

# Modelling and analysis of social network data from rank preferences



Pierre OSSELIN  
St Peter's College  
University of Oxford

A thesis submitted for the degree of  
*MSc in Mathematical Sciences*  
Trinity 2018

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Community representation</b>	<b>4</b>
2.1	Mathematic Formulation . . . . .	4
2.2	Computation of the posterior . . . . .	5
2.3	Simulations on simulated Data . . . . .	6
2.3.1	Fixed parameters for the prior . . . . .	6
2.3.2	Model with degree-correction . . . . .	8
2.3.3	Performance Comparison . . . . .	9
<b>3</b>	<b>Model with Covariates</b>	<b>10</b>
3.1	Notation . . . . .	10
3.2	Classic Graphical Models . . . . .	11
3.2.1	ERGM model . . . . .	11
3.2.2	The $p_1$ Model . . . . .	11
3.2.3	The $p_2$ Model . . . . .	11
3.2.4	The SRM Model . . . . .	12
3.3	Plackett-Luce Model with Covariates . . . . .	14
<b>4</b>	<b>EM algorithm with community representation</b>	<b>16</b>
4.1	Mathematical formulation . . . . .	16
4.2	EM formulation . . . . .	17
4.3	EM formulation with variable ranks . . . . .	18
4.4	Simulation . . . . .	20
4.4.1	Toy Example . . . . .	20
4.4.2	Performance Evaluation . . . . .	20
<b>5</b>	<b>Test on simulated Data</b>	<b>22</b>
5.1	Fixed Community representation . . . . .	22
5.2	Model With Degree Correction . . . . .	24
5.3	Model With Covariates . . . . .	25
<b>6</b>	<b>Simulations on the Sampsons monk dataset</b>	<b>26</b>

# Chapter 1

## Introduction

This paper aims at performing inference on social network through partial ranking of the individuals

We use the Plackett-Luce model:

$$p(\rho_i | (\lambda)) = \prod_{k=1}^K \frac{\lambda_{i\rho_k}}{\sum_{j \neq i} \lambda_{i,j} - \sum_{l=1}^{k-1} \lambda_{i,\rho_l}} \quad (1.1)$$

Where  $K$  is the number of ranked individuals. This model is equivalent to successively choose our favorite individual by pondering our strength of affiliation with each person, and take randomly the best remaining by taking into account the proportion of this strength among the other, which is a conceivable model.

The first thing to do is to prove that Plackett-Luce is well defined as a probability distribution. It is easy to prove it by induction over  $K \leq n$ .

We denote  $\Omega_{K,n} := \{\rho : \{1, \dots, K\} \rightarrow I_K \text{ bijective}, I_K \in [n]^{(K)}\}$ . For  $K = 1$  :

$$\sum_{\rho \in \Omega_{1,n}} p(\rho | (\lambda)) = \sum_{k \neq i} \frac{\lambda_{i,k}}{\sum_{j \neq i} \lambda_{i,j}} = 1 \quad (1.2)$$

If the result is true for  $K - 1 < n$  then, with

$\forall j \neq i, \Omega_{K,n,j} = \{\rho : \{1, \dots, K\} \rightarrow I_K \text{ bijective}, I_K \in [n]^{(K)}, \rho(1) = j\}$  and  $(\hat{\lambda})_j$  the matrix without column  $j$ :

$$\sum_{\rho \in \Omega_{K,n}} p(\rho | (\lambda)) = \sum_{j=1}^n \sum_{\rho \in \Omega_{K,n,j}} p(\rho | (\lambda)) \quad (1.3)$$

$$= \sum_{j=1}^n \frac{\lambda_{i,j}}{\sum_{l \neq i} \lambda_{i,l}} \underbrace{\sum_{\rho \in \Omega_{K-1,n-1}} p(\rho | (\hat{\lambda})_j)}_{= 1 \text{ by induction}} = 1 \quad (1.4)$$

# Chapter 2

## Community representation

### 2.1 Mathematic Formulation

Here we suppose that the strength between individuals are characterized by their similitude to certain communities.

We introduce  $p$  the number of communities and  $(w_{i,k})_{1 \leq i \leq n, 1 \leq k \leq p}$  the strength of affiliation of the individual  $i$  to the community  $p$ .

$$\lambda_{i,j} = \sum_{k=1}^p w_{i,k} w_{j,k} \quad (2.1)$$

This formulation is equivalent to the rank matrix approximation often encountered in the literature where we assume that:

$$(\lambda) = WW^T \quad (2.2)$$

where  $W$  is the matrix of the  $(w_{i,k})$

We suppose here that the strength of the relations are only determined by the proximity with a certain community hence that the relations are reciprocal. From the reordering inequality, we can observe that  $\lambda_{i,j}$  is maximized when both  $i$  and  $j$  have the same community ranking preferences. We can interpret the formula  $\lambda_{i,j} = \sum_{k=1}^p w_{i,k} w_{j,k}$  also as a scalar product in the  $p$ -dimensional space of communities. In this space, each individual is a point represented by a vector  $w_i \in (\mathbb{R}^+)^p$  the strength of affiliation is also computed by  $\lambda_{i,j} = ||w_i|| \cdot ||w_j|| \cdot \cos(\theta)$  we can interpret the norm of the vector as the individual global influence or sociability, and the  $\theta$  as the divergence in personality. Hence a given individual will appreciate another individual with great sociability but different personality or views as much as a low sociability individual with the same mindset. In this interpretation, the prior function attributed to the  $w_{i,k}$  is equivalent to putting a topography on the space, allowing points to lie more likely in a certain area. This is illustrated in the Figure 2.1

The main issue regarding this geometric approach is that having a bigger affiliation means that the community points have a greater scalar product, but doesn't mean those points are nearer. If we wanted to use a clustering approach to this problem,



Figure 2.1: Topography of the community space for the same gamma distribution

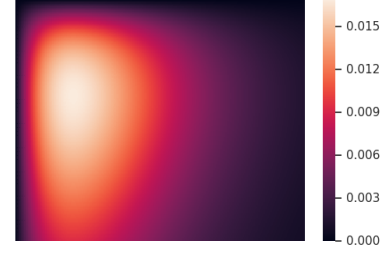


Figure 2.2: Topography of the community space for different gamma distribution

we should build a space in which the mapping to this space make it equivalent having great scalar product in the input space and being near in the output space.

## 2.2 Computation of the posterior

Let's define the log-posterior of the community representation. Which is the sum of the log priors of the parameters and the log-likelihood.

$$l((w); D) = \sum_{i=1}^n \sum_{k=1}^p (\alpha - 1) \ln(w_{i,k}) - \beta w_{i,k} + \sum_{i=1}^n \sum_{j=1}^K \ln(\lambda_{i,\rho_j}) - \ln\left(\sum_{j \neq i} \lambda_{i,j} - \sum_{l=1}^{j-1} \lambda_{i,\rho_l}\right) \quad (2.3)$$

To make the optimization easier, we can compute the gradient of the posterior according to the parameter  $w_{r,s}$

With

$$\frac{\partial \lambda_{i,j}}{\partial w_{r,s}} = 1_{i=r} w_{j,s} + 1_{j=r} w_{i,s} \quad (2.4)$$

We have :

$$\begin{aligned} \frac{\partial l((w); D)}{\partial w_{r,s}} &= \frac{a-1}{w_{r,s}} - b + \sum_{j=1}^K \frac{1}{\lambda_{r,\rho_j^{(r)}}} w_{\rho_j^{(r)},s} + \sum_{i=1; i \neq r}^n \sum_{j=1}^K \frac{1}{\lambda_{i,r}} w_{i,s} 1_{\rho_j^{(r)}=r} - \\ &\quad \sum_{j=1}^K \frac{1}{\sum_{j \neq r} \lambda_{r,j} - \sum_{l=1}^{j-1} \lambda_{r,\rho_l^{(r)}}} \sum_{j \neq r} w_{j,s} - \sum_{i=1; i \neq r}^n \sum_{j=1}^K \frac{1}{\sum_{j \neq i} \lambda_{i,j} - \sum_{l=1}^{j-1} \lambda_{i,\rho_l^{(i)}}} w_{i,s} + \\ &\quad \sum_{j=1}^K \frac{1}{\sum_{j \neq r} \lambda_{r,j} - \sum_{l=1}^{j-1} \lambda_{r,\rho_l^{(r)}}} \sum_{l=1}^{j-1} w_{\rho_l^{(r)},s} + \sum_{i=1; i \neq r}^n \sum_{j=1}^K \frac{1}{\sum_{j \neq i} \lambda_{i,j} - \sum_{l=1}^{j-1} \lambda_{i,\rho_l^{(i)}}} \sum_{l=1}^{j-1} 1_{\rho_l^{(i)}=r} w_{i,s} \end{aligned} \quad (2.5)$$

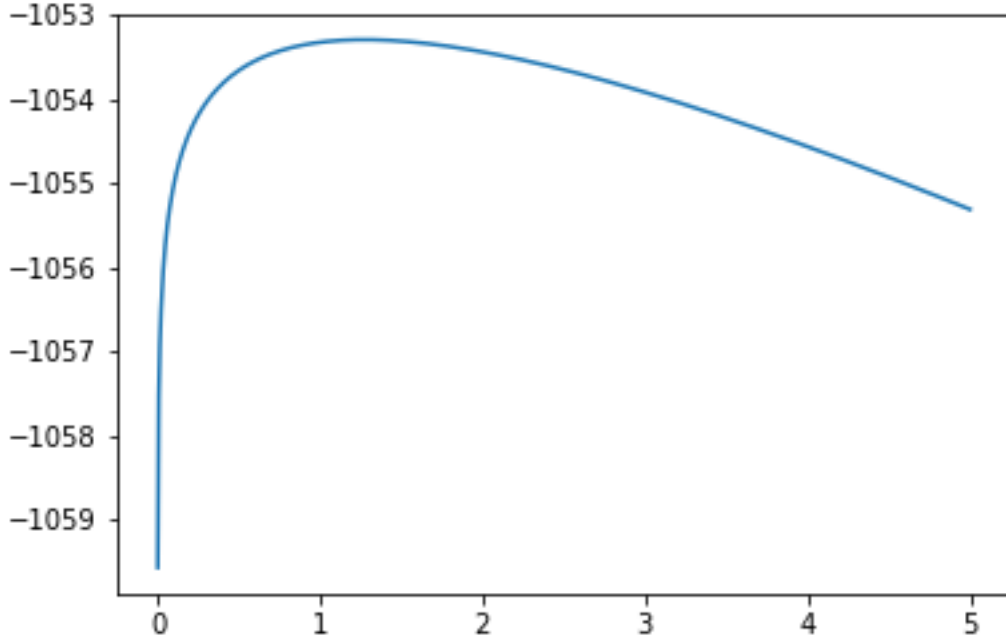


Figure 2.3: Posterior according to one coordinate

## 2.3 Simulations on simulated Data

### 2.3.1 Fixed parameters for the prior

We derived two optimization algorithm to find the MAP, one based on coordinate-wise optimization, given that for the  $a$  and  $b$  considered, the log-posterior seems coordinate-wise concave, hence assuring a more stable convergence. The other algorithm is based on the L-BFGS-B algorithm, popular in machine learning optimization.

For this simulation we took  $a = 2$ ,  $b = 1$ ,  $n = 15$ ,  $p = 3$ ,  $k = 3$ , drew a sample from  $w$ , computed  $\lambda$  and  $\rho$ . Then we maximized the posterior according to the two methods described above, in order to estimate the community strengths and then the individual mutual strength.

In figure 2.4 we can see that the coordinate and global optimization have converged to the same result. We obtain an absolute error of 0.02. Which is equivalent to 0.3 reported to the number of individual, and a ranking gap (mean absolute ranking error) of 2.13 for this individual.

It is interesting to compare the ranking gap to the expected ranking gap from a random permutation. We can manually compute with a tedious calculation that

$$\frac{1}{n} \mathbb{E}_{\sigma_n} \left[ \sum_{k=1}^n |\sigma(k) - k| \right] = \frac{n^2 - 1}{3n} \quad (2.6)$$

For  $n = 18$  we have a random expectation of the ranking gap of 6. Hence the

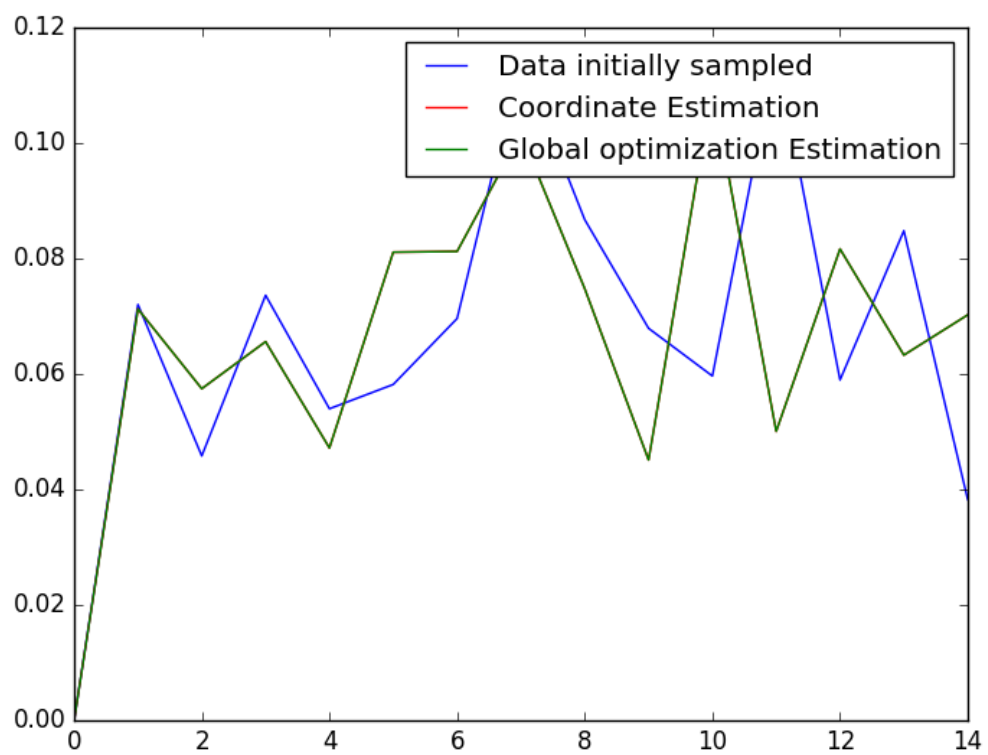


Figure 2.4: Comparison between the sampled and estimated Social affiliation of the first individual

model tends to preserve the "global" social network for each individual.

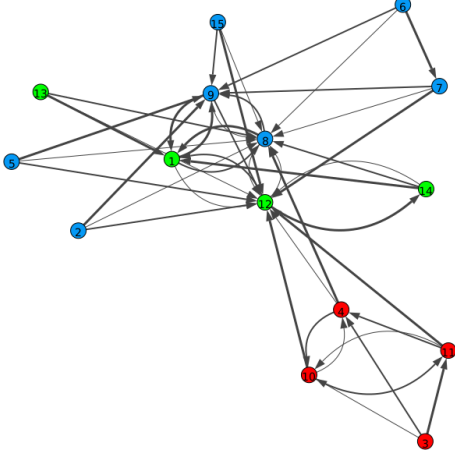


Figure 2.5: Original Community description

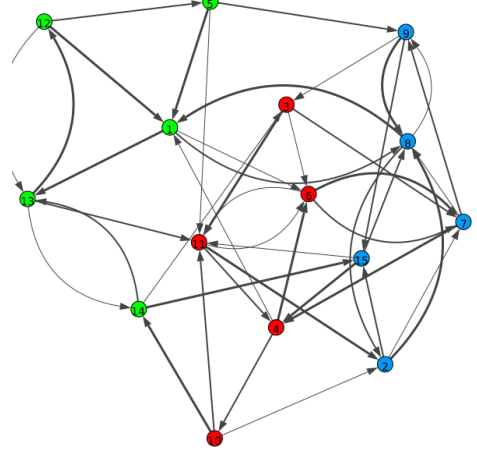


Figure 2.6: Inferred community

We can see that the inferred communities are quite similar to the original one. The original red and green communities are preserved except for the individual 6 and 5 that has passed from the blue community to the red and green respectively. Here we have assigned an individual to his preferred community, the width of an edge represents the importance the son represents for the father.

### 2.3.2 Model with degree-correction

In this subsection, we try to capture the different popularity that a community can have over the other communities. In the previous model, the regularization terms  $w_{i,j}$  in the posterior had the same weights, irrespective of the community concerned (Having a too big  $w_{i,j}$  was equally improbable regardless of the community). In a model with degree-correction, i.e in which  $w_{i,k} \sim \Gamma(\alpha, \beta_k)$ , we do not prevent the possibility of a popular community in the regularization, which occurs by the emergence of popular individuals.



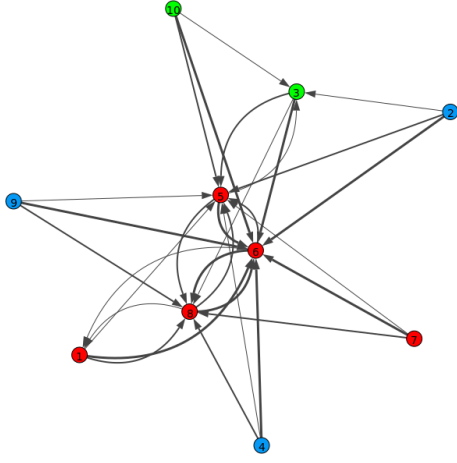


Figure 2.7: Original Community description for the degree Correction Model

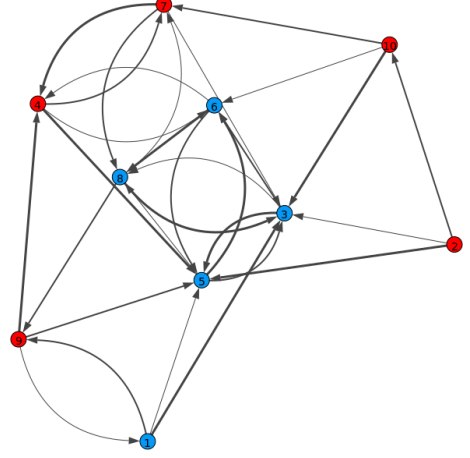


Figure 2.8: Inferred community for the degree correction Model

We can see in this example that the green community has been absorbed in an other community by the inference, and that 7 has been replaced by 3 in the other community.

### 2.3.3 Performance Comparison

To evaluate the performance of the models, we make prediction on the preferences for individuals who have not indicated their preferences, provide further preferences from individuals who did, or extract from the lists interpretable information about the structure of the underlying friendship network, e.g. in terms of latent communities. For this we will use different metrics. The ranking gap between latent appreciation, the absolute error in the strengths and the missed community. To do that, we will sample the strengths between individuals  $\lambda_{i,j}$  and then sample full rankings from it. From this ranking we then hide a certain part of it, we feed the remaining parts to our model and then compute the error in terms of metrics previously defined to evaluate the performance.

# Chapter 3

## Model with Covariates

The purpose of ERGMs, in a nutshell, is to describe parsimoniously the local selection forces that shape the global structure of a network. To this end, a network dataset may be considered like the response in a regression model, where the predictors are things like propensity for individuals of the same sex to form partnerships or propensity for individuals to form triangles of partnerships. Thus, an ERGM can help us quantify the strength of this intra-group effect. The information gleaned from use of an ERGM may then be used to understand a particular phenomenon or to simulate new random realizations of networks that retain the essential properties of the original. Handcock, Hunter, Butts, Goodreau, and Morris (2008) say more about the purpose of modeling with ERGMs; yet in this article, we focus primarily on technical details.

### 3.1 Notation

Let  $N$  be the set of nodes in the network of interest, indexed  $\{1, \dots, n\}$ . The relationships in the network may be directed (e.g., friendship nominations, messages) or undirected (e.g., sexual partnerships, conversations). In the former case, we define the set of dyads (here used to refer to potential relationships)  $\mathbb{Y}$  to be a subset of  $N \times N$ , the set of ordered pairs of nodes; in the latter case, it is a subset of  $\{\{i, j\} : (i, j) \in N \times N\}$ , the unordered pairs of nodes. (We will also use  $u(\mathbb{Y})$  to refer to an unordered version of  $\mathbb{Y}$ , i.e.,  $u(\mathbb{Y}) \equiv \{\{i, j\} : (i, j) \in \mathbb{Y}\}$ .) Usually,  $\mathbb{Y}$  is further constrained in that in most social networks studied, a node cannot have a relationship of interest with itself, excluding pairs of the form  $(i, i)$ . For binary networks, in which the relationship of interest must be either present or absent, we use  $\mathcal{Y} \in 2^{\mathbb{Y}}$ , the set of subsets of  $\mathbb{Y}$ , to refer to the set of possible networks of interest, which may be further constrained (that is,  $\mathbb{Y}$  may be a proper subset of  $2^{\mathbb{Y}}$ ). We will use  $\mathbf{Y}$  to refer to network random variables and  $y \in \mathcal{Y}$  to refer to their realizations, and  $y_{i,j}$  be an 0 - 1 indicator of whether a relationship of interest is present between  $i$  and  $j$  in a binary network context.

## 3.2 Classic Graphical Models

### 3.2.1 ERGM model

In the ERGM model (exponential-family random graph models), the distribution of  $\mathbf{Y}$  can be parameterized in the form:

$$P_{\theta}(Y = y) = \frac{\exp((\eta(\theta))^T * z(x))}{\chi(\theta)}, y \in \mathcal{Y} \quad (3.1)$$

where  $\theta$  is a q-vector of model parameters, which are mapped to a p-vector of natural parameters by  $\eta(\cdot)$ , and  $g(\cdot)$  is a p-vector of sufficient statistics, which capture network features of interest. These sufficient statistics typically embody the features of the network of interest that are believed to be significant to the social process which had produced it, such as degree distribution (e.g., propensity towards monogamy in sexual partnership networks), homophily (i.e., birds of a feather flock together), and triad-closure bias (i.e., a friend of a friend is a friend) . (Morris, Handcock and Hunter, 2008)

### 3.2.2 The $p_1$ Model

The  $p_1$  model is the first to propose log-linear models for social networks. Suppose that we take  $\mathcal{Y}$  to be the set of all directed graphs, with independent dyads [i.e., the pairs  $(Y_{i,j}, Y_{j,i})$  are independent for different choices of  $\{i, j\}$ ] and the following model for tie probabilities:

$$P(Y_{i,j} = y_1, Y_{j,i} = y_2) = \exp\{y_1(\mu + \alpha_i + \beta_j) + y_2(\mu + \alpha_j + \beta_i) + y_1 y_2 \rho\} \backslash k_{ij} \\ y_1, y_2 = 0, 1; i, j = 1, \dots, n; i \neq j \quad (3.2)$$

In this model, we can interpret  $\mu$  as a density parameter, constituting an overall mean,  $\alpha_i$  as a productivity parameter, characterizing  $i$  as a sender,  $\beta_j$  as an attractiveness parameter, characterizing  $j$  as a receiver, and  $\rho$  as a parameter indicating the force of reciprocation.

### 3.2.3 The $p_2$ Model

While ERGMs are useful for modeling global network characteristics, models based on conditional independence (given latent variables) are useful for multiple reasons. First, ERGMs are not well understood and sometimes possess undesirable properties, e.g., model degeneracy. Second, the likelihood function of ERGMs is intractable, complicating statistical computing. Third, there may be unobserved heterogeneity or unobserved structure. Models based on independence conditional on latent variables: random effects models and mixed effects models and extensions; stochastic block models and extensions, including mixed membership models; and latent space models and extensions.

The aim of the  $p_2$  model is to relate binary network data to covariates while taking into account the specific network structure. This requires a kind of bivariate logistic regression model capable of handling the dependence of network data. The  $p_2$  model is based on the same formula. The density and reciprocity parameters, however, are allowed to vary over dyads, and are therefore denoted  $\mu_{ij}$  and  $\rho_{ij}$ . Moreover, parameters  $\alpha_i$ ,  $\beta_j$ ,  $\mu_{ij}$ , and  $\rho_{ij}$  are further modelled using covariates. We first consider the node-specific parameters  $\alpha_i$  and  $\beta_j$ . Covariates and random effects are included in a linear regression model for the productivity and attractiveness parameters.

$$\begin{cases} \alpha = X_1\gamma_1 + A \\ \beta = X_2\gamma_2 + B \end{cases} \quad (3.3)$$

This formulation expresses the plausible idea that attractiveness (or popularity) and productivity (or sociability) depend on actor attributes (denoted by  $X_1$  and  $X_2$ , respectively, where the same or different attributes may be used for attractiveness and productivity) with corresponding weights  $\gamma_1$  and  $\gamma_2$ . Naturally, the attributes do not explain all variation in attractiveness and productivity parameters, as is represented by the residual terms  $A$  and  $B$ ,  $n \times 1$  vectors with components  $A_i$  and  $B_i$ . The residuals are modeled as normally distributed random variables with expectation 0 and variances  $\sigma_A^2$  and  $\sigma_B^2$ , respectively. Parameters  $\sigma_A^2$  and  $\sigma_B^2$  can be interpreted as unexplained variance, that is, the variance of the  $\alpha$ 's and  $\beta$ 's that is left after taking into account the effect of the covariates  $X_1$  and  $X_2$ . The productivity and attractiveness parameters of the same node are correlated:  $\forall i, \text{cov}(A_i, B_i) = \sigma_{AB}$ . Independence is assumed for parameters of different actors by setting  $\text{cov}(A_i, A_j) = \text{cov}(B_i, B_j) = \text{cov}(A_i, B_j) = 0$  for  $i \neq j$  (cf. Wong, 1987). If no external information on actors is available, the terms  $X_1\gamma_1$  and  $X_2\gamma_2$  vanish and  $\sigma_A^2$  and  $\sigma_B^2$  denote the variances of the  $\alpha$  and  $\beta$  parameters, respectively. Then a pure random effects model results with, apart from the density and reciprocity parameters, only two variance parameters and one covariance parameter. Thus, the  $p_2$  model without covariates is a more parsimonious model than the  $p_1$  model with well-interpretable parameters. Obviously, the fit of  $p_1$  will be better than the fit of  $p_2$  without covariates.

### 3.2.4 The SRM Model

ERGM models work very poorly, the learned parameters do not generate data that resembles the input and tend toward wholly connected or completely empty graphs. Similarly, Since a binary random graph model does not accommodate rank data, analysis of these data with such a model typically begins by reducing all positive ranked relations to ones, and all negative ranked relations to zeros, leaving unranked relations as zeros as well. Such a data analysis essentially throws away some of the information in the data. Furthermore, the support of most binary random graph models consists of all possible graphs on the node set. In contrast, the data collection scheme used by Sampson (1969) was censored, as no graph with more than three outgoing edges per node could have been observed.

We take the same formulation as for the community representation. Latent variable  $\mathbf{Y}$  representing the strength of affiliation between individuals. And  $\mathbf{S}$  the sociomatrix given in data.  $\mathbf{S} = \{s_{i,j} : i \neq j\}$ , coded so that  $s_{i,j} = 0$  if  $j$  is not nominated by  $i$ ,  $s_{i,j} = 1$  if  $j$  is  $i^{\text{th}}$  least favored nomination, and so on. Under this coding,  $s_{i,j} \leq s_{i,k}$  if  $i$  scores  $j$  more highly than  $k$ , or if  $i$  nominates  $j$  but not  $k$ . Letting  $a_i = \{1, \dots, n\} \setminus \{i\}$  be the set of individuals whom person  $i$  may potentially nominate, each observed outdegree  $d_i = \sum_{j \in a_i} 1(s_{i,j} > 0)$  satisfies  $d_i \leq m$ . We then write

$$s_{i,j} = [(m - \text{rank}_i(y_{i,j} + 1) \wedge 0)] \times 1(y_{i,j} > 0) \quad (3.4)$$

and its inverse:

$$s_{i,j} > 0 \Rightarrow y_{i,j} > 0 \quad (3.5)$$

$$s_{i,j} > s_{i,k} \Rightarrow y_{i,j} > y_{i,k} \quad (3.6)$$

$$s_{i,j} = 0 \text{ and } d_i < m \Rightarrow y_{i,j} \leq 0 \quad (3.7)$$

Given a statistical model  $\{p(Y|\theta) : \theta \in \Theta\}$  for the underlying social relations  $Y$ , inference for the parameter  $\theta$  can be based on a likelihood derived from the observed scores  $S$ . The likelihood is, as usual, the probability of the observed data  $S$  as a function of the parameter  $\theta$ . To obtain this probability, let  $F(S)$  denote the set of  $Y$ -values that are consistent with  $S$  in terms of associations (5)–(7) above. Since the entries of  $S$  are the observed scores if and only if  $Y \in F(S)$ , the likelihood is given by

$$L_F(\theta : \mathbf{S}) = \Pr(Y \in F(S)|\theta) = \int_{F(S)} p(Y|\theta) d\mu(Y) \quad (3.8)$$

where  $\mu$  is a measure that dominates the probability densities  $\{p(Y|\theta) : \theta \in \Theta\}$ .

Given current values of  $(\theta, Y)$ , one step of the Gibbs sampler for the FRN likelihood proceeds by updating the values as follows:

1. Simulate  $\theta \sim p(Y|\theta)$
2. For each  $i \neq j$ , simulate  $y_{i,j} \sim p(y_{i,j}|\theta, Y_{-(i,j)}, Y \in F(S))$  as follows:
  - (a) if  $s_{i,j} > 0$ , simulate  $y_{i,j} \sim p(y_{i,j}|Y_{-(i,j)}, \theta) \times 1(\max\{y_{i,k} : s_{i,k} < s_{i,j}\} \leq y_{i,j} \leq \min\{y_{i,k} : s_{i,k} > s_{i,j}\})$ ;
  - (b) if  $s_{i,j} = 0$ , and  $d_i < m$  simulate  $y_{i,j} \sim p(y_{i,j}|Y_{-(i,j)}, \theta) \times 1(y_{i,j} \leq 0)$ ;
  - (c) if  $s_{i,j} = 0$ , and  $d_i = m$  simulate  $y_{i,j} \sim p(y_{i,j}|Y_{-(i,j)}, \theta) \times 1(y_{i,j} \leq \min\{y_{i,k} : s_{i,k} > 0\})$ ;

In the above steps, " $y \sim f(y)$ " means "simulate  $y$  from a distribution with density proportional to  $f(y)$ ." For each ordered pair  $(i, j)$ , step 2 of this algorithm will generate a value of  $y_{i,j}$  from its full conditional distribution, constrained so that conditions (5)(7) that define the FRN likelihood are met.

To deal with the previous remarks, we apply a model in the case where  $\theta$  represents the parameters in the following standard regression model for relational data:

$$y_{i,j} = \beta^T \mathbf{x}_{i,j} + a_i + b_j + \epsilon_{i,j} \quad (3.9)$$

$$\begin{pmatrix} a_i \\ b_i \end{pmatrix}, i = 1, \dots, n \sim \text{i.i.d normal}(0, \Sigma_{ab}) \quad (3.10)$$

$$\begin{pmatrix} \epsilon_{i,j} \\ \epsilon_{j,i} \end{pmatrix}, i = 1, \dots, n \sim \text{i.i.d normal}(0, \sigma^2 \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}) \quad (3.11)$$

The additive row effect  $a_i$  is often interpreted as a measure of person  $i$ 's "sociability" whereas the additive column effect  $b_i$  is taken as a measure of  $i$ 's "popularity.". The parameter  $\rho$  represents potential correlation between  $y_{i,j}$  and  $y_{j,i}$ . In a mixed-effects version of this model, the possibility that a persons sociability  $a_i$  is correlated with their popularity  $b_i$  can be represented with a covariance matrix  $\Sigma_{ab}$ . The covariance among the elements of  $Y = y_{i,j} : i \neq j$  induced by  $\Sigma_{ab}$  and is called the SRM (Warner et al. ,1979).

### 3.3 Plackett-Luce Model with Covariates

To adapt a model with covariates to the Bradley-Terry likelihood we will model the strength matrix as follow:

$$\ln(y_{i,j}) = \beta^T \mathbf{x}_{i,j} + a_i + b_j + \epsilon_{i,j} \quad (3.12)$$

$$\begin{pmatrix} \epsilon_{i,j} \\ \epsilon_{j,i} \end{pmatrix}, i = 1, \dots, n \sim \text{i.i.d normal}(0, \sigma^2 \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}) \quad (3.13)$$

We will treat  $\beta$  ,  $\mathbf{a}$   $\mathbf{b}$  and  $\Sigma_\epsilon$  as parameters.

To perform bayesian inference, we need to compute an integral over a high dimensional, and intractable set of  $(\lambda)$ .

$$L(\theta|(\rho)) \propto p(\theta)p((\rho)|\theta) = p(\theta) \int_{(\lambda)} p((\lambda)|\theta)p((\rho)|(\lambda)) \quad (3.14)$$

Given the observed ranks  $(\rho)$  and a prior distribution  $p(\theta)$  over the parameter space, the joint posterior distribution with density  $p(\theta, (\lambda)|(\rho))$  can be approximated by generating a Markov chain whose stationary distribution is that of  $(\theta, (\lambda))|(\rho)$ . The values of  $\theta$  simulated from this chain provide an approximation to the (marginal) posterior distribution of  $\theta$  given by the information from  $(\rho)$ .

One such MCMC algorithm is the Metropolis Hasting Algorithm:

Given current values of  $(\theta, Y)$ , one step of the Metropolis - Hasting for the Plackett-Luce likelihood proceeds by updating the values as follows:

1. Define the target distribution  $\pi(\cdot)$ . As we sample from the joint distribution, we can take the following target distribution.

$$p(\theta, (\lambda)|(\rho)) \propto p(\theta) \times p((\lambda)|\theta) \times p((\rho)|(\lambda)) = \pi(\theta, (\lambda)) \quad (3.15)$$

with:

(a)  $p(\theta)$  prior over the parameters.

(b)

$$\begin{pmatrix} \log(\lambda_{i,j}) \\ \log(\lambda_{j,i}) \end{pmatrix} | \theta, i < j \sim i.i.d \mathcal{N}\left(\begin{pmatrix} \beta^T x_{i,j} + a_i + b_j \\ \beta^T x_{j,i} + a_j + b_i \end{pmatrix}, \sigma^2 \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}\right) \quad (3.16)$$

(c)

$$p(\rho_i | (\lambda)) = \prod_{k=1}^{p_i} \frac{\lambda_{i\rho_k}}{\sum_{j \neq i} \lambda_{i,j} - \sum_{l=1}^{k-1} \lambda_{i,\rho_l}} \quad (3.17)$$

2. Choose a proposal distribution  $q(\theta^*, (\lambda)^* | \theta, (\lambda))$ . It needs to be easy to simulate and determine an irreducible chain. We can take a log-normal distribution for each element of  $\log(\lambda)^*$  with median  $\log(\lambda_{i,j})$ . A normal distribution for  $\beta, a, b$ , and finally  $\mathcal{U}[\rho - a, \rho + a]$  and  $\mathcal{U}[\sigma - a, \sigma + a]$  for  $\rho^*$  and  $\sigma^*$  with  $a$  likely small.
3. Write down the algorithm. Given  $(\theta, (\lambda))$  :
  - (a) Sample  $(\theta^*, (\lambda)^*)$  from the  $q(\cdot | \theta, (\lambda))$  previous distribution.
  - (b) Sample  $u \sim \mathcal{U}[0, 1]$
  - (c) If  $u < \alpha = \min\{1, \frac{\pi(\theta^*, (\lambda)^*) q(\theta, (\lambda) | \theta^*, (\lambda)^*)}{\pi(\theta, (\lambda)) q(\theta^*, (\lambda)^* | \theta, (\lambda))}\}$  then set  $(\theta, (\lambda)) = (\theta^*, (\lambda)^*)$  else keep the same value.

## Chapter 4

# EM algorithm with community representation

We have seen that the derivation of the posterior distribution over the  $(w)$  derived in equation (2.5) is somewhat very complicated. Hence we try in this chapter to introduce latent variables in order to derive an EM algorithm.

### 4.1 Mathematical formulation

We recall the likelihood of the Plackett-Luce Model:

$$p(\rho_i | (\lambda)) = \prod_{k=1}^K \frac{\lambda_{i\rho_k}}{\sum_{j \neq i} \lambda_{i,j} - \sum_{l=1}^{k-1} \lambda_{i,\rho_l}} \quad (4.1)$$

Along with its community representation:

$$\lambda_{i,j} = \sum_{k=1}^p w_{i,k} w_{j,k} \quad (4.2)$$

With elements of  $(\lambda)$  and  $(w)$  positive.

If we add a latent variable  $Z$  with its probability distribution  $P(Z|D, \lambda)$  we obtain the following complete data-likelihood:

$$L(\lambda, Z) = P(D, Z | \lambda) = \left( \prod_{i=1}^n \prod_{k=1}^K \frac{\lambda_{i\rho_k}}{\sum_{j \neq i} \lambda_{i,j} - \sum_{l=1}^{k-1} \lambda_{i,\rho_l}} \right) P(Z|D, \lambda) \quad (4.3)$$

Hence our approach is to find latent variables that would both give sense and a nice posterior distribution over the  $(w)$  to derive an EM algorithm.

We first introduce:

$$V_{i,j,k} \sim \mathcal{E}(w_{i,k} w_{j,k}), \quad 1 \leq i \neq j \leq n, \quad 1 \leq k \leq p \quad (4.4)$$

$$V_{i,j} = \min(V_{i,j,1}, \dots, V_{i,j,p}) \sim \mathcal{E}(\lambda_{i,j}) \quad (4.5)$$



And

$$Z_{i,k}|\rho_i = \min(V_{i,j})_{j \neq i, \rho_i[1], \dots, \rho_i[k-1]} \sim \mathcal{E}\left(\sum_{j \neq i, \rho_i[1], \dots, \rho_i[k-1]} \lambda_{i,j}, 1 \leq k \leq K\right) \quad (4.6)$$

An interpretation of those latent variables is that the next person ranked by an individual is equivalent to the first arrived individual (among those remaining to be ranked) following a waiting time of parameter  $\lambda_{i,j}$ . The arriving time of this same individual is equivalent to the first arriving time of other individuals following a waiting time of parameter  $w_{i,k}w_{j,k}$  which is the contribution of the  $k^{th}$  community.

This latent variable allows us to suppress the denominator in (4.1). To suppress the nominator we introduce the following latent variables:

$$Y_{i,j} \sim \mathcal{B}\left(\frac{w_{i,1}w_{\rho_i[j],1}}{\lambda_{i,j}}, \dots, \frac{w_{i,p}w_{\rho_i[j],p}}{\lambda_{i,j}}\right), 1 \leq j \leq K \quad (4.7)$$

This variable can be interpreted as the explanation of the choice of a certain individual in terms of community.

We obtain the following log-likelihood distribution.

$$\mathcal{L}((w); Y, Z, D) = \sum_{i=1}^n \sum_{k=1}^K \left\{ -\left(\sum_{l \neq i} \lambda_{i,l} - \sum_{l=1}^{k-1} \lambda_{i,\rho_l}\right) z_{i,k} + \sum_{p=1}^P \delta_{y_{i,k},p} (\ln(w_{i,p} + \ln(w_{\rho_k,p}))) \right\} \quad (4.8)$$

## 4.2 EM formulation

E-step:

$$\mathbb{E}_{Z,Y|\rho,(w)^*}(\mathcal{L}((w); Y, Z, D)) = \sum_{i=1}^n \sum_{k=1}^K \left\{ -\frac{(\sum_{l \neq i} \lambda_{i,l} - \sum_{l=1}^{k-1} \lambda_{i,\rho_l})}{(\sum_{l \neq i} \lambda_{i,l}^* - \sum_{l=1}^{k-1} \lambda_{i,\rho_l}^*)} + \sum_{p=1}^P \frac{w_{i,p}^* w_{\rho_k,p}^*}{\lambda_{i,\rho_k}^*} (\ln(w_{i,p} + \ln(w_{\rho_k,p}))) \right\} \quad (4.9)$$

Hence, with a gamma prior on the  $(w)$  we obtain:

$$Q(w|w^*) = \sum_{i=1}^n \left[ \sum_{k=1}^K \left\{ -\frac{(\sum_{l \neq i} \lambda_{i,l} - \sum_{l=1}^{k-1} \lambda_{i,\rho_l})}{(\sum_{l \neq i} \lambda_{i,l}^* - \sum_{l=1}^{k-1} \lambda_{i,\rho_l}^*)} + \sum_{p=1}^P \frac{w_{i,p}^* w_{\rho_k,p}^*}{\lambda_{i,\rho_k}^*} (\ln(w_{i,p} + \ln(w_{\rho_k,p}))) \right\} + \sum_{p=1}^P (a-1) \ln(w_{i,p}) - b w_{i,p} \right] \quad (4.10)$$

M-step: With the same result as in chapter 2 for the full derivation:

$$\frac{\partial \lambda_{i,j}}{\partial w_{r,s}} = 1_{i=r} w_{j,s} + 1_{j=r} w_{i,s} \quad (4.11)$$

We obtain this time, after computation:

$$w_{r,s} \frac{\partial Q}{\partial w_{r,s}} = C(r, s, w^*) - A(r, s, w, w^*) w_{r,s} \quad (4.12)$$

With

$$C(r, s, w^*) = a - 1 + \sum_{k=1}^K \frac{w_{r,s}^* w_{\rho_r[k],s}^*}{\lambda_{r,\rho_r[k]}^*} + \sum_{i \neq r} \frac{w_{i,s}^* w_{r,s}^*}{\lambda_{i,r}^*} \mathbb{I}(r \in \rho_i) \quad (4.13)$$

and

$$\begin{aligned} A(r, s, w, w^*) = b + & \left( \sum_{k=1}^K \frac{1}{\Lambda^*(r, k)} \right) \left( \sum_{l \neq r} w_{l,s} \right) + \sum_{i \neq r} w_{i,s} \left( \sum_{k=1}^K \frac{1}{\Lambda^*(i, k)} \right) - \\ & \sum_{l=1}^{K-1} w_{\rho_r(l),s} \left( \sum_{k=l+1}^K \frac{1}{\Lambda^*(r, k)} \right) - \sum_{i \neq r} w_{i,s} \left( \sum_{k > \text{rank}_{\rho_i}(r)}^K \frac{1}{\Lambda^*(i, k)} \right) \mathbb{I}(r \in \rho_i) \end{aligned} \quad (4.14)$$

where

$$\Lambda^*(i, k) = \sum_{l \neq i} \lambda_{i,l}^* - \sum_{l=1}^{k-1} \lambda_{i,\rho_i[l]}^* \quad (4.15)$$

We can see that it is difficult to directly derive the maximum arguments for  $Q(w|w^*)$ , as putting the gradient to zero require to solve a set of non linear equations. However, we can notice that  $A(r, s, w, w^*)$  does not depend on  $w_{r,s}$ . Moreover, we can easily check by remounting the equations that  $A$  and  $C$  are always positive ( $C$  might be negative if  $a < 1$ ). Hence, as a function of one parameter,  $Q$  has one maximum reached by setting the derivative to zero. Hence instead of updating the whole set of variables  $w$  directly by maximizing  $Q$ , we can update the elements one by one and update the set of fixed elements  $w^*$ . As  $A$  and  $C$  are always positive, doing so will always increase the posterior likelihood. We can update  $w_{r,s}$  with

$$w_{r,s}^{(t)} = \frac{C(r, s, w^{(t-1)})}{A(r, s, w^{(t-1)})} \quad (4.16)$$

This process will converge to (at least) a local minima.

### 4.3 EM formulation with variable ranks

In this section we suppose that  $K$  depends on the individual  $i$  i.e we allows the individuals to rank as many people as they want. The resulting model is the same, with only replacing  $K$  with  $K_i$ . We obtain the following log-likelihood distribution.

$$\mathcal{L}((w); Y, Z, D) = \sum_{i=1}^n \sum_{k=1}^{K_i} \left\{ - \left( \sum_{l \neq i} \lambda_{i,l} - \sum_{l=1}^{k-1} \lambda_{i,\rho_l} \right) z_{i,k} + \sum_{p=1}^P \delta_{y_{i,k},p} (\ln(w_{i,p}) + \ln(w_{\rho_k,p})) \right\} \quad (4.17)$$

We proceed the same way as the previous section, and we obtain the same formulations by allowing different lengths for the rows of  $\rho$ . In particular, we obtain the following EM algorithm for inference on the  $w$ :

$$w_{r,s}^{(t)} = \frac{C(r, s, w^{(t-1)})}{A(r, s, w^{(t-1)})} \quad (4.18)$$

with

$$C(r, s, w^*) = a - 1 + \sum_{k=1}^{K_r} \frac{w_{r,s}^* w_{\rho_r[k],s}^*}{\lambda_{r,\rho_r[k]}^*} + \sum_{i \neq r} \frac{w_{i,s}^* w_{r,s}^*}{\lambda_{i,r}^*} \mathbb{I}(r \in \rho_i) \quad (4.19)$$

and

$$A(r, s, w, w^*) = b + \left( \sum_{k=1}^{K_r} \frac{1}{\Lambda^*(r, k)} \right) \left( \sum_{l \neq r} w_{l,s} \right) + \sum_{i \neq r} w_{i,s} \left( \sum_{k=1}^{K_i} \frac{1}{\Lambda^*(i, k)} \right) - \sum_{l=1}^{K_r-1} w_{\rho_r(l),s} \left( \sum_{k=l+1}^{K_r} \frac{1}{\Lambda^*(r, k)} \right) - \sum_{i \neq r} w_{i,s} \left( \sum_{k > \text{rank}_{\rho_i}(r)}^{K_i} \frac{1}{\Lambda^*(i, k)} \right) \mathbb{I}(r \in \rho_i) \quad (4.20)$$

where

$$\Lambda^*(i, k) = \sum_{l \neq i} \lambda_{i,l}^* - \sum_{l=1}^{k-1} \lambda_{i,\rho_i[l]}^* \quad (4.21)$$

## 4.4 Simulation

### 4.4.1 Toy Example

We drew a first simple example with a clear clustering into two groups of individuals as in figure (4.1). Running the EM algorithm implemented in R is instantaneous and return the following matrix for  $W$ :

$$W = \begin{pmatrix} 0.1294021 & 1.0349799 \\ 0.1228317 & 0.8720837 \\ 0.1187870 & 0.7219490 \\ 1.0349782 & 0.1294015 \\ 0.8720818 & 0.1228311 \\ 0.7219470 & 0.1187863 \end{pmatrix} \quad (4.22)$$

which show a clear differentiation between two groups as illustrated in figure (4.2). Furthermore, we can see that the greatest values for  $w_{i,j}$  in the respective communities are assigned to the most appreciated individuals in each group (denoted by bold arrows), which is consistent with what we would expect.

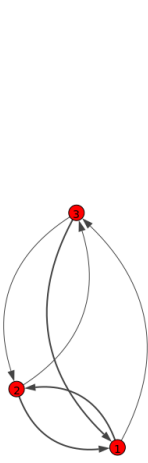


Figure 4.1: Initial Network with two independent group of individuals

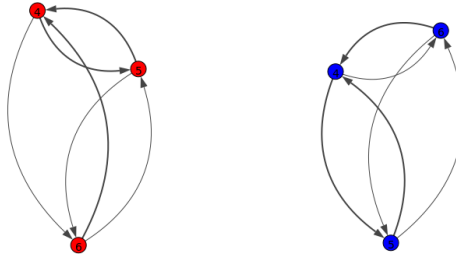


Figure 4.2: Inferred communities

### 4.4.2 Performance Evaluation

In this subsection we tested our EM algorithm on simulated random networks of size  $n = 20, K = 5, p = 4, a = 2, b = 2$  and plotted the log posterior for each update in the matrix ( $w$ ), one epoch being one full update of the matrix. We also performed the EM algorithm on the same network with different initial values of ( $w$ ).

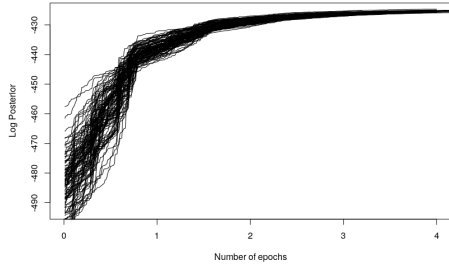


Figure 4.3: Performance on the first 4 epochs

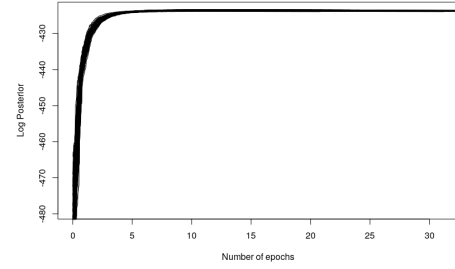


Figure 4.4: Performance on 30 epochs

We can see that the algorithm performs quite well, with three updates for each parameters  $w_{i,j}$  we have a very good estimation of the maximum a posteriori.

# Chapter 5

## Test on simulated Data

### 5.1 Fixed Community representation

We draw a simulation of a network with  $n = 25$ ,  $p = 3$ ,  $k = 4$ ,  $a = 2$ ,  $b = 2$ , optimized it. We obtain the graph with ranking gap of 5.83 for a mean of 8.32 and AE of 0.013

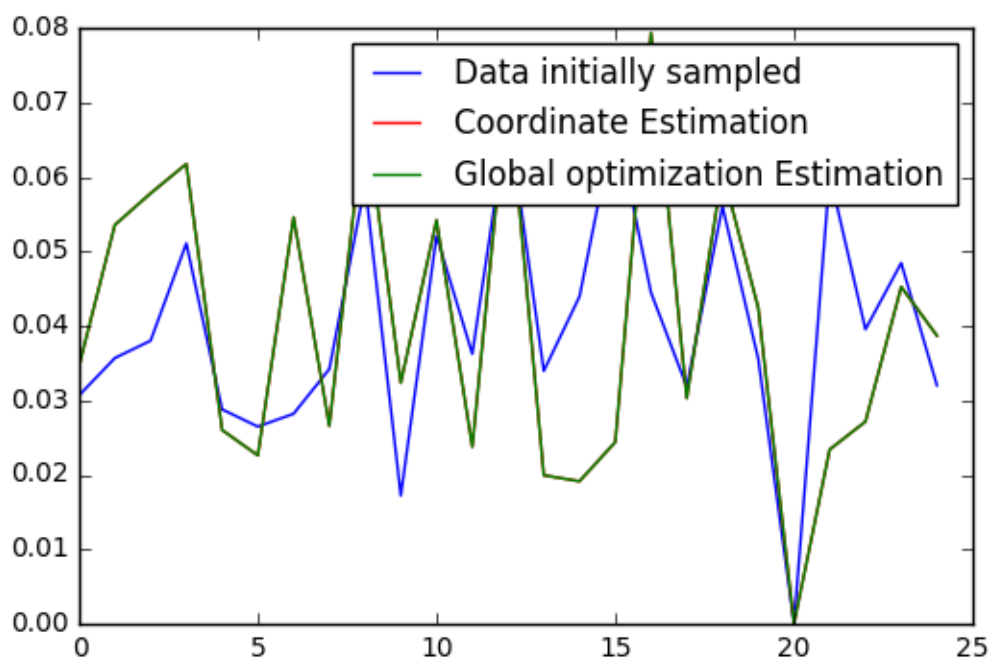
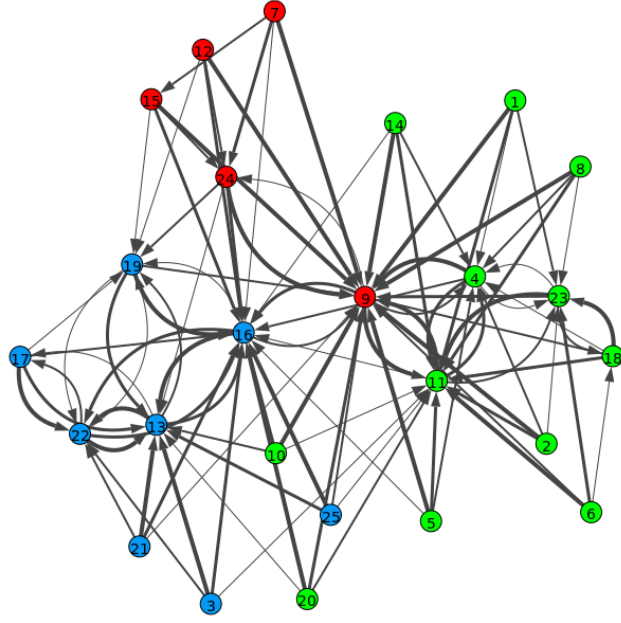
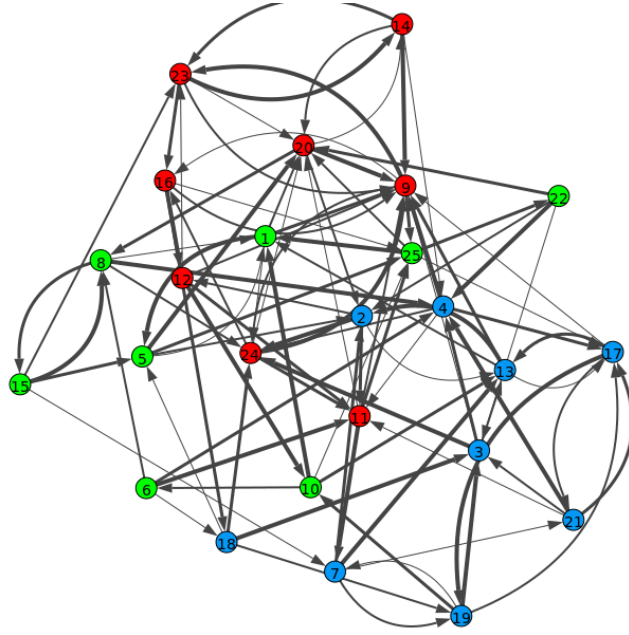


Figure 5.1: Lambda Estimation Fixed Prior



(a)



(b)

Figure 5.2: (a) Original Graph. (b) Inferred Graph

Now we made a simulation of 10 networks and measured the mean Absolute Error and Ranking Gap. We have a mean ranking gap of 2.887111111111114 compared to a mean of 5. A a mean AE of 0.022484955968703896

We run for 100 networks with the same parameters. We have mean ranking gap of 2.9304888888888896 and mean AE of 0.02122157183931606

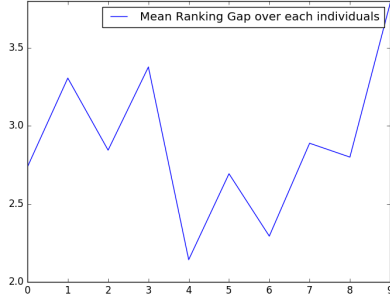


Figure 5.3: Ranking Gap Community Representation for 10 networks

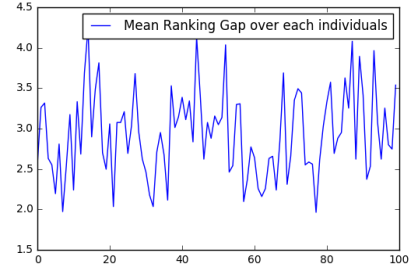


Figure 5.4: Ranking Gap Community Representation for 100 networks

## 5.2 Model With Degree Correction

We draw a simulation of 100 networks with  $n = 100$ ,  $p = 3$ ,  $k = 4$ ,  $a = 2$ ,  $b = 2$ , optimized it. We obtain the graph with ranking gap of 2.9426666666666666 for a mean of 5 and AE of 0.03798207013695899

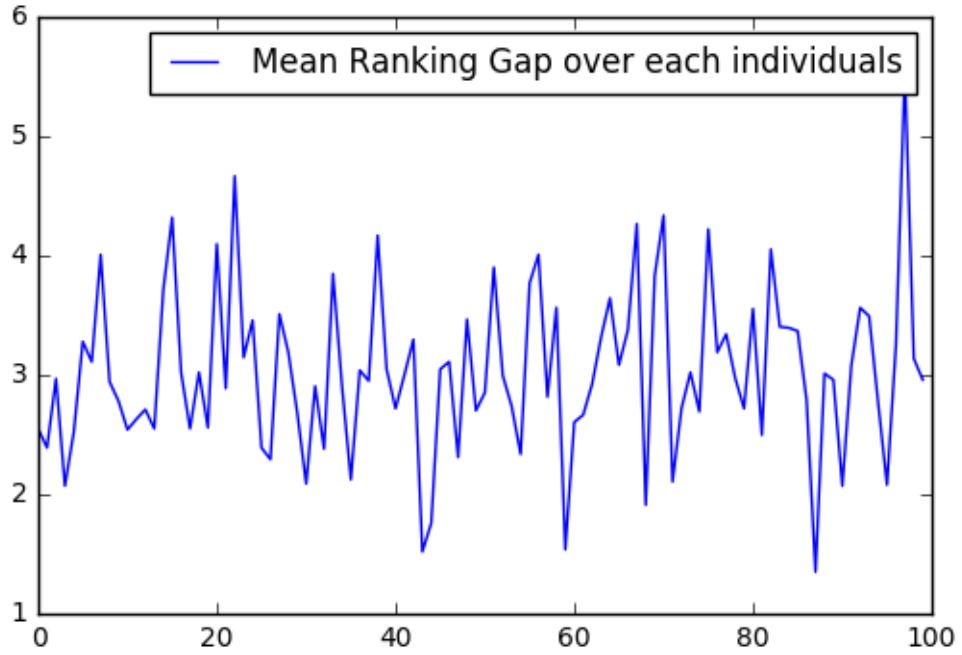


Figure 5.5: Ranking Gap Degree Correction



### 5.3 Model With Covariates

We simulated 100 networks of 100 individuals with the rank model:

$$\begin{aligned} y_{i,j} &= \beta_0 + a_i + b_j + \epsilon_{i,j} \\ \begin{pmatrix} a_i \\ b_i \end{pmatrix}, i = 1, \dots, n &\sim \text{i.i.d normal}(0, \Sigma_{ab}) \\ \begin{pmatrix} \epsilon_{i,j} \\ \epsilon_{j,i} \end{pmatrix}, i = 1, \dots, n &\sim \text{i.i.d normal}(0, \sigma^2 \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}) \end{aligned}$$

## Chapter 6

# Simulations on the Sampsons monk dataset

In this part we use the The data set collected by Sampson (1968) this is a classic data set in social network analysis. The data set summarizes relationships, observed at three distinct time points, among 18 monks who were about to enter a monastery when a conflict erupted. We use here the directed network where  $y_{i,j} = 1$  denotes that monk  $i$  ranked monk  $j$  among his top three preferred individuals and  $y_{i,j} = 0$  otherwise. The directed network is shown in Figure 6.1 , where circles represent monks and directed edges are oriented from  $i$  to  $j$  whenever  $y_{i,j} = 1$ . The monks are divided by Sampson into three groups: Loyal Opposition, Turks, and Outcasts.

We applied the same method as for the simulated data and obtained the following community. We can observe that our model identifies exactly the same clustering of communities as the ones observed by Sampson. Namely the Loyal Opposition (Red community), Turks (Green Community), and Outcasts (Blue Community). In this graph the individuals are labeled according to their number provided with their names in the data set package and not according to their rank in the data-frame.

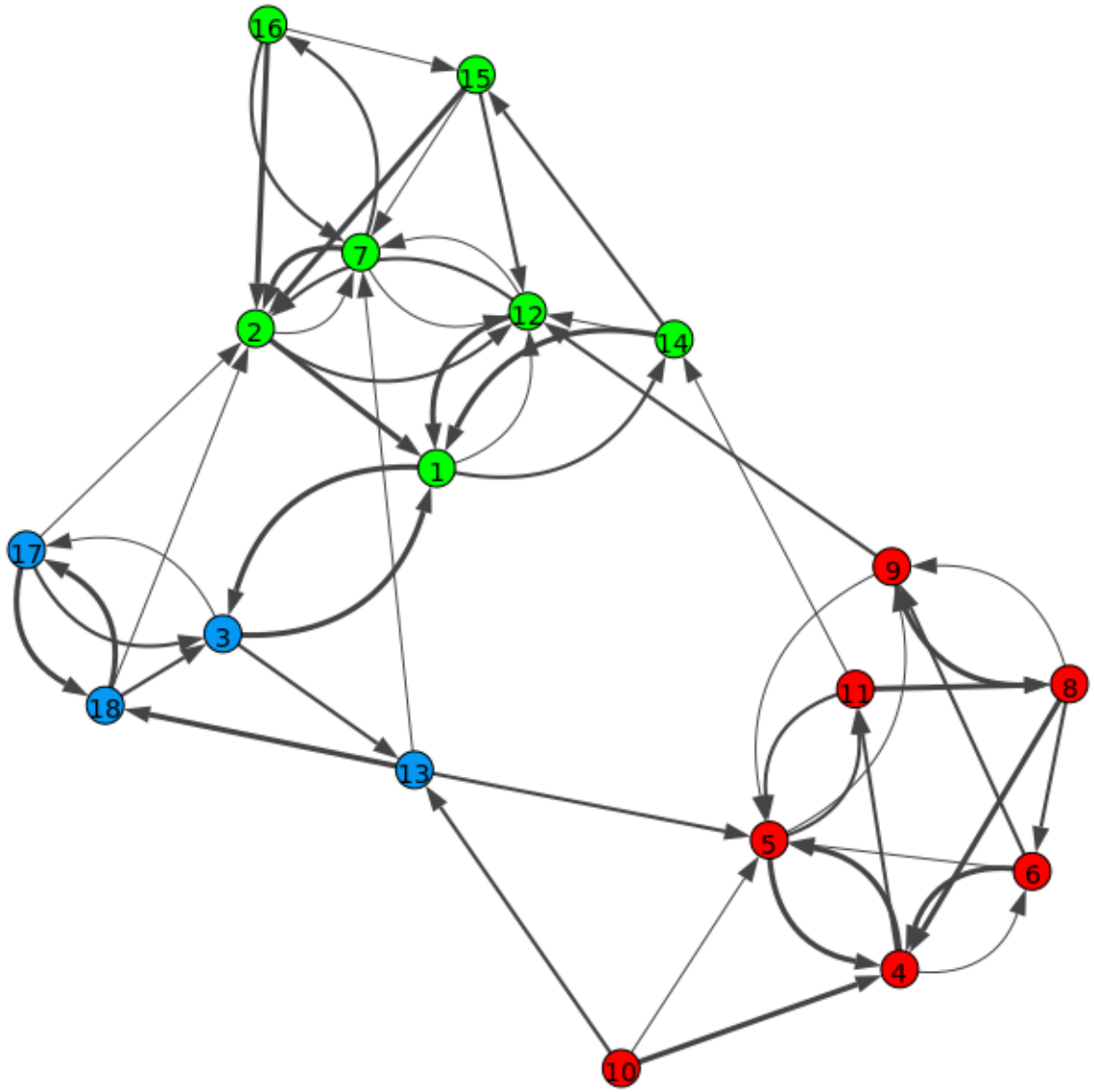


Figure 6.1: Inferred community for the Sampson DataSet