

**Statistica Sinica Preprint No: SS-2016-0397.R1**

<b>Title</b>	Network Inference From Grouped Observations Using Hub Models
<b>Manuscript ID</b>	SS-2016-0397.R1
<b>URL</b>	<a href="http://www.stat.sinica.edu.tw/statistica/">http://www.stat.sinica.edu.tw/statistica/</a>
<b>DOI</b>	10.5705/ss.202016.0397
<b>Complete List of Authors</b>	Yunpeng Zhao and Charles Weko
<b>Corresponding Author</b>	Yunpeng Zhao
<b>E-mail</b>	yzhao15@gmu.edu

---

# Network Inference From Grouped Observations Using Hub Models

YUNPENG ZHAO<sup>1</sup> AND CHARLES WEKO<sup>2</sup>

*George Mason University<sup>1</sup> and United States Army<sup>2</sup>*

2     *Abstract:* In medical research, economics, and the social sciences data frequently  
3     appear as subsets of a set of objects. Over the past century a number of descriptive  
4     statistics have been developed to infer network structure from such data. However,  
5     these measures lack a generating mechanism that links the inferred network struc-  
6     ture to the observed groups. To address this issue, we propose a model-based  
7     approach called the *Hub Model* which assumes that every observed group has a  
8     leader and that the leader has brought together the other members of the group.  
9     The performance of Hub Models is demonstrated by simulation studies. We apply  
10    this model to the characters in a famous 18<sup>th</sup> century Chinese novel.

11    *Key words and phrases:* Social network analysis, affiliation network, expectation-  
12    maximization algorithm, half weight index, Dream of the Red Chamber.

---

## 13 1 INTRODUCTION

14 A network can be denoted by  $N = (V, E)$ , where  $V = \{v_1, v_2, \dots, v_n\}$  is the set  
15 of  $n$  nodes, and  $E$  is the set of edges between nodes. In this article, we focus  
16 on symmetric weighted networks represented by an  $n \times n$  adjacency matrix,  
17  $A$ , where the element  $A_{ij}$  measures the relationship strength between nodes  
18  $v_i$  and  $v_j$ .

19 Traditionally, statistical network analysis focuses on modeling *observed*  
20 network structure (e.g., highway systems or electrical transmission grids). In  
21 this situation, nodes are well defined and the physical links between nodes  
22 is observable (Hiller and Lieberman (2001); Newman (2011)). In some fields  
23 of research (e.g., the social sciences) network structure is not explicit, the  
24 observable data are groups of individuals and a model is presumed to produce  
25 the groups. The fundamental task is to estimate model parameters from such  
26 data.

27 Wasserman and Faust (1994) introduce inference of relationships with the  
28 example of children attending birthday parties. In their example, the children  
29 act as nodes in the network and the birthday parties represent subsets of  
30 children.

31 In this paper, a collection of nodes observed in the same sample is called

---

32 a *group* and a dataset is called *grouped data*. In Wasserman and Faust's  
33 example, each party defines a group and the set of all parties is the grouped  
34 data. Two individuals are said to *co-occur* if they appear in the same group.

35 One common technique used to estimate an adjacency matrix from grouped  
36 data is to count the number of times that a pair of nodes appears in the  
37 same group (Zachary (1977);Freeman et al. (1989);Wasserman and Faust  
38 (1994);Kolaczyk (2009);Brent et al. (2011)). Frequently, a threshold is ap-  
39 plied to this count to create an unweighted adjacency matrix; however,  
40 Choudhury et al. (2010) show that the characteristics of networks inferred by  
41 this technique are sensitive to the threshold. We adopt a generalized version  
42 of the inter-citation frequency (Kolaczyk (2009)) which measures the number  
43 of times a pair of nodes is observed to co-occur in the dataset. We refer to  
44 this measure as the *co-occurrence matrix*.

45 An alternative technique, called the *half weight index* (Cairns and Schwa-  
46 ger (1987)), estimates an adjacency matrix by the frequency that two nodes  
47 co-occur given that one of them is observed. This addresses a shortcoming  
48 of the co-occurrence matrix in which nodes that appear rarely can be es-  
49 timated to have a weak relationship even though the relationship is quite  
50 strong (Voelkl et al. (2011)).

51 The co-occurrence matrix and half weight index both have probabilis-

---

52 tic interpretations. The co-occurrence matrix estimates the probability that  
53 two nodes will be observed together. The half weight index estimates the  
54 probability that two nodes will be observed together given that one of them  
55 is observed. These are not equivalent to the probability of an active rela-  
56 tionship between nodes, and neither of these techniques describe the process  
57 which leads to the generation of the observed groups. It is unclear how these  
58 descriptive statistics relate to the grouped data in these methods.

59 We propose a model-based approach for grouped data generation which  
60 we refer to as the *Hub Model* because each observed group is assumed to be  
61 brought together by a hub node (see Figure 1).

62 The Hub Model differs from such classical network models as the stochas-  
63 tic blockmodel and its variants (Holland et al. (1983);Airoldi et al. (2008)),  
64 the exponential random graph models (Frank and Strauss (1986);Robins et al.  
65 (2007)), the latent space model and its variants (Hoff et al. (2002);Handcock  
66 et al. (2007)), among others (see Goldenberg et al. (2010) for a comprehen-  
67 sive review). These models focus on modeling the statistical behavior of the  
68 network, treating the network as the observed data, while the Hub Model  
69 treats the network as latent governing the grouping behavior of a popula-  
70 tion. Our task is to estimate the latent network, the adjacency matrix, from  
71 the observed group data. In this article, we treat the adjacency matrix as

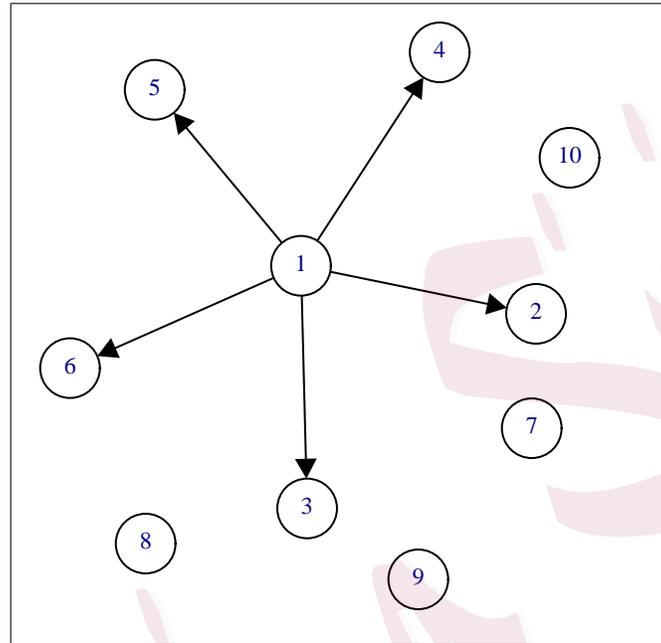


Figure 1: The generating mechanism of the Hub Model is demonstrated on a group of 10 nodes. In the observed sample, nodes  $v_1, \dots, v_6$  are members of the group while nodes  $v_7, \dots, v_{10}$  are not members of the group. The observed group is the result of the hub node,  $v_1$ , bringing together nodes  $v_2, \dots, v_6$ .

72 fixed parameters and make no structural assumption about it. If there were  
73 *a priori* information about the latent network, such as that it follows the  
74 stochastic blockmodel or the exponential random graph model, then one  
75 could take a Bayesian approach and use this model as *a priori*. For more  
76 discussion, refer to Section 7.

77 The Hub Model belongs to the family of finite mixture models which  
78 has been applied in many situations, including text classification (Carreira-

---

79 Perpinan and Renals (2000)), topic models (Anandkumar et al. (2015)), fin-  
80 gerprint identification (Vretos et al. (2012)), and product recommendation  
81 (Colace et al. (2015)).

82 Hub Models have the advantage that relationship strength is both math-  
83 ematically well defined and practical to researchers. In the Hub Model,  $A_{ij}$ ,  
84 is defined as the probability that node  $v_i$  will include node  $v_j$  when  $v_i$  is the  
85 hub node of a group. The formal definition of the Hub Model is given in  
86 Section 3.

87 As an introduction, consider the hypothetical relationships in Figure 2a.  
88 In this example there is a pair of nodes,  $v_1$  and  $v_2$ , that never directly pair to  
89 each other, but have an 80% chance of interacting with five nodes:  $A_{ij} = 0.8$   
90 for all  $i \leq 2$  and  $j \geq 3$ , while  $A_{ij} = 0$  otherwise. In Figure 2b, the co-  
91 occurrence matrix mistakenly assigns a relatively strong relationship to nodes  
92  $v_1$  and  $v_2$  because they often co-occur. In Figure 2c, the half weight index  
93 arrives at a similar conclusion. In both Figures 2b and 2c, the non-existent  
94 relationship between nodes  $v_1$  and  $v_2$  is estimated to be stronger than all other  
95 relationships. By contrast, the Hub Model in Figure 2d clearly captures the  
96 relationships of the population.

97 To the best of our knowledge, there have been limited attempts to apply  
98 model-based approaches to these data. Rabbat et al. (2008) provide an

---

99 application for telecommunication networks. They model group formation  
100 as a random walk from a source node to a terminal node. This model assumes  
101 a distinctly different process of group formation than do Hub Models. The  
102 nodes along the path are subjected to an unknown permutation to account  
103 for the lack of order information. Treating permutations as missing data,  
104 they employ a *Monte Carlo EM* algorithm based on importance sampling to  
105 estimate the parameters of the model.

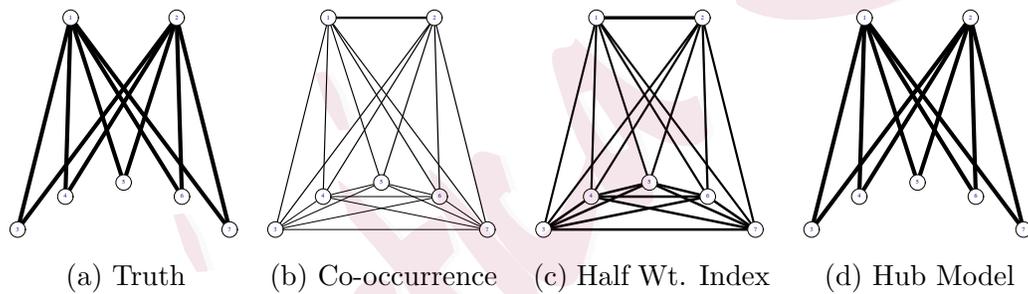


Figure 2: Comparison of Estimation Techniques

106 In the following sections we present a formal description of the grouped  
107 data structure, review existing techniques, and define Hub Models. Then  
108 we address Hub Model identifiability and provide a theorem that proves  
109 that a symmetric adjacency matrix is a sufficient condition for identifiability.  
110 We propose an EM algorithm to solve the maximum likelihood estimator of  
111 the Hub Model. We have evaluated the model performance by simulation  
112 studies. We have applied the Hub Model to infer the relationships among

---

113 the characters of the 18<sup>th</sup> century Chinese novel, *Dream of the Red Chamber*.  
114 We close with a discussion of how the size of the population impacts model  
115 efficiency and ways to incorporate network structure assumptions to simplify  
116 the model.

## 117 2 GROUPED DATA

### 118 2.1 Data Structure

119 For a population of  $n$  individuals,  $V = \{v_1, \dots, v_n\}$ , we observe  $T$  subsets  
120 of the global population,  $\{V^{(t)} | V^{(t)} \subseteq V, t = 1, \dots, T\}$ . Each observed subset  
121 can be coded as an  $n$  length row vector  $G^{(t)}$  where

$$G_i^{(t)} = \begin{cases} 1 & \text{if } v_i \in V^{(t)} \\ 0 & \text{if } v_i \notin V^{(t)} \end{cases}$$

122 The full set of observations is denoted by a  $T \times n$  matrix,  $G$ . The  $t^{\text{th}}$  row  
123 of  $G$  is  $G^{(t)}$ .

124 **2.2 Existing Methods**

125 Inferring relationships from grouped data relies on descriptive statistics that  
126 count the number of times that two nodes are observed together. We focus on  
127 two popular techniques which estimate probabilities of individual behavior.

128 A simple measure of grouped data is the *co-occurrence matrix*. Versions  
129 of this technique appear throughout the literature under many names and  
130 notations including: *capacity matrix* (Zachary (1977)), *sociomatrix* (Wasser-  
131 man and Faust (1994)), *inter-citation frequency* (Kolaczyk (2009)), *cocitation*  
132 *matrix* (Newman (2011)), and *strength* (Brent et al. (2011)).

133 A co-occurrence matrix,  $O$ , is the  $n \times n$  symmetric matrix

$$O = \frac{G'G}{T}, \quad (2.1)$$

134 which estimates the frequency that the nodes  $v_i$  and  $v_j$  are observed in  
135 the same group.

136 One shortcoming of the co-occurrence matrix is that it estimates the prob-  
137 ability that two nodes *will be observed* to co-occur in a given observation. If  
138 two nodes have a strong relationship, but appear in the dataset infrequently,  
139 the co-occurrence matrix estimates a low probability that the two nodes *will*  
140 *be observed* to co-occur.

141 As an example, consider four nodes  $v_1, \dots, v_4$  and the grouped data represented in Table 1. For this dataset,  $O_{1,2} = \frac{2}{5}$  and  $O_{3,4} = \frac{2}{5}$ , but every

Event	Node			
	$v_1$	$v_2$	$v_3$	$v_4$
1	1	0	0	0
2	1	1	0	0
3	1	1	0	0
4	1	0	1	1
5	0	1	1	1

Table 1: Notional Grouped Data

142

143 time node  $v_3$  is present node  $v_4$  is also present. A researcher might conclude  
144 that there is some aspect of the relationship between nodes  $v_3$  and  $v_4$  which  
145 has been understated.

146 As an alternative, the *half weight index* estimates the probability that  
147 two nodes will be observed to co-occur given that one of them is observed  
148 (Cairns and Schwager (1987)).

149 The half weight index has been introduced in a number of equivalent  
150 forms (Dice (1945)). Computationally, the most direct form is

$$H_{ij} = \frac{2 \sum_t G_i^{(t)} G_j^{(t)}}{\sum_t G_i^{(t)} + \sum_t G_j^{(t)}}. \quad (2.2)$$

151 Returning to the example in Table 1,  $H_{1,2} = \frac{4}{7}$  while  $H_{3,4} = \frac{4}{4}$ . Therefore,

---

152 the half weight index infers a different network than the co-occurrence matrix.

### 153 3 HUB MODELS

#### 154 3.1 Generating Mechanism

155 Hub Models (HM) assume that each group is a star subgraph on the global  
156 population. The hub node connecting the observed group is represented by  
157 an  $n$  length row vector,  $S^{(t)}$ , where

$$S_i^{(t)} = \begin{cases} 1 & \text{if } v_i \text{ the hub node of sample } t, \\ 0 & \text{otherwise.} \end{cases}$$

158 There is one and only one element of  $S^{(t)}$  that is equal to 1, and each group  
159 is independently generated by a two step process: we take the hub node to  
160 be drawn from a multinomial distribution with parameter  $\rho = (\rho_1, \dots, \rho_n)$ ,  
161 and we suppose the hub node,  $v_i$ , chooses to include  $v_j$  in the group with  
162 probability  $A_{ij} = \mathbb{P}(G_j^{(t)} = 1 | S_i^{(t)} = 1)$ , with  $A_{ii} = 1$  for all  $i$ .

163 In most practical applications, the hub node of each group is unknown,  
164 and we focus on this case. We refer to the model where leaders are known  
165 as the Known Hub Model (KHM).

166 Since the co-occurrence matrix and half weight index produce a symmetric

167 adjacency matrix, we assume the Hub Model adjacency matrix is symmetric.  
168 Symmetry ensures the identifiability of the Hub Model when group leaders  
169 are unobserved (Supplemental Material S1.2).

170 This generating mechanism implies that each observed group is indepen-  
171 dent of every other. In particular,  $G^{(t)}$  is not a transformation of  $G^{(t-1)}$  and  
172 the order in which groups are observed contains no information about the  
173 relationships between group members. Researchers often collect data in such  
174 a way as to ensure this property (Bejder et al. (1998)).

### 175 **3.2 Likelihood of the Hub Model**

176 Under the HM, the probability of an observation has the form of a finite  
177 mixture model with  $n$  components

$$\mathbb{P}(G^{(t)}|A, \rho) = \sum_{i=1}^n \rho_i G_i^{(t)} \prod_j A_{ij}^{G_j^{(t)}} (1 - A_{ij})^{1-G_j^{(t)}}. \quad (3.1)$$

178 By taking the log of the product of individual observed groups, the log  
179 likelihood function for the full set of observations is

$$\mathcal{L}(G|A, \rho) = \sum_t \log \left[ \sum_{i=1}^n \rho_i G_i^{(t)} \prod_j A_{ij}^{G_j^{(t)}} (1 - A_{ij})^{1-G_j^{(t)}} \right]. \quad (3.2)$$

### 3.2 Likelihood of the Hub Model

---

180 Solving for the MLE of HM is an optimization problem with the con-  
181 straints  $\sum_i \rho_i = 1$ , and  $A_{ij} = A_{ji}$  for all  $i$  and  $j$ . This gives the Lagrange  
182 function

$$\Lambda(G|A, \rho) = \mathcal{L}(G|A, \rho) - \lambda_o[(\sum_i \rho_i) - 1] - \sum_{i < j} \lambda_{ij}(A_{ij} - A_{ji}). \quad (3.3)$$

183 The log likelihood does not have a closed-form solution for the MLE.  
184 Instead we derive estimating equations that can be incorporated into an  
185 Expectation Maximization algorithm. Before doing so we investigate the  
186 identifiability of the Hub Model.

A basic requirement for any model is *identifiability*. For Hub Models, this means that, for any two sets of parameters  $\{A, \rho\}$  and  $\{A^*, \rho^*\}$ ,

$$\mathbb{P}(G = g|A, \rho) = \mathbb{P}(G = g|A^*, \rho^*) \quad \forall g \implies A = A^*, \rho = \rho^*. \quad (3.4)$$

187 The generating mechanism for Hub Models is equivalent to a finite mix-  
188 ture model of multivariate Bernoulli random variables. In general, such a  
189 model is not identifiable (Teicher (1961)). This shortcoming does not pre-  
190 vent such models from being useful in many applications. For example, when  
191 dealing with classification problems where the researcher only has to identify

192 which component density an observation came from, this type of mixture can  
193 be effectively used (Carreira-Perpinan and Renals (2000)). In such a situ-  
194 ation, the individual parameters of the multivariate Bernoulli random vari-  
195 ables are not of interest, but identifiability presents a challenge here because  
196 we are specifically interested in the individual parameters of the adjacency  
197 matrix.

198 If no constraint is put on the adjacency matrix, the model is unidenti-  
199 fiable. We have a sufficient condition for identifiability, see Supplemental  
200 Material S1 for more details.

201 **Theorem 1.** Let  $A$  and  $A^*$  be symmetric adjacency matrices with  $A_{ii} =$   
202  $A_{ii}^* = 1$  for all  $i$ ,  $A_{ij} < 1$  and  $A_{ij}^* < 1$  for all  $i \neq j$ . If  $\mathbb{P}(g|A, \rho) = \mathbb{P}(g|A^*, \rho^*)$   
203 for all  $g$ , then  $\{A, \rho\} = \{A^*, \rho^*\}$ .

204 Even though symmetry of the adjacency matrix is a natural assumption,  
205 it is only a sufficient condition for identifiability. For future work, we will  
206 explore other assumptions to ensure identifiability.

### 207 3.3 Estimating Equations

208 In Supplemental Materials S2, we derive (3.5) and (3.6) as estimating equa-  
209 tions that the MLE must satisfy. The maximum likelihood estimator does

---

210 not have a closed-form solution for the parameters as the right side of the  
211 estimating equations includes the estimated parameters. We will show that  
212 solving these equations iteratively is equivalent to an EM algorithm.

$$\hat{A}_{xy} = \frac{\sum_t G_y^{(t)} \mathbb{P}(S_x = 1|G^{(t)}) + \sum_t G_x^{(t)} \mathbb{P}(S_y = 1|G^{(t)})}{\sum_t [\mathbb{P}(S_x = 1|G^{(t)}) + \mathbb{P}(S_y = 1|G^{(t)})]}. \quad (3.5)$$

$$\hat{\rho}_x = \frac{\sum_{t=1}^T \mathbb{P}(S_x^{(t)} = 1|G^{(t)})}{T}. \quad (3.6)$$

## 213 4 EM ALGORITHM

214 These estimating equations depend on the probability  $\mathbb{P}(S_x^{(t)} = 1|G^{(t)})$ . This  
215 suggests an algorithm updating  $\{\hat{A}, \hat{\rho}\}$  and  $\mathbb{P}(S_x^{(t)} = 1|G^{(t)})$  iteratively, which  
216 can be fitted into the general framework of an EM algorithm.

217 EM algorithms formulate a complete data model, then solve the model as  
218 if some data is observed and other data is missing. In this case, the Known  
219 Hub Model serves as the complete data model,  $G$  is the observed data, and  
220  $S$  is the missing data. Each iteration of the EM algorithm consists of an  
221 expectation step followed by a maximization step (McLachlan and Krishnan  
222 (2008)).

223

## E-Step

Since the log likelihood function of the complete data model is linear in the unobserved data, the E-Step (on the  $(m + 1)^{th}$  iteration) simply requires calculating the current conditional expectation of  $S_i^{(t)}$  given the observed data (see McLachlan and Krishnan (2008) for a detailed explanation).

$$\begin{aligned} E[S_x^{(t)}|G^{(t)}] &= \mathbb{P}(S_x^{(t)} = 1|G^{(t)}) \\ &= \frac{\rho_x G_x^{(t)} \prod_j A_{xj}^{G_j^{(t)}} (1 - A_{xj})^{1-G_j^{(t)}}}{\sum_{i=1}^n \rho_i G_i^{(t)} \prod_j A_{ij}^{G_j^{(t)}} (1 - A_{ij})^{1-G_j^{(t)}}} \end{aligned} \quad (4.1)$$

224

## M-Step

225 The M-Step replaces  $\mathbb{P}(S_x^{(t)} = 1|G^{(t)})$  on the right hand side of (3.5) and  
226 (3.6) with  $E[S_x^{(t)}|G^{(t)}]$  from (4.1).

227

## Algorithm

228 Several standard techniques are used to improve the performance of the  
229 EM algorithm. We first run the EM algorithm ten times with different start-  
230 ing points and choose the solution with the highest likelihood. We limit the  
231 number of iterations applied to a starting point on the grounds that with a

232 bad starting point, it takes a long time to converge to a point not close to  
 233 the maximum. As a final step, we treat any  $\hat{A}_{xy} \leq 10^{-4}$  as  $\hat{A}_{xy} = 0$ . We  
 234 apply this finishing step to remove clutter from the returned solutions.

```

Data: G
Result:  $\hat{A}, \hat{\rho}$ 
Initialize:
 $\mathcal{L}(G|\hat{A}) = -\infty$ 
for rep=1 to 10 do
  Initialize:
   $\hat{A}_{ij}^{(0)} = \text{unif}(0, 1) \quad \forall \{i, j\}$ 
   $X_i = \text{unif}(0, 1) \quad \forall i$ 
   $\hat{\rho}_i^{(0)} = \frac{X_i}{\sum_k X_k}$ 
   $\Delta\mathcal{L}(G|A^{(0)}) = 10^4$ 
  counter=1
  while  $\frac{|\Delta\mathcal{L}(G|A^{(m+1)})|}{\mathcal{L}(G|A^{(m)})} > 10^{-4}$  and counter < 100 do
    E-Step
    Update  $\mathbb{P}(S_k^{(t)} = 1|G^{(t)})$  by Equation 4.1
    M-Step
    Update  $A^{(m+1)}$  by Equation S2.10
    Update  $\rho^{(m+1)}$  by Equation S2.13
     $\Delta\mathcal{L}(G|A^{(m+1)}) = \mathcal{L}(G|A^{(m+1)}) - \mathcal{L}(G|A^{(m)})$ 
    counter=counter+1
  end
  if  $\mathcal{L}(G|A^{(m+1)}) > \mathcal{L}(G|\hat{A})$  then
    if  $\hat{A}_{ij} \leq 10^{-4}$  then
       $\hat{A}_{ij} = 0$ 
    else
       $\hat{A}_{ij} = A_{ij}^{(m+1)}$ 
    end
  end
end
end

```

**Algorithm 1:** Expectation Maximization Algorithm for the Hub Model

---

## 235 5 SIMULATION

236 To perform simulations, we generated parameters  $\{A, \rho\}$  as follows.

237 For  $\rho$ , we selected  $n$  random numbers,  $X_i$ , uniformly and divided each  
238 random number by the sum of all  $X_i$ 's,  $\rho_i = \frac{X_i}{\sum_i X_i}$ .

239 We used a two-step process to generate the adjacency matrix. First, we  
240 created a symmetric unweighted undirected random graph on  $n$  nodes using  
241 the configuration model (Jackson (2010)) with the power law distribution  
242  $\mathbb{P}(k) \propto k^{-\eta}$ , where  $k$  is the possible value of the node degree. We assumed  
243 a power law degree distribution because it is commonly believed that many  
244 social networks have such a property (Newman (2011)). In all simulations,  
245 we chose  $\eta = 2$ ; many networks are reported to have a power between 2 and 3  
246 and a power of 2 generates the densest of them. We refer to this unweighted  
247 graph as the *structure* of the network.

248 Each edge in the graph was assigned a relationship strength with a beta  
249 distribution,

$$A_{ij} = \begin{cases} \text{Beta}(\alpha, \beta) & \text{if there is an edge between } v_i \text{ and } v_j \\ 0 & \text{otherwise} \end{cases}$$

250 We let  $A_{ji} = A_{ij}$  to ensure symmetry. We set  $\alpha = 1$  and  $\beta = 4$  in the beta

---

251 distribution so that the average relationship strength is less than 0.5, which  
252 we believe is realistic.

253 In Tables 2 and 3, we consider five different network sizes  $n = 10, 20, 50, 100, 150$ .  
254 For the first two cases, we set the minimum node degree to be 1 in the power  
255 law distribution; for the last three cases, we set the minimum degree to  
256 be 5 in order to make sure the networks were not too sparse. For each  
257 size, we generated 100 sets of parameters  $\{A, \rho\}$  using the setup described  
258 above. For each  $\{A, \rho\}$ , we generated a dataset with  $T$  groups. Each aver-  
259 age and standard deviation was calculated over this 100 datasets. We took  
260  $T = 100, 200, 500, 1000, 2000, 5000, 10000, 20000, 50000$ .

261 We first measured the ability of the estimated adjacency matrix  $\hat{A}$  to  
262 correctly identify the structure. To do this we defined true positives and  
263 true negatives as

$$TP = \sum_{i < j} \mathbf{1}_{(A_{ij} > 0)} \mathbf{1}_{(\hat{A}_{ij} > 10^{-4})},$$
$$TN = \sum_{i < j} \mathbf{1}_{(A_{ij} = 0)} \mathbf{1}_{(\hat{A}_{ij} \leq 10^{-4})}.$$

264 Here,  $v_i$  and  $v_j$  were considered to have no relationship if the estimated  
265 link strength was below  $10^{-4}$ . False positives and false negatives were cal-

---

266 culated similarly. We used the Matthews correlation coefficient (MCC) to  
267 measure the identification of the structure because it is a binary classification  
268 measure that accounts for situations where the number of ones is significantly  
269 different than the number of zeros (Liu et al. (2015)). Based on our setup,  
270 our simulated structures had many more zeros than ones.

271 For the non-zero elements  $A_{ij}$ , we further evaluated the difference between  
272 the numerical values of  $A_{ij}$  and  $\hat{A}_{ij}$  by calculating the mean absolute error  
273 (MAE) of non-zero  $A_{ij}$ ,

$$MAE(A) = \frac{\sum_{i < j} |\hat{A}_{ij} - A_{ij}| \mathbf{1}_{(A_{ij} > 0)}}{\sum_{i < j} \mathbf{1}_{(A_{ij} > 0)}}.$$

274 We also report the average run time and the average number of iterations  
275 for the EM algorithm when the simulation is run on an Intel Pentium CPU  
276 G2030 at 3.00 GHz with 4.00GB of RAM.

277 The first observation from Tables 2 and 3 is that for a fixed value of  $n$   
278 the average error of both the MCC and the MAE decline as the number of  
279 observations increases. For a fixed number of observations, the average error  
280 increases as the number of nodes increases.

---

281 The standard deviation of estimates generally improves once the number  
282 of observations exceeds the number of parameters in the model. For example,  
283 with 100 nodes there are roughly 10,000 parameters to estimate, thus samples  
284 of only 2,000 or 5,000 observations demonstrate high standard deviations.

285 Average run time generally increases as the number of observations and  
286 the number of nodes increase. An important factor affecting the run time  
287 is the number of iterations the EM algorithm performs before converging.  
288 In Table 2 the number of iterations declines as observations increase until it  
289 appears to approach a minimum number. Table 3 provides further insight as  
290 the number of iterations generally increases until the number of observations  
291 is roughly equal to the number of parameters in the model, after which the  
292 iterations declines. Up to that point, the algorithm quickly converges to  
293 an adjacency matrix which is sparser than the true adjacency matrix due  
294 to insufficient sample size. The implication of these declining iterations is  
295 that run time is not strictly a function of the size of the dataset, but the  
296 relationship between the number of nodes and the number of observations.

	$n = 10$					
Obs	Avg MCC	StDev MCC	Avg MAE(A)	StDev MAE(A)	Avg Run Time (sec)	Avg Iterations
100	0.8010	0.0977	0.0533	0.0219	0.0472	20.258
200	0.8929	0.0903	0.0349	0.0128	0.0431	16.670
500	0.9487	0.0530	0.0212	0.0071	0.0411	13.618
1000	0.9770	0.0364	0.0147	0.0047	0.0369	12.011
2000	0.9865	0.0279	0.0102	0.0030	0.0353	10.613
5000	0.9984	0.0115	0.0067	0.0019	0.0298	9.604
10000	0.9988	0.0086	0.0045	0.0014	0.0295	9.416
20000	0.9994	0.0060	0.0035	0.0009	0.0305	9.327
50000	1	0	0.0020	0.0006	0.0316	9.210
	$n = 20$					
100	0.6727	0.0972	0.0833	0.0210	0.1005	21.007
200	0.7984	0.0756	0.0599	0.0154	0.0992	19.961
500	0.8781	0.0576	0.0340	0.0079	0.1039	17.793
1000	0.9147	0.0594	0.0225	0.0056	0.1131	15.418
2000	0.9360	0.0612	0.0150	0.0033	0.1473	13.803
5000	0.9734	0.0367	0.0099	0.0024	0.1653	11.571
10000	0.9842	0.0393	0.0069	0.0019	0.1806	10.662
20000	0.9937	0.0187	0.0048	0.0013	0.2052	10.260
50000	0.9989	0.0070	0.0031	0.0006	0.2320	9.888

Table 2: Average and Standard Deviation of Mean Absolute Error as Observations Increase

	<i>n</i> = 50					
Obs	Avg MCC	StDev MCC	Avg MAE(A)	StDev MAE(A)	Avg Run Time (sec)	Avg Iterations
100	0.3454	0.0503	0.1680	0.0139	0.2272	5.261
200	0.3987	0.0622	0.1368	0.0081	0.9216	16.237
500	0.5815	0.0668	0.0936	0.0085	2.7233	36.148
1000	0.8499	0.0302	0.0526	0.0049	2.6903	38.222
2000	0.9013	0.0176	0.0345	0.0030	2.3761	24.713
5000	0.9127	0.0193	0.0212	0.0017	2.8953	17.802
10000	0.9074	0.0259	0.0145	0.0012	5.1788	15.343
20000	0.9080	0.0327	0.0104	0.0008	7.1548	13.932
50000	0.9142	0.0383	0.0065	0.0006	12.190	12.866
	<i>n</i> = 100					
100	0.2620	0.0352	0.1955	0.0096	0.2058	2.040
200	0.3187	0.0346	0.1756	0.0109	0.2922	2.533
500	0.3495	0.0519	0.1359	0.0070	1.8683	9.151
1000	0.3857	0.0498	0.1109	0.0074	6.9431	25.852
2000	0.5343	0.1055	0.0748	0.0100	14.6644	44.035
5000	0.8236	0.1469	0.0351	0.0080	17.5031	34.544
10000	0.9128	0.0826	0.0219	0.0028	19.4031	23.370
20000	0.9355	0.0579	0.0148	0.0015	22.4366	17.494
50000	0.9484	0.0282	0.0092	0.0006	33.8123	13.905
	<i>n</i> = 150					
100	0.2247	0.0366	0.1994	0.0105	0.3373	1.536
200	0.2674	0.0316	0.1909	0.0081	0.3705	1.547
500	0.2965	0.0431	0.1632	0.0091	0.8822	2.623
1000	0.2625	0.0600	0.1363	0.0067	7.4969	11.65
2000	0.2354	0.0628	0.1247	0.0089	42.4597	47.525
5000	0.2700	0.1402	0.1075	0.0144	98.8080	75.973
10000	0.4276	0.2247	0.0822	0.0252	150.6061	72.416
20000	0.6025	0.2601	0.0532	0.0280	184.3534	60.144
50000	0.7602	0.2441	0.0275	0.0230	217.9005	41.975

Table 3: Average and Standard Deviation of Mean Absolute Error as Observations Increase (continued)

---

## 297 6 DATA ANALYSIS

298 We performed data analysis on the 18<sup>th</sup> century Chinese novel, *Dream of*  
299 *the Red Chamber*. The observed groups in this dataset do not necessarily  
300 conform to the Hub Model assumption, but we found that, even without this  
301 assumption being explicitly valid, important information about the relation-  
302 ships can be estimated.

303 The Supplemental Materials S3 include two additional data sets estimat-  
304 ing co-sponsorship of legislation in the Senate of the 110<sup>th</sup> United States  
305 Congress and the dispersion of plant species across North America.

306 As noted by Kolaczyk (2009), a significant challenge with estimating the  
307 parameters of implicit networks is that for a given dataset there is usually no  
308 way to verify the extent to which the estimate matches reality. Hence, there  
309 is no “ground truth” or “golden standard” to compare the estimated results  
310 against. Therefore, it is useful to analyze data about which there is some  
311 qualitative knowledge of the relationships between nodes. To this end, we  
312 constructed a dataset of characters from *Dream of the Red Chamber*. Since  
313 novels contain a qualitative social structure that is familiar to readers, the  
314 results of quantitative analysis can be compared to this standard.

315 This novel was chosen for two reasons: the relationships between the

---

316 characters are subtle and complex, and the novel has been carefully studied  
317 by scholars. The story then presents a challenge to estimating relationships  
318 and without a body of knowledge to compare the estimates against.

319 Traditionally datasets are built from novels by carefully reading the text  
320 and identifying dyadic interactions between characters based on criteria es-  
321 tablished by the researchers, e.g., characters  $A$  and  $B$  have a conversation  
322 (MacCarron and Kenna (2013)). This method may construct high quality  
323 datasets, but to identify interactions requires readers who have time to build  
324 them. Since *Dream of the Red Chamber* is written in classical Chinese and  
325 the English translation runs over 2,600 pages, directly generating the dataset  
326 would be excessively time consuming.

327 We built our dataset using text mining and defining a group as characters  
328 who co-occur in the same paragraph. Paragraphs with no characters named  
329 in them were ignored. For a complete description of the text mining protocol,  
330 see Supplemental Materials S5.

331 We analyzed the relationships of 29 important characters. The character  
332 names presented here are based on the original pinyin pronunciations and  
333 the David Hawkes translation (Hawkes (1974)). A Chinese version of the  
334 novel was used for text-mining. The complete novel contains 120 chapters,  
335 but we focused on the first 80 because it is commonly believed that the last

---

336 40 chapters are written by a different author and may not reflect the original  
337 themes of the novel (Hsueh-Chin (2016)). The resulting dataset had 1,389  
338 observations of groups containing at least one of the 29 characters.

339 In Figure 3, the adjacency matrix is represented as an  $n \times n$  grid where  
340 the  $i^{th} \times j^{th}$  cell represents the relationship between nodes  $v_i$  and  $v_j$ . The  
341 relationship strength is represented by the cell's color: nodes with weak  
342 relationships have light cells while nodes with strong relationships have dark  
343 cells. Cells representing relationships of intermediate strength are shaded  
344 along the gray scale.

345 This visualization demonstrates another difference in the performance of  
346 the techniques. The co-occurrence matrix estimates all relationships as being  
347 very weak and it is difficult to differentiate strong relationships from the  
348 absence of a relationship. The half-weight index presents a much stronger  
349 set of relationships but there is evidence of relationships which have been  
350 imputed transitively. In general, HM returns a much sparser network where  
351 relationship strengths demonstrate higher contrast. This tendency towards  
352 sparsity is discussed in more detail in the Supplemental Materials S4.2.

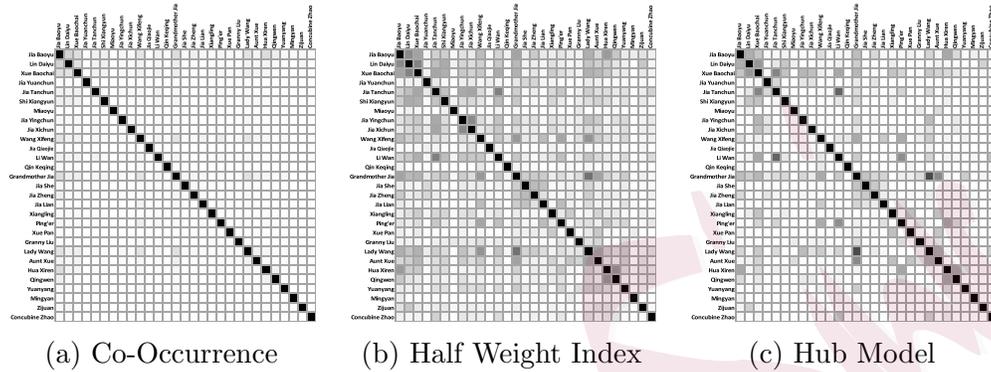


Figure 3: Comparison of Results for *Dream of the Red Chamber*

353 The EM algorithm of HM provides stable solutions. By selecting multiple  
 354 starting points, we find that the adjacency matrix (Figure 3c) is repeatedly  
 355 returned as the most likely parameter of the observed data.

356 The Hub Model parameter's standard deviation was estimated using the  
 357 bootstrap technique. In general, the standard deviation was low. This was  
 358 particularly true for  $\hat{\rho}$  where the maximum standard deviation was 0.0173.  
 359 Table 4 presents the standard deviation of the estimated adjacency matrix  
 at different percentiles.

Percentile	Max	95 %	75 %	Med	25 %	5 %	Min
StDev	0.2696	0.1025	0.0374	0.0100	0.0000	0.0000	0.0000

Table 4: Percentiles of Standard Deviation in  $\hat{A}$  estimated by HM for *Dream of the Red Chamber*

360

361 One of the main themes of *Dream of the Red Chamber* is the love story

---

362 surrounding the protagonist Jia Baoyu (1st character in Figure 3c) and two  
363 potential fiances, the sickly Lin Daiyu (2nd character) and the “ideal” Xue  
364 Baochai (3rd character). Although Jia Baoyu shares a special bond with  
365 Lin Daiyu and has no significant emotional connection to Xue Baochai, he  
366 is ultimately tricked into marrying Xue Baochai (Hsueh-Chin (2016)). In  
367 Table 5, we present the relationships between these two girls and the other  
368 characters as estimated by the co-occurrence matrix, half weight index, and  
369 HM.

370 From the novel, Lin Daiyu is a sensitive girl who prefers to be alone. By  
371 contrast, Xue Baochai is a social and calculating girl. She is extremely good  
372 at interpersonal communication especially with the protagonist’s mother  
373 (Lady Wang) and grandmother (Grandmother Jia) (Hsueh-Chin (2016)).  
374 These different personalities are clearly represented by the HM estimator  
375 while the other estimators do not identify this difference.

	Co-Occurrence Matrix ( $O$ )		Half Weight Index ( $H$ )		Hub ( $\hat{A}$ )	
	Lin Daiyu	Xue Baochai	Lin Daiyu	Xue Baochai	Lin Daiyu	Xue Baochai
Jia Baoyu	0.1728	0.1274	0.4563	0.3587	0.3113	0.2258
Lin Daiyu	1.0000	0.1109	1.0000	0.4866	1.0000	0.4072
Xue Baochai	0.1109	1.0000	0.4866	1.0000	0.4072	1.0000
Jia Yuanchun	0.0072	0.0050	0.0531	0.0449	0.0156	0.0228
Jia Tanchun	0.0439	0.0533	0.2490	0.3482	0.0915	0.4848
Shi Xiangyun	0.0590	0.0490	0.3273	0.3119	0.2194	0.2365
Miaoyu	0.0072	0.0036	0.0552	0.0337	0.0597	0
Jia Yingchun	0.0252	0.0274	0.1667	0.2141	0	0.2846
Jia Xichun	0.0187	0.0202	0.1313	0.1692	0.0102	0.2461
Wang Xifeng	0.0497	0.0526	0.1840	0.2131	0.0317	0.0697
Jia Qiaojie	0.0022	0.0022	0.0170	0.0208	0	0.0348
Li Wan	0.0367	0.0482	0.2086	0.3160	0.0580	0.3384
Qin Keqing	0.0007	0.0007	0.0052	0.0062	0	0
Grandmother Jia	0.0655	0.0648	0.2725	0.2985	0.1925	0.2820
Jia She	0.0065	0.0043	0.0449	0.0357	0	0
Jia Zheng	0.0122	0.0144	0.0701	0.0952	0.0143	0.0174
Jia Lian	0.0072	0.0036	0.0423	0.0245	0.0002	0.0073
Xiangling	0.0180	0.0252	0.1185	0.1961	0.0741	0.2344
Ping'er	0.0122	0.0209	0.0668	0.1306	0.0016	0.1643
Xue Pan	0.0043	0.0101	0.0292	0.0809	0	0
Granny Liu	0.0072	0.0050	0.0493	0.0411	0.0101	0.0113
Lady Wang	0.0490	0.0590	0.2248	0.3037	0.0224	0.2065
Aunt Xue	0.0302	0.0396	0.1806	0.2750	0.0479	0.1657
Hua Xiren	0.0403	0.0389	0.1938	0.2105	0.0283	0.1469
Qingwen	0.0166	0.0115	0.1020	0.0829	0.0155	0.0886
Yuanyang	0.0086	0.0101	0.0556	0.0763	0	0.0430
Mingyan	0.0007	0.0007	0.0053	0.0064	0	0
Zijuan	0.0317	0.0108	0.2184	0.0888	0.1775	0.0376
Concubine Zhao	0.0050	0.0058	0.0361	0.0495	0	0.0338

Table 5: Relationships of Lin Daiyu and Xue Baochai to other characters in *Dream of the Red Chamber*

---

## 376 7 CONCLUSION

377 To the best of our knowledge, Hub Models introduce an innovative approach  
378 to the task of implicit network inference. By defining a model-based gener-  
379 ating mechanism to link the latent network to observed grouped data and  
380 applying an EM algorithm, we are able to estimate the network.

381 Not only are the estimators easy to calculate in a reasonable amount of  
382 time, but they have a practical interpretation. The parameter  $\rho_i$  measures  
383 the probability that node  $v_i$  will form a group.  $A_{ij}$  measures the probability  
384 that a member of the population will be included in a group formed by node  
385  $v_i$ .

386 The Hub Models compare favorably against existing techniques. Since  
387 the co-occurrence matrix and half weight index lack a generating mechanism  
388 to connect them to the observed grouped data, these measures often cannot  
389 detect important features of a network. By applying the Hub Model to the  
390 18<sup>th</sup> century Chinese novel *Dream of the Red Chamber*, we demonstrate that  
391 the HM is able to detect important features in the relationships between  
392 nodes in complex situations.

393 By the standards of statistical network analysis, the size of the adjacency  
394 matrices presented in this paper are small. An important question is how

---

395 the Hub Model would perform with 10,000 or even 1,000,000 nodes. While  
396 it is computationally feasible to apply the Hub Model to populations of this  
397 size, there is a practical challenge of collecting enough observations to have  
398 sufficient statistical power.

399 We observe that how “small” or “large” a dataset is depends on the rela-  
400 tionship between the number of nodes and the number of observed groups. In  
401 principle, if there are  $n$  nodes, the Hub Model must estimate  $n^2$  parameters.  
402 If the number of observations is less than the number of nodes, multiple sets  
403 of parameters have the same likelihood and parameter estimation is unstable.  
404 In general, it is only when the number of observations exceeds the square of  
405 the number of nodes, that we have stable estimates.

406 This means that to estimate the Hub Model parameters of a popula-  
407 tion with hundreds of thousands of nodes, would require tens of billions of  
408 observations. Therefore, applying Hub Models directly to text or even a  
409 recommender system would be impractical.

410 In order to make the Hub Model useful for such large populations, some  
411 technique must be applied to reduce the number of parameters in the model.  
412 In this paper, we have placed no restrictions on the adjacency matrix. How-  
413 ever, there are a number of restrictions which could be applied to enable us  
414 to handle populations with “small” datasets.

---

415 One way is to make an assumption about the structure of the underlying  
416 network. For example, one might assume that the latent network is itself  
417 the result of a block model or exponential random graph model. Such an  
418 approach would create a hierarchical model for group formation.

419 A second way that assumptions about the structure of the underlying  
420 network could be applied is to change the dimensions of the adjacency matrix.  
421 In doing this, researchers may limit the number of nodes which can act as  
422 leaders or treat some nodes as having the same behavior.

423 The Hub Model can potentially be useful to model the term-document  
424 matrix in text mining. Such a matrix describes the frequency of terms that  
425 occur in a collection of documents, which is similar to the format of group  
426 data. Many text mining techniques are based on a co-occurrence matrix  
427 created from the term-document matrix. The Hub Model may provide more  
428 meaningful estimates of the relations between terms.

## 429 **Supplementary Materials**

430 The supplemental materials contain additional details regarding the proof  
431 of Theorem 1, calculation of the estimating equations 3.5 and 3.6. Addition-  
432 ally, we provide data analysis for co-sponsorship of the 110<sup>th</sup> Congress and  
433 a dataset of North American flora. We conclude with a discussion of iden-

434 tifiability, self-sparsity, and the protocol for text mining *Dream of the Red*  
435 *Chamber*.

### 436 **Acknowledgements**

437 This work is partially supported by NSF DMS 1513004.

### 438 **Author's Statement**

439 The views expressed in this paper are those of the authors and do not  
440 reflect the official policy or position of the US Army, the Department of  
441 Defense, or the US Government.

### 442 **References**

443 Airoldi, E. M., Blei, D. M., Fienberg, S. E., and Xing, E. P. (2008). Mixed membership stochastic  
444 blockmodels. *Journal of Machine Learning Research*, 9(Sep):1981–2014.

445 Anandkumar, A., Foster, D. P., Hsu, D., Kakade, S. M., and Liu, Y. (2015). A spectral algorithm  
446 from latent dirichlet allocation. *Algorithmica*, 72(1):193–214.

447 Bejder, L., Fletcher, D., and Brager, S. (1998). A method for testing association patterns of social  
448 animals. *Animal Behavior*, 56:719–725.

449 Brent, L. J. N., Lehmann, J., and Ramos-Fernandez, G. (2011). Social network analysis in the  
450 study of nonhuman primates: A historical perspective. *American Journal of Primatology*,

## REFERENCES

---

- 451       73:720–730.
- 452 Cairns, S. J. and Schwager, S. J. (1987). A comparison of association indices. *Animal Behavior*,
- 453       35.
- 454 Carreira-Perpinan, M. A. and Renals, S. (2000). Practical identifiability of finite mixtures of
- 455       multivariate bernoulli distributions. *Neural Computation*, 12:141–152.
- 456 Choudhury, M., M., W. A., Hofman, J. M., and Watts, D. J. (2010). Inferring relevant social
- 457       networks from interpersonal communication. *International World Wide Web Conference*
- 458       *Committee*.
- 459 Colace, F., De Santo, M., Greco, L., Moscato, V., and Picariello, A. (2015). A collaborative
- 460       user-centered framework for recommending items in online social networks. *Computers in*
- 461       *Human Behavior*.
- 462 Dice, L. R. (1945). Measures of the amount of ecological association between species. *Ecology*,
- 463       26:297–302.
- 464 Frank, O. and Strauss, D. (1986). Markov graphs. *Journal of the American Statistical Associa-*
- 465       *tion*, 81:832–842.
- 466 Freeman, L. C., White, D. R., and Romney, A. K. (1989). *Research Methods in Social Network*
- 467       *Analysis*. George Mason University Press.
- 468 Goldenberg, A., Zheng, A. X., Fienberg, S. E., and Airoldi, E. M. (2010). A survey of statistical
- 469       network models. *Foundations and Trends in Machine Learning*, 2:129–233.

---

## REFERENCES

- 470 Handcock, M. D., Raftery, A. E., and Tantrum, J. M. (2007). Model-based clustering for social  
471 networks. *J. R. Statist. Soc. A*, 170:301–354.
- 472 Hawkes, D. (1974). *The Story of the Stone, or The Dream of the Red Chamber, Vol. 1: The*  
473 *Golden Days*. Penguin Classics.
- 474 Hiller, F. S. and Lieberman, G. L. (2001). *Introduction to Operations Research*. McGraw-Hill.
- 475 Hoff, P. D., Raftery, A. E., and Handcock, M. S. (2002). Latent space approaches to social  
476 network analysis. *Journal of the American Statistical Association*, 97:1090–1098.
- 477 Holland, P. W., Laskey, K. B., and Leinhardt, S. (1983). Stochastic blockmodels: first steps.  
478 *Social Networks*, 5(2):109–137.
- 479 Hsueh-Chin, T. (2016). *CliffsNotes: Dream of the Red Chamber*. Houghton Mifflin Harcourt.
- 480 Jackson, M. O. (2010). *Social and economic networks*. Princeton university press.
- 481 Kolaczyk, E. D. (2009). *Statistical Analysis of Network Data: Methods and Models*. Springer.
- 482 Liu, Y., Cheng, J., Yan, C., Wu, X., and Chen, F. (2015). Research on the matthews correla-  
483 tion coefficients metrics of personalized recommendation algorithm evaluation. *International*  
484 *Journal of Hybrid Information Technology*, 8(1):163–172.
- 485 MacCarron, P. and Kenna, R. (2013). Viking sagas: Six degrees of icelandic separation-social  
486 networks from the viking era. *Significance*, pages 12–17.
- 487 McLachlan, G. J. and Krishnan, T. (2008). *The EM Algorithm and Extensions*. John Wiley and  
488 Sons, Inc.

## REFERENCES

---

- 489 Newman, M. E. J. (2011). *Networks: An Introduction*. Oxford University Press.
- 490 Rabbat, M., Figueiredo, M., and Nowak, R. (2008). Network inference from co-occurrences. *IEEE*  
491 *Transactions on Information Technology*, 54(9):4053–4068.
- 492 Robins, G., Pattison, P., Kalish, Y., and Lusher, D. (2007). An introduction to exponential  
493 random graph ( $p^*$ ) models for social networks. *Social networks*, 29(2):173–191.
- 494 Teicher, H. (1961). Identifiability of mixtures. *The Annals of Mathematical Statistics*, 32(1):244–  
495 248.
- 496 Voelkl, B., Kasper, C., and Schwab, C. (2011). Network measures for dyadic interactions: Stability  
497 and reliability. *American Journal of Primatology*, 73:731–740.
- 498 Vretos, N., Nikolaidis, N., and Pitas, I. (2012). Video fingerprinting using latent dirichlet alloca-  
499 tion and facial images. *Pattern Recognition*, 45(7):2489–2498.
- 500 Wasserman, S. and Faust, C. (1994). *Social Network Analysis: Methods and Applications*. Cam-  
501 bridge University Press.
- 502 Zachary, W. W. (1977). An information flow model for conflicts and fission in small groups.  
503 *Journal of Anthropological Research*, 33:452–473.
- 504 Department of Statistics  
505 George Mason University  
506 4400 University Drive, MS 4A7  
507 Fairfax, VA 22030-4444

---

## REFERENCES

- 508 E-mail: (yzhao15@gmu.edu)
- 509 United States Army
- 510 1400 Defense Pentagon
- 511 Washington, DC 20301
- 512 E-mail: (charles.w.weko.mil@mail.mil)