APPROXIMATE VARIATIONAL ESTIMATION FOR A MODEL OF NETWORK FORMATION

ANGELO MELE AND LINGJIONG ZHU

ABSTRACT. We develop approximate estimation methods for exponential random graph models (ERGMs), whose likelihood is proportional to an intractable normalizing constant. The usual approach approximates this constant with Monte Carlo simulations, however convergence may be exponentially slow. We propose a deterministic method, based on a variational mean-field approximation of the ERGM's normalizing constant. We compute lower and upper bounds for the approximation error for any network size, adapting nonlinear large deviations results. This translates into bounds on the distance between true likelihood and mean-field likelihood. Monte Carlo simulations suggest that in practice our deterministic method performs better than our conservative theoretical approximation bounds imply, for a large class of models.

Keywords: Networks, Microeconometrics, Large Networks, Variational Inference, Large deviations, Mean-Field Approximations

1. INTRODUCTION

This paper studies variational mean-field methods to approximate the likelihood of exponential random graph models (ERGMs), a class of statistical network formation models that has become popular in sociology, machine learning, statistics and more recently economics. While a large part of the statistical network literature is devoted to models with unconditionally or conditionally independent links (Graham, 2017; Airoldi et al., 2008; Bickel et al., 2013), ERGMs allow for conditional and unconditional dependence among links (Snijders, 2002; Wasserman and Pattison, 1996). These models have recently gained attention in economics, because several works have shown that ERGMs have a microeconomic foundation. In fact, the ERGM likelihood naturally

Date: First version: June 15, 2015. This version: May 9, 2020.

We are grateful to the editor and three excellent referees for their suggestions. We thank Anton Badev, Vincent Boucher, Aureo DePaula, Bryan Graham, Mert Gürbüzbalaban, Matt Jackson, Hiro Kaido, Michael Leung, Xiaodong Liu and Demian Pouzo for comments on previous versions of this paper. The second author is partially supported by NSF Grant DMS-1613164.

emerges as the stationary equilibrium of a potential game, where players engage in a myopic bestresponse dynamics of link formation (Blume, 1993; Mele, 2017; Badev, 2013; Chandrasekhar, 2016; Chandrasekhar and Jackson, 2014; Boucher and Mourifie, 2017), and in a large class of evolutionary games and social interactions models (Blume, 1993; Durlauf and Ioannides, 2010).

Estimation and inference for ERGMs are challenging, because the likelihood of the observed network is proportional to an intractable normalizing constant, that cannot be computed exactly, even in small networks. Therefore, exact Maximum Likelihood estimation (MLE) is infeasible. The usual estimation approach, the Markov Chain Monte Carlo MLE (MCMC-MLE), consists of simulating many networks using the model's conditional link probabilities and approximating the constant and the likelihood with Monte Carlo methods (Snijders, 2002; Koskinen, 2004; Chatterjee and Diaconis, 2013; Mele, 2017). Estimates of the MCMC-MLE converge almost surely to the MLE if the likelihoods are well-behaved (Gever and Thompson, 1992). However, a recent literature has shown that the simulation methods used to compute the MCMC-MLE may have exponential slow convergence, making estimation and approximation of the likelihood impractical or infeasible for a large class of ERGMs (Bhamidi et al., 2011; Chatterjee and Diaconis, 2013; Mele, 2017). An alternative is the Maximum Pseudo-likelihood estimator (MPLE), that finds the parameters that maximize the product of the conditional link probabilities of the model. While MPLE is simple and computationally fast, the properties of the estimator are not well understood, except in special cases, when some regularity conditions are satisfied (Boucher and Mourifie, 2017; Besag, 1974); in practice MPLE may give misleading estimates when the dependence among links is strong (Gever and Thompson, 1992). Furthermore, since the ERGMs are exponential families, networks with the same sufficient statistics will produce the same MLE, but may have different MPLE.

Our work departs from the standard methods of estimation, proposing deterministic approximations of the likelihood, based on the approximated solution of a variational problem. Our strategy is to use a mean-field algorithm to approximate the normalizing constant of the ERGM, at any given parameter value (Wainwright and Jordan, 2008; Bishop, 2006; Chatterjee and Diaconis, 2013). We then maximize the resulting approximate log-likelihood, with respect to the parameters. To be concrete, our approximation consists of using the likelihood of a simpler model with independent links to approximate the constant of the ERGM. The mean-field approximation algorithm finds the likelihood with independent links that minimizes the Kullback-Leibler divergence from the ERGM likelihood. Using this likelihood with independent links, we compute an approximate normalizing constant. We then evaluate the log-likelihood of our model, where the exact normalizing constant is replaced by its mean-field approximation.

Our main contribution is the computation of exact bounds for the approximation error of the normalizing constant's mean-field estimate. Our proofs use the theoretical machinery of Chatterjee and Dembo (2016) for non-linear large deviations in models with intractable normalizing constants. Using this powerful tool, we provide explicit lower and upper bounds to the error of approximation for the mean-field normalizing constant. The bounds depend on the magnitude of the parameters of our model and the size of link externalities (Mele, 2017; Boucher and Mourifie, 2017; Chandrasekhar, 2016; DePaula, 2017). The result holds for dense and moderately sparse networks. Remarkably and conveniently the mean-field error converges to zero as the network becomes large. This guarantees that for large networks, the log-normalizing constant of an ERGM is well approximated by our mean-field log-normalizing constant.

The main implication of our main result is that we can compute bounds to the distance between the log-likelihood of the ERGM and our approximate log-likelihood; these also converge in supnorm as the network grows large. As a consequence, we can use the approximated likelihood for estimation in large networks. If the likelihood is strictly concave, it is possible to show that our approximate estimator converges to the maximum likelihood estimator as long as the network grows large. Furthermore, because our bounds may not be sharp, in practice convergence could be faster than what is implied in these results.

While our method is guaranteed to perform well in large graphs, many applications involve small networks. For example, the school networks data in the National Longitudinal Study of Adoloscent Health (Add Health) (Boucher and Mourifie, 2017; Moody, 2001; Badev, 2013) or the Indian villages in Banerjee et al. (2013) include on average about 200-300 nodes. To understand the performance of our estimator in practice, we perform simple Monte Carlo exercises in

networks with few hundreds nodes, comparing mean-field estimates to MCMC-MLE and MPLE. Our Monte Carlo results show that in practice our estimator works better than the theoretical results suggest, for networks with 50 to 1000 nodes. The median mean-field approximation point estimates are close to the true parameters, but exhibit a small bias. Both MCMC-MLE and MPLE show a larger variability of point estimates for the two-stars and triangle parameters, measured as median absolute deviation. When we increase the network size, all three estimators improve, as expected. We conclude that our method's performance is comparable to available estimators in small networks. While our code can be made faster by exploiting efficient matrix algebra libraries and parallelization, the CPU time for estimation is comparable to the estimators implemented in the ergm package in R for networks with less than 200 nodes.

The main message of our theoretical results and Monte Carlo simulations is that the approximate mean-field approach is a valid alternative to existing methods for estimation of a large class of ERGMs. We note that our theoretical bounds may not be sharp, and in practice the meanfield algorithm may have better performance than what is implied by our conservative results, as confirmed by our Monte Carlo experiments.

To the best of our knowledge, this paper is one of the first works in economics to use meanfield approximations for approximate estimation of complex models. We show that our application of variational approximations has theoretical guarantees, and we can bound the error of approximation. While similar deterministic methods have been used to provide an approximation to the normalizing constant of the ERGM model (Chatterjee and Diaconis, 2013; Amir et al., 2012; Mele, 2017; He and Zheng, 2013; Aristoff and Zhu, 2018), we are the first to characterize the variational approximation error for a model with covariates and its computational feasibility.

Our technique can be applied to other models in economics and social sciences. For example, models of social interactions with binary decisions like in Blume (1993), Badev (2013), Durlauf and Ioannides (2010), models for bundles (Fox and Lazzati, 2017), and models of choices from menus (Kosyakova et al., 2018) have similar likelihoods with intractable normalizing constants . Therefore our method of approximation may allow estimation of these models for large sets of bundles or menu choices.

The rest of the paper is organized as follows. Section 2 presents the theoretical model and variational approximations. Section 3 contains the main theoretical results and the error bounds. Section 4 presents the Monte Carlo results and Section 5 concludes. All the proofs and additional Monte Carlo simulations are in the Appendix. Additional results and discussions are presented in the Online Appendix.

2. NETWORK FORMATION MODEL AND VARIATIONAL METHODS

2.1. Exponential random graph models. The class of exponential random graphs is an important generative model for networks and has been extensively used in applications in many disciplines (Wasserman and Pattison, 1996; Jackson, 2010; DePaula, 2017; Mele, 2017; Moody, 2001; Wimmer and Lewis, 2010; Amir et al., 2012). In this paper we consider a model with nodal covariates, two-stars and triangles.

Our model assumes that the network consists of n heterogeneous nodes, indexed by i = 1, ..., n; each node is characterized by a S-dimensional vector of observed attributes $\tau_i \in \mathcal{X} := \bigotimes_{j=1}^S \mathcal{X}_j$, i = 1, ..., n. The sets \mathcal{X}_j can represent age, race, gender, income, etc.¹ Let α be a $n \times n$ symmetric matrix with elements $\alpha_{ij} := \nu(\tau_i, \tau_j)$, where $\nu : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is a symmetric function and let β and γ be scalars. For ease of exposition we focus on the case in which the attributes are discrete and finite, but our results hold when this assumption is relaxed and the number of attributes is allowed to increase with the size of the network.

The likelihood $\pi_n(g, \alpha, \beta, \gamma)$ of observing the adjacency matrix g depends on the composition of links, the number of two-stars and the number of triangles

(2.1)
$$\pi_n(g;\alpha,\beta,\gamma) = \frac{\exp\left[Q_n(g;\alpha,\beta,\gamma)\right]}{\sum_{\omega\in\mathcal{G}_n}\exp\left[Q_n(\omega;\alpha,\beta,\gamma)\right]}$$

where the function Q is called a *potential function* and takes the form

$$(2.2) \qquad Q_n(g;\alpha,\beta,\gamma) = \sum_{i=1}^n \sum_{j=1}^n \alpha_{ij} g_{ij} + \frac{\beta}{2n} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n g_{ij} g_{jk} + \frac{2\gamma}{3n} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n g_{ij} g_{jk} g_{ki}.$$

¹For instance, if we consider gender and income, then S = 2, and we can take $\bigotimes_{j=1}^{2} \mathcal{X}_{j} = \{\text{male,female}\} \times \{\text{low, medium, high}\}$. The sets \mathcal{X}_{j} can be both discrete and continuous. For example, if we consider gender and income, we can also take $\bigotimes_{j=1}^{2} \mathcal{X}_{j} = \{\text{male,female}\} \times [\$50,000,\$200,000]$. Below we restrict the covariates to be discrete, but we allow the number of types to grow with the size of the network.

and $c(\alpha, \beta, \gamma) := \sum_{\omega \in \mathcal{G}_n} \exp [Q_n(\omega; \alpha, \beta, \gamma)]$ is a normalizing constant that guarantees that likelihood (2.1) is a proper distribution. The second and third term of the potential function (2.2) are the counts of two-stars and triangles in the network, rescaled by *n*. We rewrite (2.1) as

(2.3)
$$\pi_n(g;\alpha,\beta,\gamma) = \exp\left\{n^2 \left[T_n(g;\alpha,\beta,\gamma) - \psi_n(\alpha,\beta,\gamma)\right]\right\},$$

where $T_n(g; \alpha, \beta, \gamma) = Q_n(g; \alpha, \beta, \gamma)n^{-2}$ is the potential scaled by n^2 and the log-normalizing constant (scaled by n^2) is,

(2.4)
$$\psi_n(\alpha,\beta,\gamma) = \frac{1}{n^2} \log \sum_{\omega \in \mathcal{G}_n} \exp\left[n^2 T_n(\omega;\alpha,\beta,\gamma)\right],$$

and $\mathcal{G}_n := \{\omega = (\omega_{ij})_{1 \le i,j \le n} : \omega_{ij} = \omega_{ji} \in \{0,1\}, \omega_{ii} = 0, 1 \le i, j \le n\}$ is the set of all binary matrices with *n* nodes. The re-scaling of the potential and the log-normalizing constant is necessary for the asymptotic results, to avoid the explosion of the potential function as the size of the network grows large.

2.2. **Microeconomic equilibrium foundations.** ERGMs caught the attention of economists because recent works proves a behavioral and equilibrium interpretation of likelihood (2.3).² In fact, these likelihoods naturally arise as equilibrium of best-response dynamics in potential games (Blume, 1993; Monderer and Shapley, 1996; Butts, 2009; Mele, 2011).

To be concrete, let's consider the following game. Players' payoffs are a function of the composition of direct links, friends' popularity and the number of common friends. The utility of network g for player i is given by

(2.5)
$$u_i(g,\tau) = \sum_{j=1}^n \alpha_{ij} g_{ij} + \frac{\beta}{n} \sum_{j=1}^n \sum_{k=1}^n g_{ij} g_{jk} + \frac{\gamma}{n} \sum_{j=1}^n \sum_{k=1}^n g_{ij} g_{jk} g_{ki},$$

Each player forms links with other nodes, maximizing utility (2.5), but taking into account the strategies of other players. We can show that this game of network formation converges to an exponential random graph in a stationary equilibrium, under the following assumptions:³ (1) the

²Butts (2009), Mele (2017), Chandrasekhar and Jackson (2014), Boucher and Mourifie (2017), Badev (2013), DePaula (2017).

³See Mele (2017) or Badev (2013) for more technical details and variants of these assumptions. See also Chandrasekhar (2016), DePaula (2017), Chandrasekhar and Jackson (2014), Boucher and Mourifie (2017).

network formation is sequential, with only two active players in each period; (2) two players meet over time with probability $\rho_{ij} := \rho(\tau_i, \tau_j, g_{-ij}) > 0$, where g_{-ij} indicate the network g but link g_{ij} ; and these meetings are i.i.d. over time; (3) before choosing whether to form or delete a link, players receive an i.i.d. logistic shock $(\varepsilon_{ij1}, \varepsilon_{ij0})$. At time t, the link g_{ij}^t is formed if

$$u_i(g_{ij}^t = 1, g_{-ij}^{t-1}, \tau) + u_j(g_{ij}^t = 1, g_{-ij}^{t-1}, \tau) + \varepsilon_{ij1}^t \ge u_i(g_{ij}^t = 0, g_{-ij}^{t-1}, \tau) + u_j(g_{ij}^t = 0, g_{-ij}^{t-1}, \tau) + \varepsilon_{ij0}^t$$

Mele (2017) shows that such a model is a potential game (Monderer and Shapley, 1996) with potential function given by equation (2.2). The probability of observing network g in the long run is given by (2.3) (Theorem 1 in Mele (2017)), thus (2.3) describes the stationary behavior of the model. In the long-run we observe with high probability the pairwise stable networks, where no pair of players want to form or delete a link.⁴

2.3. Variational Approximations. The constant $\psi_n(\alpha, \beta, \gamma)$ in (2.4) is intractable because it is a sum over all $2^{\binom{n}{2}}$ possible networks with *n* nodes; if there are n = 10 nodes, the sum involves computation of 2^{45} potential functions, which is infeasible.⁵ In the literature on exponential family likelihoods with intractable normalizing constant, this problem is solved by approximating the normalizing constant using Markov Chain Monte Carlo (Snijders, 2002; Mele, 2017; Goodreau et al., 2009; Koskinen, 2004; Caimo and Friel, 2011; Murray et al., 2006). However, Bhamidi et al. (2011) has shown that such methods may have exponentially slow convergence for many ERGMs specifications.

We propose methods that avoid simulations and we find an approximate likelihood $q_n(g)$ that minimizes the Kullback-Leibler divergence $KL(q_n|\pi_n)$ between q_n and the true likelihood π_n :

(2.6)

$$KL(q_n|\pi_n) = \sum_{\omega \in \mathcal{G}_n} q_n(\omega) \log \left[\frac{q_n(\omega)}{\pi_n(\omega; \alpha, \beta)} \right]$$

$$= \sum_{\omega \in \mathcal{G}_n} q_n(\omega) \left[\log q_n(\omega) - n^2 T_n(\omega; \alpha, \beta, \gamma) + n^2 \psi_n(\alpha, \beta, \gamma) \right] \ge 0.$$

⁴In the Online Appendix \mathbf{E} we provide more details about the microeconomic foundation of the model for the interested reader.

⁵See Geyer and Thompson (1992), Murray et al. (2006), Snijders (2002) for examples.

With some algebra we obtain a lower-bound for the constant $\psi_n(\alpha, \beta, \gamma)$

$$\psi_n(\alpha,\beta,\gamma) \ge \mathbb{E}_{q_n}\left[T_n(\omega;\alpha,\beta,\gamma)\right] + \frac{1}{n^2}\mathcal{H}(q_n) := \mathcal{L}(q_n),$$

where $\mathcal{H}(q_n) = -\sum_{\omega \in \mathcal{G}_n} q_n(\omega) \log q_n(\omega)$ is the entropy of distribution q_n , and $\mathbb{E}_{q_n}[T_n(\omega; \alpha, \beta, \gamma)]$ is the expected value of the re-scaled potential, computed according to the distribution q_n .

To find the best likelihood approximation we minimize $KL(q_n|\pi_n)$ with respect to q_n , which is equivalent to finding the supremum of the lower-bound $\mathcal{L}(q_n)$, i.e.

(2.7)
$$\psi_n(\alpha,\beta,\gamma) = \sup_{q_n \in \mathcal{Q}_n} \mathcal{L}(q_n) = \sup_{q_n \in \mathcal{Q}_n} \left\{ \mathbb{E}_{q_n} \left[T_n(\omega;\alpha,\beta,\gamma) \right] + \frac{1}{n^2} \mathcal{H}(q_n) \right\}$$

where Q_n is the set of all the probability distributions on G_n . We have transformed the problem of computing an intractable sum into a variational problem, i.e. a maximization problem.

In general, problem (2.7) has no closed-form solution, thus the literature suggests to restrict Q_n to be the set of all completely factorized distribution⁶

(2.8)
$$q_n(g) = \prod_{i,j} \mu_{ij}^{g_{ij}} (1 - \mu_{ij})^{1 - g_{ij}},$$

where $\mu_{ij} = \mathbb{E}_{q_n}(g_{ij}) = \mathbb{P}_{q_n}(g_{ij} = 1)$. This approximation is called a *mean-field approximation* of the discrete exponential family. Straightforward algebra shows that the entropy of q_n is additive

$$\frac{1}{n^2}\mathcal{H}(q_n) = -\frac{1}{2n^2}\sum_{i=1}^n\sum_{j=1}^n \left[\mu_{ij}\log\mu_{ij} + (1-\mu_{ij})\log(1-\mu_{ij})\right],$$

and the expected potential can be computed as

$$\mathbb{E}_{q_n}\left[T_n\left(\omega;\alpha,\beta,\gamma\right)\right] = \frac{\sum_i \sum_j \alpha_{ij}\mu_{ij}}{n^2} + \beta \frac{\sum_i \sum_j \sum_k \mu_{ij}\mu_{jk}}{2n^3} + \gamma \frac{2\sum_i \sum_j \sum_k \mu_{ij}\mu_{jk}\mu_{ki}}{3n^3}.$$

⁶See Wainwright and Jordan (2008), Bishop (2006)

The mean-field approximation leads to a *lower bound of* $\psi_n(\alpha, \beta, \gamma)$, because we restricted Q_n , and the simpler variational problem is to find a $n \times n$ symmetric matrix $\mu(\alpha, \beta, \gamma)$ that solves

$$\psi_{n}(\alpha,\beta,\gamma) \geq \psi_{n}^{MF}(\boldsymbol{\mu}(\alpha,\beta,\gamma))$$

$$= \sup_{\boldsymbol{\mu}\in[0,1]^{n^{2}}:\mu_{ij}=\mu_{ji},\forall i,j} \left\{ \frac{1}{n^{2}} \sum_{i,j} \alpha_{ij}\mu_{ij} + \frac{\beta}{2n^{3}} \sum_{i,j,k} \mu_{ij}\mu_{jk} + \frac{2\gamma}{3n^{3}} \sum_{i,j,k} \mu_{ij}\mu_{jk}\mu_{ki} - \frac{1}{2n^{2}} \sum_{i=1}^{n} \sum_{j=1}^{n} [\mu_{ij}\log\mu_{ij} + (1-\mu_{ij})\log(1-\mu_{ij})] \right\}.$$
(2.9)

The mean-field problem is in general *nonconvex* and the maximization can be performed using any global optimization method, e.g. simulated annealing or Nelder-Mead.⁷

3. THEORETICAL RESULTS

3.1. Convergence of the variational mean-field approximation. For finite n, the variational mean-field approximation contains an error of approximation. In the next theorem, we provide a lower and upper bound to the error of approximation for our model.

THEOREM 3.1. For fixed network size *n*, the approximation error of the variational mean-field problem is bounded as

$$(3.1) \quad \frac{C_3(\beta,\gamma)}{n} \le \psi_n(\alpha,\beta,\gamma) - \psi_n^{MF}(\boldsymbol{\mu}(\alpha,\beta,\gamma)) \le C_1(\alpha,\beta,\gamma) \left(\frac{\log n}{n}\right)^{1/5} + \frac{C_2(\alpha,\beta,\gamma)}{n^{1/2}},$$

where $C_1(\alpha, \beta, \gamma)$, $C_2(\alpha, \beta, \gamma)$ are constants depending on α , β and γ and $C_3(\beta, \gamma)$ is a constant depending only on β, γ :

$$C_{1}(\alpha, \beta, \gamma) := c_{1} \cdot \left(\max_{i,j} |\alpha_{ij}| + |\beta|^{4} + |\gamma|^{4} + 1 \right),$$

$$C_{2}(\alpha, \beta, \gamma) := c_{2} \cdot \left(\max_{i,j} |\alpha_{ij}| + |\beta| + |\gamma| + 1 \right)^{1/2} \cdot (1 + |\beta|^{2} + |\gamma|^{2})^{1/2},$$

$$C_{3}(\beta, \gamma) := |\beta| + 4|\gamma|,$$

where $c_1, c_2 > 0$ are some universal constants.

⁷See Wainwright and Jordan (2008) and Bishop (2006) for more details.

The constants in Theorem 3.1 are functions of the parameters α , β and γ . The upper bound depends on the maximum payoff from direct links and the intensity of payoff from indirect connections. The lower bound only depends on the strength of indirect connections payoffs (popularity and common friends, that is β and γ). One consequence is that our result holds when the network is dense, but also when it is moderately sparse, in the sense that $|\alpha_{ij}|$, $|\beta|$ and $|\gamma|$ can have moderate growth in n instead of being bounded, and the difference of ψ_n and ψ_n^{MF} goes to zero if $C_1(\alpha, \beta, \gamma)$ grows slower than $n^{1/5}/(\log n)^{1/5}$ and $C_2(\alpha, \beta, \gamma)$ grows slower than $n^{1/2}$ as $n \to \infty$. For example, if $\max_{i,j} |\alpha_{ij}| = O(n^{\delta_1})$, $|\beta| = O(n^{\delta_2})$, $|\gamma| = O(n^{\delta_3})$ where $\delta_1 < \frac{1}{5}$ and $\delta_2, \delta_3 < \frac{1}{20}$, then $\psi_n - \psi_n^{MF}$ goes to zero as $n \to \infty$. On the other hand, if the graph is too sparse, e.g. $|\beta| = \Omega(n)$, $|\gamma| = \Omega(n)$, then ψ_n cannot be approximated by ψ_n^{MF} .

Our main Theorem 3.1 implies that *we can approximate the log-likelihood of the ERGM* using the mean-field approximated constant.

PROPOSITION 3.1. Let $\ell_n(g_n, \alpha, \beta, \gamma)$ be the log-likelihood of the ERGM

$$\ell_n(g_n, \alpha, \beta, \gamma) := n^{-2} \log \left(\pi_n(g_n, \alpha, \beta, \gamma) \right) = T_n(g_n, \alpha, \beta, \gamma) - \psi_n(\alpha, \beta, \gamma),$$

and $\ell_n^{MF}(g_n, \alpha, \beta, \gamma)$ be the "mean-field log-likelihood" obtained by approximating ψ_n with ψ_n^{MF} :

$$\ell_n^{MF}(g_n, \alpha, \beta, \gamma) := T_n(g_n, \alpha, \beta, \gamma) - \psi_n^{MF}(\alpha, \beta, \gamma).$$

Then for any compact parameter space Θ *,*

$$(3.2) \quad 0 \le \sup_{\alpha,\beta,\gamma\in\Theta} \left[\ell_n^{MF} - \ell_n\right] \le \sup_{\alpha,\beta,\gamma\in\Theta} C_1(\alpha,\beta,\gamma) n^{-1/5} (\log n)^{1/5} + \sup_{\alpha,\beta,\gamma\in\Theta} C_2(\alpha,\beta,\gamma) n^{-1/2}.$$

Proposition 3.1 shows that as the network size grows large, the mean-field approximation of the log-likelihood ℓ_n^{MF} is arbitrarily close to the ERGM log-likelihood ℓ_n . This approximation is similar in spirit to the MCMC-MLE method, where the log-normalizing constant is approximated via MCMC to obtain an approximated log-likelihood (Geyer and Thompson, 1992; Snijders, 2002; DePaula, 2017; Moller and Waagepetersen, 2004). The main difference is that our approximation is *deterministic* and does not require any simulation.

Note that $\ell_n^{MF} = T_n - \psi_n^{MF}$ and $\ell_n = T_n - \psi_n$. If ℓ_n converges to ℓ_∞ uniformly on a compact parameter space Θ , then so does ℓ_n^{MF} . If ℓ_n, ℓ_n^{MF} and ℓ_∞ are continuous and strictly concave, $\hat{\theta}_n$, $\hat{\theta}_n^{MF}$, the unique maximizers of ℓ_n and ℓ_n^{MF} will converge to the unique maximizer of ℓ_∞ and hence $\hat{\theta}_n - \hat{\theta}_n^{MF}$ will go to zero as $n \to \infty$. In the Online Appendix we provide further results on the behavior of the mean-field approximation as $n \to \infty$, where we discuss the convergence of the log-constant.⁸

The result in Proposition 3.1 can be used to bound the distance between the mean-field estimate and the maximum likelihood estimate, for any network size rather than for large n. However, such bounds require additional and stronger assumptions on the shape of the likelihood. Indeed, in Appendix **B**, we show that a sufficient conditions for computing the bound is a strongly concave likelihood. Under such assumption, we can use the bound in Proposition 3.1 for the log-likelihood to provide a bound on the distance between MLE and mean-field estimator for any network size n. However, these bounds may not be sharp, and therefore we consider them very conservative. In the next section we show via Monte Carlo simulation that in many cases our estimator performs better than the bounds would imply.

4. ESTIMATION EXPERIMENTS

To understand the performance of the variational approximation in smaller networks, we perform some Monte Carlo experiments. We compare the mean-field approximation with the standard simulation-based MCMC-MLE Geyer and Thompson (1992); Snijders (2002) and the MPLE (Besag, 1974). Our method converges in n^2 steps, while the MCMC-MLE may converge in e^{n^2} steps. The MPLE usually converges faster.

4.1. **Approximation algorithm for the normalizing constant.** We implemented our variational approximation for few models in the R package mfergm, available in Github.⁹ We follow the statistical machine learning literature and use an iterative algorithm that is guaranteed to converge

⁸The strict concavity of the likelihood is closely related to the identification of parameters in ERGM models, for which there is a lack of general results (see Mele (2017), Chatterjee and Diaconis (2013), Aristoff and Zhu (2018) for examples in special cases).

⁹See https://github.com/meleangelo/mfergm, with instructions for installation and few examples.

to a local maximum of the mean-field problem (Wainwright and Jordan, 2008; Bishop, 2006). The algorithm is derived from first-order conditions of the variational mean-field problem.

Let μ^* be the matrix that solves the variational problem (2.9). If we take the derivative with respect to μ_{ij} and equate to zero, we get

(4.1)
$$\mu_{ij}^* = \left\{ 1 + \exp\left[-2\alpha_{ij} - \beta n^{-1} \sum_{k=1}^n \left(\mu_{jk}^* + \mu_{ki}^* \right) - 4\gamma n^{-1} \sum_{k=1}^n \mu_{jk}^* \mu_{ki}^* \right] \right\}^{-1}$$

The logit equation in (4.1) characterizes a system of equations, whose fixed point is a solution of the mean-field problem. We can therefore start from a matrix μ and iterate the updates in (4.1) until we reach a fixed point, as described in the following algorithm.

ALGORITHM 1. Approximation of log-normalizing constant. Fix parameters α, β, γ and a relatively small tolerance value ϵ_{tol} . Initialize the $n \times n$ matrix $\mu^{(0)}$ as $\mu_{ij}^{(0)} \stackrel{iid}{\sim} U[0, 1]$, for all i, j. Fix the maximum number of iterations as T. Then for each t = 0, ..., T:

Step 1. Update the entries of matrix $\mu^{(t)}$ for all i, j = 1, ..., n

(4.2)
$$\mu_{ij}^{(t+1)} = \left\{ 1 + \exp\left[-2\alpha_{ij} - \beta n^{-1} \sum_{k=1}^{n} \left(\mu_{jk}^{(t)} + \mu_{ki}^{(t)} \right) - 4\gamma n^{-1} \sum_{k=1}^{n} \mu_{jk}^{(t)} \mu_{ki}^{(t)} \right] \right\}^{-1}$$

Step 2. Compute the value of the variational mean-field log-constant $\psi_n^{MF(t)}$ as

$$\psi_n^{MF(t)} = \frac{\sum_i \sum_j \alpha_{ij} \mu_{ij}^{(t)}}{n^2} + \beta \frac{\sum_i \sum_j \sum_k \mu_{ij}^{(t)} \mu_{jk}^{(t)}}{2n^3} + \gamma \frac{2 \sum_i \sum_j \sum_k \mu_{ij}^{(t)} \mu_{jk}^{(t)} \mu_{ki}^{(t)}}{3n^3} - \frac{1}{2n^2} \sum_{i=1}^n \sum_{j=1}^n \left[\mu_{ij}^{(t)} \log \mu_{ij}^{(t)} + (1 - \mu_{ij}^{(t)}) \log(1 - \mu_{ij}^{(t)}) \right].$$

Step 3. Stop at $t^* \leq T$ if: $\psi_n^{MF(t^*)} - \psi_n^{MF(t^*-1)} \leq \epsilon_{tol}$. Otherwise go back to Step 1.

The algorithm is initialized at a random uniform matrix $\mu^{(0)}$ and iteratively applies the update (4.1) to each entry of the matrix, until the increase in the objective function is less than a tolerance level. Since the problem is concave in each μ_{ij} , this iterative method is guaranteed to find a local maximum of (2.9).¹⁰ In our simulations we use a tolerance level of $\epsilon_{tol} = 0.0001$. To improve

¹⁰There are other alternatives to the random uniform matrix. Indeed a simple starting value could be the set of conditional probabilities of the model at parameters α , β , γ . We did not experiment with this alternative method.

convergence we can re-start the algorithm from different random matrices, as usually done with local optimizers.¹¹ This step is easily parallelizable, thus preserving the order n^2 convergence; while the standard MCMC-MLE is an intrinsically sequential algorithm and cannot be parallelized.

4.2. Monte Carlo design. All the computations in this section are performed on a PC Dell T6610 with 6 Quad-core Intel i7 (48 threads) and 64GB RAM. We test our approximation using 1000 simulated networks. Each node *i* has a binary attribute x_i , i.e. $x_i \stackrel{iid}{\sim} Bernoulli(0.5)$. Let $z_{ij} = 1$ if $x_i = x_j$ and $z_{ij} = 0$ otherwise.

$$(4.3) \quad t_z(g) := \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n g_{ij} z_{ij}; \ t_{-z}(g) := \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n g_{ij} (1 - z_{ij}),$$
$$t_e(g) := \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n g_{ij}; \ t_s(g) := \frac{1}{n^3} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n g_{ij} g_{jk}; \ t_t(g) := \frac{1}{n^3} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n g_{ij} g_{jk} g_{ki},$$

where $t_e(g)$, $t_s(g)$ and $t_t(g)$ are the fraction of links, two-stars and triangles respectively. And $t_z(g)$ and $t_{-z}(g)$ are the fractions of links of the same type and different type, respectively. The log-likelihood of the model $\ell_n(g; \alpha, \beta, \gamma)$ is

(4.4)
$$\ell_n(g, x; \alpha, \beta, \gamma) = \alpha_1 t_z(g) + \alpha_2 t_{-z}(g) + (\beta/2) t_s(g) + (2\gamma/3) t_t(g) - \psi_n(\alpha_1, \alpha_2, \beta, \gamma).$$

For computational convenience we rewrite model (4.4) in a slightly different but equivalent way

(4.5)
$$\ell_n(g, x; \tilde{\alpha}, \beta, \gamma) = \tilde{\alpha}_1 t_e(g) + \tilde{\alpha}_2 t_z(g) + (\beta/2) t_s(g) + (2\gamma/3) t_t(g) - \psi_n(\alpha_1, \alpha_2, \beta, \gamma) + (\beta/2) t_s(g) + (\beta/2) t_s$$

where we have defined $\tilde{\alpha}_1 := \alpha_2$ and $\tilde{\alpha}_2 := \alpha_1 - \alpha_2$. We use specification (4.5) in our simulations.¹²

To generate the artificial networks, we draw i.i.d. attributes $x_i \sim Bernoulli(0.5)$, initialize a network with n nodes as an Erdos-Renyi graph with probability $p = e^{\tilde{\alpha}_1}/(1 + e^{\tilde{\alpha}_1})$, and then run

¹¹In the Monte Carlo exercises we have experimented with different numbers of re-starts of the iterative algorithm. However, it is not clear what would be the optimal number of re-starts. A fixed number of restarts could be suboptimal. It seems reasonable to increase this number as the network grows larger.

¹²There are other small differences in how we have specified the model and how we have setup computations using the statnet package in R, that can affect the comparability of the simulation results, in particular the normalizations of the sufficient statistics. This is handled by our mfergm package, to guarantee comparability of the estimates obtained with MCMC-MLE, MPLE and Mean-field approximate inference.

the Metropolis-Hastings network sampler using the simulate.ergm command in the R package ergm to sample 1000 networks, each separated by 10,000 iterations, and after a burn-in of 10 million iterations.¹³ The MCMC-MLE estimator is solved using the Stochastic approximation method of Snijders (2002), where each simulation has a burnin of 100,000 iterations of the Metropolis-Hastings sampler and networks are sampled every 1000 iterations. The other convergence parameters are kept at default of the ergm package. The MPLE estimate is obtained using the default parameters in ergm. To be sure that our results do not depend on the initialization of the parameters, we start each estimator at the true parameter value, thus decreasing the computational time required for convergence. All the code is available in Github for replication.

TABLE 4.1. Monte Carlo estimates, comparison of three methods. True parameter vector is $(\tilde{\alpha}_1, \tilde{\alpha}_2, \beta, \gamma) = (-2, 1, 1, 1)$

n = 50		MCM	C-MLE			MEAN-	FIELD		MPLE			
	$\tilde{\alpha}_1$	\tilde{lpha}_2	β	γ	$\tilde{\alpha}_1$	\tilde{lpha}_2	β	γ	$\tilde{\alpha}_1$	\tilde{lpha}_2	β	γ
median	-2.002	1.024	0.716	-2.042	-2.000	0.998	1.000	0.999	-1.957	1.016	0.118	-0.584
mad	0.295	0.238	3.412	26.132	0.044	0.040	0.012	0.012	0.268	0.179	3.261	16.540
n = 100	MCMC-MLE					MEAN-	FIELD		MPLE			
median	-1.991	0.991	0.886	1.183	-2.002	0.995	1.001	0.999	-1.974	0.991	0.713	1.020
mad	0.197	0.117	2.324	16.150	0.020	0.017	0.005	0.005	0.178	0.085	2.237	10.478
n = 200		MCM	C-MLE		MEAN-FIELD				MPLE			
median	-2.000	1.000	1.043	0.438	-2.003	0.995	1.001	0.999	-1.990	1.000	0.853	0.657
mad	0.127	0.064	1.686	10.627	0.009	0.009	0.002	0.002	0.125	0.046	1.613	7.950
n = 500	MCMC-MLE				MEAN-FIELD				MPLE			
median	-2.000	1.001	1.000	0.706	-2.002	0.994	1.016	0.992	-1.994	1.001	0.912	0.762
mad	0.084	0.033	1.090	6.962	0.007	0.008	0.023	0.011	0.074	0.023	0.945	4.691

Results of 1000 Monte Carlo estimates using three methods. MCMC-MLE is the Monte Carlo Maximum Likelihood estimator of Geyer and Thompson (1992), as implemented in ergm in R, with a stochastic approximation algorithm Snijders (2002). MEAN-FIELD is our method. MPLE is the Maximum Pseudo-Likelihood Estimate. Each network is generated with a 10 million run of the Metropolis-Hastings sampler of the ergm command in R, sampling every 10000 iterations. mad is the median absolute deviation.

4.3. **Results.** The first model has true parameter vector $(\tilde{\alpha}_1, \tilde{\alpha}_2, \beta, \gamma) = (-2, 1, 1, 1)$ and the summaries of point estimates are reported in Table 4.1. We show results for n = 50, 100, 200 and 500; reporting median and median absolute deviation (mad) of point estimates for each parameter.

¹³The code is available in the Github package mfergm, and the function is simulate.model#, where # stands for the model number: 2 is the model with $\gamma = 0$, 3 is the model with $\beta = 0$, and 4 is the model with $\beta \neq 0$ and $\gamma \neq 0$.

The median estimates of the mean-field approximation are quite stable and exhibit a small bias, as is well known in the literature (Wainwright and Jordan, 2008; Bishop, 2006). The median results for MCMC-MLE and MPLE are quite precise for $\tilde{\alpha}_1$ and $\tilde{\alpha}_2$, but vary a lot for β and γ , as shown by the large median absolute deviation. Nonetheless the median point estimates of β and γ are slowly converging to the true parameter vector as n increases.¹⁴ Therefore, the mean-field approximation provides estimates in line with MPLE and MCMC-MLE, with more reliability for β and γ in these small sample estimation exercises.

TABLE 4.2. Monte Carlo estimates, comparison of three methods. True parameter vector is $(\tilde{\alpha}_1, \tilde{\alpha}_2, \beta, \gamma) = (-3, 2, 1, 3)$

n = 50		MCMC-MLE				MEAN-	FIELD			MPLE			
	\tilde{lpha}_1	\tilde{lpha}_2	β	γ	\tilde{lpha}_1	\tilde{lpha}_2	β	γ	\tilde{lpha}_1	\tilde{lpha}_2	β	γ	
median	-3.041	2.064	0.743	-0.512	-3.007	1.993	1.000	3.000	-3.026	2.083	0.215	1.764	
mad	0.476	0.424	3.811	25.109	0.026	0.026	0.013	0.014	0.514	0.401	3.593	16.538	
n = 100		MCM	C-MLE			MEAN-	FIELD			MI	PLE		
median	-3.006 2.015 0.932 0.587				-3.011	1.989	1.000	2.999	-2.991	2.018	0.682	1.773	
mad	0.261	0.206	2.538	17.905	0.016	0.016	0.008	0.008	0.259	0.194	2.364	12.123	
n = 200		MCM	C-MLE		MEAN-FIELD				MPLE				
median	-3.012	2.007	1.069	2.807	-3.011	1.988	1.000	2.999	-3.005	2.011	0.932	2.988	
mad	0.158	0.117	1.822	11.360	0.008	0.008	0.004	0.004	0.156	0.109	1.714	8.144	
n = 500		MCM	C-MLE		MEAN-FIELD			MPLE					
median	-2.998	2.000	0.951	3.047	-3.011	1.988	1.002	2.999	-2.998	2.001	0.921	3.117	
mad	0.096	0.061	1.276	7.191	0.003	0.003	0.002	0.002	0.083	0.049	1.077	5.378	

Notes: see notes for Table 4.1.

TABLE 4.3. Monte Carlo estimates, comparison of three methods. True parameter vector is $(\tilde{\alpha}_1, \tilde{\alpha}_2, \beta, \gamma) = (-3, 1, 2, 1)$

n = 500		MCM	C-MLE		MEAN-FIELD				MPLE				
	\tilde{lpha}_1	\tilde{lpha}_2	β	γ	\tilde{lpha}_1	\tilde{lpha}_2	β	γ	\tilde{lpha}_1	\tilde{lpha}_2	β	γ	
median	-3.001	0.998	2.028	-19.034	-3.000	1.000	2.000	1.000	-2.996	1.000	1.488	-7.923	
mad	0.086	0.065	7.205	165.600	0.011	0.011	0.0001	0.0001	0.078	0.044	6.345	84.681	
n = 1000		MCM	C-MLE		MEAN-FIELD				MPLE				
	\tilde{lpha}_1	\tilde{lpha}_2	β	γ	\tilde{lpha}_1	\tilde{lpha}_2	β	γ	$\tilde{\alpha}_1$	\tilde{lpha}_2	β	γ	
median	-2.999	1.004	1.809	-0.716	-3.000	1.000	2.000	1.000	-2.999	1.002	1.757	0.540	
mad	0.057	0.037	4.891	125.293	0.005	0.005	0.0001	0.0001	0.049	0.022	4.113	61.328	

Notes: see notes for Table 4.1. The case with n = 1000 contains only 500 monte carlo replications.

¹⁴Some of the bias in the mean-field approximation may be due to the fact that we only initialize μ once in these simulations.

The second set of results is for a model with parameters $(\tilde{\alpha}_1, \tilde{\alpha}_2, \beta, \gamma) = (-3, 2, 1, 3)$, see Table 4.2. The pattern is similar to Table 4.1. Indeed the mean-field estimator seems to work relatively well in most cases, especially for the estimates of β and γ . For parameters $\tilde{\alpha}_1, \tilde{\alpha}_2$ our mean-field estimator (median) bias persists as n increases. Finally, we also report a simulation with a larger network with n = 500, 1000 in Table 4.3. The results are the same as the other tables and the mean-field approximation is robustly close to the true parameter values in most simulations.

These Monte Carlo experiments suggests that our approximation method performs well in practice. We conclude that in most cases the mean-field approximation algorithm works better than our conservative theoretical results suggest.¹⁵

5. CONCLUSIONS AND FUTURE WORK

We have shown that for a large class of exponential random graph models (ERGM), we can approximate the normalizing constant of the likelihood using a mean-field variational approximation algorithm (Wainwright and Jordan, 2008; Bishop, 2006; Chatterjee and Diaconis, 2013; Mele, 2017). Our theoretical results use nonlinear large deviations methods (Chatterjee and Dembo, 2016) to bound the error of approximation, showing that it converges to zero as the network grows.

Our estimation method consists of replacing the log-normalizing constant in the log-likelihood of the ERGM with the value approximated by the mean-field algorithm; we then find the parameters that maximize such approximate log-likelihood. Since our approximated constant converges to the true constant in large networks, the approximate log-likelihood converges to the correct loglikelihood in sup-norm, as the network becomes large. If the likelihoods are well-behaved and not too flat around the maximizers, we can also show that our estimate converges to MLE.

Using an iterative procedure to find the approximate mean-field constant, we compare our method to MCMC-MLE and MPLE (Snijders, 2002; Boucher, 2015; Besag, 1974; DePaula, 2017) in a simple Monte Carlo study for small networks. The mean-field approximation exhibits a small bias, but the median estimates are similar to MCMC-MLE and MPLE. Theoretically, our method

¹⁵While these results are encouraging, in Appendix we report some example of non-convergence of the mean-field algorithm, mostly due to our iterative algorithm getting trapped in a local maximum in some simulations.

converges in a number of steps proportional to the number of potential links of a network, while MCMC-MLE could be exponentially slow.

While these results are encouraging, there are several open problems and possible research directions. First, it is not clear that the mean-field estimates are consistent. Our small Monte Carlo seem to indicate that there is a persistent bias term, but there is no general proof in this setting along the lines of Bickel et al. (2013) for stochastic block models. Second, it is not clear that the ERGM model is identified for all parameter values. Indeed some results in this literature suggest otherwise (Chatterjee and Diaconis, 2013; Mele, 2017; Boucher and Mourifie, 2017). A promising research avenue for the future is the use of the large n mean-field approximation to understand identification, similarly to what has been done with graph limits in Chatterjee and Diaconis (2013). Third, while the mean-field approximation is simple and we are able to compute the approximation errors, our lower and upper bounds may not be sharp. This raises the question of whether there is another factorization (like in structured mean-field) that leads to better approximations and faster convergence (Xing et al., 2003). We hope that our work will stimulate additional research and more applications of this class of approximations.

REFERENCES

- Airoldi, Edoardo M., David Blei, Stephen E. Fienberg and Eric P. Xing (2008), 'Mixed membership stochastic blockmodels', *Journal of Machine Learning Research* **9**, 1981–2014.
- Amir, Eyal, Wen Pu and Dorothy Espelage (2012), Approximating partition functions for exponential-family random graph models, *in* 'Advances in Neural Information Processing Systems (NIPS)'.
- Aristoff, David and Lingjiong Zhu (2018), 'On the phase transition curve in a directed exponential random graph model', *Advances in Applied Probability* **50**, 272–301.
- Badev, Anton (2013), Discrete games in endogenous networks: Theory and policy.
- Banerjee, Abhijit, Arun G. Chandrasekhar, Esther Duflo and Matthew O. Jackson (2013), 'The diffusion of microfinance', *Science* **341**(6144).
- Besag, Julian (1974), 'Spatial interaction and the statistical analysis of lattice systems', *Journal of the Royal Statistical Society Series B (Methodological)* **36**(2), 192–236.
- Bhamidi, Shankar, Guy Bresler and Allan Sly (2011), 'Mixing time of exponential random graphs', *The Annals of Applied Probability* **21**(6), 2146–2170.
- Bickel, Peter, David Choi, Xiangyu Chang and Hai Zhang (2013), 'Asymptotic normality of maximum likelihood and its variational approximation for stochastic blockmodels', *Ann. Statist.*41(4), 1922–1943.
- Bishop, Christopher (2006), Pattern Recognition and Machine Learning, Springer, New York.
- Blume, Lawrence E. (1993), 'The statistical mechanics of strategic interaction', *Games and Economic Behavior* **5**(3), 387–424.
- Borgs, C., J.T. Chayes, L. Lovász, V.T. Sós and K. Vesztergombi (2008), 'Convergent sequences of dense graphs i: Subgraph frequencies, metric properties and testing', *Advances in Mathematics* 219(6), 1801 – 1851.
- Boucher, Vincent (2015), 'Structural homophily', International Economic Review 56(1), 235–264.
- Boucher, Vincent and Ismael Mourifie (2017), 'My friends far far away: A random field approach to exponential random graph models', *Econometrics Journal* **20**(3), S14–S46.

Butts, Carter (2009), Using potential games to parameterize ERG models. working paper.

- Caimo, Alberto and Nial Friel (2011), 'Bayesian inference for exponential random graph models', *Social Networks* **33**(1), 41–55.
- Chandrasekhar, Arun (2016), *in* Y.Bramoulle, A.Galeotti and B. W.Rogers, eds, 'The Oxford Handbook of the Economics of Networks', Oxford University Press, chapter Econometrics of Network Formation.
- Chandrasekhar, Arun and Matthew Jackson (2014), Tractable and consistent exponential random graph models. working paper.
- Chatterjee, Sourav and Amir Dembo (2016), 'Nonlinear large deviations', *Advances in Mathematics* **299**, 396–450.
- Chatterjee, Sourav and Persi Diaconis (2013), 'Estimating and understanding exponential random graph models', *The Annals of Statistics* **41**(5).
- Chatterjee, Sourav and S. R. S. Varadhan (2011), 'The large deviation principle for the Erdos-Rényi random graph', *European Journal of Combinatorics* **32**(7), 1000 1017.
- DePaula, Aureo (2017), Econometrics of network models, *in* B.Honore, A.Pakes, M.Piazzesi and L.Samuelson, eds, 'Advances in Economics and Econometrics: Eleventh World Congress', Cambridge University Press.
- DePaula, Aureo, Seth Richards-Shubik and Elie Tamer (2018), 'Identifying preferences in networks with bounded degree', *Econometrica* **86**(1), 263–288.
- Durlauf, Steven N. and Yannis M. Ioannides (2010), 'Social interactions', Annual Review of Economics 2(1), 451–478.
- Fox, Jeremy T. and Natalia Lazzati (2017), 'A note on identification of discrete choice models for bundles and binary games', *Quantitative Economics* **8**(3), 1021–1036.
- Geyer, Charles and Elizabeth Thompson (1992), 'Constrained Monte Carlo maximum likelihood for dependent data', *Journal of the Royal Statistical Society, Series B (Methodological)* **54**(3), 657–699.
- Goodreau, S. M., Kitts J. A. and Morris M. (2009), 'Birds of a feather, or friend of a friend? using exponential random graph models to investigate adolescent social networks', *Demography* 46(1), 103–125.

- Graham, Bryan (2017), 'An empirical model of network formation: with degree heterogeneity', *Econometrica* **85**(4), 1033–1063.
- He, Ran and Tian Zheng (2013), Estimation of exponential random graph models for large social networks via graph limits, *in* 'Proceedings of the 2013 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining', ASONAM '13, ACM, New York, NY, USA, pp. 248–255.
- Iijima, Ryota and Yuichiro Kamada (2014), Social distance and network structures. Working Paper. Jackson, Matthew O. (2010), *Social and Economics Networks*, Princeton University Press.
- Koskinen, Johan (2004), Bayesian analysis of exponential random graphs estimation of parameters and model selection, Research report 2004:2, Department of Statistics, Stockholm University.
- Kosyakova, Tetyana, Thomas Otter, Sanjog Misra and Christian Neuerburg (2018), Exact MCMC for choices from menus measuring substitution and complementarity among menu items.
- Lovasz, L. (2012), *Large Networks and Graph Limits*, American Mathematical Society colloquium publications, American Mathematical Society.
- Mele, Angelo (2011), Segregation in social networks: Monte Carlo maximum likelihood estimation. Working Paper.
- Mele, Angelo (2017), 'A structural model of dense network formation', *Econometrica* **85**, 825–850.
- Moller, Jesper and Rasmus Plenge Waagepetersen (2004), *Statistical inference and simulation for spatial point processes*, Monographs on Statistics and Applied Probability 100, Chapman and Hall.
- Monderer, Dov and Lloyd Shapley (1996), 'Potential games', *Games and Economic Behavior* **14**, 124–143.
- Moody, James (2001), 'Race, school integration, and friendship segregation in America', *American Journal of Sociology* **103**(7), 679–716.
- Murray, Iain A., Zoubin Ghahramani and David J. C. MacKay (2006), MCMC for doublyintractable distributions, *in* 'Proceedings of the Twenty-Second Conference on Uncertainty in

Artificial Intelligence', pp. 359–366.

- Radin, Charles and Mei Yin (2013), 'Phase transitions in exponential random graphs', *The Annals of Applied Probability* **23**(6), 2458–2471.
- Sheng, Shuyang (2012), Identification and estimation of network formation games. working paper.
- Snijders, Tom A.B (2002), 'Markov chain Monte Carlo estimation of exponential random graph models', *Journal of Social Structure* **3**(2).
- Train, Kenneth (2009), Discrete Choice Methods with Simulation, Cambridge University Press.
- Wainwright, M.J. and M.I. Jordan (2008), 'Graphical models, exponential families, and variational inference', *Foundations and Trends@ in Machine Learning* **1**(1-2), 1–305.
- Wasserman, Stanley and Philippa Pattison (1996), 'Logit models and logistic regressions for social networks: I. An introduction to Markov graphs and p*', *Psychometrika* **61**(3), 401–425.
- Wimmer, Andreas and Kevin Lewis (2010), 'Beyond and below racial homophily: ERG models of a friendship network documented on Facebook', *American Journal of Sociology* **116**(2), 583– 642.
- Xing, Eric P., Michael I. Jordan and Stuart Russell (2003), A generalized mean field algorithm for variational inference in exponential families, *in* 'Proceedings of the Nineteenth Conference on Uncertainty in Artificial Intelligence', UAI'03, pp. 583–591.

APPENDIX

A.1. **Proof of Theorem 3.1.** In this proof we will try to follow closely the notation in Chatterjee and Dembo (2016). Suppose that $f : [0, 1]^N \to \mathbb{R}$ is twice continuously differentiable in $(0, 1)^N$, so that f and all its first and second order derivatives extend continuously to the boundary. Let ||f||denote the supremum norm of $f : [0, 1]^N \to \mathbb{R}$. For each i and j, denote

(A.1)
$$f_i := \frac{\partial f}{\partial x_i}, \qquad f_{ij} := \frac{\partial^2 f}{\partial x_i \partial x_j},$$

and let

(A.2)
$$a := ||f||, \quad b_i := ||f_i||, \quad c_{ij} := ||f_{ij}||.$$

Given $\epsilon > 0$, $\mathcal{D}(\epsilon)$ is the finite subset of \mathbb{R}^N so that for any $x \in \{0,1\}^N$, there exists $d = (d_1, \ldots, d_N) \in \mathcal{D}(\epsilon)$ such that

(A.3)
$$\sum_{i=1}^{N} (f_i(x) - d_i)^2 \le N\epsilon^2.$$

Let us define

(A.4)
$$F := \log \sum_{x \in \{0,1\}^N} e^{f(x)},$$

and for any $x = (x_1, ..., x_N) \in [0, 1]^N$,

(A.5)
$$I(x) := \sum_{i=1}^{N} [x_i \log x_i + (1 - x_i) \log(1 - x_i)]$$

In the proof we rely on Theorem 1.5 in Chatterjee and Dembo (2016) that we reproduce in Theorem A.1 to help the reader:

THEOREM A.1 (Chatterjee and Dembo (2016)). For any $\epsilon > 0$,

(A.6)
$$\sup_{x \in [0,1]^N} \{f(x) - I(x)\} - \frac{1}{2} \sum_{i=1}^N c_{ii} \le F \le \sup_{x \in [0,1]^N} \{f(x) - I(x)\} + \mathcal{E}_1 + \mathcal{E}_2,$$

where

(A.7)
$$\mathcal{E}_1 := \frac{1}{4} \left(N \sum_{i=1}^N b_i^2 \right)^{1/2} \epsilon + 3N\epsilon + \log |\mathcal{D}(\epsilon)|,$$

and

(A.8)
$$\mathcal{E}_{2} := 4 \left(\sum_{i=1}^{N} (ac_{ii} + b_{i}^{2}) + \frac{1}{4} \sum_{i,j=1}^{N} (ac_{ij}^{2} + b_{i}b_{j}c_{ij} + 4b_{i}c_{ij}) \right)^{1/2} + \frac{1}{4} \left(\sum_{i=1}^{N} b_{i}^{2} \right)^{1/2} \left(\sum_{i=1}^{N} c_{ii}^{2} \right)^{1/2} + 3 \sum_{i=1}^{N} c_{ii} + \log 2.$$

We will use the Theorem A.1 to derive the lower and upper bound of the mean-field approximation problem. Our results extend Theorem 1.7. in Chatterjee and Dembo (2016) from the ERGM with two-stars and triangles to the model that allows nodal covariates. Notice that in our case the N of the theorem is the number of links, i.e. $N = \binom{n}{2}$. Let

(A.9)
$$Z_n := \sum_{x_{ij} \in \{0,1\}, x_{ij} = x_{ji}, 1 \le i < j \le n} e^{\sum_{1 \le i, j \le n} \alpha_{ij} x_{ij} + \frac{\beta}{2n} \sum_{1 \le i, j, k \le n} x_{ij} x_{jk} + \frac{2\gamma}{3n} \sum_{1 \le i, j, k \le n} x_{ij} x_{jk} x_{ki}},$$

be the normalizing factor and also define

(A.10)
$$L_n := \sup_{x_{ij} \in [0,1], x_{ij} = x_{ji}, 1 \le i < j \le n} \left\{ \frac{1}{n^2} \sum_{i,j} \alpha_{ij} x_{ij} + \frac{\beta}{2n^3} \sum_{i,j,k} x_{ij} x_{jk} + \frac{2\gamma}{3n^3} \sum_{i,j,k} x_{ij} x_{jk} x_{ki} - \frac{1}{n^2} \sum_{1 \le i < j \le n} [x_{ij} \log x_{ij} + (1 - x_{ij}) \log(1 - x_{ij})] \right\}.$$

Notice that $n^{-2}Z_n = \psi_n$ and $L_n = \psi_n^{MF}$. For our model, the function $f : [0, 1]^{\binom{n}{2}} \to \mathbb{R}$ is defined as

(A.11)
$$f(x) = \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_{ij} x_{ij} + \frac{\beta}{2n} \sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{k=1}^{n} x_{ij} x_{jk} + \frac{2\gamma}{3n} \sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{k=1}^{n} x_{ij} x_{jk} x_{ki}.$$

Then, we can compute that,

(A.12)
$$a = \|f\| \le \sum_{i=1}^{n} \sum_{j=1}^{n} |\alpha_{ij}| + \frac{1}{2} |\beta| n^{2} + \frac{2}{3} |\gamma| n^{2}$$
$$\le n^{2} \left[\max_{i,j} |\alpha_{i,j}| + \frac{1}{2} |\beta| + \frac{2}{3} |\gamma| \right].$$

Let $k \in \mathbb{N}$, and H be a finite simple graph on the vertex set $[k] := \{1, \ldots, k\}$. Let E be the set of edges of H and |E| be its cardinality. For a function $T : [0, 1]^{\binom{n}{2}} \to \mathbb{R}$

(A.13)
$$T(x) := \frac{1}{n^{k-2}} \sum_{q \in [n]^k} \prod_{\{\ell, \ell'\} \in E} x_{q_\ell q_{\ell'}}$$

Chatterjee and Dembo (2016) (Lemma 5.1.) showed that, for any i < j, i' < j',

(A.14)
$$\left\|\frac{\partial T}{\partial x_{ij}}\right\| \le 2|E|,$$

and

(A.15)
$$\left\|\frac{\partial^2 T}{\partial x_{ij}\partial x_{i'j'}}\right\| \le \begin{cases} 4|E|(|E|-1)n^{-1} & \text{if } |\{i,j,i',j'\}| = 2 \text{ or } 3, \\ 4|E|(|E|-1)n^{-2} & \text{if } |\{i,j,i',j'\}| = 4. \end{cases}$$

Therefore, by (A.14), we can compute that

(A.16)
$$b_{(ij)} = \left\| \frac{\partial f}{\partial x_{ij}} \right\| \le 2 \max_{i,j} |\alpha_{ij}| + 2|\beta| + 8|\gamma|.$$

By (A.15), we can also compute that

$$(A.17) c_{(i,j)(i'j')} = \left\| \frac{\partial^2 f}{\partial x_{ij} \partial x_{i'j'}} \right\| \\ \leq \begin{cases} 4\left(\frac{1}{2}|\beta|2(2-1) + \frac{2}{3}|\gamma|3(3-1)\right)n^{-1} & \text{if } |\{i,j,i',j'\}| = 2 \text{ or } 3, \\ 4\left(\frac{1}{2}|\beta|2(2-1) + \frac{2}{3}|\gamma|3(3-1)\right)n^{-2} & \text{if } |\{i,j,i',j'\}| = 4, \end{cases} \\ = \begin{cases} 4\left(|\beta| + 4|\gamma|\right)n^{-1} & \text{if } |\{i,j,i',j'\}| = 2 \text{ or } 3, \\ 4\left(|\beta| + 4|\gamma|\right)n^{-2} & \text{if } |\{i,j,i',j'\}| = 4. \end{cases}$$

Next, we compute that

(A.18)
$$\frac{\partial f}{\partial x_{ij}} = 2\alpha_{ij} + \frac{\partial}{\partial x_{ij}} \left[\frac{\beta}{2n} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n x_{ij} x_{jk} + \frac{2\gamma}{3n} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n x_{ij} x_{jk} x_{ki} \right].$$

Let T_1 and T_2 be defined as

(A.19)
$$T_1(x) := \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n x_{ij} x_{jk}, \qquad T_2(x) := \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n x_{ij} x_{jk} x_{ki}.$$

Then, we have

(A.20)
$$\frac{\partial f}{\partial x_{ij}} = 2\alpha_{ij} + \frac{\beta}{2}\frac{\partial T_1}{\partial x_{ij}} + \frac{2\gamma}{3}\frac{\partial T_2}{\partial x_{ij}}.$$

Chatterjee and Dembo (2016) (Lemma 5.2.) showed that for the T_1 and T_2 defined above, there exist a set $\mathcal{D}_1(\epsilon)$ and $\mathcal{D}_2(\epsilon)$ satisfying the criterion (A.3) (with $f = T_1$ and $f = T_2$) so that

(A.21)
$$|\mathcal{D}_1(\epsilon)| \le \exp\left\{\frac{\tilde{C}_1 2^4 3^4 n}{\epsilon^4} \log \frac{\tilde{C}_2 2^4 3^4}{\epsilon^4}\right\} = \exp\left\{\frac{\tilde{C}_1 6^4 n}{\epsilon^4} \log \frac{\tilde{C}_2 6^4}{\epsilon^4}\right\},$$

(A.22)
$$|\mathcal{D}_2(\epsilon)| \le \exp\left\{\frac{\tilde{C}_1 3^4 3^4 n}{\epsilon^4} \log \frac{\tilde{C}_2 3^4 3^4}{\epsilon^4}\right\} = \exp\left\{\frac{\tilde{C}_1 3^8 n}{\epsilon^4} \log \frac{\tilde{C}_2 3^8}{\epsilon^4}\right\},$$

where \tilde{C}_1 and \tilde{C}_2 are universal constants. Let us define

(A.23)

$$\mathcal{D}(\epsilon) := \left\{ 2\alpha_{ij} + \frac{\beta}{2}d_1 + \frac{2\gamma}{3}d_2 : d_1 \in \mathcal{D}_1\left(\frac{2}{\beta} \cdot \frac{\epsilon}{\sqrt{2}}\right), d_2 \in \mathcal{D}_2\left(\frac{3}{2\gamma} \cdot \frac{\epsilon}{\sqrt{2}}\right), 1 \le i \le j \le n \right\}.$$

Hence, $\mathcal{D}(\epsilon)$ satisfies the criterion (A.3) and

$$\begin{aligned} (\mathbf{A}.24) \quad |\mathcal{D}(\epsilon)| &\leq \frac{1}{2}n(n+1) \left| \mathcal{D}_1\left(\sqrt{2}\epsilon/\beta\right) \right| \cdot \left| \mathcal{D}_2\left(3\epsilon/2\sqrt{2}\gamma\right) \right| \\ &\leq \frac{1}{2}n(n+1) \exp\left\{ \frac{\tilde{C}_1 6^4\beta^4 n}{4\epsilon^4} \log \frac{\tilde{C}_2 6^4\beta^4}{4\epsilon^4} \right\} \exp\left\{ \frac{\tilde{C}_1 3^8 2^6\gamma^4 n}{3^4\epsilon^4} \log \frac{\tilde{C}_2 3^8 2^6\gamma^4}{3^4\epsilon^4} \right\}. \end{aligned}$$

Therefore, by recalling \mathcal{E}_1 from (A.7), we get

$$(A.25) \quad \mathcal{E}_{1} = \frac{1}{4} \left(\binom{n}{2} \sum_{1 \le i < j \le n} b_{(ij)}^{2} \right)^{1/2} \epsilon + 3\binom{n}{2} \epsilon + \log |\mathcal{D}(\epsilon)|$$

$$\leq \left[\frac{1}{4} \left(2 \max_{i,j} |\alpha_{ij}| + 2|\beta| + 8|\gamma| \right) + 3 \right] \binom{n}{2} \epsilon$$

$$+ \log \left(\frac{1}{2} n(n+1) \right) + \frac{\tilde{C}_{1} 6^{4} \beta^{4} n}{4\epsilon^{4}} \log \frac{\tilde{C}_{2} 6^{4} \beta^{4}}{4\epsilon^{4}} + \frac{\tilde{C}_{1} 3^{4} 2^{6} \gamma^{4} n}{\epsilon^{4}} \log \frac{\tilde{C}_{2} 3^{4} 2^{6} \gamma^{4}}{\epsilon^{4}} \right]$$

$$\leq C_{1}(\alpha, \beta, \gamma) n^{2} \epsilon + \frac{C_{1}(\alpha, \beta, \gamma) n}{\epsilon^{4}} \log \frac{C_{1}(\alpha, \beta, \gamma)}{\epsilon^{4}}$$

$$= C_{1}(\alpha, \beta, \gamma) n^{9/5} (\log n)^{1/5},$$

by choosing $\epsilon = (\frac{\log n}{n})^{1/5}$, where $C_1(\alpha, \beta, \gamma)$ is a constant depending only on α, β, γ :

(A.26)
$$C_1(\alpha, \beta, \gamma) := c_1 \left(\max_{i,j} |\alpha_{ij}| + |\beta|^4 + |\gamma|^4 + 1 \right),$$

where $c_1 > 0$ is some universal constant. To see why we can choose $C_1(\alpha, \beta, \gamma)$ as in (A.26) so that (A.25) holds, we first notice that it follows from (A.25) that we can choose $C_1(\alpha, \beta, \gamma)$ such that $C_1(\alpha, \beta, \gamma) \ge \max\{\tilde{c}_1 \max_{ij} |\alpha_{ij}| + \tilde{c}_2 |\beta| + \tilde{c}_3 |\gamma| + \tilde{c}_4, \tilde{c}_5 \beta^4, \tilde{c}_6 \gamma^4\}$, where $\tilde{c}_1, \tilde{c}_2, \tilde{c}_3, \tilde{c}_4, \tilde{c}_5, \tilde{c}_6 > 0$ are some universal constants. Note that $\max\{\tilde{c}_1 \max_{ij} |\alpha_{ij}| + \tilde{c}_2 |\beta| + \tilde{c}_3 |\gamma| + \tilde{c}_4, \tilde{c}_5 \beta^4, \tilde{c}_6 \gamma^4\} \le$ $\tilde{c}_1 \max_{ij} |\alpha_{ij}| + \tilde{c}_2 |\beta| + \tilde{c}_3 |\gamma| + \tilde{c}_4 + \tilde{c}_5 \beta^4 + \tilde{c}_6 \gamma^4 \le c_1 (\max_{i,j} |\alpha_{ij}| + |\beta|^4 + |\gamma|^4 + 1)$ for some universal constant $c_1 > 0$. Thus, we can take $C_1(\alpha, \beta, \gamma)$ as in (A.26).

We can also compute from (A.8) that

$$\mathcal{E}_{2} = 4 \bigg(\sum_{1 \le i < j \le n} (ac_{(ij)(ij)} + b_{(ij)}^{2}) + \frac{1}{4} \sum_{1 \le i < j \le n, 1 \le i' < j' \le n} (ac_{(ij)(i'j')}^{2} + b_{(ij)}b_{(i'j')}c_{(ij)(i'j')} + 4b_{(ij)}c_{(ij)(i'j')}) \bigg)^{1/2} + \frac{1}{4} \bigg(\sum_{1 \le i < j \le n} b_{(ij)}^{2} \bigg)^{1/2} \bigg(\sum_{1 \le i < j \le n} c_{(ij)(ij)}^{2} \bigg)^{1/2} + 3 \sum_{1 \le i < j \le n} c_{(ij)(ij)} + \log 2,$$

so that

$$\begin{split} \mathcal{E}_{2} &\leq 4 \left\{ \binom{n}{2} \left(n \left(\max_{i,j} |\alpha_{ij}| + \frac{1}{2} |\beta| + \frac{2}{3} |\gamma| \right) 4(|\beta| + 4|\gamma|) + \left(2 \max_{i,j} |\alpha_{ij}| + 2|\beta| + 8|\gamma| \right)^{2} \right) \\ &\quad + \frac{1}{4} n^{2} \left[\max_{i,j} |\alpha_{ij}| + \frac{1}{2} |\beta| + \frac{2}{3} |\gamma| \right] \\ &\quad \cdot \left[\binom{n}{2} \binom{n-2}{2} 4^{2} (|\beta| + 4|\gamma|)^{2} n^{-4} + \left(\binom{n}{2}^{2} - \binom{n}{2} \binom{n-2}{2} \right) 4^{2} (|\beta| + 4|\gamma|)^{2} n^{-2} \right] \\ &\quad + \left(2 \max_{i,j} |\alpha_{ij}| + 2|\beta| + 8|\gamma| \right) \cdot \left(\max_{i,j} |\alpha_{ij}| + \frac{1}{2} |\beta| + \frac{2}{3} |\gamma| \right) \\ &\quad \cdot \left[\binom{n}{2} \binom{n-2}{2} 4(|\beta| + 4|\gamma|) n^{-2} + \left(\binom{n}{2}^{2} - \binom{n}{2} \binom{n-2}{2} \right) 4(|\beta| + 4|\gamma|) n^{-1} \right] \right\}^{1/2} \\ &\quad + \frac{1}{4} \binom{n}{2} \left(2 \max_{i,j} |\alpha_{ij}| + 2|\beta| + 8|\gamma| \right) 4(|\beta| + 4|\gamma|) n^{-1} + 3\binom{n}{2} 4(|\beta| + 4|\gamma|) n^{-1} + \log 2 \\ &\leq C_{2}(\alpha, \beta, \gamma) n^{3/2}, \end{split}$$

where we used the formulas for a, $b_{(ij)}$, and $c_{(ij)(i'j')}$ that we derived earlier and the combinatorics identities:

$$\sum_{1 \le i < j \le n, 1 \le i' < j' \le n, |\{i, j, i', j'\}| = 4} 1 = \sum_{1 \le i < j \le n} \sum_{1 \le i' < j' \le n, |\{i, j, i', j'\}| = 4} 1 = \binom{n}{2} \binom{n-2}{2},$$
$$\sum_{1 \le i < j \le n, 1 \le i' < j' \le n, |\{i, j, i', j'\}| = 2 \text{ or } 3} 1 = \binom{n}{2}^2 - \binom{n}{2}\binom{n-2}{2},$$

and $C_2(\alpha,\beta,\gamma)$ is a constant depending only on α,β,γ that can be chosen as:

(A.27)
$$C_2(\alpha,\beta,\gamma) := c_2 \left(\max_{i,j} |\alpha_{ij}| + |\beta| + |\gamma| + 1 \right)^{1/2} (1 + |\beta|^2 + |\gamma|^2)^{1/2},$$

where $c_2 > 0$ is some universal constant.

Finally, to get lower bound, notice that

(A.28)
$$\frac{1}{2} \sum_{1 \le i < j \le n} c_{(ij)(ij)} \le \frac{1}{2} \binom{n}{2} 4(|\beta| + 4|\gamma|) n^{-1} \le C_3(\beta, \gamma) n,$$

where $C_3(\beta, \gamma)$ is a constant depending only on β, γ and we can simply take $C_3(\beta, \gamma) = |\beta| + 4|\gamma|$.

A.2. **Proof of Proposition 3.1.** We can approximate ψ_n by ψ_n^{MF} as seen in Theorem 3.1, and as a result, we can approximate the log-likelihood as follows.

$$\ell_n(g,\alpha,\beta,\gamma) := \frac{1}{n^2} \log(\pi_n(g,\alpha,\beta,\gamma)) = T_n(g,\alpha,\beta,\gamma) - \psi_n(\alpha,\beta,\gamma),$$

by the mean-field log-likelihood:

$$\ell_n^{MF}(g,\alpha,\beta,\gamma) := T_n(g,\alpha,\beta,\gamma) - \psi_n^{MF}(\alpha,\beta,\gamma),$$

Then the difference between the mean-field likelihood and the ERGM likelihood is bounded uniformly over $g \in \mathcal{G}$, for any α, β, γ :

$$0 \leq \ell_n^{MF}(g,\alpha,\beta,\gamma) - \ell_n(g,\alpha,\beta,\gamma) \leq C_1(\alpha,\beta,\gamma)n^{-1/5}(\log n)^{1/5} + C_2(\alpha,\beta,\gamma)n^{-1/2}.$$

Therefore, for any compact Θ , we have

$$0 \leq \sup_{\alpha,\beta,\gamma\in\Theta} \left[\ell_n^{MF}(g,\alpha,\beta,\gamma) - \ell_n(g,\alpha,\beta,\gamma) \right]$$

$$\leq \sup_{\alpha,\beta,\gamma\in\Theta} \left[C_1(\alpha,\beta,\gamma) n^{-1/5} (\log n)^{1/5} + C_2(\alpha,\beta,\gamma) n^{-1/2} \right]$$

$$\leq \sup_{\alpha,\beta,\gamma\in\Theta} C_1(\alpha,\beta,\gamma) n^{-1/5} (\log n)^{1/5} + \sup_{\alpha,\beta,\gamma\in\Theta} C_2(\alpha,\beta,\gamma) n^{-1/2}.$$

This proves the result.

APPENDIX B. A BOUND BETWEEN MLE AND MEAN-FIELD ESTIMATOR

We use the bounds on the likelihoods to also derive a bound on the distance between the MLE and our mean-field estimator, when the MLE exists and it is well-behaved. Because our bounds may not be sharp, this proves to be quite hard. We therefore, consider a *local* version of this convergence. We know that the ERGM likelihood is concave in parameters because it is an exponential family. We also know that the mean-field log-constant is convex in parameters¹⁶, therefore the approximate log-likelihood is also concave. However, to get a bound on the distance between

 $[\]overline{{}^{16}\psi_n^{MF}}$ is convex in (α, β, γ) by its definition in (2.9) since the expression inside the supremum in (2.9) is affine in (α, β, γ) and supremum over any affine function is convex.

estimates we need well-behaved objective functions, with enough curvature at least close to their maximizers. If the objective functions is too flat, the distance between the estimator may be too large in terms of our upper bounds.¹⁷ Therefore we assume that the likelihood and its mean-field approximation have enough curvature.

PROPOSITION B.1. Assume (α, β, γ) lives on a compact set Θ . Let $\hat{\theta}_n := (\hat{\alpha}_n, \hat{\beta}_n, \hat{\gamma}_n)$ and $\hat{\theta}_n^{MF} := (\hat{\alpha}_n^{MF}, \hat{\beta}_n^{MF}, \hat{\gamma}_n^{MF})$ be the maximizers of ℓ_n and ℓ_n^{MF} , respectively, in the interior of Θ . Moreover, we assume that ψ_n and ψ_n^{MF} are differentiable and μ_n - and μ_n^{MF} -strongly convex in (α, β, γ) , respectively, on Θ , where $\mu_n > 0$ and $\mu_n^{MF} > 0$. Then

(**B**.1)

$$\|\hat{\theta}_n - \hat{\theta}_n^{MF}\| \le \frac{2}{(\mu_n + \mu_n^{MF})^{\frac{1}{2}}} \left[\sup_{\alpha, \beta, \gamma \in \Theta} C_1^{\frac{1}{2}}(\alpha, \beta, \gamma) \left(\frac{\log n}{n} \right)^{\frac{1}{10}} + \sup_{\alpha, \beta, \gamma \in \Theta} C_2^{\frac{1}{2}}(\alpha, \beta, \gamma) n^{-\frac{1}{4}} \right],$$

where C_1 and C_2 are defined in Theorem 3.1 and $\|\cdot\|$ denotes the Euclidean norm.

In Proposition B.1, if μ_n and μ_n^{MF} goes to zero sufficiently fast as n goes zero, then the bound in (B.1) may not go to zero as n goes to zero. If for example μ_n, μ_n^{MF} are uniformly bounded from below, and both $\sup_{\alpha,\beta,\gamma\in\Theta} C_1(\alpha,\beta,\gamma)$ and $\sup_{\alpha,\beta,\gamma\in\Theta} C_2(\alpha,\beta,\gamma)$ are O(1), then $\|\hat{\theta}_n - \hat{\theta}_n^{MF}\| = O(n^{-1/10}(\log n)^{1/10})$.

B.1. **Proof of Proposition B.1.** We assume that ψ_n (resp. ψ_n^{MF}) is differentiable and μ_n -strongly convex (resp. μ_n^{MF} -strongly convex) in $\theta := (\alpha, \beta, \gamma) \in \Theta$. Note that

$$\ell_n = T_n - \psi_n, \qquad \ell_n^{MF} = T_n - \psi_n^{MF},$$

and T_n is linear in $\theta = (\alpha, \beta, \gamma)$, we have that ℓ_n (resp. ℓ_n^{MF}) is differentiable and μ_n -strongly concave in $\theta := (\alpha, \beta, \gamma) \in \Theta$ so that for any $x, y \in \Theta$,

(B.2)
$$\ell_n(y) \le \ell_n(x) + \nabla \ell_n(x)^T (y - x) - \frac{\mu_n}{2} \|y - x\|^2,$$

¹⁷Geyer and Thompson (1992) mentions similar problems arise for the MCMC-MLE. Indeed, as mentioned above, the MLE may not exist. For example, if the number of triangles is zero in the data, it will be impossible to estimate γ and the MCMC-MLE may give an approximation with solution that tends to infinity.

and in particular,

(B.3)
$$\ell_n(\hat{\theta}_n^{MF}) \le \ell_n(\hat{\theta}_n) + \nabla \ell_n(\hat{\theta}_n)^T (\hat{\theta}_n^{MF} - \hat{\theta}_n) - \frac{\mu_n}{2} \|\hat{\theta}_n^{MF} - \hat{\theta}_n\|^2$$
$$= \ell_n(\hat{\theta}_n) - \frac{\mu_n}{2} \|\hat{\theta}_n^{MF} - \hat{\theta}_n\|^2,$$

and similarly, for any $x, y \in \Theta$,

(B.4)
$$\ell_n^{MF}(y) \le \ell_n^{MF}(x) + \nabla \ell_n^{MF}(x)^T (y-x) - \frac{\mu_n}{2} \|y-x\|^2,$$

and in particular,

(B.5)
$$\ell_n^{MF}(\hat{\theta}_n) \le \ell_n^{MF}(\hat{\theta}_n^{MF}) + \nabla \ell_n^{MF}(\hat{\theta}_n^{MF})^T(\hat{\theta}_n - \hat{\theta}_n^{MF}) - \frac{\mu_n^{MF}}{2} \|\hat{\theta}_n - \hat{\theta}_n^{MF}\|^2 = \ell_n^{MF}(\hat{\theta}_n^{MF}) - \frac{\mu_n^{MF}}{2} \|\hat{\theta}_n - \hat{\theta}_n^{MF}\|^2.$$

Adding the inequalities (B.3) and (B.5), we get

$$\begin{aligned} \|\hat{\theta}_n - \hat{\theta}_n^{MF}\|^2 &\leq \frac{2}{\mu_n^{MF} + \mu_n} \left[\left(\ell_n^{MF}(\hat{\theta}_n^{MF}) - \ell_n(\hat{\theta}_n^{MF}) \right) + \left(\ell_n(\hat{\theta}_n) - \ell_n^{MF}(\hat{\theta}_n) \right) \right] \\ &\leq \frac{4}{\mu_n^{MF} + \mu_n} \sup_{\theta \in \Theta} |\ell_n^{MF}(\theta) - \ell_n(\theta)|. \end{aligned}$$

By applying Theorem 3.1, we get

$$\begin{split} \|\hat{\theta}_{n} - \hat{\theta}_{n}^{MF}\| &\leq \frac{2}{(\mu_{n} + \mu_{n}^{MF})^{\frac{1}{2}}} \left[\sup_{\alpha,\beta,\gamma\in\Theta} C_{1}(\alpha,\beta,\gamma) n^{-\frac{1}{5}} (\log n)^{\frac{1}{5}} + \sup_{\alpha,\beta,\gamma\in\Theta} C_{2}(\alpha,\beta,\gamma) n^{-\frac{1}{2}} \right]^{\frac{1}{2}} \\ &\leq \frac{2}{(\mu_{n} + \mu_{n}^{MF})^{\frac{1}{2}}} \left[\sup_{\alpha,\beta,\gamma\in\Theta} C_{1}^{\frac{1}{2}}(\alpha,\beta,\gamma) n^{-\frac{1}{10}} (\log n)^{\frac{1}{10}} + \sup_{\alpha,\beta,\gamma\in\Theta} C_{2}^{\frac{1}{2}}(\alpha,\beta,\gamma) n^{-\frac{1}{4}} \right], \end{split}$$

where the last step is due to the inequality $\sqrt{x+y} \le \sqrt{x} + \sqrt{y}$ for any $x, y \ge 0$. The proof is complete.

APPENDIX C. ADDITIONAL SIMULATION RESULTS

C.1. No covariates, edges and two-stars model. We have estimated a model with no covariates. This corresponds to a model in which $\tilde{\alpha}_2 = 0$ or $\alpha_1 = \alpha_2 = \alpha$. The results of our simulations for small networks are in Table C.1. Our method performs relatively well in this simpler case. Indeed in this case there are results that would allow us to solve the variational problem in closed form for large n (Chatterjee and Diaconis, 2013; Mele, 2017; Aristoff and Zhu, 2018; Radin and Yin, 2013). The MPLE and MCMC-MLE median estimate seems to converge to the true value as we increase n, but our approximation seems to perform slightly better here.

TABLE C.1. Monte Carlo estimates, comparison of three methods. True parameter vector is $(\tilde{\alpha}_1, \tilde{\alpha}_2, \beta) = (-2, 0, 1)$

n = 50	М	CMC-M	ILE	ME	AN-FIE	LD	MPLE			
	$\tilde{\alpha}_1$	\tilde{lpha}_2	β	$\tilde{\alpha}_1$	\tilde{lpha}_2	β	$\tilde{\alpha}_1$	\tilde{lpha}_2	β	
median	-2.063	0.016	-0.324	-2.021	0.007	0.999	-1.983	0.018	-1.006	
0.05	-2.692	-0.614	-23.828	-2.412	-0.372	0.975	-2.439	-0.368	-34.177	
0.95	-1.363	0.657	22.738	-1.783	0.413	1.015	-1.449	0.401	14.465	
n = 100	М	CMC-M	ILE	ME	AN-FIE	LD		MPLE		
	$\tilde{\alpha}_1$	\tilde{lpha}_2	β	$\tilde{\alpha}_1$	\tilde{lpha}_2	β	$\tilde{\alpha}_1$	\tilde{lpha}_2	β	
median	-1.970	-0.042	0.221	-1.981	-0.017	1.000	-1.949	-0.023	-1.231	
0.05	-2.241	-0.333	-13.226	-2.101	-0.194	0.993	-2.168	-0.196	-14.402	
0.95	-1.602	0.249	16.316	-1.874	0.134	1.012	-1.643	0.142	9.328	
n = 200	М	CMC-M	ILE	ME	AN-FIE	LD		MPLE		
	$\tilde{\alpha}_1$	\tilde{lpha}_2	β	$\tilde{\alpha}_1$	\tilde{lpha}_2	β	\tilde{lpha}_1	\tilde{lpha}_2	β	
median	-2.012	-0.005	1.483	-1.998	0.002	1.000	-2.003	-0.001	1.225	
0.05	-2.214	-0.184	-9.515	-2.067	-0.093	0.997	-2.160	-0.095	-9.682	
0.95	-1.796	0.161	12.179	-1.935	0.091	1.003	-1.790	0.095	8.784	

Notes. See notes for Table 4.1.

C.2. Model with 2-stars. In this subsection we report estimates of a model where the triangle term is excluded from the specification ($\gamma = 0$ in log-likelihood (4.5)). In Table C.2 we report results for 100 simulations of a model with $(\tilde{\alpha}_1, \tilde{\alpha}_2, \beta) = (-2, 1, 2)$. We run simulations for networks of size n = 50, 100, 200, to show how our method compares to MCMC-MLE and MPLE when the size of the network grows. In general, we expect more precise results as n grows large.

The results are encouraging and the mean-field approximation seems to behave as expected. Indeed, the median estimate is very close to the true parameters that generate the data. As the size of the network grows from n = 50 to n = 200, both MCMC-MLE and MPLE also improve in precision. The fastest method in terms of computational time is the MPLE. This is because the MPLE's speed depends on the number of parameters. Our mean-field approximation is as fast as the MCMC-MLE.

n = 50	MC	CMC-M	LE	ME	AN-FIE	LD		MPLE			
	\tilde{lpha}_1	\tilde{lpha}_2	β	$\tilde{\alpha}_1$	\tilde{lpha}_2	β	$\tilde{\alpha}_1$	\tilde{lpha}_2	β		
median	-2.015	0.999	2.303	-1.993	1.000	2.004	-1.996	0.998	1.820		
0.05	-2.433	0.641	-1.085	-2.060	0.885	1.916	-2.325	0.780	-2.556		
0.95	-1.666	1.337	6.118	-1.905	1.090	2.087	-1.573	1.273	4.783		
n = 100	MC	CMC-M	LE	ME	AN-FIE	LD		MPLE			
	\tilde{lpha}_1	\tilde{lpha}_2	β	$\tilde{\alpha}_1$	\tilde{lpha}_2	β	$\tilde{\alpha}_1$	\tilde{lpha}_2	β		
median	-1.995	1.012	1.932	-1.980	1.011	2.011	-1.980	1.010	1.783		
0.05	-2.189	0.861	0.701	-2.032	0.969	1.992	-2.175	0.901	0.329		
0.95	-1.833	1.157	3.314	-1.944	1.044	2.088	-1.816	1.141	2.867		
n = 200	MC	CMC-M	LE	ME	AN-FIE	LD		MPLE			
	$\tilde{\alpha}_1$	\tilde{lpha}_2	β	$\tilde{\alpha}_1$	\tilde{lpha}_2	β	$\tilde{\alpha}_1$	\tilde{lpha}_2	β		
median	-2.000	1.009	1.938	-1.986	1.005	2.016	-1.997	1.007	1.930		
0.05	-2.182	0.925	0.843	-2.004	0.932	1.999	-2.176	0.950	0.592		
0.95	-1.882	1.087	4.119	-1.935	1.028	2.214	-1.847	1.069	3.541		

TABLE C.2. Monte Carlo estimates, comparison of three methods. True parameter vector is $(\tilde{\alpha}_1, \tilde{\alpha}_2, \beta) = (-2, 1, 2)$

Notes. See notes for Table 4.1.

The second set of Monte Carlo experiments is reported in Table C.3, where the data are generated by parameter vector $(\tilde{\alpha}_1, \tilde{\alpha}_2, \beta) = (-2, 1, 3)$. The pattern is similar to the previous table, but the mean field estimates exhibit a little more bias.

TABLE C.3. Monte Carlo estimates, comparison of three methods. True parameter vector is $(\tilde{\alpha}_1, \tilde{\alpha}_2, \beta) = (-2, 1, 3)$

n = 50	MC	CMC-M	LE	ME	AN-FIE	LD		MPLE		
	$\tilde{\alpha}_1$	\tilde{lpha}_2	β	$\tilde{\alpha}_1$	\tilde{lpha}_2	β	$\tilde{\alpha}_1$	\tilde{lpha}_2	β	
median	-1.978	1.010	2.742	-1.958	1.026	3.025	-1.921	1.016	2.357	
0.05	-2.308	0.745	1.342	-2.045	0.878	2.938	-2.201	0.823	-0.742	
0.95	-1.689	1.229	4.466	-1.811	1.141	3.468	-1.547	1.202	4.288	
n = 100	MC	CMC-M	LE	ME	AN-FIE	LD		MPLE		
	\tilde{lpha}_1	\tilde{lpha}_2	β	$\tilde{\alpha}_1$	\tilde{lpha}_2	β	\tilde{lpha}_1	\tilde{lpha}_2	β	
median	-2.005	1.002	3.022	-1.851	1.091	3.166	-1.997	1.001	3.009	
0.05	-2.116	0.892	2.665	-2.274	0.866	2.998	-2.098	0.924	2.514	
0.95	-1.902	1.110	3.414	-1.670	1.861	4.092	-1.895	1.096	3.425	
n = 200	MC	CMC-M	LE	ME	AN-FIE	LD	MPLE			
	$\tilde{\alpha}_1$	\tilde{lpha}_2	β	$\tilde{\alpha}_1$	\tilde{lpha}_2	β	$\tilde{\alpha}_1$	\tilde{lpha}_2	β	
median	-2.003	1.000	2.959	-1.923	1.030	3.107	-1.984	1.000	2.847	
0.05	-2.151	0.934	2.314	-2.059	0.922	3.000	-2.104	0.951	2.096	
0.95	-1.902	1.064	3.944	-1.836	1.164	4.222	-1.861	1.039	3.666	

Notes. See notes for Table 4.1.

C.3. Model with triangles. The second set of simulations involves a model with no two-stars, that is $\beta = 0$, in Table C.4. In this specification our mean-field approximation seems to do better than the other estimators, at least for this small networks.

TABLE C.4. Monte Carlo estimates, comparison of three methods. True parameter vector is $(\tilde{\alpha}_1, \tilde{\alpha}_2, \gamma) = (-2, 1, -2)$

n = 50	Μ	CMC-N	1LE	ME	AN-FIE	ELD	MPLE			
	\tilde{lpha}_1	\tilde{lpha}_2	γ	\tilde{lpha}_1	\tilde{lpha}_2	γ	\tilde{lpha}_1	\tilde{lpha}_2	γ	
median	-2.024	1.026	-13.959	-2.000	1.005	-2.000	-2.031	1.012	-9.804	
0.05	-2.384	0.622	-60.419	-2.321	0.168	-6.425	-2.398	0.758	-45.881	
0.95	-1.689	1.457	49.585	-0.739	2.246	-1.777	-1.809	1.394	21.696	
n = 100	Μ	CMC-N	1LE	ME	AN-FIE	ELD	MPLE			
	\tilde{lpha}_1	\tilde{lpha}_2	γ	\tilde{lpha}_1	\tilde{lpha}_2	γ	\tilde{lpha}_1	\tilde{lpha}_2	γ	
median	-2.006	1.019	-6.053	-1.967	1.035	-2.007	-2.002	1.015	-4.980	
0.05	-2.164	0.832	-35.171	-3.472	0.951	-7.368	-2.124	0.876	-23.937	
0.95	-1.824	1.183	27.361	-1.388	3.763	-1.910	-1.890	1.153	13.519	
n = 200	Μ	CMC-N	1LE	ME	AN-FIE	ELD		MPLE		
	$\tilde{\alpha}_1$	\tilde{lpha}_2	γ	$\tilde{\alpha}_1$	\tilde{lpha}_2	γ	$\tilde{\alpha}_1$	\tilde{lpha}_2	γ	
median	-2.007	1.001	-1.002	-1.972	1.031	-2.006	-2.003	1.000	-1.913	
0.05	-2.083	0.901	-23.049	-2.014	1.008	-2.115	-2.061	0.929	-15.721	
0.95	-1.931	1.095	16.760	-1.473	1.636	-1.983	-1.952	1.072	9.153	

Notes. See notes for Table 4.1.

n = 50		MCM	C-MLE			MEAN	-FIELD		MPLE			
	\tilde{lpha}_1	\tilde{lpha}_2	β	γ	\tilde{lpha}_1	\tilde{lpha}_2	β	γ	\tilde{lpha}_1	\tilde{lpha}_2	β	γ
median	-2.008	1.023	-1.256	-4.943	-1.977	1.030	-1.018	-1.002	-1.959	1.015	-2.032	-3.296
mad	0.320	0.267	4.898	43.074	0.153	0.165	0.144	0.154	0.307	0.191	4.532	24.826
n = 100		MCM	C-MLE			-FIELD		MPLE				
median	-1.996	1.004	-1.138	-3.173	-1.932	1.177	-1.057	-1.021	-1.974	1.006	-1.566	-1.489
mad	0.219	0.133	3.364	28.410	0.567	0.553	0.335	0.346	0.207	0.093	3.119	16.695
n = 200		MCM	C-MLE		MEAN-FIELD					M	PLE	
median	-1.995	1.007	-1.155	-0.980	-1.603	1.645	-1.317	-1.078	-1.987	1.003	-1.340	-1.308
mad	0.133	0.069	2.098	18.167	0.559	0.794	0.656	0.558	0.127	0.047	2.064	11.196
n = 500		MCM	C-MLE		MEAN-FIELD				MPLE			
median	-1.998	1.002	-1.070	-1.315	-1.682	1.836	-1.431	-1.155	-1.991	1.000	-1.113	-1.227
mad	0.084	0.033	1.496	10.897	0.805	0.849	0.776	0.883	0.079	0.020	1.340	7.036

TABLE C.5. Monte Carlo estimates, comparison of three methods. True parameter vector is $(\tilde{\alpha}_1, \tilde{\alpha}_2, \beta, \gamma) = (-2, 1, -1, -1)$

Notes: see notes for Table 4.1.

TABLE C.6. Monte Carlo estimates, comparison of three methods. True parameter vector is $(\tilde{\alpha}_1, \tilde{\alpha}_2, \beta, \gamma) = (-2, 1, -2, 3)$

n = 50		MCM	C-MLE			MEAN	-FIELD		MPLE			
	\tilde{lpha}_1	\tilde{lpha}_2	β	γ	\tilde{lpha}_1	\tilde{lpha}_2	β	γ	\tilde{lpha}_1	\tilde{lpha}_2	β	γ
median	-2.005	1.024	-2.368	-4.197	-1.955	1.037	-2.022	2.998	-1.958	1.017	-3.006	-0.198
mad	0.349	0.292	5.767	46.688	0.095	0.085	0.088	0.082	0.307	0.196	4.707	26.076
n = 100		MCM	C-MLE		MEAN-FIELD				MPLE			
median	-2.000	0.995	-2.333	1.560	-1.909	1.082	-2.100	2.983	-1.972	0.997	-2.708	2.617
mad	0.216	0.145	3.429	31.810	0.151	0.147	0.199	0.130	0.195	0.099	3.221	17.184
n = 200		MCM	C-MLE		MEAN-FIELD					M	PLE	
median	-1.998	0.997	-2.062	1.847	-1.593	1.512	-2.849	2.711	-1.985	0.999	-2.321	2.326
mad	0.129	0.073	2.302	22.032	0.565	0.677	1.195	0.594	0.124	0.049	2.167	13.057
n = 500		MCM	C-MLE		MEAN-FIELD				MPLE			
median	-2.004	1.002	-1.944	2.531	-1.523	1.605	-3.493	2.557	-2.002	1.002	-2.059	2.786
mad	0.091	0.038	1.579	11.813	0.782	0.726	1.472	0.982	0.080	0.024	1.472	8.068

Notes: see notes for Table 4.1.

C.4. Some examples of nonconvergence. In Tables C.5 and C.6 we show examples in which our mean-field approximation performs worse than the alternative estimators. There are several possible explanations for this poor convergence. First, it may be that we are not finding the maximizer of the approximation variational problem (2.9), given the local nature of updates (4.1). In these simulations we do not start the matrix $\mu^{(0)}$ at different initial values, therefore we converge to a local maximum that may not be global. Our package mfergm allows the researcher to initialize

 $\mu^{(0)}$ at different random starting points. This can improve convergence. In principle we should increase the number of re-starts as *n* grows, as it is known that these models may have multiple modes. Ideally, one can use a Nelder-Mead or Simulated Annealing algorithm to find the maximizer of the variational problem, but this is more time-consuming. All these ideas lead to simple parallelization of our package's functions that are beyond the scope of the present work. Second, the tolerance level that we use $\epsilon_{tol} = 0.0001$ may be too large. Third, the likelihood may exhibit a phase transition and thus a small difference in parameters may cause a large change in the behavior of the model. We conjecture that some of these issues are related to identification and we plan to explore this in future work.

C.5. A note on computational speed. In our Monte Carlo exercises, we note that the computational speed of the three estimators is similar for small networks. For n = 100, the mean-field approximation takes about 3.5s to estimate the model, while an MCMC-MLE with a burnin of 100,000 and sampling every 1000 iterations takes approximately 5.5s and the MPLE takes about 1.7s. For n = 50 the estimates take 1.6s for mean-field, 4s for MCMC-MLE and 1.2s for MPLE.

However, for larger networks, our code is computationally inefficient and results in much larger computational time than using the built-in functions in the ergm package in R for MCMC-MLE and MPLE. We have experimented with faster iterative routines that could speed up the approximate solution of the variational mean-field problem, but these are not fully stable. Additionally our code does not make efficient use of the memory, as the matrix μ is dense and we are not using efficient matrix algebra libraries to speed up the computation. We believe that such improvement in our benchmark code will make computational time comparable to MPLE.

ONLINE APPENDIX - NOT FOR PUBLICATION

APPENDIX D. ASYMPTOTIC RESULTS

In this section we consider the model as $n \to \infty$. We have seen previously that the log normalizing constant $\psi_n(\alpha, \beta, \gamma)$ can be approximated by $\psi_n^{MF}(\mu(\alpha, \beta, \gamma))$ by the mean-field approximation, where $\mu(\alpha, \beta, \gamma)$ solves the optimization problem in (2.9) and $\psi_n^{MF}(\mu(\alpha, \beta, \gamma))$ is its optimal value, where we recall that

$$\psi_{n}^{MF}(\boldsymbol{\mu}(\alpha,\beta,\gamma)) = \sup_{\boldsymbol{\mu}\in[0,1]^{n^{2}}:\mu_{ij}=\mu_{ji},\forall i,j} \left\{ \frac{1}{n^{2}} \sum_{i,j} \alpha_{ij}\mu_{ij} + \frac{\beta}{2n^{3}} \sum_{i,j,k} \mu_{ij}\mu_{jk} + \frac{2\gamma}{3n^{3}} \sum_{i,j,k} \mu_{ij}\mu_{jk}\mu_{ki} - \frac{1}{2n^{2}} \sum_{i,j} [\mu_{ij}\log\mu_{ij} + (1-\mu_{ij})\log(1-\mu_{ij})] \right\},$$

We will study the limit as $n \to \infty$. Before we proceed, we need a representation of the vector α in the infinite network. The following assumption guarantee that we can switch from the discrete to the continuum.

ASSUMPTION D.1. Assume that

$$\alpha_{ij} = \alpha \left(i/n, j/n \right),$$

where $\alpha(x, y) : [0, 1]^2 \to \mathbb{R}$ is a deterministic exogenous function that is symmetric, i.e., $\alpha(x, y) = \alpha(y, x)$.¹⁸

Since we have n players, the number of types for the players must be finite, although it may grow as n grows. α_{ij} are symmetric, and can take at most $\frac{n(n+1)}{2}$ values. As $n \to \infty$, the number of types can become infinite and $\alpha(x, y)$ may take infinitely many values. On the other hand, in terms of practical applications, finitely many values often suffice ¹⁹.

¹⁸To ease the notations, we project $\otimes_{j=1}^{S} \mathcal{X}_j$ onto [0, 1] and the function $\alpha(\tau_i, \tau_j)$ defined previously is now re-defined from $[0, 1]^2$ to \mathbb{R} .

¹⁹If an entry of the vector τ_i is continuous, we can always transform the variable in a discrete vector using thresholds. For example, if $\mathcal{X}_j = [\$50,000,\$200,000]$, we can bucket the incomes into three levels, low: [\$50,000,\$100,000), medium [\$100,000,\$150,000) and high: [\$150,000, \$200,000].



FIGURE D.1. Examples of function $\alpha(x, y)$.

The figure provides several examples of possible partitions of the net benefit function $\alpha(x, y)$ with finite covariates. The asymptotic version of this function is defined over the unit square.

ASSUMPTION D.2. We assume that $\alpha(x, y)$ is uniformly bounded in x and y:

(D.1)
$$\sup_{(x,y)\in[0,1]^2} |\alpha(x,y)| < \infty.$$

As a simple example, let us consider gender: the population consists of males and female agents. For example, half of the nodes (population) are males, say $i = 1, 2, ..., \frac{n}{2}$ and the other half are females, $i = \frac{n}{2} + 1, \frac{n}{2} + 2, ..., n$.²⁰ That means, $\alpha(x, y)$ takes three values according to the three

 $^{^{20}\}mbox{Here,}$ we assume without loss of generality that n is an even number.

regions:

$$\begin{split} \left\{ (x,y) : 0 < x, y < \frac{1}{2} \right\}, \\ \left\{ (x,y) : \frac{1}{2} < x, y < 1 \right\}, \\ \left\{ (x,y) : 0 < x < \frac{1}{2} < y < 1 \right\} \bigcup \left\{ (x,y) : 0 < y < \frac{1}{2} < x < 1 \right\}, \end{split}$$

and these three regions correspond precisely to pairs: male-male, female-female, and male-female. This example is represented in Figure D.1(C).

The work of Chatterjee and Diaconis (2013) show that the variational problem in (2.7) translates into an analogous variational problem for the graph limit.²¹ In the special case $\alpha(x, y) \equiv \alpha$, it is shown in Chatterjee and Diaconis (2013) that as $n \to \infty$ the log-constant of the ERGM converges to the solution of the variational problem (D.3), that is

(D.2)
$$\psi_n(\alpha,\beta,\gamma) \to \psi(\alpha,\beta,\gamma),$$

where

(D.3)

$$\begin{split} \psi(\alpha,\beta,\gamma) &= \sup_{h\in\mathcal{W}} \bigg\{ \alpha \int_0^1 \int_0^1 h(x,y) dx dy + \frac{\beta}{2} \int_0^1 \int_0^1 \int_0^1 h(x,y) h(y,z) dx dy dz \\ &+ \frac{2\gamma}{3} \int_0^1 \int_0^1 \int_0^1 h(x,y) h(y,z) h(z,x) dx dy dz - \frac{1}{2} \int_0^1 \int_0^1 I(h(x,y)) dx dy \bigg\}, \end{split}$$

where

(D.4)
$$\mathcal{W} := \left\{ h : [0,1]^2 \to [0,1], h(x,y) = h(y,x), 0 \le x, y \le 1 \right\},$$

and we define the entropy function:

$$I(x) := x \log x + (1 - x) \log(1 - x), \qquad 0 \le x \le 1,$$

with I(0) = I(1) = 0.

²¹See also Mele (2017) for similar results in a directed network.

In essence the first three terms in (D.3) correspond to the expected potential function in the continuum, while the last term in (D.3) corresponds to the entropy of the graph limit.

We will show that (D.2) holds with

(D.5)

$$\begin{split} \psi(\alpha,\beta,\gamma) &= \sup_{h\in\mathcal{W}} \bigg\{ \int_0^1 \int_0^1 \alpha(x,y)h(x,y)dxdy + \frac{\beta}{2} \int_0^1 \int_0^1 \int_0^1 h(x,y)h(y,z)dxdydz \\ &+ \frac{2\gamma}{3} \int_0^1 \int_0^1 \int_0^1 h(x,y)h(y,z)h(z,x)dxdydz - \frac{1}{2} \int_0^1 \int_0^1 I(h(x,y))dxdy \bigg\}, \end{split}$$

The function h in the expressions above is known as the graphon from the graph limits literature ²², large deviations literature for random graphs²³ and analysis of the resulting variational problem.²⁴ and it is a representation of an infinite network, where h is a simple symmetric function $h : [0, 1]^2 \rightarrow [0, 1]$, and h(x, y) = h(y, x). Note that our goal is to approximate ψ_n^{MF} and hence ψ_n by ψ , whose definition involves the function h, and we call such a function a graphon in the rest of the paper, to be consistent with the literature, while we are not attempting here to establish a theory of graph limits to allow nodal covariates. That is an interesting research direction worth investigating in the future, but is out of the scope of the current paper.

The following proposition shows that for a model with finitely many types the variational approximation is asymptotically exact.

PROPOSITION D.1. Under Assumptions D.1 and D.2, as $n \to \infty$

$$\psi_n(\alpha,\beta,\gamma) \to \psi(\alpha,\beta,\gamma),$$

where $\psi(\alpha, \beta, \gamma)$ is defined in (D.5).

Proof. It follows directly from Theorem 3.1 and $\psi_n^{MF}(\mu(\alpha, \beta, \gamma)) \to \psi(\alpha, \beta, \gamma)$, as $n \to \infty$. \Box

The proposition states that as n becomes large, we can approximate the exponential random graph using a model with independent links (conditional on finitely many types). This is a very

²²See Lovasz (2012), Borgs et al. (2008)

²³See Chatterjee and Varadhan (2011), Chatterjee and Diaconis (2013)

²⁴See Aristoff and Zhu (2018), Radin and Yin (2013) among others.

useful result because the latter approximation is simple and tractable, while the exponential random graph model contains complex dependence patterns that make estimation computationally expensive.

D.1. Approximation of the limit log normalizing constant. We can analyze and provide an approximation of the log-constant in the large network limit. The variational formula for $\psi(\alpha, \beta, \gamma)$ is an infinite-dimensional problem which is intractable in most cases. Nevertheless, we can always bound the infinite dimensional problem with finite dimensional ones (both lower and upper bounds), at least in the absence of transitivity. For details, see Proposition F.2 in the Online Appendix. The lower-bound in Proposition F.2 coincides with the structured mean-field approach of Xing et al. (2003). In a model with homogeneous players, the lower-bound corresponds to the computational approximation of graph limits implemented in He and Zheng (2013).

In the case of extreme homophily, we can also obtain finite-dimensional approximation, see Proposition F.1 in the Online Appendix.

D.2. Characterization of the variational problem. We recall that the log normalizing constant in the $n \to \infty$ limit is given by the variational problem:

$$(D.6) \quad \psi(\alpha,\beta,\gamma) = \sup_{h \in \mathcal{W}} \left\{ \int_0^1 \int_0^1 \alpha(x,y)h(x,y)dxdy + \frac{\beta}{2} \int_0^1 \int_0^1 \int_0^1 h(x,y)h(y,z)dxdydz + \frac{2\gamma}{3} \int_0^1 \int_0^1 \int_0^1 h(x,y)h(y,z)h(z,x)dxdydz - \frac{1}{2} \int_0^1 \int_0^1 \left[h(x,y)\log h(x,y) + (1-h(x,y))\log(1-h(x,y))\right]dxdy \right\}.$$

PROPOSITION D.2. *The optimal graphon h that solves the variational problem* (D.6) *satisfies the Euler-Lagrange equation:*

(D.7)

$$2\alpha(x,y) + \beta \int_0^1 h(x,y)dx + \beta \int_0^1 h(x,y)dy + 4\gamma \int_0^1 h(x,z)h(y,z)dz = \log\left(\frac{h(x,y)}{1-h(x,y)}\right).$$

Proof. The proof follows from the same argument as in Theorem 6.1. in Chatterjee and Diaconis (2013). □

COROLLARY 1. If $\alpha(x, y)$ is not a constant function, then the optimal graphon h that solves the variational problem (D.6) is not a constant function.

Proof. If the optimal graphon h is a constant function, then (D.7) implies that α is a constant function. Contradiction.

In general, if a graphon satisfies the Euler-Lagrange equation, that only indicates that the graphon is a stationary point, and it is not clear if the graphon is the local maximizer, local minimizer or neither. In the next result, we will show that when β is negative, any graphon satisfying the Euler-Lagrange equation in our model is indeed a local maximizer.

PROPOSITION D.3. Assume that $\beta < 0$ and $\gamma = 0$. If h is a graphon that satisfies the Euler-Lagrange equation (D.7), then h is a local maximizer of the variational problem (D.6).

Proof. Let us define

(D.8)
$$\Lambda[h] := \int_0^1 \int_0^1 \alpha(x, y) h(x, y) dx dy + \frac{\beta}{2} \int_0^1 \int_0^1 \int_0^1 h(x, y) h(y, z) dx dy dz \\ - \frac{1}{2} \int_0^1 \int_0^1 \left[h(x, y) \log h(x, y) + (1 - h(x, y)) \log(1 - h(x, y)) \right] dx dy.$$

Let h satisfy (D.7) and for any symmetric function g and $\epsilon > 0$ sufficiently small, we have

$$\begin{aligned} \text{(D.9)} \quad & \Lambda[h+\epsilon g] - \Lambda[h] \\ &= \epsilon^2 \left[\frac{\beta}{2} \int_0^1 \left(\int_0^1 g(x,y) dy \right)^2 dx - \frac{1}{4} \int_0^1 \int_0^1 I''(h(x,y)) g^2(x,y) dx dy \right] + O(\epsilon^3) \\ &= \epsilon^2 \left[\frac{\beta}{2} \int_0^1 \left(\int_0^1 g(x,y) dy \right)^2 dx - \frac{1}{4} \int_0^1 \int_0^1 \frac{g^2(x,y)}{h(x,y)(1-h(x,y))} dx dy \right] + O(\epsilon^3), \end{aligned}$$

and since $\beta < 0$, we conclude that h is a local maximizer in (D.6).

Remark D.1. In general, the variational problem for the graphons and the corresponding Euler-Lagrange equation (D.7) does not yield a closed form solution. In the special case $\beta = \gamma = 0$,

(D.10)
$$\psi(\alpha, 0, 0) = \sup_{h \in \mathcal{W}} \left\{ \iint_{[0,1]^2} \alpha(x, y) h(x, y) dx dy - \frac{1}{2} \iint_{[0,1]^2} I(h(x, y)) dx dy \right\},$$

where $I(x) := x \log x + (1-x) \log(1-x)$ and it is easy to see that the optimal graphon h(x, y) is given by $h(x, y) = \frac{e^{2\alpha(x,y)}}{e^{2\alpha(x,y)}+1}$, and therefore, $\psi(\alpha, 0, 0) = \frac{1}{2} \iint_{[0,1]^2} \log(1 + e^{2\alpha(x,y)}) dxdy$.

APPENDIX E. DETAILS OF EQUILIBRIUM ECONOMIC FOUNDATIONS

E.1. Setup and preferences. Consider a population of n heterogeneous players (the nodes), each characterized by an exogenous type $\tau_i \in \bigotimes_{j=1}^S \mathcal{X}_j$, i = 1, ..., n. The attribute τ_i is an S-dimensional vector and the sets \mathcal{X}_j can represent age, race, gender, income, etc. ²⁵ We collect all τ_i 's in an $n \times S$ matrix τ . The network's adjacency matrix g has entries $g_{ij} = 1$ if i and j are linked; and $g_{ij} = 0$ otherwise. The network is undirected, i.e. $g_{ij} = g_{ji}$, and $g_{ii} = 0$, for all i's.²⁶ The utility of player i is

(E.1)
$$u_i(g,\tau) = \sum_{j=1}^n \alpha_{ij} g_{ij} + \frac{\beta}{n} \sum_{j=1}^n \sum_{k=1}^n g_{ij} g_{jk},$$

where $\alpha_{ij} := \nu(\tau_i, \tau_j)$ are symmetric functions $\nu : \bigotimes_{j=1}^S \mathcal{X}_j \times \bigotimes_{j=1}^S \mathcal{X}_j \to \mathbb{R}$ and $\nu(\tau_i, \tau_j) = \nu(\tau_j, \tau_i)$ for all i, j; and β is a scalar. The utility of player i depends on the number of direct links, each weighted according to a function ν of the types τ . This payoff structure implies that the net benefit of forming a direct connection depends on the characteristics of the two individuals involved in the link.

Players also care about the number of links that each of their direct contacts have formed.²⁷ For example, when $\beta > 0$, there is an incentive to form links to people that have many friends, e.g. popular kids in school. On the other hand, when $\beta < 0$ the incentive is reversed. For example, one

²⁶Extensions to directed networks are straightforward (see Mele (2017)).

²⁵For instance, if we consider gender and income, then S = 2, and we can take $\bigotimes_{j=1}^2 \mathcal{X}_j = \{\text{male,female}\} \times \{\text{low, medium, high}\}$. The sets \mathcal{X}_j can be both discrete and continuous. For example, if we consider gender and income, we can also take $\bigotimes_{j=1}^2 \mathcal{X}_j = \{\text{male,female}\} \times [\$50,000,\$200,000]$. Below we restrict the covariates to be discrete, but we allow the number of types to grow with the size of the network.

²⁷The normalization of β by *n* is necessary for the asymptotic analysis.

can think that forming links to a person with many connections could decrease our visibility and decrease the effectiveness of interactions. Similar utility functions have been used extensively in the empirical network formation literature.²⁸

The preferences in (E.1) include only direct links and friends' populatity. However, we can also include other types of link externalities. For example, in many applications the researcher is interested in estimating preferences for common neighbors. This is an important network statistics to measure transitity and clustering in networks. In our model we can easily add an utility component to capture these effects.

(E.2)
$$u_i(g,\tau) = \sum_{j=1}^n \alpha_{ij}g_{ij} + \frac{\beta}{n}\sum_{j=1}^n \sum_{k=1}^n g_{ij}g_{jk} + \frac{\gamma}{n}\sum_{j=1}^n \sum_{k=1}^n g_{ij}g_{jk}g_{ki},$$

These preferences include an additional parameter γ that measures the effect of common neighbors. The potential function for this model is

(E.3)
$$Q_n(g;\alpha,\beta) = \sum_{i=1}^n \sum_{j=1}^n \alpha_{ij} g_{ij} + \frac{\beta}{2n} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n g_{ij} g_{jk} + \frac{2\gamma}{3n} \sum_{j=1}^n \sum_{k=1}^n g_{ij} g_{jk} g_{ki}.$$

In general, all the results that we show below extend to more general utility functions that include payoffs for link externalities similar to (2.5).

The probability that i and j meet can depend on their networks: it could be a function of their common neighbors, or a function of their degrees and centralities, for example. In Assumption E.1, we assume that the existence of a link between i and j does not affect their probability of meeting. This is because we prove the existence and functional form of the stationary distribution (2.3) using the detailed balance condition, which is not satisfied if we allow the meeting probabilities to depend on the link between i and j.

The model can easily be extended to directed networks and the results on equilibria and longrun stationary distribution will hold. The results about the approximations of the likelihood shown below will also hold for directed networks, with minimal modifications of the proofs.

²⁸See Mele (2017), Sheng (2012), DePaula et al. (2018), Chandrasekhar and Jackson (2014), Badev (2013), Butts (2009).

Finally, while our model generates dense graphs, the approximations using variational methods and nonlinear large deviations that we develop in the rest of the paper also work in moderately sparse graphs. More precisely, the utility of player i is given by

(E.4)
$$u_i(g,\tau) = \sum_{j=1}^n \alpha_{ij}^{(n)} g_{ij} + \frac{\beta^{(n)}}{n} \sum_{j=1}^n \sum_{k=1}^n g_{ij} g_{jk} + \frac{\gamma^{(n)}}{n} \sum_{j=1}^n \sum_{k=1}^n g_{ij} g_{jk} g_{ki},$$

where $|\alpha_{ij}^{(n)}|$, $|\beta^{(n)}|$ and $|\gamma^{(n)}|$ can have moderate growth in *n* instead of being bounded. We will give more details later in our paper.²⁹

Example E.1. (Homophily) Consider a model with $\nu(\tau_i, \tau_j) = V - c(\tau_i, \tau_j)$, where V > 0 is the benefit of a link and $c(\tau_i, \tau_j)$ (= $c(\tau_j, \tau_i)$) is the cost of the link between i and j. To model homophily in this framework let the cost function be

(E.5)
$$c(\tau_i, \tau_j) = \begin{cases} c & \text{if } \tau_i = \tau_j, \\ C & \text{if } \tau_i \neq \tau_j. \end{cases}$$

For example, consider the parameterization 0 < c < V < C and $\beta = 0$, $\gamma = 0$. In this case the players have no incentive to form links with agents of other groups. On the other hand, if we have 0 < c < V < C and $\beta, \gamma > 0$, also links across groups will be formed, as long as β, γ are sufficiently large.

Example E.2. (Social Distance Model) Let the payoff from direct links be a function of the social distance among the individuals. Formally, let $\nu(\tau_i, \tau_j) := \eta d(\tau_i, \tau_j) - c$, where $d(\tau_i, \tau_j)$ is a distance function, η is a parameter that determines the sensitivity to the social distance and c > 0 is the cost of forming a link.³⁰ The case with $\eta < 0$ represents a world where individuals prefer linking to similar agents and $\eta > 0$ represents a world where individuals prefer linking with people at larger social distance. Note that even when $\eta < 0$, if we have $\beta, \gamma > 0$ sufficiently large, individuals may still have an incentive to form links with people at larger social distance.

 $[\]overline{^{29}$ See Chatterjee and Dembo (2016) for additional applications of nonlinear large deviations.

³⁰See Iijima and Kamada (2014) for a more general example of such model.

E.2. Meetings and equilibrium. The network formation process follows a stochastic best-response dynamics:³¹ in each period t, two random players meet with probability ρ_{ij} ; upon meeting they have the opportunity to form a link (or delete it, if already in place). Players are myopic: when they form a new link, they do not consider how the new link will affect the incentives of the other player in the future evolution of the network.

ASSUMPTION E.1. The meeting process is a function of types and the network. Let g_{-ij} indicate the network g without considering the link g_{ij} . Then the probability that i and j meet is

(E.6)
$$\rho_{ij} := \rho(\tau_i, \tau_j, g_{-ij}) > 0$$

for all pairs i and j, and i.i.d. over time.

Assumption E.1 implies that the meeting process can depend on covariates and the state of the network. For example, if two players have many friends in common they may meet with high probability; or people that share some demographics may meet more often. Crucially, every pair of players has a strictly positive probability of meeting. This guarantees that each link of the network has the opportunity of being revised.

Upon meetings, players decide whether to form or delete a link by maximizing the sum of their current utilities, i.e. the total surplus generated by the relationship. We are implicitly assuming that individuals can transfer utilities. When deciding whether to form a new link or deleting an existing link, players receive a random matching shock ε_{ij} that shifts their preferences.

At time t, the links g_{ij} is formed if

$$u_i(g_{ij} = 1, g_{-ij}, \tau) + u_j(g_{ij} = 1, g_{-ij}, \tau) + \varepsilon_{ij}(1) \ge u_i(g_{ij} = 0, g_{-ij}, \tau) + u_j(g_{ij} = 0, g_{-ij}, \tau) + \varepsilon_{ij}(0).$$

We make the following assumptions on the matching value.

ASSUMPTION E.2. Individuals receive a logistic shock before they decide whether to form a link (i.i.d. over time and players).

³¹See Blume (1993), Mele (2017), Badev (2013).

The logistic assumption is standard in many discrete choice models in economics and statistics (Train (2009)).

We can now characterize the equilibria of the model, following Mele (2017) and Chandrasekhar and Jackson (2014). In particular, we can show that the network formation is a potential game (Monderer and Shapley (1996)).

PROPOSITION E.1. The network formation is a potential game, and there exists a potential function $Q_n(g; \alpha, \beta)$ that characterizes the incentives of all the players in any state of the network

(E.7)
$$Q_n(g;\alpha,\beta) = \sum_{i=1}^n \sum_{j=1}^n \alpha_{ij} g_{ij} + \frac{\beta}{2n} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n g_{ij} g_{jk} + \frac{2\gamma}{3n} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n g_{ij} g_{jk} g_{ki}.$$

Proof. The proposition follows the same lines as Proposition 1 in Mele (2017) and it is omitted for brevity. \Box

The potential function $Q_n(g; \alpha, \beta)$ is such that, for any g_{ij}

$$Q_n(g;\alpha,\beta) - Q_n(g-ij;\alpha,\beta) = u_i(g) + u_j(g) - \left[u_i(g-ij) + u_j(g-ij)\right].$$

Thus we can keep track of all players' incentives using the scalar $Q_n(g; \alpha, \beta)$. It is easy to show that all the pairwise stable (with transfers) networks are the local maxima of the potential function.³² The sequential network formation follows a *Glauber* dynamics, therefore converging to a unique stationary distribution.

THEOREM E.1. In the long run, the model converges to the stationary distribution π_n , defined as

(E.8)
$$\pi_n(g;\alpha,\beta) = \frac{\exp\left[Q_n(g;\alpha,\beta)\right]}{\sum_{\omega\in\mathcal{G}}\exp\left[Q_n(\omega;\alpha,\beta)\right]} = \exp\left\{n^2\left[T_n(g;\alpha,\beta) - \psi_n(\alpha,\beta)\right]\right\}$$

where $T_n(g; \alpha, \beta) = n^{-2}Q_n(g; \alpha, \beta)$,

(E.9)
$$\psi_n(\alpha,\beta) = \frac{1}{n^2} \log \sum_{\omega \in \mathcal{G}} \exp\left[n^2 T_n(\omega;\alpha,\beta)\right],$$

³²A network g is pairwise stable with transfers if: (1) $g_{ij} = 1 \Rightarrow u_i(g,\tau) + u_j(g,\tau) \ge u_i(g-ij,\tau) + u_j(g-ij,\tau)$ and (2) $g_{ij} = 0 \Rightarrow u_i(g,\tau) + u_j(g,\tau) \ge u_i(g+ij,\tau) + u_j(g+ij,\tau)$; where g+ij represents network g with the addition of link g_{ij} and network g-ij represents network g without link g_{ij} . See Jackson (2010) for more details.

and
$$\mathcal{G} := \{ \omega = (\omega_{ij})_{1 \le i,j \le n} : \omega_{ij} = \omega_{ji} \in \{0,1\}, \omega_{ii} = 0, 1 \le i, j \le n \}.$$

Proof. The proof is an extension of Theorem 1 in Mele (2017). See also Chandrasekhar and Jackson (2014) and Butts (2009).

Notice that the likelihood (2.3) corresponds to an ERGM model with heterogeneous nodes and two-stars. As a consequence our model inherits all the estimation and identification challenges of the ERGM model.

APPENDIX F. SPECIAL CASE: THE EDGE-STAR MODEL

The general solution of the variational problem (D.3) is complicated. However, there are some special cases where we can characterize the solution with extreme detail. These examples show how we can solve the variational approximation in stylized settings, and we use them to explain how the method works in practice. In this section, we consider the special case in the absence of transitivity, i.e. $\gamma = 0$ and we get further results for the edge-star model.

F.1. Extreme homophily. We can exploit homophily to obtain a tractable approximation. Suppose that there are M types in the population. The cost of forming links among individuals of the same group is finite, but there is a large cost of forming links among people of different groups (potentially infinite). We show that in this case the normalizing constant can be approximated by solving M independent univariate maximization problems. In the special case of extreme homophily, our model converges to a block-diagonal model.

PROPOSITION F.1. Let $0 = a_0 < a_1 < \cdots < a_M = 1$ be a given sequence. Assume that

(F.1)
$$\alpha(x,y) = \alpha_{mm}, \quad \text{if } a_{m-1} < x, y < a_m, \quad m = 1, 2, \dots, M.$$

and $\alpha(x, y) \leq -K$ otherwise is a given function. Let $\psi(\alpha, \beta, 0; -K)$ be the variational problem for the graphons and $\psi(\alpha, \beta, 0; -\infty) = \lim_{K \to \infty} \psi(\alpha, \beta, 0; -K)$. Then, we have

(F.2)
$$\psi(\alpha,\beta,0;-\infty) = \sum_{m=1}^{M} (a_m - a_{m-1})^2 \sup_{0 \le x \le 1} \left\{ \alpha_{mm} x + \frac{\beta}{2} x^2 - \frac{1}{2} I(x) \right\}.$$

Proof. First, observe that

$$\begin{split} (\text{F.3}) \qquad \psi(\alpha,\beta,0;-\infty) \\ &= \sup_{h\in\mathcal{W}^-} \Big\{ \sum_{i=1}^M \alpha_i \iint_{[a_{i-1},a_i]^2} h(x,y) dx dy + \frac{\beta}{2} \int_0^1 \int_0^1 h(x,y) h(y,z) dx dy dz \\ &\quad -\frac{1}{2} \sum_{i=1}^M \iint_{[a_{i-1},a_i]^2} I(h(x,y)) dx dy \Big\} \\ &= \sup_{h\in\mathcal{W}^-} \Big\{ \sum_{i=1}^M \alpha_i \iint_{[a_{i-1},a_i]^2} h(x,y) dx dy + \frac{\beta}{2} \sum_{i=1}^M \int_{a_{i-1}}^{a_i} \left(\int_{a_{i-1}}^{a_i} h(x,y) dy \right)^2 dx \\ &\quad -\frac{1}{2} \sum_{i=1}^M \iint_{[a_{i-1},a_i]^2} I(h(x,y)) dx dy \Big\} \\ &= \sum_{i=1}^M \sup_{\substack{h:[a_{i-1},a_i]^2 \to [0,1]\\h(x,y) = h(y,x)}} \Big\{ \alpha_i \iint_{[a_{i-1},a_i]^2} h(x,y) dx dy + \frac{\beta}{2} \int_{a_{i-1}}^{a_i} \left(\int_{a_{i-1}}^{a_i} h(x,y) dy \right)^2 dx \\ &\quad -\frac{1}{2} \iint_{[a_{i-1},a_i]^2} I(h(x,y)) dx dy \Big\}, \end{split}$$

where

(F.4)
$$\mathcal{W}^- := \left\{ h \in \mathcal{W} : h(x, y) = 0 \text{ for any } (x, y) \notin \bigcup_{i=1}^M [a_{i-1}, a_i]^2 \right\}.$$

By taking h to be a constant on $[a_{i-1}, a_i]^2$, it is clear that

(F.5)
$$\psi(\alpha, \beta, 0; -\infty) \ge \sum_{i=1}^{M} (a_i - a_{i-1})^2 \sup_{0 \le x \le 1} \left\{ \alpha_i x + \frac{\beta}{2} x^2 - \frac{1}{2} I(x) \right\}.$$

By Jensen's inequality

$$\begin{aligned} \text{(F.6)} \quad \psi(\alpha,\beta,0;-\infty) &\leq \sum_{i=1}^{M} \sup_{\substack{h:[a_{i-1},a_i]^2 \to [0,1] \\ h(x,y) = h(y,x)}} \left\{ \alpha_i \int_{a_{i-1}}^{a_i} \left(\int_{a_{i-1}}^{a_i} h(x,y) dy \right)^2 dx \\ &\quad + \frac{\beta}{2} \int_{a_{i-1}}^{a_i} \left(\int_{a_{i-1}}^{a_i} h(x,y) dy \right)^2 dx \\ &\quad - \frac{1}{2} (a_i - a_{i-1}) \int_{a_{i-1}}^{a_i} I\left(\frac{1}{a_i - a_{i-1}} \int_{a_{i-1}}^{a_i} h(x,y) dy \right) dx \right\} \\ &\leq \sum_{i=1}^{M} (a_i - a_{i-1})^2 \sup_{0 \leq x \leq 1} \left\{ \alpha_i x + \frac{\beta}{2} x^2 - \frac{1}{2} I(x) \right\}. \end{aligned}$$

The net benefit function $\alpha(x, y)$ assumed in the Proposition is shown in Figure D.1(D). Essentially this result means that with extreme homophily, we can approximate the model, assuming perfect segregation: thus we can independently solve the variational problem of each type. This approach is computationally very simple, since each variational problem becomes a univariate maximization problem.

The solution of such univariate problem has been studied and characterized in previous work by Chatterjee and Diaconis (2013), Radin and Yin (2013), Aristoff and Zhu (2018) and Mele (2017). It can be shown that the solutions μ_m^* , where m = 1, ..., M, are the fixed point of equations

(F.7)
$$\mu_m = \frac{\exp\left[\alpha_{mm} + \beta\mu_m\right]}{1 + \exp\left[\alpha_{mm} + \beta\mu_m\right]},$$

for each group *m*, and $\beta \mu_m^*(1 - \mu_m^*) < 1$. The global maximizer μ_m^* is unique except on a phase transition curve $\{(\alpha_{mm}, \beta) : \alpha_{mm} + \beta = 0, \alpha_{mm} < -1\}$, see e.g. Radin and Yin (2013); Aristoff and Zhu (2018). It is shown in Chatterjee and Diaconis (2013) that the network of each group corresponds to an Erdős-Rényi graph with probability of a link equal to μ_m^* .

F.2. Analytically Tractable Bounds. In this section, for the edge-star model, we provide analytically tractable bounds for $\psi(\alpha, \beta, \gamma)$ when $\gamma = 0$. **PROPOSITION F.2.** Let $\gamma = 0$ and $0 = a_0 < a_1 < \cdots < a_{M-1} < a_M = 1$ be a given sequence.

Let us assume that

$$\alpha(x, y) = \alpha_{ml},$$
 if $a_{m-1} < x < a_m$ and $a_{l-1} < y < a_l$, where $1 \le m, l \le M$.

Then, we have

$$\begin{split} \sup_{\substack{0 \le u_{ml} \le 1\\ u_{ml} = u_{lm}, 1 \le m, l \le M}} \sum_{m=1}^{M} (a_m - a_{m-1}) \bigg\{ \sum_{l=1}^{M} (a_l - a_{l-1}) \alpha_{ml} u_{ml} \\ &+ \frac{\beta}{2} \left(\sum_{l=1}^{M} (a_l - a_{l-1}) u_{ml} \right)^2 - \frac{1}{2} \sum_{l=1}^{M} (a_l - a_{l-1}) I(u_{ml}) \bigg\} \\ &\le \psi(\alpha, \beta, 0) \le \sum_{m=1}^{M} (a_m - a_{m-1}) \sup_{\substack{0 \le u_{ml} \le 1\\ 1 \le l \le M}} \bigg\{ \sum_{l=1}^{M} (a_l - a_{l-1}) \alpha_{ml} u_{ml} + \frac{\beta}{2} \left(\sum_{l=1}^{M} (a_l - a_{l-1}) u_{ml} \right)^2 \\ &- \frac{1}{2} \sum_{l=1}^{M} (a_l - a_{l-1}) I(u_{ml}) \bigg\}. \end{split}$$

Proof. To compute the lower and upper bounds, let us define

(F.8)
$$u_{ij}(x) = \frac{1}{a_j - a_{j-1}} \int_{a_{j-1}}^{a_j} h(x, y) dy, \quad \text{for any } a_{i-1} < x < a_i.$$

We can compute that

(F.9)
$$\iint_{[0,1]^2} \alpha(x,y) h(x,y) dx dy = \sum_{i=1}^M \sum_{j=1}^M (a_j - a_{j-1}) \int_{a_{i-1}}^{a_i} \alpha_{ij} u_{ij}(x) dx.$$

Moreover,

(F.10)
$$\frac{\beta}{2} \int_0^1 \int_0^1 \int_0^1 h(x, y) h(y, z) dx dy dz = \frac{\beta}{2} \int_0^1 \left(\int_0^1 h(x, y) dy \right)^2 dx$$
$$= \frac{\beta}{2} \sum_{i=1}^M \int_{a_{i-1}}^{a_i} \left(\sum_{j=1}^M (a_j - a_{j-1}) u_{ij}(x) \right)^2 dx.$$

By Jensen's inequality, we can also compute that

$$\begin{aligned} \text{(F.11)} \\ &\frac{1}{2} \int_0^1 \int_0^1 I(h(x,y)) dx dy = \frac{1}{2} \sum_{i=1}^M \int_{a_{i-1}}^{a_i} \left[\sum_{j=1}^M \int_{a_{j-1}}^{a_j} I(h(x,y)) dy \right] dx \\ &= \frac{1}{2} \sum_{i=1}^M \int_{a_{i-1}}^{a_i} \left[\sum_{j=1}^M (a_j - a_{j-1}) \frac{1}{a_j - a_{j-1}} \int_{a_{j-1}}^{a_j} I(h(x,y)) dy \right] dx \\ &\geq \frac{1}{2} \sum_{i=1}^M \int_{a_{i-1}}^{a_i} \left[\sum_{j=1}^M (a_j - a_{j-1}) I\left(\frac{1}{a_j - a_{j-1}} \int_{a_{j-1}}^{a_j} h(x,y) dy \right) \right] dx \\ &= \frac{1}{2} \sum_{i=1}^M \int_{a_{i-1}}^{a_i} \sum_{j=1}^M (a_j - a_{j-1}) I(u_{ij}(x)) dx \end{aligned}$$

Hence, by (F.9), (F.10), (F.11), we get

$$\begin{split} \psi(\alpha,\beta,0) &\leq \sum_{i=1}^{M} \sum_{j=1}^{M} (a_{j}-a_{j-1}) \int_{a_{i-1}}^{a_{i}} \alpha_{ij} u_{ij}(x) dx + \frac{\beta}{2} \sum_{i=1}^{M} \int_{a_{i-1}}^{a_{i}} \left(\sum_{j=1}^{M} (a_{j}-a_{j-1}) u_{ij}(x) \right)^{2} dx \\ &- \frac{1}{2} \sum_{i=1}^{M} \int_{a_{i-1}}^{a_{i}} \sum_{j=1}^{M} (a_{j}-a_{j-1}) I(u_{ij}(x)) dx \\ &\leq \sum_{i=1}^{M} (a_{i}-a_{i-1}) \sup_{\substack{0 \leq u_{ij} \leq 1\\ 1 \leq j \leq M}} \left\{ \sum_{j=1}^{M} (a_{j}-a_{j-1}) \alpha_{ij} u_{ij} + \frac{\beta}{2} \left(\sum_{j=1}^{M} (a_{j}-a_{j-1}) u_{ij} \right)^{2} \\ &- \frac{1}{2} \sum_{j=1}^{M} (a_{j}-a_{j-1}) I(u_{ij}) \right\} \end{split}$$

On the other hand, by restricting the supremum over the graphons h(x, y)

(F.12)
$$h(x,y) = u_{ij},$$
 if $a_{i-1} < x < a_i$ and $a_{j-1} < y < a_j$, where $1 \le i, j \le M$,

where $(u_{ij})_{1 \le i,j \le M}$ is a symmetric matrix of the constants, and optimize over all the possible values $0 \le u_{ij} \le 1$, we get the lower bound:

(F.13)
$$\psi(\alpha, \beta, 0) \geq \sup_{\substack{0 \leq u_{ij} \leq 1 \\ u_{ij} = u_{ji}, 1 \leq i, j \leq M}} \sum_{i=1}^{M} (a_i - a_{i-1}) \left\{ \sum_{j=1}^{M} (a_j - a_{j-1}) \alpha_{ij} u_{ij} + \frac{\beta}{2} \left(\sum_{j=1}^{M} (a_j - a_{j-1}) u_{ij} \right)^2 - \frac{1}{2} \sum_{j=1}^{M} (a_j - a_{j-1}) I(u_{ij}) \right\}.$$

CAREY BUSINESS SCHOOL, JOHNS HOPKINS UNIVERSITY, 100 INTERNATIONAL DR, BALTIMORE, MD 21202

DEPARTMENT OF MATHEMATICS, FLORIDA STATE UNIVERSITY, 208 LOVE BUILDING, 1017 ACADEMIC WAY, TALLAHASSEE, FL 32306