

# ESTIMATING PARAMETERS OF A DIRECTED WEIGHTED GRAPH MODEL WITH BETA-DISTRIBUTED EDGE-WEIGHTS

M. Bolla<sup>1</sup>, J. Mala<sup>1,2</sup>, and A. Elbanna<sup>1,3</sup>

We introduce a directed, weighted random graph model, where the edge-weights are independent and beta distributed with parameters depending on their endpoints. We will show that the row- and column-sums of the transformed edge-weight matrix are sufficient statistics for the parameters, and use the theory of exponential families to prove that the ML estimate of the parameters exists and is unique. Then an algorithm to find this estimate is introduced together with convergence proof that uses properties of the digamma function. Simulation results and applications are also presented.

## 1. Introduction

The theory of ML estimation in the following types of exponential family random graph models has frequently been investigated in the last decade, see, e.g., [5, 6, 8, 10, 11]. The graph has  $n$  vertices, and the adjacency relations between them are given by the  $n \times n$  random edge-weight matrix  $\mathbf{W} = (w_{ij})$  of zero diagonal. If  $\mathbf{W}$  is symmetric, then we have an undirected graph; otherwise, our graph is directed, where  $w_{ij}$  is the nonnegative weight assigned to the  $i \rightarrow j$  edge according to the model. We assume that the edge-weights (above or out of the main diagonal) are completely independent (but their distribution usually depends on different parameters), and have an exponential family distribution  $P_{\theta}$ . So the likelihood function has the general form

$$L_{\theta}(\mathbf{W}) = e^{\langle \theta, \mathbf{t}(\mathbf{W}) \rangle - Z(\theta)} \cdot h(\mathbf{W}), \quad (1)$$

with the *canonical parameter*  $\theta$ , *log-partition (cumulant) function*  $Z(\theta)$ , and *canonical sufficient statistic*  $\mathbf{t}$ . In these random graph models, the components of  $\mathbf{t} = \mathbf{t}(\mathbf{W})$  are the row-sums and/or column-sums of  $\mathbf{W}$  or some  $\mathbf{W}$ -related matrix, i.e., they are vertex-degrees or in- and out-degrees of the observed undirected or directed, weighted or unweighted graph (in the weighted case, the edge-weights may undergo a suitable transformation). Also,  $h(\mathbf{W})$  is usually 1 over the support of the likelihood function, indicating that given the canonical sufficient statistics, the joint distribution of the entries is uniform (microcanonical) in these models.

To make inference on the parameters, typically we have only one observation for the graph. It may seem that it is a one-element sample, but there are the adjacencies that form the sample; the number of them is  $\binom{n}{2}$  in the undirected, and  $n(n-1)$  in the directed case. The number of parameters contained in  $\theta$  is  $n$  in the undirected and  $2n$  in the directed case. The parameters can be considered as affinities or potentials of the vertices to make ties in the undirected, and to emanate or adsorb edges in the directed case. It is important that we divide the components of the canonical parameter  $\theta$  of the underlying distribution of the  $ij$  or  $i \rightarrow j$  edge between the connected vertices, like  $\alpha_i + \alpha_j$  in the undirected and  $\alpha_i + \beta_j$  in the directed case ( $i \neq j$ ), see [5, 6, 10].

<sup>1</sup>Institute of Mathematics, Budapest University of Technology and Economics, Budapest, Hungary, e-mail: marib@math.bme.hu, jmala@math.bme.hu, ahmed@math.bme.hu

<sup>2</sup>ELTE Eötvös Loránd University, Institute of Mathematics, Budapest, Hungary, e-mail: mala.jozsef@sek.elte.hu

<sup>3</sup>Tanta University, Faculty of Science, Mathematics Department, Tanta, Egypt

In *regular exponential families* ( $\Theta$  is open), the ML equation  $\nabla_{\theta} \ln L_{\theta}(\mathbf{W}) = \mathbf{0}$  is equivalent to

$$\nabla_{\theta} Z(\theta) = \mathbf{t}. \quad (2)$$

Since  $\nabla_{\theta} Z(\theta) = \mathbf{E}_{\theta} \mathbf{t}$ , the ML Eq. (2) means that the canonical sufficient statistic is made equal to its expectation. But when is it possible? Now we briefly summarize existing theoretical results on this issue.

Let  $\mathcal{M} = \{\mathbf{E}_{\theta} \mathbf{t} : \theta \in \Theta\}$  denote the so-called *mean parameter space* in the model; it is necessarily convex. Let  $\mathcal{M}^0$  denote its interior. When the canonical statistic is also complete, and hence, minimal sufficient, the representation (1) is *minimal* (i.e., the model is not overparametrized).

**Proposition 1** [11, Proposition 3.2]. *In exponential family, the gradient mapping  $\nabla Z : \Theta \rightarrow \mathcal{M}$  is one-to-one, if and only if the exponential family representation is minimal.*

**Proposition 2** [11, Theorem 3.3]. *In a minimal exponential family, the gradient mapping  $\nabla Z$  is onto  $\mathcal{M}^0$ .*

By Propositions 1 and 2, any parameter in  $\mathcal{M}^0$  is uniquely realized by the  $\mathbf{P}_{\theta}$  distribution for some  $\theta \in \Theta$ . Also, in a regular and minimal exponential family,  $\mathcal{M}$  is an open set and is identical to  $\mathcal{M}^0$ .

As the ML estimate of  $\theta$  is the solution of (2), we have the following.

**Proposition 3** [10, Proposition 5]. *Assume that the (canonical) parameter space  $\Theta$  is open. Then there exists a solution  $\hat{\theta} \in \Theta$  to the ML equation  $\nabla_{\theta} Z(\theta) = \mathbf{t}$  if and only if  $\mathbf{t} \in \mathcal{M}^0$ ; further, if such a solution exists, it is also unique.*

Note that in regular and minimal exponential families,  $\mathcal{M}^0$  is also the interior of  $\mathcal{T}$ , which is the convex hull of all possible values of  $\mathbf{t}$ , see, e.g., [6, 9]. In the case of discrete distributions, it frequently happens that the boundary of  $\mathcal{T}$  has positive measure. For instance, the so-called threshold graphs are located on the boundary of the polyhedron, determined by the Erdős–Gallai conditions, in the model of [6] which uses Bernoulli distributed entries. However, in the case of an absolutely continuous  $\mathbf{P}_{\theta}$  distribution, the boundary of  $\mathcal{T}$  has zero Lebesgue measure, and so, probability zero with respect to the  $\mathbf{P}_{\theta}$  measure. Therefore, in view of Proposition 3, the ML equation has a unique solution with probability 1.

The organization of the paper is as follows. In Section 2, we introduce a model for directed edge-weighted graphs and prove that a unique ML estimate of the parameters exists. In Section 3, we define an iterative algorithm to find this solution, and prove its convergence with a convenient starting. In Section 4, the algorithm is applied to randomly generated and real-word data. In Appendix A, properties of the digamma function, whereas in Appendix B, the boundary of our  $\mathcal{M}$ , is discussed. The long proof of the main convergence theorem of the iteration algorithm, introduced in Section 3, is presented in Appendix C.

We remark that edge-weighted graphs of uniformly bounded edge-weights are prototypes of real-world networks, see, e.g., [4]. Without loss of generality, if the edge-weights are transformed into the  $[0,1]$  interval, the *beta distribution* for them, with varying parameters, is capable to model a wide range of possible probability densities on them. This indicates the soundness of the model to be introduced in Section 2.

## 2. A random graph model with beta distributed edge-weights

Let  $\mathbf{W} = (w_{ij})$  be the  $n \times n$  (usually not symmetric) edge-weight matrix of a random directed graph on  $n$  vertices:  $w_{ii} = 0$  ( $i = 1, \dots, n$ ) and  $w_{ij} \in [0, 1]$  is the weight of the  $i \rightarrow j$  edge ( $i \neq j$ ). Our model is the following: the  $i \neq j$  weight obeys a beta distribution with parameters  $a_i > 0$  and  $b_j > 0$ . The parameters are collected in  $\mathbf{a} = (a_1, \dots, a_n)$  and  $\mathbf{b} = (b_1, \dots, b_n)$ , or briefly, in  $\theta = (\mathbf{a}, \mathbf{b})$ . Here  $a_i$  can be thought of as the potential of the vertex  $i$  to send messages out, and  $b_i$  is its resistance to receive messages in.

The likelihood function is factorized as

$$\begin{aligned}
 L_{\mathbf{a},\mathbf{b}}(\mathbf{W}) &= \prod_{i \neq j} \frac{\Gamma(a_i + b_j)}{\Gamma(a_i)\Gamma(b_j)} w_{ij}^{a_i-1} (1 - w_{ij})^{b_j-1} = \\
 &= C(\mathbf{a}, \mathbf{b}) \prod_{i \neq j} \exp[(a_i - 1) \ln w_{ij} + (b_j - 1) \ln(1 - w_{ij})] = \\
 &= \exp \left[ \sum_{i=1}^n (a_i - 1) \sum_{j \neq i} \ln w_{ij} + \sum_{j=1}^n (b_j - 1) \sum_{i \neq j} \ln(1 - w_{ij}) - Z(\mathbf{a}, \mathbf{b}) \right],
 \end{aligned}$$

where  $C(\mathbf{a}, \mathbf{b})$  is the normalizing constant, and  $Z(\mathbf{a}, \mathbf{b}) = -\ln C(\mathbf{a}, \mathbf{b})$  is the log-partition (cumulant) function. Since the likelihood function depends on  $\mathbf{W}$  only through the row-sums of the  $n \times n$  matrix  $\mathbf{U} = \mathbf{U}(\mathbf{W})$  of general entry  $\ln w_{ij}$  and the column-sums of the  $n \times n$  matrix  $\mathbf{V} = \mathbf{V}(\mathbf{W})$  of general entry  $\ln(1 - w_{ij})$ , by the Neyman–Fisher factorization theorem, the row-sums  $R_1, \dots, R_n$  of  $\mathbf{U}$  and column-sums  $C_1, \dots, C_n$  of  $\mathbf{V}$  are sufficient statistics for the parameters. Moreover,  $\mathbf{t} = (\mathbf{R}, \mathbf{C}) = (R_1, \dots, R_n, C_1, \dots, C_n)$  is the canonical sufficient statistic, which is also minimal. Note that  $\mathbf{U}$  contains the log-weights of the original graph, while  $\mathbf{V}$  contains the log-weights of the complement graph of edge-weight matrix  $\overline{\mathbf{W}}$  with entries  $1 - w_{ij}$  ( $i \neq j$ ). The first factor in the Neyman–Fisher factorization (including gamma functions) depends only on the parameters and on the sample through these sufficient statistics, whereas the seemingly not present other factor – which would merely depend on  $\mathbf{W}$  – is constantly 1, indicating that the conditional joint distribution of the entries, given the row- and column-sums of the log-weight and log-complement matrix is uniform (microcanonical) in this model. So under the conditions on the margins of  $\mathbf{U}$  and  $\mathbf{V}$ , the directed graphs coming from the above model are uniformly distributed.

The system of likelihood equations is obtained by making the derivatives of  $L_{\mathbf{a},\mathbf{b}}(\mathbf{W})$  with respect to the parameters equal to 0:

$$\begin{aligned}
 \frac{\partial \ln L_{\mathbf{a},\mathbf{b}}(\mathbf{W})}{\partial a_i} &= \sum_{j \neq i} \psi(a_i + b_j) - (n - 1)\psi(a_i) + R_i = 0, \quad i = 1, \dots, n, \\
 \frac{\partial \ln L_{\mathbf{a},\mathbf{b}}(\mathbf{W})}{\partial b_j} &= \sum_{i \neq j} \psi(a_i + b_j) - (n - 1)\psi(b_j) + C_j = 0, \quad j = 1, \dots, n.
 \end{aligned} \tag{3}$$

Here  $\psi(x) = \frac{\partial \ln \Gamma(x)}{\partial x} = \frac{\Gamma'(x)}{\Gamma(x)}$  for  $x > 0$  is the *digamma function*. For its properties, see Appendix A.

To apply the theory of Section 1, we utilize that the parameter space  $\Theta \subset \mathbb{R}_+^{2n}$  is open, akin to the canonical parameter space,  $(-1, \infty)^{2n}$ . Note that the canonical parameter is, in fact,  $(\mathbf{a}', \mathbf{b}') = \boldsymbol{\theta}' = \boldsymbol{\theta} - \mathbf{1}$ , where  $\mathbf{1} \in \mathbb{R}^{2n}$  is the vector of all 1 coordinates. With it, the log-partition function is

$$Z(\mathbf{a}', \mathbf{b}') = - \sum_{j \neq i} \Gamma(a_i' + b_j' + 2) + \sum_{j \neq i} \Gamma(a_i' + 1) + \sum_{i \neq j} \Gamma(b_j' + 1).$$

In view of (2), the ML equation is equivalent to

$$\begin{aligned}
 \frac{\partial Z(\mathbf{a}', \mathbf{b}')}{\partial a_i'} &= - \sum_{j \neq i} \psi(a_i' + b_j' + 2) + (n - 1)\psi(a_i' + 1) = R_i, \quad i = 1, \dots, n, \\
 \frac{\partial Z(\mathbf{a}', \mathbf{b}')}{\partial b_j'} &= - \sum_{i \neq j} \psi(a_i' + b_j' + 2) + (n - 1)\psi(b_j' + 1) = C_j, \quad i = 1, \dots, n.
 \end{aligned}$$

But this system of equations is the same as (3), in terms of the parameter  $\boldsymbol{\theta}'$  instead of  $\boldsymbol{\theta}$ .

In view of Section 1, the mean parameter space  $\mathcal{M}$  consists of parameters  $(A_1, \dots, A_n, B_1, \dots, B_n)$  obtained by the gradient mapping, that is,

$$\begin{aligned} A_i &= A_i(\mathbf{a}, \mathbf{b}) = - \sum_{j \neq i} [\psi(a_i + b_j) - \psi(a_i)], \quad i = 1, \dots, n, \\ B_j &= B_j(\mathbf{a}, \mathbf{b}) = - \sum_{i \neq j} [\psi(a_i + b_j) - \psi(b_j)], \quad j = 1, \dots, n. \end{aligned} \quad (4)$$

Space  $\mathcal{M}$  is an open set, whose boundary is determined by the limit properties between the digamma and the log functions; see Appendix B for details. There we also find a correspondence between the points on the boundary of  $\mathcal{M}$  and those on the boundary of the convex hull  $\mathcal{T}$  of the possible sufficient statistics  $\mathbf{t} = (\mathbf{R}, \mathbf{C})$  within  $\mathbb{R}^{2n}$ . It is interesting that while the boundary points of  $\mathcal{M}$  do not belong to the open set  $\mathcal{M}$ , the boundary points of  $\mathcal{T}$  do belong to  $\mathcal{T}$  and can be realized as row- and column-sums of the  $\mathbf{U}(\mathbf{W})$  and  $\mathbf{V}(\mathbf{W})$  matrices with a  $\mathbf{W}$  of off-diagonal entries in  $(0, 1)$ . However, this boundary has 0 probability, and so, any canonical sufficient statistic  $\mathbf{t}$  of the observed graph is in  $\mathcal{M}$ , with probability 1. Therefore, by Proposition 3, we can state the following

**Theorem 1.** *The system of the ML Eq. (3) has a unique solution  $\hat{\boldsymbol{\theta}} = (\hat{\mathbf{a}}, \hat{\mathbf{b}})$ , with probability 1.*

Later we will use the following trivial upper bound for the sum of row- and column-sums (of the  $\mathbf{U}$  and  $\mathbf{V}$  matrices):

$$\sum_{i=1}^n R_i + \sum_{j=1}^n C_j = \sum_{i=1}^n \sum_{j \neq i} \ln w_{ij} + \sum_{j=1}^n \sum_{i \neq j} \ln(1 - w_{ij}) = \sum_{i \neq j} \ln[w_{ij}(1 - w_{ij})] \leq -2 \ln 2 n(n-1) \quad (5)$$

due to rearranging the terms and the relation  $w_{ij}(1 - w_{ij}) \leq 1/4$  for  $w_{ij} \in [0, 1]$  with equality if and only if  $w_{ij} = \frac{1}{2}$  ( $i \neq j$ ). For finer estimates see Appendix B.

Also note that the Hessian of the system of ML equations (consisting of the second order partial derivatives of  $L_{\boldsymbol{\theta}}$  at  $\hat{\boldsymbol{\theta}}$ ) does not contain the sufficient statistics anymore; therefore the negative of it is the Fisher-information matrix at  $\hat{\boldsymbol{\theta}}$ . Because of the regularity conditions, the information matrix is positive, and so, the Hessian is negative definite. This is also an indication of the existence of a unique ML estimate.

### 3. Iteration algorithm to find the parameters

To use a fixed point iteration, now we rewrite the system of likelihood equations in the form  $\boldsymbol{\theta} = f(\boldsymbol{\theta})$ , where  $\boldsymbol{\theta} = (\mathbf{a}, \mathbf{b})$ , as follows:

$$\begin{aligned} a_i &= \psi^{-1} \left[ \frac{1}{n-1} R_i + \frac{1}{n-1} \sum_{j \neq i} \psi(a_i + b_j) \right] =: g_i(\mathbf{a}, \mathbf{b}), \quad i = 1, \dots, n, \\ b_j &= \psi^{-1} \left[ \frac{1}{n-1} C_j + \frac{1}{n-1} \sum_{i \neq j} \psi(a_i + b_j) \right] =: h_j(\mathbf{a}, \mathbf{b}), \quad j = 1, \dots, n. \end{aligned} \quad (6)$$

Here  $g_i$ 's and  $h_j$ 's are the coordinate functions of  $f = (g, h) : \mathbb{R}^{2n} \rightarrow \mathbb{R}^{2n}$ . Then, starting at  $\boldsymbol{\theta}^{(0)}$ , we use the successive approximation  $\boldsymbol{\theta}^{(it)} := f(\boldsymbol{\theta}^{(it-1)})$  for  $it = 1, 2, \dots$ , until convergence. Now the statement of convergence of the above iteration to the theoretically guaranteed unique  $\hat{\boldsymbol{\theta}}$  (see Theorem 1) follows.

**Theorem 2.** *Let  $\hat{\boldsymbol{\theta}} = (\hat{\mathbf{a}}, \hat{\mathbf{b}})$  be the unique solution of the ML Eq. (3). Then the above mapping  $f = (g, h)$  is a contraction in some closed neighborhood  $K$  of  $\hat{\boldsymbol{\theta}}$ , and so, starting at any  $\boldsymbol{\theta}^{(0)} \in K$ , the fixed point of the iteration  $\boldsymbol{\theta}^{(it)} = f(\boldsymbol{\theta}^{(it-1)})$  exists and is  $\hat{\boldsymbol{\theta}}$ .*

The proof of this theorem is to be found in Appendix C.

Since  $K$  is only theoretically guaranteed, we need some practical considerations about the choice of  $\boldsymbol{\theta}^{(0)}$ , which should be adapted to the sufficient statistics. In the sequel, for two vectors  $\mathbf{x} = (x_1, \dots, x_n)$ ,  $\mathbf{y} = (y_1, \dots, y_n)$  we use the notation  $\mathbf{x} > \mathbf{y}$  if  $x_i > y_i$  for each  $i = 1, \dots, n$ . Likewise,  $\mathbf{x} \geq \mathbf{y}$  is the shorthand for  $x_i \geq y_i$  for each  $i = 1, \dots, n$ .

Recall that  $f = (g, h)$  is the mapping (6) of the fixed point iteration, and  $\hat{\boldsymbol{\theta}} = (\hat{\mathbf{a}}, \hat{\mathbf{b}}) > \mathbf{0}$  is the (only) solution of the equation  $f(\boldsymbol{\theta}) = \boldsymbol{\theta}$ , where  $\mathbf{0} \in \mathbb{R}^{2n}$  is the vector of all 0 coordinates.

**Proposition 4.** *Let*

$$M := \max \left\{ \max_{i \in \{1, \dots, n\}} \left( -\frac{R_i}{n-1} \right), \max_{i \in \{1, \dots, n\}} \left( -\frac{C_i}{n-1} \right) \right\} \quad (7)$$

and  $\varepsilon > 0$  be the (only) solution of the equation  $\psi(2x) - \psi(x) = M$ . Then  $(\hat{\mathbf{a}}, \hat{\mathbf{b}}) \geq \varepsilon \mathbf{1}$ .

**Proof.** In view of (5) we have that  $M \geq \ln 2$ . Since equality in (5) is attained with probability 0, we have that  $M > \ln 2$  with probability 1. Therefore, by Lemma 3 of Appendix A, there exists an  $\varepsilon$ , with probability 1, such that  $\psi(2\varepsilon) - \psi(\varepsilon) = M$ .

Without loss of generality we can assume that

$$a_{i_0} = \min \left\{ \min_{i \in \{1, \dots, n\}} \hat{a}_i, \min_{i \in \{1, \dots, n\}} \hat{b}_i \right\}.$$

Then by the ML equation, the monotonicity of  $\psi$ , and Lemma 3 of Appendix A, we get

$$(n-1)M \geq -R_{i_0} = \sum_{j \neq i_0} \psi(\hat{a}_{i_0} + \hat{b}_j) - (n-1)\psi(\hat{a}_{i_0}) \geq (n-1)[\psi(2\hat{a}_{i_0}) - \psi(\hat{a}_{i_0})].$$

Therefore,  $a_{i_0} \geq \varepsilon$ , whence  $\hat{a}_i, \hat{b}_i \geq \varepsilon$  holds for every  $i = 1, \dots, n$ .

**Proposition 5.** *With the solution  $\varepsilon$  of  $\psi(2x) - \psi(x) = M$  of (7), we have  $f(\varepsilon \mathbf{1}) \geq \varepsilon \mathbf{1}$ .*

**Proof.** We have

$$g_i(\varepsilon \mathbf{1}) = \psi^{-1} \left( \psi(2\varepsilon) + \frac{R_i}{n-1} \right) \geq \psi^{-1}(\psi(2\varepsilon) - M) = \varepsilon.$$

Likewise,

$$h_i(\varepsilon \mathbf{1}) = \psi^{-1} \left( \psi(2\varepsilon) + \frac{C_i}{n-1} \right) \geq \psi^{-1}(\psi(2\varepsilon) - M) = \varepsilon.$$

It is also clear that we have the following.

**Proposition 6.** *If  $(\mathbf{a}, \mathbf{b}) \geq (\mathbf{x}, \mathbf{y}) > \mathbf{0}$ , then  $f(\mathbf{a}, \mathbf{b}) \geq f(\mathbf{x}, \mathbf{y})$ .*

**Theorem 3.** *With  $\varepsilon$  satisfying  $\psi(2\varepsilon) - \psi(\varepsilon) = M$  of (7), and starting at  $\boldsymbol{\theta}^{(0)} = \varepsilon \mathbf{1}$ , the sequence  $\boldsymbol{\theta}^{(it)}$  of the iteration  $\boldsymbol{\theta}^{(it)} = f(\boldsymbol{\theta}^{(it-1)})$  for  $it \rightarrow \infty$  converges at a geometric rate to the unique solution  $(\hat{\mathbf{a}}, \hat{\mathbf{b}})$  of the ML equation.*

**Proof.** From Propositions 4 and 5 we obtain that the sequence  $\boldsymbol{\theta}^{(it)}$  is coordinate-wise increasing. Moreover, it is clear that  $(\boldsymbol{\theta}^{(it)})$  is bounded from above by  $(\hat{\mathbf{a}}, \hat{\mathbf{b}})$ , due to Proposition 6. Therefore, the convergence of  $\boldsymbol{\theta}^{(it)}$  follows, and by the continuity of  $f$ , the limit is clearly a fixed point of  $f$ . However, in view of Section 2, the solution of the ML equation is a fixed point of  $f$ , and it cannot be anything but the unique solution  $(\hat{\mathbf{a}}, \hat{\mathbf{b}})$ , guaranteed by Theorem 1. Further, from Theorem 2 we get that the rate of convergence is (at least) geometric.

Therefore, a good starting can be chosen by these considerations. Also note that at the above  $\boldsymbol{\theta}^{(0)}$  and possibly at its first (finitely many) iterates,  $f$  is usually not a contraction. It becomes a contraction

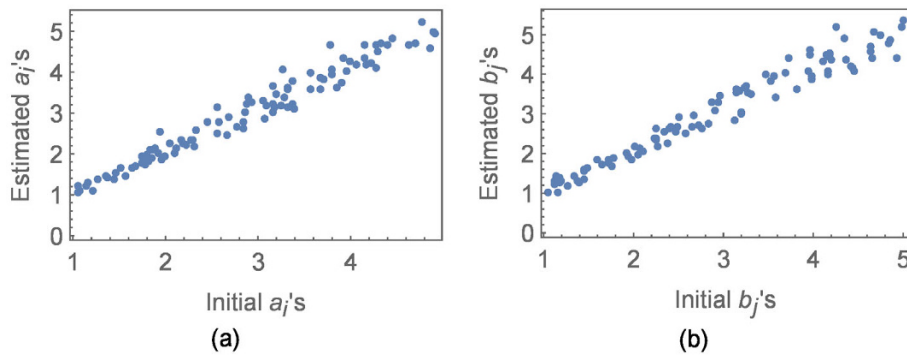
only when some iterate  $\theta^{(it_0)}$  gets into the neighborhood  $K$  of  $\hat{\theta}$  of Theorem 2, which is inevitable in view of the convergence of the sequence  $\theta^{(it)}$ . So, Theorem 2 is literally applicable only if we start the iteration at  $\theta^{(it_0)}$ . In practice, however, we do not know the theoretically guaranteed neighborhood  $K$ . The practical merit of Theorem 3 is just that it offers a realizable starting.

#### 4. Applications

First we generated a random directed edge-weighted graph on  $n = 100$  vertices. The edge-weight matrix  $\mathbf{W}$  had zero diagonal, and the off-diagonal entries  $w_{ij}$ 's were independent. Further, for  $i \neq j$ , the weight  $w_{ij}$  was generated according to beta distribution with parameters  $a_i > 0$  and  $b_j > 0$ , where  $a_i$ 's and  $b_j$ 's were chosen randomly in the interval  $[1,5]$ .

Then we estimated the parameters based on  $\mathbf{W}$ , and plotted the  $a_i, \hat{a}_i$  ( $i = 1, \dots, n$ ) and  $b_j, \hat{b}_j$  ( $j = 1, \dots, n$ ) pairs.

Figure 1 shows a good fit between them.



**Fig. 1.** Panel (a) shows the original versus the estimated parameters  $a_i$ 's with  $MSE = 0.0628806$ , while Panel (b) shows the original versus the estimated parameters  $b_j$ 's with  $MSE = 0.0768382$ .

We also applied the algorithm to migration data among 34 countries. Here  $w_{ij}$  is proportional to the number of people in thousands who moved from country  $i$  to country  $j$  (to find jobs) during the year 2011, and it is normalized so that be in the interval  $(0,1)$ . The estimated parameters are in Table 1.

In this context,  $a_i$ 's are related to the emigration and  $b_i$ 's to the counter-immigration potentials. When  $a_i$  is large, country  $i$  has a relatively large potential for emigration. On the contrary, when  $b_i$  is large, country  $i$  tends to have a relatively large resistance against immigration.

It should be noted again that edge-weighted graphs of this type very frequently model real-world directed networks.

## Appendix

### A. Properties of the digamma function

Though we do not use it explicitly, the following approximation of the digamma function  $\psi(x) = \frac{\partial \ln \Gamma(x)}{\partial x} = \frac{\Gamma'(x)}{\Gamma(x)}$  ( $x > 0$ ) is interesting in its own right.

**Lemma 1.** *The equation  $\psi(x) = \ln(x - \frac{1}{2}) + \mathcal{O}(\frac{1}{x^2})$  holds for  $x > 1$ .*

The statement of the lemma easily follows by Taylor expansion.

**Lemma 2.** *The equation  $\frac{1}{\psi'(x+y)} > \frac{1}{\psi'(x)} + \frac{1}{\psi'(y)}$  holds for  $x, y > 0$ .*

**Table 1.** Estimated parameters for migration data, 2011

i	Country	$a_i$	$b_i$	i	Country	$a_i$	$b_i$
1	Australia	0.26931	1475.75242	18	Japan	0.23211	9926.91644
2	Austria	0.27403	632.81653	19	Korea	0.22310	4199.25005
3	Belgium	0.33380	46.01197	20	Luxembourg	0.17543	107.91399
4	Canada	0.27383	2363.23435	21	Mexico	0.26706	4655.95370
5	Chile	0.21236	28940.59777	22	Netherlands	0.37754	39.52320
6	Czech Rep.	0.31188	470.28651	23	New Zealand	0.20542	2568.00582
7	Denmark	0.26514	847.34887	24	Norway	0.22646	519.12451
8	Estonia	0.23235	25602.33371	25	Poland	0.62846	1106.55946
9	Finland	0.29357	1100.00568	26	Portugal	0.31011	1606.59979
10	France	0.52721	37.92122	27	Slovak Rep.	0.27871	42451.19093
11	Germany	0.62020	1.64064	28	Slovenia	0.19720	6824.54028
12	Greece	0.29708	6319.19184	29	Spain	0.39732	182.47160
13	Hungary	0.31443	32750.88310	30	Sweden	0.39627	57.34509
14	Iceland	0.18051	2950.72653	31	Switzerland	0.33611	4524.67821
15	Ireland	0.27555	364.52781	32	Turkey	0.25900	146175.82805
16	Israel	0.25854	1926.04551	33	United Kingdom	0.49301	48.61626
17	Italy	0.50522	135.14076	34	United States	0.38019	2433.78269

**Proof.** First we prove that the function  $u(x) = \frac{1}{\psi'(x)}$ ,  $x \in (0, \infty)$ , is strictly convex. Indeed, one can easily see that

$$u''(x) = \frac{-\psi'''(x)[\psi'(x)]^2 + 2[\psi''(x)]^2\psi'(x)}{[\psi'(x)]^4},$$

and this is positive due to  $\psi'(x) > 0$  and the fact that

$$\frac{[\psi''(x)]^2}{\psi'''(x)\psi'(x)} > \frac{1}{2}.$$

The latter is a particular case of Corollary 2.3 in [2].

Now, in view of  $\lim_{x \rightarrow 0} \psi'(x) = \infty$ , we can extend  $u$  continuously to 0 by setting  $u(0) = 0$ . Then  $u$  is still strictly convex, and therefore, for every  $x, y > 0$  we have

$$u(x) = u\left(\frac{y}{x+y} \cdot 0 + \frac{x}{x+y} \cdot (x+y)\right) < \frac{y}{x+y}u(0) + \frac{x}{x+y}u(x+y).$$

Consequently,

$$u(x) < \frac{x}{x+y}u(x+y), \tag{8}$$

and likewise,

$$u(y) < \frac{y}{x+y}u(x+y). \tag{9}$$

Adding (8) and (9) together, we get the statement of the lemma.

**Lemma 3.** *The function  $\psi(2x) - \psi(x)$ ,  $x \in (0, \infty)$ , is decreasing and its range is  $(\ln 2, \infty)$ .*

**Proof.** This is easily seen from the identity  $\psi(2x) = \frac{1}{2}\psi(x) + \frac{1}{2}\psi(x + \frac{1}{2}) + \ln 2$ , which can be found in [1].

In the last lemma we collect some limiting properties of the digamma function and its derivative, see, e.g., [1–3] for details.

**Lemma 4.** *The digamma function  $\psi$  is a strictly concave, smooth function on  $(0, \infty)$  that satisfies the following limit relations:*

$$\begin{aligned} \lim_{x \rightarrow 0^+} \psi(x) &= -\infty, & \lim_{x \rightarrow \infty} \psi(x) &= \infty, & \lim_{x \rightarrow \infty} \psi'(x) &= 0, \\ \lim_{x \rightarrow \infty} (\psi(x) - \ln x) &= 0, & \lim_{x \rightarrow 0^+} (\psi(2x) - \psi(x)) &= \infty. \end{aligned}$$

## B. Considerations on the boundary of the mean parameter space

In Section 2, we saw that the mean parameter space  $\mathcal{M}$  consists of  $2n$ -tuples  $(A_1, \dots, A_n, B_1, \dots, B_n)$  obtained from the parameters  $(\mathbf{a}, \mathbf{b}) = (a_1, \dots, a_n, b_1, \dots, b_n)$  of the underlying beta distributions by Eq. (4).

Denoting by  $L(\mathbf{a}, \mathbf{b}) = (A_1(\mathbf{a}, \mathbf{b}), \dots, A_n(\mathbf{a}, \mathbf{b}), B_1(\mathbf{a}, \mathbf{b}), \dots, B_n(\mathbf{a}, \mathbf{b}))$  this dependence, i.e., the  $\Theta \rightarrow \mathcal{M}$  (one-to-one) mapping, a boundary point  $\bar{L} = (\bar{A}_1, \dots, \bar{A}_n, \bar{B}_1, \dots, \bar{B}_n)$  of  $\mathcal{M}$  can be obtained as  $\bar{L} = \lim_{k \rightarrow \infty} L(\mathbf{a}^k, \mathbf{b}^k)$ , where  $\mathbf{a}^k = (a_1^k, \dots, a_n^k)$ ,  $\mathbf{b}^k = (b_1^k, \dots, b_n^k)$ , and

$$\begin{aligned} \lim_{k \rightarrow \infty} a_i^k &= \bar{a}_i \in [0, \infty], \quad i = 1, \dots, n, \\ \lim_{k \rightarrow \infty} b_j^k &= \bar{b}_j \in [0, \infty], \quad j = 1, \dots, n. \end{aligned}$$

In view of Lemma 4, only the  $\bar{a}_i, \bar{b}_j = \infty$  cases have relevance. The sequence  $(\mathbf{a}^k, \mathbf{b}^k)$  can be chosen such that

$$\lim_{k \rightarrow \infty} \frac{a_i^k}{b_j^k} = x_{ij} \quad \text{with} \quad 0 < x_{ij} < \infty \quad \text{for} \quad i \neq j. \quad (10)$$

Then, using (4),

$$\begin{aligned} \bar{A}_i &= - \lim_{k \rightarrow \infty} \sum_{j \neq i} \left[ \psi(a_i^k + b_j^k) - \psi(a_i^k) \right] = \lim_{k \rightarrow \infty} \sum_{j \neq i} \left[ \ln(a_i^k) - \ln(a_i^k + b_j^k) \right] = \sum_{j \neq i} \ln \frac{x_{ij}}{1 + x_{ij}}, \quad i = 1, \dots, n, \\ \bar{B}_j &= - \lim_{k \rightarrow \infty} \sum_{i \neq j} \left[ \psi(a_i^k + b_j^k) - \psi(b_j^k) \right] = \lim_{k \rightarrow \infty} \sum_{i \neq j} \left[ \ln(b_j^k) - \ln(a_i^k + b_j^k) \right] = \sum_{i \neq j} \ln \frac{1}{1 + x_{ij}}, \quad j = 1, \dots, n. \end{aligned}$$

These equations show that the boundary point  $\bar{L}$  of  $\mathcal{M}$  contains – in its coordinates – the row- and column-sums of the matrices  $\mathbf{U}(\mathbf{W})$  and  $\mathbf{V}(\mathbf{W})$  respectively (see Section 2), where the general off-diagonal entry of the  $n \times n$  edge-weight matrix  $\mathbf{W}$  is  $\frac{x_{ij}}{1+x_{ij}}$ .

Observe that  $2n - 1$   $x_{ij}$ 's can be chosen freely, and all the others are obtainable from them. To see this, consider the complete bipartite graph on vertex classes  $(a_1, \dots, a_n)$  and  $(b_1, \dots, b_n)$ , where to the edge connecting  $a_i$  and  $b_j$  we assign  $\frac{a_i}{b_j}$ . Choose a minimal spanning tree of this graph (it contains  $2n - 1$  edges), and consider the sequence of  $(\mathbf{a}^k, \mathbf{b}^k)$ 's satisfying condition (10). Then, as  $k \rightarrow \infty$ , the  $x_{ij}$ 's of the edges not included in the spanning tree can be obtained from the  $x_{ij}$ 's of the  $2n - 1$  edges included in the spanning tree. Therefore the row- and column-sums of the edge-weight matrix  $\mathbf{W}$  of entries  $\frac{x_{ij}}{1+x_{ij}}$  ( $i \neq j$ ) are on a  $(2n - 1)$ -dimensional manifold in  $\mathbb{R}_+^{2n}$ , so they are on the boundary of the convex hull  $\mathcal{T}$  of the possible sufficient statistics  $(\mathbf{R}, \mathbf{C})$ . However, this boundary has zero Lebesgue measure, and so, zero probability with respect to the underlying absolutely continuous distribution.

## C. Proof of Theorem 2

It suffices to prove that some induced matrix norm of the matrix of the first derivatives  $\mathbf{J}$  of  $f$  at  $\hat{\theta}$  is strictly less than 1. We prove this for the  $L_1$ -norm. From (6) we obtain that

$$\frac{\partial g_i}{\partial a_i}(\hat{\mathbf{a}}, \hat{\mathbf{b}}) = \frac{\frac{1}{n-1} \sum_{j \neq i} \psi'(\hat{a}_i + \hat{b}_j)}{\psi' \left[ \psi^{-1} \left( \frac{1}{n-1} \sum_{j \neq i} \psi(\hat{a}_i + \hat{b}_j) + \frac{R_i}{n-1} \right) \right]}. \quad (11)$$



From (3) we have  $\frac{1}{n-1} \sum_{j \neq i} \psi(\hat{a}_i + \hat{b}_j) + \frac{R_i}{n-1} = \psi(\hat{a}_i)$ . Substituting it into (11), we get

$$\frac{\partial g_i}{\partial a_j}(\hat{\mathbf{a}}, \hat{\mathbf{b}}) = \begin{cases} \frac{\sum_{s \neq i} \frac{1}{n-1} \psi'(\hat{a}_i + \hat{b}_s)}{\psi'(\hat{a}_i)}, & \text{if } j = i, \\ 0, & \text{if } j \neq i. \end{cases}$$

Likewise,

$$\frac{\partial g_i}{\partial b_j}(\hat{\mathbf{a}}, \hat{\mathbf{b}}) = \begin{cases} 0, & \text{if } j = i, \\ \frac{\frac{1}{n-1} \psi'(\hat{a}_i + \hat{b}_j)}{\psi'(\hat{a}_i)}, & \text{if } j \neq i. \end{cases}$$

Further,

$$\frac{\partial h_i}{\partial a_j}(\hat{\mathbf{a}}, \hat{\mathbf{b}}) = \begin{cases} 0, & \text{if } j = i, \\ \frac{\frac{1}{n-1} \psi'(\hat{a}_j + \hat{b}_i)}{\psi'(\hat{b}_i)}, & \text{if } j \neq i \end{cases}$$

and

$$\frac{\partial h_i}{\partial b_j}(\hat{\mathbf{a}}, \hat{\mathbf{b}}) = \begin{cases} \frac{\sum_{s \neq i} \frac{1}{n-1} \psi'(\hat{a}_s + \hat{b}_i)}{\psi'(\hat{b}_i)}, & \text{if } j = i, \\ 0, & \text{if } j \neq i. \end{cases}$$

Observe that  $J(\hat{\mathbf{a}}, \hat{\mathbf{b}})$  has nonnegative entries. Therefore, its  $L_1$ -norm is the maximum of its column-sums. The  $j$ th column-sum of  $J(\hat{\mathbf{a}}, \hat{\mathbf{b}})$  is equal to

$$\frac{\sum_{s \neq j} \frac{1}{n-1} \psi'(\hat{a}_j + \hat{b}_s)}{\psi'(\hat{a}_j)} + \sum_{s \neq j} \frac{\frac{1}{n-1} \psi'(\hat{a}_j + \hat{b}_s)}{\psi'(\hat{b}_s)} = \frac{1}{n-1} \sum_{s \neq j} \psi'(\hat{a}_j + \hat{b}_s) \left( \frac{1}{\psi'(\hat{a}_j)} + \frac{1}{\psi'(\hat{b}_s)} \right) \quad (12)$$

for  $j = 1, \dots, n$ ; and likewise, the  $(n+j)$ th column-sum of  $J(\hat{\mathbf{a}}, \hat{\mathbf{b}})$  is

$$\frac{1}{n-1} \sum_{s \neq j} \psi'(\hat{a}_s + \hat{b}_j) \left( \frac{1}{\psi'(\hat{a}_s)} + \frac{1}{\psi'(\hat{b}_j)} \right) \quad (13)$$

for  $j = 1, \dots, n$ . As (12) and (13) are of similar in appearance, it suffices to prove that the right-hand side of (12) is less than 1. But

$$\psi'(\hat{a}_j + \hat{b}_s) \left( \frac{1}{\psi'(\hat{a}_j)} + \frac{1}{\psi'(\hat{b}_s)} \right) < 1$$

holds by Lemma 2, and we have  $n-1$  terms in the summation.

Since  $f: \mathbb{R}^{2n} \rightarrow \mathbb{R}^{2n}$  is continuously differentiable in a neighborhood of  $\hat{\boldsymbol{\theta}} = (\hat{\mathbf{a}}, \hat{\mathbf{b}})$ , Theorem 3 of [7] implies that there is a closed neighborhood  $K$  of  $\hat{\boldsymbol{\theta}}$  such that  $f$  is a contraction on  $K$ . In particular, the fixed point iteration  $f(\boldsymbol{\theta}^{(it-1)}) = \boldsymbol{\theta}^{(it)}$  ( $it \rightarrow \infty$ ) converges for every  $\boldsymbol{\theta}^{(0)} \in K$  to  $\hat{\boldsymbol{\theta}}$ , which is the unique solution of (6).

## Acknowledgments

The first author was supported by BME — Artificial Intelligence FIKP Grant of EMMI (BME FIKP-MI/SC).

## REFERENCES

1. M. Abramowitz and I. A. Stegun, *Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables*, Dover Publishing, New York (1972).
2. H. Alzer and J. Wells, “Inequalities for the polygamma functions,” *SIAM J. Math. Anal.*, **29**, No. 6, 1459–1466 (1998).
3. J. Bernardo, “Psi (digamma) function. Algorithm AS 103,” *Appl. Stat.*, **25**, 315–317 (1976).
4. M. Bolla, *Spectral Clustering and Biclustering*, Wiley, New York (2013).
5. M. Bolla and A. Elbanna, “Estimating parameters of a probabilistic heterogeneous block model via the EM algorithm,” *J. Probab. Stat.*, Article ID 657965 (2015).
6. S. Chatterjee, P. Diaconis, and A. Sly, “Random graphs with a given degree sequence,” *Ann. Stat.*, **21**, 1400–1435 (2010).
7. M. Grasmair, *Fixed point iterations*,  
[https://wiki.math.ntnu.no/\\_media/ma2501/2014v/fixedpoint.pdf](https://wiki.math.ntnu.no/_media/ma2501/2014v/fixedpoint.pdf)
8. C. J. Hillar and A. Wibisono, “Maximum entropy distributions on graphs,” *ArXiv*, arXiv:1301.3321v2 (2013).
9. S. L. Lauritzen, *Graphical Models*, Oxford University Press (1995).
10. T. Yan, C. Leng, and J. Zhu, “Asymptotics in directed exponential random graph models with an increasing bi-degree sequence,” *Ann. Stat.*, **44**, 31–57 (2016).
11. M. Wainwright and M. I. Jordan, “Graphical models, exponential families, and variational inference,” *Found. Trend. Mach. Learn.*, **1**, No. 1–2, 1–305 (2008).