}(\bar{\mathbf{y}}^{\theta^t} | \mathbf{y}^{\text{obs}})$ towards $\bar{\mathbf{g}}(\bar{\mathbf{y}}^{\theta^t})$ until $\bar{\mathbf{g}}(\bar{\mathbf{y}}^{\theta^t} | \mathbf{y}^{\text{obs}})$ is sufficiently deep in $\text{Conv}\{\mathbf{g}(\bar{\mathbf{y}}^{\theta^t})\}$.

For CD, this creates a problem: while $\bar{\mathbf{y}}^{\theta^t} | \mathbf{y}^{\text{obs}}$ depends on \mathbf{y}^{obs} only through the observed dyads and information about which dyads are missing due to the ergodic property of MCMC, sampling from $\text{ERGM}_{\text{CD}_{(m,k)}}(\boldsymbol{\theta}^t)$ requires a specific initial network and depends on it strongly. In the context of CD, these problems can be partially addressed by using higher k s: the longer the MCMC chain, the less important \mathbf{y}^{obs} , but more efficient and stable approaches are subject for research. (A similar issue exists for the MPLE: the composite likelihood is a sum of (5) over possible imputations of missing dyads in \mathbf{y}^{obs} , and simply excluding the unobserved dyads from the product (5) still conditions on them.)

Alternatively, a network might not be observed at all, only its sufficient statistic vector \mathbf{g}^{obs} along with its sample space \mathcal{Y} . By sufficiency, MLE is unaffected by this. (Hummel et al., 2012) MPLE, MCLE, and CD are, however. A simple practical solution is to use simulated annealing to construct a network \mathbf{y}^{sim} such that $\mathbf{g}(\mathbf{y}^{\text{sim}}) \approx \mathbf{g}^{\text{obs}}$ and use it as a surrogate for \mathbf{y}^{obs} . It may not be possible to obtain a perfectly matched network, but this can be addressed in the same way as with missing data.

Lastly, we have focused on ERGMs in particular, but these methods are agnostic to the nature of the data, operating only on sufficient statistics, so this development is equally applicable to other domains. In particular, the problem of a complex $\Theta_{\mathbb{N}}$ is present in Strauss and related point processes as well (Geyer and Thompson, 1992, for example).

References

Arthur U. Asuncion, Qiang Liu, Alexander T. Ihler, and Padhraic Smyth. Learning with blocks: Composite likelihood and contrastive divergence. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics (AISTATS-10)*,

2010. URL http://machinelearning.wustl.edu/mlpapers/papers/AISTATS2010_AsuncionLIS10.
- Julian Besag. Spatial interaction and the statistical analysis of lattice systems (with discussion). *Journal of the Royal Statistical Society, Series B*, 36: 192–236, 1974. ISSN 0035-9246.
- Miguel Á. Carreira-Perpiñ and Geoffrey Hinton. On contrastive divergence learning. In Robert G. Cowell and Zoubin Ghahramani, editors, *Proceedings of the Tenth International Workshop on Artificial Intelligence and Statistics, Jan 6-8, 2005, Savannah Hotel, Barbados*, pages 33–40. Society for Artificial Intelligence and Statistics, 2005. URL <http://www.gatsby.ucl.ac.uk/aistats/>.
- Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory*. Wiley Series in Telecommunications. John Wiley & Sons, Inc., 1991. ISBN 0-471-06259-6.
- Ian E. Fellows. Why (and when and how) contrastive divergence works. *arXiv preprint arXiv:1405.0602*, 2014.
- Charles J. Geyer and Elizabeth A. Thompson. Constrained Monte Carlo maximum likelihood for dependent data (with discussion). *Journal of the Royal Statistical Society. Series B*, 54(3):657–699, 1992. ISSN 0035-9246.
- Mark S. Handcock and Krista J. Gile. Modeling social networks from sampled data. *Annals of Applied Statistics*, 4(1):5–25, 2010. ISSN 1932-6157. doi:10.1214/08-AOAS221.
- Mark S. Handcock, David R. Hunter, Carter T. Butts, Steven M. Goodreau, Pavel N. Krivitsky, and Martina Morris. *ergm: Fit, Simulate and Diagnose Exponential-Family Models for Networks*. The Statnet Project (<http://www.statnet.org>), 2015. URL <http://CRAN.R-project.org/package=ergm>. R package version 3.4.0.
- Ran He and Tian Zheng. GLMLE: graph-limit enabled fast computation for fitting exponential random graph models to large social networks. *Social Network Analysis and Mining*, 5(1):8, 2015. ISSN 1869-5450. doi:10.1007/s13278-015-0247-3.

- Geoffrey E. Hinton. Training products of experts by minimizing contrastive divergence. *Neural computation*, 14(8):1771–1800, 2002.
- Harold Hotelling. The generalization of student’s ratio. *Annals of Mathematical Statistics*, 2(3):360–378, 08 1931. doi:10.1214/aoms/1177732979.
- Ruth M. Hummel. *Improving Estimation for Exponential-Family Random Graph Models*. PhD thesis, The Pennsylvania State University, may 2011. URL <https://etda.libraries.psu.edu/paper/11493/>.
- Ruth M. Hummel, David R. Hunter, and Mark S. Handcock. Improving simulation-based algorithms for fitting ergms. *Journal of Computational and Graphical Statistics*, 21(4):920–939, 2012. doi:10.1080/10618600.2012.679224.
- David R. Hunter and Mark S. Handcock. Inference in curved exponential family models for networks. *Journal of Computational and Graphical Statistics*, 15(3):565–583, 2006. ISSN 1061-8600.
- David R. Hunter, Mark S. Handcock, Carter T. Butts, Steven M. Goodreau, and Martina Morris. `ergm`: A package to fit, simulate and diagnose exponential-family models for networks. *Journal of Statistical Software*, 24(3):1–29, May 2008. ISSN 1548-7660. URL <http://www.jstatsoft.org/v24/i03>.
- David R. Hunter, Pavel N. Krivitsky, and Michael Schweinberger. Computational statistical methods for social network models. *Journal of Computational and Graphical Statistics*, 21(4):856–882, 2012. doi:10.1080/10618600.2012.732921.
- Aapo Hyvriinen. Consistency of pseudolikelihood estimation of fully visible boltzmann machines. *Neural Computation*, 18(10):2283–2292, October 2006. ISSN 0899-7667. doi:10.1162/neco.2006.18.10.2283.
- Bruce Kapferer. *Strategy and Transaction in an African Factory: African Workers and Indian Management in a Zambian Town*. Manchester University Press, 1972.
- Pavel N. Krivitsky. Exponential-family random graph models for valued networks. *Electronic Journal of Statistics*, 6:1100–1128, 2012. doi:10.1214/12-EJS696.

- Emmanuel Lazega and Philippa E. Pattison. Multiplexity, generalized exchange and cooperation in organizations: a case study. *Social Networks*, 21(1):67–90, 1999. ISSN 0378-8733. doi:10.1016/S0378-8733(99)00002-7.
- Bruce G. Lindsay. Composite likelihood methods. *Contemporary Mathematics*, 80:221–239, 1988. doi:10.1090/conm/080/999014.
- Martina Morris, Mark S. Handcock, and David R. Hunter. Specification of exponential-family random graph models: Terms and computational aspects. *Journal of Statistical Software*, 24(4):1–24, May 2008. ISSN 1548-7660. URL <http://www.jstatsoft.org/v24/i04>.
- Saisuke Okabayashi and Charles J. Geyer. Long range search for maximum likelihood in exponential families. *Electronic Journal of Statistics*, 6:123–147, 2012. doi:10.1214/11-EJS664.
- R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2015. URL <http://www.R-project.org/>.
- A. Ramachandra Rao, Rabindranath Jana, and Suraj Bandyopadhyay. A markov chain Monte Carlo method for generating random $(0, 1)$ -matrices with given marginals. *Sankhyā: The Indian Journal of Statistics, Series A*, 58(2):225–242, June 1996. ISSN 0581-572X.
- Alessandro Rinaldo, Stephen E. Fienberg, and Yi Zhou. On the geometry of discrete exponential families with application to exponential random graph models. *Electronic Journal of Statistics*, 3:446–484, 2009. ISSN 1935-7524. doi:10.1214/08-EJS350.
- Herbert Robbins and Sutton Monro. A stochastic approximation method. *The Annals of Mathematical Statistics*, 22(3):400–407, September 1951. ISSN 00034851.
- Garry Robins, Philippa Pattison, and Stanley S. Wasserman. Logit models and logistic regressions for social networks: III. Valued relations. *Psychometrika*, 64(3):371–394, 1999. ISSN 0033-3123.
- Shai S. Shen-Orr, Ron Milo, Shmoolik Mangan, and Uri Alon. Network motifs in the transcriptional regulation network of *Escherichia coli*. *Nature Genetics*, 31(1):64–68, May 2002. ISSN 1061-4036. doi:10.1038/ng881.

- Galit Shmueli, Thomas P. Minka, Joseph B. Kadane, Sharad Borle, and Peter Boatwright. A useful distribution for fitting discrete data: Revival of the Conway–Maxwell–Poisson distribution. *Journal of the Royal Statistical Society: Series C*, 54(1):127–142, January 2005. ISSN 1467-9876. doi:10.1111/j.1467-9876.2005.00474.x.
- Tom A. B. Snijders. Markov chain Monte Carlo estimation of exponential random graph models. *Journal of Social Structure*, 3(2), 2002.
- Tom A. B. Snijders, Philippa E. Pattison, Garry L. Robins, and Mark S. Handcock. New specifications for exponential random graph models. *Sociological Methodology*, 36(1):99–153, 2006.
- David Strauss and Michael Ikeda. Pseudolikelihood estimation for social networks. *Journal of the American Statistical Association*, 85(409):204–212, 1990. ISSN 0162-1459.
- Marijtje A. J. van Duijn, Krista J. Gile, and Mark S. Handcock. A framework for the comparison of maximum pseudo-likelihood and maximum likelihood estimation of exponential family random graph models. *Social Networks*, 31(1):52–62, 2009. ISSN 0378-8733. doi:10.1016/j.socnet.2008.10.003.
- Peng Wang, Garry Robins, Philippa Pattison, and Johan Koskinen. *MPNet User Manual*. Melbourne School of Psychological Sciences, University of Melbourne, Melbourne, Australia, June 2014. URL <http://sna.unimelb.edu.au/PNet>.
- Stanley S. Wasserman and Philippa Pattison. Logit models and logistic regressions for social networks: I. An introduction to Markov graphs and p^* . *Psychometrika*, 61(3):401–425, 1996. ISSN 0033-3123.
- Wayne W. Zachary. An information flow model for conflict and fission in small groups. *Journal of Anthropological Research*, 33(4):452–473, 1977. ISSN 0091-7710.

A Details of the Examples

A.1 Lazega and Pattison’s Law firm

Hunter and Handcock (2006), in their development of inference for curved

ERGMs, used data collected by Lazega and Pattison (1999), describing patterns of collaboration of lawyers in a firm. The model they fit included covariates such as the effect of seniority, type of practice, whether the two lawyers had the same practice, were of the same gender, and worked in the same office; and it modeled triadic closure using Alternating k -triangles (also known as Geometrically-Weighted Edgewise Shared Partners (GWESP)), a curved ERGM term. (See the article in question for the details.)

We fit two variants of their Model 2 to these data: a variant whose sample space was restricted to have the same edge count as the observed network (which is what was fit by Hunter and Handcock) and a variant not conditioned on edge count, but using edge count as an additional model statistic.

A.2 *E. coli* transcriptional regulation network

Hummel et al. (2012), in illustrating their computational methods on a difficult model, used the *E. coli* transcriptional regulation network of Shen-Orr et al. (2002). Here, we fit two variants demonstrated by Hummel et al., referred to as “Model 2”: edge count, counts of actors with degree 2–5 (separately), and Geometrically-Weighted Degree (GWD) term with decay coefficient fixed at 0.25) and “Model 2 plus self-edges”, contains all of the above terms and, in addition, nodal covariates indicating whether a node has a non-self-edge and whether it has a self-edge.

A.3 Kapferer’s sociational data

Hummel et al. (2012) also demonstrated their approach on a well-known dataset collected by Kapferer (1972) on workers in a tailor shop in Zambia, and we reproduce the two models they had fit. The first model had, as its terms, count of edges, and the GWD, the GWESP, and the Geometrically-Weighted Dyadwise Shared Partners (GWDSP) statistics, the latter three having their decay coefficient fixed at 0.25. The second model dropped the GWD term.

A.4 Valued ties in a Zachary’s Karate club

Besides being particularly difficult for pseudolikelihood calculation, ERGMs for networks whose relationships have values present the additional challenge

that, if the set of possible relationship states \mathbb{S} is modeled without an *a priori* bound, and therefore the sample space of networks \mathcal{Y} is infinite, the natural parameter space $\Theta_{\mathbb{N}} = \{\boldsymbol{\theta} \in \mathbb{R}^q : \kappa(\boldsymbol{\theta}) < \infty\}$ may not equal \mathbb{R}^q , and if it turns out that $\boldsymbol{\theta}^0 \in \Theta_{\mathbb{N}}^c$, MCMC of the very first optimization step will not converge.

Contrastive divergence offers a solution to this problem: by limiting the number of MCMC steps, $\text{ERGM}_{\text{CD}(m,k)}(\boldsymbol{\theta})$ can avoid the runaway simulation, while providing the direction for the optimization to reach the parameter space.

We illustrate this on two examples of valued ERGMs, both using data collected by Zachary (1977), who reported observations of social relations in a university karate club with membership that varied between 50 and 100. The actors—32 ordinary club members and officers, the club president (“John A.”), and the part-time instructor (“Mr. Hi”)—were the ones who consistently interacted outside of the club. Over the course of the study, the club divided into two factions, and, ultimately, split into two clubs, one led by Hi and the other by John and the original club’s officers. The split was driven by a disagreement over whether Hi could unilaterally change the level of compensation for his services.

Zachary reported, for each pair of actors, the count of social contexts in which they interacted. The 8 contexts considered were academic classes at the university; Hi’s private karate studio in his night classes; Hi’s private karate studio where he taught on weekends; student-teaching at Hi’s studio; the university rathskeller (bar) located near the karate club; a bar located near the university campus; open karate tournaments in the area; and inter-collegiate karate tournaments. The highest number of contexts of interaction for a pair of individuals that was observed was 7.

In Model 1, we model the distribution of counts as a binomial-reference ERGM, i.e., $\mathbb{S} = 0..8$ and $h(\mathbf{y}) = \prod_{(i,j) \in \mathbb{Y}} \binom{8}{y_{i,j}}$, zero-modified by adding a term of the form $g_{\text{nonzero}}(\mathbf{y}) = \sum_{(i,j) \in \mathbb{Y}} 1_{y_{i,j} \neq 0}$.

In Model 2, we instead use a Poisson-reference ERGM (i.e., having dyad-wise sample space of $\mathbb{S} = \{0, 1, 2, \dots\}$) with $h(\mathbf{y}) \equiv 1 / \prod_{(i,j) \in \mathbb{Y}} y_{i,j}!$, and we include two statistics to affect the dyadwise distribution of counts: g_{nonzero} to control the overall propensity to have ties (i.e., have interactions in more than 0 contexts) and a statistic of the form $g_{\text{CMP}}(\mathbf{y}) = \sum_{(i,j) \in \mathbb{Y}} \log(y_{i,j}!)$, which, added to a geometric- or Poisson-reference ERGM models each relationship value as distributed Conway–Maxwell–Poisson (CMP) (Shmueli

et al., 2005; Krivitsky, 2012). A linear ERGM with this term—for example, with sufficient statistic $\mathbf{g}(\mathbf{y}) = (\sum_{(i,j) \in \mathbb{Y}} y_{i,j}, \sum_{(i,j) \in \mathbb{Y}} \log(y_{i,j}!))$, has a constrained natural parameter space $\Theta_N = \{\boldsymbol{\theta} \in \mathbb{R}^2 : \theta_2 = 1 \wedge \theta_1 < 0 \vee \theta_2 < 1\}$, making it neither regular nor steep (Krivitsky, 2012, App. B). For this reference, we use a Tie-Non-Tie (TNT) (Morris et al., 2008) augmentation of the zero-inflated Poisson algorithm of Krivitsky (2012, Alg. 1).

We model the structure of the network using two more terms: the faction leader effects, $\sum_{(i,j) \in \mathbb{Y}} y_{i,j} \mathbf{1}_{i=\text{Mr. Hi} \vee j=\text{Mr. Hi}}$ and $\sum_{(i,j) \in \mathbb{Y}} y_{i,j} \mathbf{1}_{i=\text{John A.} \vee j=\text{John A.}}$, and transitivity, the statistic described by Krivitsky (2012, eq. 12).

Unlike other fits, where we start the optimization at $\boldsymbol{\theta}^0 = \mathbf{0}_q$, in Model 2, we start the optimization at $\theta_{\text{CMP}}^0 = +2$, deliberately outside the parameter space. Also, we change a few non-CD-specific tuning parameters to accommodate non-binary data.