

Homophily and Community Structure at Scale: An Application to a Large Professional Network

Juan Nelson Martínez Dahbura martinez.dahbura@sansan.com
Shota Komatsu shota.komatsu@sansan.com
Takanori Nishida nishida@sansan.com
Sansan, Inc.

Angelo Mele angelo.mele@jhu.edu
Johns Hopkins University

ASSA Meeting 2023, New Orleans

web: <http://meleangelo.github.io>

R package: <https://github.com/sansan-inc/lighthergm>

Motivation

- ① Professional networks:
 - how do they form?
 - how do they affect labor markets?
- ② Observable and unobservable heterogeneity; network externalities
- ③ Structural models of network formation
 - Econometric challenges

De Paula (2018), Chandrasekhar (2016), Graham (2017), Mele (2017,2021), Menzel (2016), Sheng (2021), De Paula et Al (2018), Graham (2016), Boucher and Mourifie (2017), Gao (2020), Badev (2021), Leung and Moon (2021) and many more...
 - Computational challenges

Dahbura et al 2021, Gaonkar and Mele (2018), Babkin et al (2021), Mele et al (2021), Vu et al (2013)
- ④ Potential uses: recommendation system, effect of policy or shocks on professional networks, key player analysis

R package



<https://github.com/sansan-inc/lighthergm>

- Open Source library for scalable estimation of this class of models
- Solves memory problems of original `hergm` library
- Allows the inclusion of (discrete) covariates in the model
- Improvements in speed by a factor of 14000 for some operations

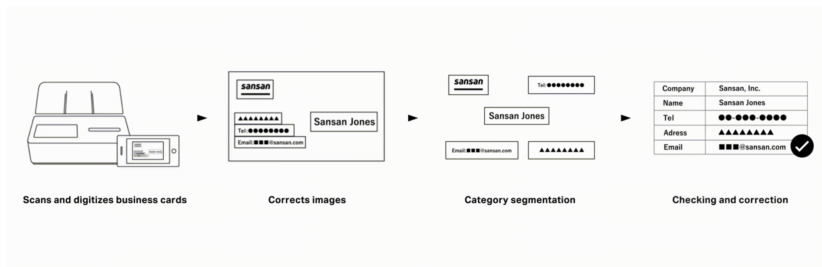
Estimation in the paper uses a Ubuntu Linux machine with 32GB of memory and 8 cores. The computation is performed with about 20 GB of memory for the block recovery step accounting for node covariates. All cores are in use during most of the calculation time

Eight

- In this work we use anonymized data from Eight on connections formed Jan-Dec 2019
- We include only users located in Tokyo who have uploaded a profile card at least once by the end of 2019 and have accepted terms of service
- We keep only nodes for which all covariates used in the analysis have non-missing values and that belong to the largest connected component of the network's 10-core.

Eight

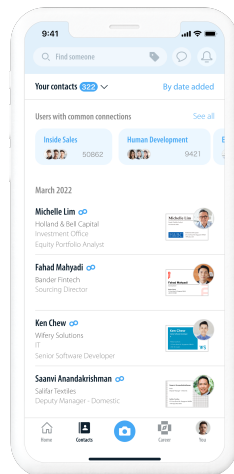
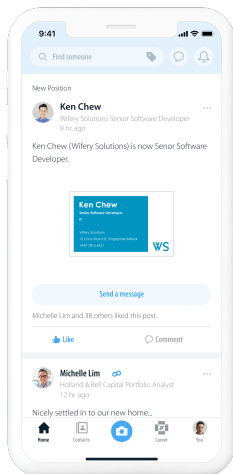
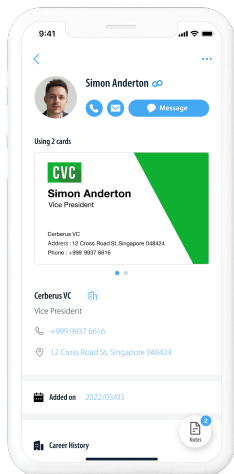
- Contact and career management app¹ with over 3 million users in Japan
- Allows users to scan physical business cards employing the smartphone's camera
- High quality digitization is achieved through the usage of advanced OCR algorithms and the help of human operators²



¹<https://8card.net/en/>

²More about the digitization process at: <https://bit.ly/3CqM0xp>

Eight



Data

- The resulting network has 30,323 nodes and 321,188 edges.
- The network is very sparse, with a density of roughly 0.0007.
- The data is highly geographically concentrated. About 84% of the nodes are located in just five districts of Tokyo.
- We also observe industrial concentration, especially in the Technology (22%) and Consulting (14%) industries.

Model

Communities and sequential network formation

- Time is discrete: $t = 0, 1, 2, 3, \dots$
- At $t = 0$ Nature assigns types

$$\mathbf{Z}_i \stackrel{iid}{\sim} \text{Multinomial}(1; \eta_1, \dots, \eta_K) \quad (1)$$

Remark: types not too large wrt network

Remark: each node belongs to one type only

(extensions to multiple communities possible as in Airoldi et al 2008)

- Conditional on $\mathbf{Z} = \mathbf{z}$, **network g is formed sequentially.**
- In each period t
 - 1 Two users i and j meet
 - 2 Users receive random matching shock ε_{ij}
 - 3 Users decide whether to form/cut/keep link g_{ij}
 → **maximize surplus generated by g_{ij}**

Assumptions

- 1 Users have positive probability of meeting any user

$$\text{Prob. } i \text{ and } j \text{ meet} := \rho(g_{-ij}, z_i, z_j, x_i, x_j, n) > 0 \quad (2)$$

- 2 Payoff of user i

$$U_i(\mathbf{g}, \mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) = \sum_{j=1}^n g_{ij} \left(u_{ij}(\boldsymbol{\alpha}, \boldsymbol{\beta}) + \sum_{r \neq i, j} g_{jr} g_{ri} v_{ijr}(\boldsymbol{\gamma}) \right)$$

- $u_{ij}(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \begin{cases} \alpha_w + \sum_{p=1}^P \beta_{wp} \mathbf{1}\{\mathbf{x}_{ip} = \mathbf{x}_{jp}\} & \text{if } i, j \text{ belong to same } k \\ \alpha_b + \sum_{p=1}^P \beta_{bp} \mathbf{1}\{\mathbf{x}_{ip} = \mathbf{x}_{jp}\} & \text{otherwise} \end{cases}$
- **Local transitivity:** $v_{ijr}(\boldsymbol{\gamma}) = \begin{cases} \gamma & \text{if } i, j, r \text{ belong to same } k \\ 0 & \text{otherwise} \end{cases}$

- 3 Matching shock ε_{ij} is logistic iid

Equilibrium

PROPOSITION. Conditional on \mathbf{z} , the long-run network distribution factorizes into **WITHIN**- and **BETWEEN**-types components

$$\begin{aligned} \pi(\mathbf{g}, \mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) &= \prod_{k=1}^K \frac{\exp [Q_{k,k}(\mathbf{g}_{k,k}, \mathbf{x}^{(k)}, \mathbf{z}; \boldsymbol{\alpha}_w, \boldsymbol{\beta}_w, \gamma)]}{c_{k,k}(\mathcal{G}_{k,k}, \mathbf{x}^{(k)}; \boldsymbol{\theta})} \\ &\times \prod_{l>k}^K \prod_{i \in \mathcal{C}_k} \prod_{j \in \mathcal{C}_l} \frac{\exp [g_{ij} (u_{ij}(\alpha_b, \beta_b) + u_{ji}(\alpha_b, \beta_b))]}{1 + \exp [(u_{ij}(\alpha_b, \beta_b) + u_{ji}(\alpha_b, \beta_b))]} \end{aligned}$$

where

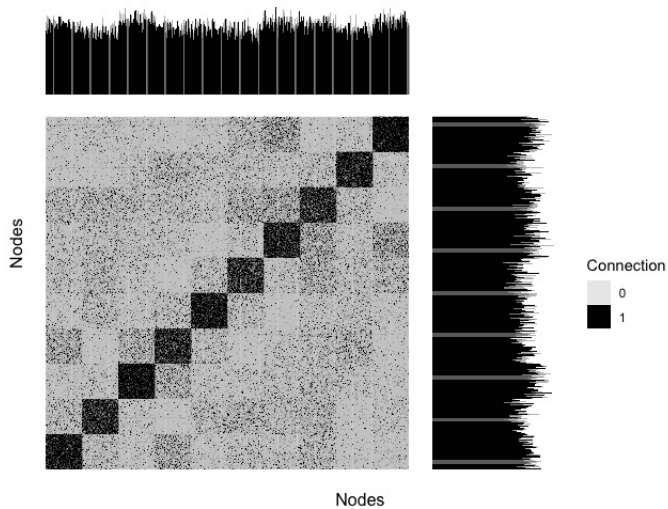
$$\begin{aligned} Q_{kk} &= \sum_{i=1}^n \sum_{j=1}^n z_{ik} z_{jk} g_{ij} \left(u_{ij}(\alpha_w, \beta_w) + \frac{2\gamma}{3} \sum_{r \neq i, j} z_{rk} g_{jr} g_{ri} \right) \\ c_{kk} &= \sum_{\omega \in \mathcal{G}} e^{Q_{kk}} \end{aligned}$$

REMARK. In long-run \rightarrow HERGM (Schweinberger-Handcock 2015)

Approximate Estimation

Approximate Maximum Likelihood

- **State-of-the-art:** Bayesian estimation (Mele JBES, forthcoming)
- For large networks, the Bayesian approach is *impractical or infeasible*
- On the other hand, if some conditions are satisfied:
 - ① communities small enough and
 - ② network large
 - ⇒ most probability mass is **across** blocks
 - ⇒ **conditionally independent links**
- Network resembles stochastic blockmodel *except within blocks*



Approximate Maximum Likelihood

Step 1: Compute approximate \hat{z} using SBM ($\gamma = 0$)

- **Conditions for good approximation:** Schweinberger-Stewart 2021, Babkin et al 2020
- **Variational approximations:**
Jordan and Wainwright 2003, Mele and Zhu 2023, Bickel and Chen 2013
- **Minorization-maximization:**
Vu et al 2013
- **Spectral methods:**
Athreya et al 2018, Mele et al 2022, Cong et al 2022, Hao et al 2022

Step 2: Approximate likelihood, given \hat{z}

- **Monte Carlo Maximum Likelihood (MCMC-MLE)**
Geyer and Thompson 1992, Snijders 2002, Mele 2011
- **Maximum Pseudolikelihood (MPLE)**
Snijders 2002, Boucher and Mourifie 2017, Babkin et al 2020

STEP 1: Variational Approximation

The full log-likelihood of our model can be written as follows

$$\begin{aligned} \mathcal{L}(\mathbf{g}, \mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\eta}) &= \log \sum_{\mathbf{z} \in \mathcal{Z}} P_{\boldsymbol{\eta}}(\mathbf{Z} = \mathbf{z}) \pi(\mathbf{g}, \mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) = \log \sum_{\mathbf{z} \in \mathcal{Z}} L(\mathbf{g}, \mathbf{x}, \mathbf{z}; \boldsymbol{\theta}, \boldsymbol{\eta}) \\ &\approx \log \sum_{\mathbf{z} \in \mathcal{Z}} L(\mathbf{g}, \mathbf{x}, \mathbf{z}; \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\gamma} = \mathbf{0}, \boldsymbol{\eta}) \end{aligned} \quad (3)$$

$$\begin{aligned} &= \log \sum_{\mathbf{z} \in \mathcal{Z}} \frac{q_{\boldsymbol{\xi}}(\mathbf{z})}{q_{\boldsymbol{\xi}}(\mathbf{z})} L_0(\mathbf{g}, \mathbf{x}, \mathbf{z}; \boldsymbol{\theta}, \boldsymbol{\eta}) \\ &\geq \ell_B(\mathbf{g}, \mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\eta}; \boldsymbol{\xi}) \end{aligned} \quad (4)$$

$$\begin{aligned} &= \sum_{i < j}^n \sum_{k=1}^K \sum_{l=1}^K \xi_{ik} \xi_{jl} \log \pi_{ij,kl}(g_{ij}, \mathbf{x}, \mathbf{z}) \\ &+ \sum_{i=1}^n \sum_{k=1}^K \xi_{ik} (\log \eta_k - \log \xi_{ik}) \end{aligned} \quad (5)$$

where

$$\pi_{ij,kl}(g_{ij}, \mathbf{x}, \mathbf{z}) = \text{Prob } i \text{ and } j \text{ of type } k \text{ and } l \text{ are connected} \quad (6)$$

STEP 1: Minorization-Maximization

Find function approximating $\ell_B(\mathbf{g}, \mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\eta}; \boldsymbol{\xi})$, but simpler to maximize.

$M(\boldsymbol{\xi}; \mathbf{g}, \mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\eta}, \boldsymbol{\xi}^{(s)})$ minorizes $\ell_B(\mathbf{g}, \mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\eta}; \boldsymbol{\xi})$ at $\boldsymbol{\xi}^{(s)}$ and iteration s if

$$M(\boldsymbol{\xi}; \mathbf{g}, \mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\eta}, \boldsymbol{\xi}^{(s)}) \leq \ell_B(\mathbf{g}, \mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\eta}; \boldsymbol{\xi}) \quad \text{for all } \boldsymbol{\xi} \quad (7)$$

$$M(\boldsymbol{\xi}^{(s)}; \mathbf{g}, \mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\eta}, \boldsymbol{\xi}^{(s)}) = \ell_B(\mathbf{g}, \mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\eta}; \boldsymbol{\xi}^{(s)}) \quad (8)$$

For stochastic blockmodels, Vu et al 2013 suggest

$$\begin{aligned} M(\boldsymbol{\xi}; \mathbf{g}, \mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\eta}, \boldsymbol{\xi}^{(s)}) &:= \sum_{i < j}^n \sum_{k=1}^K \sum_{l=1}^K \left(\xi_{ik}^2 \frac{\xi_{jl}^{(s)}}{2\xi_{ik}^{(s)}} + \xi_{jl}^2 \frac{\xi_{ik}^{(s)}}{2\xi_{jl}^{(s)}} \right) \log \pi_{ij;kl}^{(s)}(\mathbf{g}_{ij}, \mathbf{x}, \mathbf{z}) \\ &+ \sum_{i=1}^n \sum_{k=1}^K \xi_{ik} \left(\log \eta_k^{(s)} - \log \xi_{ik}^{(s)} - \frac{\xi_{ik}}{\xi_{ik}^{(s)}} + 1 \right). \end{aligned} \quad (9)$$

Parallelizes to n independent maximization problems

Variational Updates with discrete covariates

The update rules for ξ , η , and $\pi_{ij;kl}(g_{ij}, \mathbf{x}, \mathbf{z})$ follow

$$\xi^{(s+1)} := \arg \max_{\xi} M \left(\xi; \mathbf{g}, \mathbf{x}, \alpha^{(s)}, \beta^{(s)}, \eta^{(s)}, \xi^{(s)} \right),$$

$$\eta_k^{(s+1)} := \frac{1}{n} \sum_{i=1}^n \xi_{ik}^{(s+1)}, \quad k = 1, \dots, K,$$

and

$$\pi_{ij;kl}^{(s+1)}(d, \chi_1, \dots, \chi_p, \mathbf{z}) := \frac{\sum_{i=1}^n \sum_{j \neq i} \xi_{ik}^{(s+1)} \xi_{jl}^{(s+1)} \mathbf{1}\{g_{ij} = d, \chi_{1,ij} = \chi_1, \dots, \chi_{p,ij} = \chi_p\}}{\sum_{i=1}^n \sum_{j \neq i} \xi_{ik}^{(s+1)} \xi_{jl}^{(s+1)} \mathbf{1}\{\chi_{1,ij} = \chi_1, \dots, \chi_{p,ij} = \chi_p\}},$$

for $k, l = 1, \dots, K$ and $d, \chi_1, \dots, \chi_p \in \{0, 1\}$, respectively.

$\chi_{p,ij} = \mathbf{1}\{x_{ip} = x_{jp}\}$. Generalizations of this specification are allowed.

STEP 2: Maximum Pseudolikelihood Estimator

Given estimated $\hat{\mathbf{z}}$, compute conditional prob of link

WITHIN BLOCKS

$$p_{ij}(\mathbf{g}, \mathbf{x}, \boldsymbol{\theta}; \hat{\mathbf{z}}) = \Lambda \left(u_{ij}(\boldsymbol{\alpha}_w, \boldsymbol{\beta}_w) + u_{ji}(\boldsymbol{\alpha}_w, \boldsymbol{\beta}_w) + 4\gamma \sum_{r \neq i, j} g_{jr} g_{ir} \right)$$

BETWEEN BLOCKS

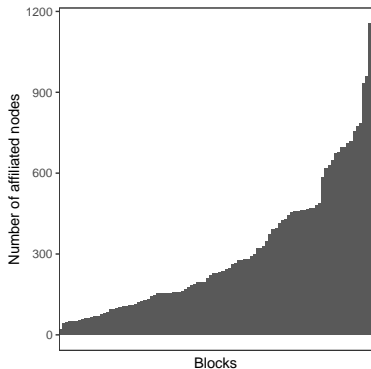
$$p_{ij}(\mathbf{g}, \mathbf{x}, \boldsymbol{\theta}; \hat{\mathbf{z}}) = \Lambda (u_{ij}(\boldsymbol{\alpha}_b, \boldsymbol{\beta}_b) + u_{ji}(\boldsymbol{\alpha}_b, \boldsymbol{\beta}_b))$$

where $\Lambda(u) = e^u / (1 + e^u)$ is the logistic function.

The **pseudolikelihood estimator** solves

$$\begin{aligned} \hat{\boldsymbol{\theta}}_{PL} &= \arg \max_{\boldsymbol{\theta}} \ell_{PL}(\mathbf{g}, \mathbf{x}, \boldsymbol{\theta}; \hat{\mathbf{z}}) \\ &= \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^n \sum_{j>i}^n [g_{ij} \log p_{ij}(\mathbf{g}, \mathbf{x}, \boldsymbol{\theta}) + (1 - g_{ij}) \log(1 - p_{ij}(\mathbf{g}, \mathbf{x}, \boldsymbol{\theta}))] \end{aligned}$$

Empirical results: block size



	Between	Within
	(1)	(2)
Intercept (α)	-7.709*** (0.002)	-4.754*** (0.005)
Shared Contacts (γ)		0.736*** (0.004)
Same Location (β_1) (H3 Tile)	0.333*** (0.007)	0.006 (0.012)
Same Industry (β_2)	0.694*** (0.005)	0.034*** (0.009)
Same Occupation (β_3)	0.409*** (0.006)	0.041*** (0.010)
Bayesian Inf. Crit.	4,171,768	808,597

Notes: * $p < 0.1$; ** $p < 0.05$; *** $p < 0.01$.

All estimates are obtained using a maximum pseudolikelihood estimator, conditioning on the estimated block structure. Block recovery was performed for a total of 100 blocks. We employ 20,000 EM iterations without employing node covariates, and a final run of 100 iterations employing covariates.

Summary

Summary

- Equilibrium model with community structure
- Approximate maximum likelihood
 - Use SBM likelihood
 - Variational Approximations for SBM
 Bickel et al 2013; Jordan and Wainwright 2003; Mele and Zhu 2020
 - Use Minorization algorithm to speed up computation
 Vu et al 2013; Babkin et al 2021
- Find evidence of homophily and transitivity (see also Dahbura et al 2021)

In progress

- Improve package speed: MPLE vs. MC-MLE; Spectral Methods instead of Variational approximations; initialization with InfoMap; Other clustering methods: Louvain, etc.
- Estimation with more covariates (discrete)
- Goodness of fit, counterfactual simulations
- Effect of networks on outcomes: key player simulations

THANK YOU!

More of this at:

arxiv: <https://arxiv.org/abs/2105.12704>

Contact:

Juan Nelson Martínez Dahbura martinez.dahbura@sansan.com

Shota Komatsu shota.komatsu@sansan.com

Takanori Nishida nishida@sansan.com

Angelo Mele angelo.mele@jhu.edu

R package: <https://github.com/sansan-inc/lighthergm>

Does it work?

Monte Carlo experiments

Monte Carlo simulation to check approximate ML

Specification of payoff - no covariates

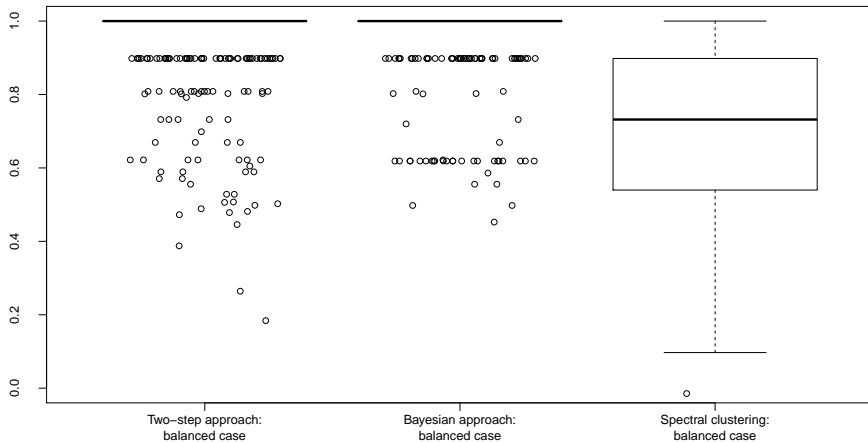
$$U_i(g, x, z; \theta) = \theta_1 \sum_{k=1}^K \sum_{j=1}^n z_{ik} z_{jk} g_{ij} + \theta_2 \sum_{j=1}^n z_{ik} z_{jk} z_{rk} g_{ij} \mathbf{1}_{ij} \\ + \theta_B \sum_{k=1}^K \sum_{\ell > k} \sum_{j=1}^n z_{ik} z_{j\ell} g_{ij}$$

where $\mathbf{1}_{ij} = 1$ if i and j have at least 1 partner in common

- $n = \{30, 150, 2500\}$
- $K = \{3, 100\}$
- Model with no covariates, only edges and transitive triples
- Simulate 500 networks and estimate model
- Parameters $(\theta_1, \theta_2) = (-1, .5) / \log(n_k)$ for $n = \{30, 150\}$
- Parameters $(\theta_1, \theta_2) = (-2, 1) / \log(n_k)$ for $n = 2500$
- Parameter $\theta_B = -3 / \log(n)$

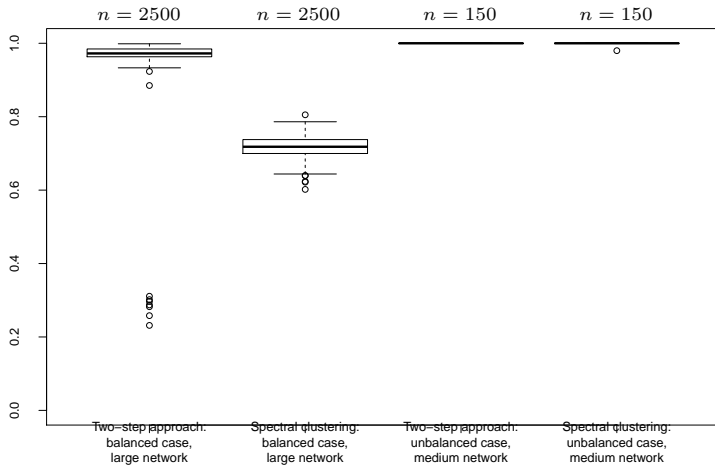
Results: Approximate ML vs Bayesian vs Spectral

Estimation of \hat{z} , $n = 30$, $K = 3$



Results: Approximate ML

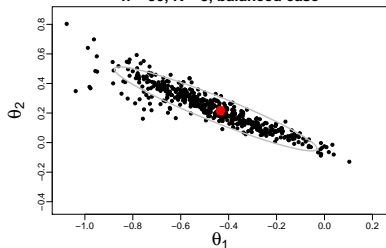
Estimation of \hat{z} , large and medium size networks



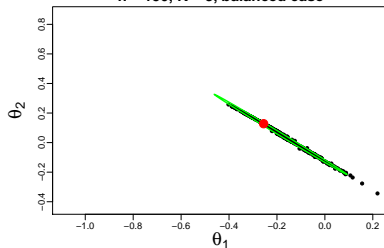
Results: Approximate ML

Parameters estimates, point estimates and 95% ellipses

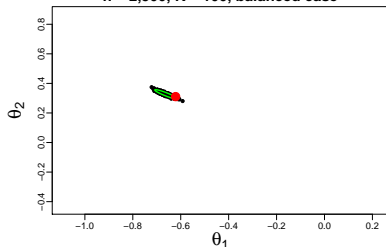
n = 30, K = 3, balanced case



n = 150, K = 3, balanced case



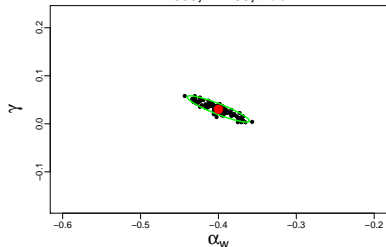
n = 2,500, K = 100, balanced case



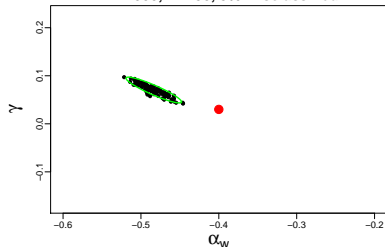
Results: Approximate ML

Parameters estimates as a function of misclassification rate for z

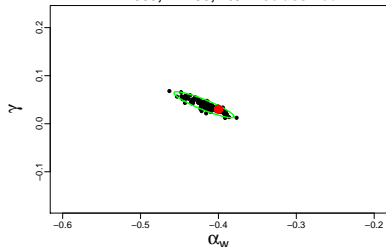
$n = 1000, K = 50, \text{true } z$



$n = 1000, K = 50, 5\% \text{ misclassified}$



$n = 1000, K = 50, 1\% \text{ misclassified}$



$n = 1000, K = 50, 10\% \text{ misclassified}$

