
Leveraging Node Attributes for Incomplete Relational Data

He Zhao¹ Lan Du¹ Wray Buntine¹

Abstract

Relational data are usually highly incomplete in practice, which inspires us to leverage side information to improve the performance of community detection and link prediction. This paper presents a Bayesian probabilistic approach that incorporates various kinds of node attributes encoded in binary form in relational models with Poisson likelihood. Our method works flexibly with both directed and undirected relational networks. The inference can be done by efficient Gibbs sampling which leverages sparsity of both networks and node attributes. Extensive experiments show that our models achieve the state-of-the-art link prediction results, especially with highly incomplete relational data.

1. Introduction

Relational learning from network data, particularly with probabilistic methods, has gained a wide range of applications such as social network analysis (Xiang et al., 2010), recommender systems (Gopalan et al., 2014b), knowledge graph completion (Hu et al., 2016b), and bioinformatics (Huopaniemi et al., 2010). Generally speaking, the goal of relational learning is to discover and analyse latent clusters of entities (i.e., community detection), and predict missing links (i.e., link prediction).

The standard approach for modelling relational data is latent factor analysis via matrix factorisation and its variations. Among the existing approaches, Non-negative Matrix Factorisation (NMF) and the Stochastic Block Model (SBM) are prominent foundational methods. NMF is usually used to model relationships between two sets of entities such as users and movies in collaborative filtering (Mnih & Salakhutdinov, 2008). While developed independently, SBM (Wang & Wong, 1987; Nowicki & Snijders, 2001) can be viewed as an extension of NMF that introduces

a block matrix to capture the interactions between latent factors. There have been many Bayesian extensions of these two methods, relaxing the assumptions and/or introducing extra components, such as the Infinite Relational Model (IRM) (Kemp et al., 2006), the mixture membership stochastic block model (MMSB) (Airoldi et al., 2008), and the non-parametric latent feature models (NLFM) (Miller et al., 2009). Poisson Factorisation (PF) (Dunson & Herzing, 2005; Zhou et al., 2012), is a popular version of NMF which models count data with convenient statistical properties (Gopalan et al., 2014b; 2015). Combining the ideas of PF and SBM, the infinite Edge Partition Model (EPM) (Zhou, 2015) and its extensions (Hu et al., 2016b) have proven successful for relational networks.

When a network has less data, relational learning becomes more difficult. One extreme case is the *cold-start* problem (Lin et al., 2013; Sedhain et al., 2014; Zhang & Wang, 2015), where a node has no observed links, making suggestion of links for that node even more challenging. In such cases, it is natural to appeal to side information such as node attributes or features. For instance, papers in citation networks are often associated with categories and authors, and users in Facebook or Twitter are often asked to provide information such as age, gender and interests. It is reasonable to assume that nodes having similar attributes are more likely to relate to each other (i.e., homophily, Nickel et al., 2016). Thus, node attributes serve as important complementary information to relational data.

There are few Bayesian probabilistic relational models that are able to leverage side information. For example, NLFM uses a linear regression model to transform the features of each node into a single number, which contributes to link probabilities. However, side information in NLFM cannot directly influence the latent factors, which gives little support for community detection. As an extension of MMSB, the Non-parametric Meta-data Dependent Relational (NMDR) model (Kim et al., 2012) incorporates attributes into the mixed-membership distribution of each node with the logistic-normal transform, which results in non-conjugacy for inference. Fan et al. (2016) further developed this idea in the Node information Involved Mixture Membership model (niMM), where side information is integrated in a conjugate way. Although these models demonstrate improvement using side information, they

¹Faculty of Information Technology, Monash University, Australia. Correspondence to: He Zhao <he.zhao@monash.edu>.

scale quadratically in the number of nodes and the incorporation of side information is often complicated.

Several recent methods (Gopalan et al., 2014a; Acharya et al., 2015; Hu et al., 2016a) extend PF with side information using the additivity of the Poisson and gamma distributions/processes. With improved scalability, the Structural Side Information Poisson Factorisation (SSI-PF) (Hu et al., 2016a) models directed unweighted networks with node labels, such as citation networks with papers labelled with one of several categories. However, its performance remains untested when a node has multiple attributes. Moreover, undirected networks are not handled by SSI-PF.

In this paper we present the Node Attribute Relational Model (NARM)¹, a fully Bayesian approach that models large, sparse, and unweighted relational networks with arbitrary node attributes encoded in binary form. It works with Poisson gamma relational models to incorporate side information. Specifically, we propose the Symmetric NARM (Sym-NARM) for undirected networks, an extension of EPM (Zhou, 2015) and the Asymmetric NARM (Asym-NARM) for directed networks, an extension of PF (Zhou et al., 2012). The proposed models have several key properties: **(1) Effectively modelling node attributes:** the proposed models are able to achieve improved link prediction performance, especially where training data are limited. **(2) Fully Bayesian and conjugate:** the inference is done by efficient, closed-form Gibbs sampling which scales linearly in the number of observed links and takes advantage of the sparsity of node attributes. It makes our models scalable for large but sparse relational networks with large sets of node attributes. **(3) Flexibility:** the proposed models work on directed and undirected relational networks with flat and hierarchical node attributes.

2. The Node Attribute Relational Model

Here we focus on modelling unweighted networks that can be either directed (i.e., the relationship is asymmetric) or undirected. Assume a relational network with N nodes is stored in a binary adjacency matrix $\mathbf{Y} \in \{0, 1\}^{N \times N}$ where $y_{i,j} = 1$ indicates the presence of a link between nodes i and j . If the relationship described in the network is symmetric, then $y_{i,j} = y_{j,i}$, and if asymmetric, possibly $y_{i,j} \neq y_{j,i}$. Node attributes are encoded in a binary matrix $\mathbf{F} \in \{0, 1\}^{N \times L}$, where L is the total number of attributes. Attribute $f_{i,l} = 1$ indicates attribute l is active with node i and vice versa. Although our models incorporate binary attributes, categorical attributes and real-valued attributes can be converted into binary values with proper transformations (Kim et al., 2012; Fan et al., 2016; Hu et al., 2016a).

¹Code available at <https://github.com/ethanhezha0/NARM/>

2.1. The Symmetric Node Attribute Relational Model

Sym-NARM works with undirected networks. Its generative process is shown in Figure 1. Instead of modelling the binary matrix \mathbf{Y} directly, it applies the Bernoulli-Poisson link (BPL) function (Zhou, 2015) using an underlying latent count matrix \mathbf{X} . One first draws a latent count $x_{i,j}$ from the Poisson distribution and then thresholds it at 1 to generate a binary value $y_{i,j}$. This is shown in Eqs. (1)-(3). Analysed in (Zhou, 2015; Hu et al., 2016b;a), BPL has the appealing property that if $y_{i,j} = 0$, then $x_{i,j} = 0$ with probability one. Thus, only non-zeros in \mathbf{Y} need to be sampled, giving huge computational savings for large sparse networks, illustrated in Section 3 and Section 5.4.

The latent matrix \mathbf{X} is further factorised into K latent factors with a non-negative bilinear model: $\mathbf{X} \sim \text{Poi}(\Phi \mathbf{\Lambda} \Phi^T)$ where $\Phi \in \mathbb{R}_+^{N \times K}$ and $\mathbf{\Lambda} \in \mathbb{R}_+^{K \times K}$. Φ is referred to as the *node factor loading matrix* where $\phi_{i,k}$ models the strength of the connection between node i and latent factor k . As in SBM, the correlations of the latent factors are modelled in a symmetric matrix $\mathbf{\Lambda}$, referred to as the *block matrix*. Following (Zhou, 2015), we draw $\mathbf{\Lambda}$ from a hierarchical relational gamma process (implemented with truncation as a vector of gamma variables), shown in Eqs. (8) and (9).

One appealing aspect of our model is the incorporation of node attributes on the prior of $\phi_{i,k}$ (i.e., $g_{i,k}$). Shown in Eq. (5), $g_{i,k}$ is constructed with a log linear combination of $f_{i,l}$. $h_{l,k}$ is referred to as the k^{th} *attribute factor loading* of attribute l , which influences $g_{i,k}$ iff attribute l is active with node i (i.e., $f_{i,l} = 1$). b_k acts as an attribute-free bias for each latent factor k . $h_{l,k}$ and b_k are gamma distributed with mean 1, hence if attribute l does not contribute to latent factor k or is less useful, $h_{l,k}$ is expected to be near 1 and to have little influence on $g_{i,k}$. The hyper-parameter μ_0 controls the variation of $h_{l,k}$.

The intuition of our model is: if two nodes have more common attributes, their gamma shape parameters will be more similar, with similar node factor loadings, resulting in a larger probability that they relate to each other. Moreover, instead of incorporating the node attributes directly into the node factor loadings, Sym-NARM uses them as the prior information using Eq. (4), which results in a principled way of balancing the side information and the network data. In addition, different attributes can contribute differently to the latent factors. For example, the gender of an author may be much less important to co-authorship with others than the research fields. This is controlled by the attribute factor loading $h_{l,k}$ in our model.

2.2. The Asymmetric Node Attribute Relational Model

Extending the Beta Gamma Gamma Poisson factorisation (BGGPF) (Zhou et al., 2012), Asym-NARM works on di-

$$y_{i,j} = \mathbf{1}_{(x_{i,j} > 0)} \quad (1)$$

$$x_{i,j} = \sum_{k_1, k_2=1}^K x_{i, k_1, k_2, j} \quad (2)$$

$$x_{i, k_1, k_2, j} \sim \text{Poi}(\phi_{i, k_1} \lambda_{k_1, k_2} \phi_{j, k_2}) \quad (3)$$

$$\phi_{i, k} \sim \text{Ga}(g_{i, k}, 1/c_i) \quad (4)$$

$$g_{i, k} = b_k \prod_{l=1}^L h_{l, k}^{f_{i, l}} \quad (5)$$

$$h_{l, k} \sim \text{Ga}(\mu_0, 1/(1/\mu_0)) \quad (6)$$

$$b_k \sim \text{Ga}(\mu_0, 1/(1/\mu_0)) \quad (7)$$

$$\lambda_{k_1, k_2} \sim \begin{cases} \text{Ga}(\epsilon r_k, 1/a_0), & \text{if } k_1 = k_2 = k \\ \text{Ga}(r_{k_1} r_{k_2}, 1/a_0), & \text{otherwise} \end{cases} \quad (8)$$

$$r_k \sim \text{Ga}(\gamma_0/K, 1/c_0) \quad (9)$$

Figure 1. The generative model of Sym-NARM. $\mathbf{1}_{(\cdot)}$ is the indicator function. $\text{Poi}(\cdot)$ and $\text{Ga}(\cdot, \cdot)$ stand for the Poisson distribution and the gamma distribution respectively. Conjugate gamma priors are imposed on the hyper-parameters: γ_0 , ϵ , c_0 , c_i , and a_0 .

rected relational networks with node attributes incorporated in a similar way to Sym-NARM. Figure 2 shows its generative process. Here the latent count matrix \mathbf{X} is factorised as $\mathbf{X} \sim \text{Poi}(\Phi \Theta)$, where $\Phi \in \mathbb{R}_+^{N \times K}$ and $\Theta \in \mathbb{R}_+^{K \times N}$ are referred to as the *factor loading matrix* and the *factor score matrix* respectively. Similar to SSI-PF, the node attributes are incorporated on the prior of Φ .

2.3. Incorporating Hierarchical Node Attributes

Relational networks can be associated with hierarchical side information (Hu et al., 2016a). For example, in a patent citation network, patents can be labelled with the International Patent Classification (IPC) code, which is a hierarchy of patent categories and sub-categories. Suppose the second level attributes are stored in a binary matrix $\mathbf{F}' \in \{0, 1\}^{L \times M}$ where M is the number of attributes in the second level. Our models can be used to incorporate hierarchical node attributes via a straightforward extension: replace hyper-parameter μ_0 in Eq. (6) with $\mu_{l, k} = \prod_{m=1}^M \delta_{m, k}^{f'_{l, m}}$. This extension mirrors what is done for first level attributes.

3. Inference with Gibbs Sampling

Both Sym-NARM and Asym-NARM enjoy local conjugacy so the inference of all latent variables can be done by closed-form Gibbs sampling. Moreover, the inference only needs to be conducted on the non-zero entries in \mathbf{Y} and \mathbf{F} . This section focuses on the sampling of $h_{l, k}$ (b_k), the key variable in the proposed incorporation of node attributes. The sampling of the other latent variables is similar to those in EPM and BGGPF, detailed in (Zhou, 2015;

$$y_{i,j} = \mathbf{1}_{(x_{i,j} > 0)} \quad (10)$$

$$x_{i,j} \sim \sum_k x_{i,j,k} \quad (11)$$

$$x_{i,j,k} \sim \text{Poi}(\phi_{i,k} \theta_{j,k}) \quad (12)$$

$$\phi_{i,k} \sim \text{Ga}\left(g_{i,k}, \frac{q_k}{1-q_k}\right) \quad (13)$$

$$q_k \sim \text{Be}(c_0 \epsilon, c_0(1-\epsilon)) \quad (14)$$

$$g_{i,k} = b_k \prod_{l=1}^L h_{l,k}^{f_{i,l}} \quad (15)$$

$$h_{l,k} \sim \text{Ga}(\mu_0, 1/(1/\mu_0)) \quad (16)$$

$$b_k \sim \text{Ga}(\mu_0, 1/(1/\mu_0)) \quad (17)$$

$$\theta_{:,k} \sim \text{Dir}_N(a_0 \mathbf{1}) \quad (18)$$

Figure 2. The generative model of Asym-NARM. $\text{Dir}_N(\cdot)$ and $\text{Be}(\cdot, \cdot)$ stand for the N dimensional Dirichlet distribution and the beta distribution respectively. μ_0 , ν_0 , a_0 , e_0 , f_0 , c_0 , ϵ are the hyper-parameters.

Zhou et al., 2012). As the sampling for $h_{l, k}$ is analogous in Sym-NARM and Asym-NARM, our introduction will be based on Asym-NARM alone.

With the Poisson gamma conjugacy, the likelihood for $g_{i, k}$ with $\phi_{i, k}$ marginalised out is:

$$p(g_{i, k} | x_{i, \cdot, k}) \propto (1 - q_k)^{g_{i, k}} \frac{\Gamma(g_{i, k} + x_{i, \cdot, k})}{\Gamma(g_{i, k})} \quad (19)$$

where $x_{i, \cdot, k} = \sum_j x_{i, j, k}$ and $x_{i, j, k}$ is the latent count. The gamma ratio in Eq. (19), i.e., the Pochhammer symbol for a rising factorial, can be augmented with an auxiliary variable $t_{i, k}$: $\frac{\Gamma(g_{i, k} + x_{i, \cdot, k})}{\Gamma(g_{i, k})} = \sum_{t_{i, k}=0}^{x_{i, \cdot, k}} S_{t_{i, k}}^{x_{i, \cdot, k}} g_{i, k}^{t_{i, k}}$ where S_t^x indicates an unsigned Stirling number of the first kind (Chen et al., 2011; Teh et al., 2012; Zhou & Carin, 2015).

Taking $\mathcal{O}(x_{i, \cdot, k})$, $t_{i, k}$ can be directly sampled by a Chinese Restaurant Process with $g_{i, k}$ as the concentration and $x_{i, \cdot, k}$ as the number of customers:

$$t_{i, k} \leftarrow t_{i, k} + \text{Bern}\left(\frac{g_{i, k}}{g_{i, k} + i'}\right) \text{ for } i' = 1 : x_{i, \cdot, k} \quad (20)$$

where $\text{Bern}(\cdot)$ is the Bernoulli distribution. Alternatively, for large $x_{i, \cdot, k}$, because the standard deviation of $t_{i, k}$ is $\mathcal{O}(\sqrt{\log x_{i, \cdot, k}})$ (Buntine & Hutter, 2012), one can sample $t_{i, k}$ in a small window around the current value (Du et al., 2010).

With the above augmentation and Eq. (15), we get:

$$p(\mathbf{G}, \mathbf{H} | x_{:, \cdot, :}, \mathbf{T}, \mathbf{F}) \propto \quad (21)$$

$$\prod_{i=1}^N \prod_{k=1}^K S_{t_{i, k}}^{x_{i, \cdot, k}} e^{-\log\left(\frac{1}{1-q_k}\right) g_{i, k}} \cdot \prod_{l=1}^L \prod_{k=1}^K h_{l, k}^{\sum_{i=1}^N f_{i, l} t_{i, k}}$$

Recall that all the attributes are binary and $h_{l, k}$ influences $g_{i, k}$ only when $f_{i, l} = 1$. Extracting all the terms related to

$h_{l,k}$ in Eq. (21), we get the likelihood of $h_{l,k}$:

$$p\left(h_{l,k} \mid \frac{g_{i,k}}{h_{l,k}}, t_{:,k}, f_{:,l}\right) \propto e^{-h_{l,k} \log\left(\frac{1}{1-q_k}\right) \sum_{i=1:f_{i,l}=1}^N \frac{g_{i,k}}{h_{l,k}} h_{l,k} \sum_{i=1}^N f_{i,l} t_{i,k}} \quad (22)$$

where $\frac{g_{i,k}}{h_{l,k}}$ is the value of $g_{i,k}$ with $h_{l,k}$ removed when $f_{i,l} = 1$. The likelihood function above is in a form that is conjugate to the gamma prior. Therefore, it is straightforward to yield the following sampling strategy for $h_{l,k}$:

$$h_{l,k} \sim \text{Ga}(\mu', 1/\nu') \quad (23)$$

$$\mu' = \mu_0 + \sum_{i=1:f_{i,l}=1}^N t_{i,k} \quad (24)$$

$$\nu' = 1/\mu_0 - \log(1 - q_k) \sum_{i=1:f_{i,l}=1}^N \frac{g_{i,k}}{h_{l,k}} \quad (25)$$

Precomputed with Eq. (15), $g_{i,k}$ can be updated with Eq. (26), after $h_{l,k}$ is sampled.

$$g_{i,k} \leftarrow \frac{g_{i,k} h'_{l,k}}{h_{l,k}} \text{ for } i = 1 : N \text{ and } f_{i,l} = 1 \quad (26)$$

where $h'_{l,k}$ is the newly sampled value of $h_{l,k}$.

To compute Eqs. (24)-(26), we only need to iterate over the nodes that attribute l is active with (i.e., $f_{i,l} = 1$). Thus, the sampling for \mathbf{H} takes $\mathcal{O}(D'KL)$ where D' is the average number of nodes that an attribute is active with. This demonstrates how the sparsity of node attributes is leveraged. As the mean of $x_{i,\cdot,k}$ is D/K , sampling the tables $\mathbf{T} \in \mathbb{N}^{N \times K}$ takes $\mathcal{O}(ND)$ which can be accelerated with the window sampling technique explained above.

We show the computational complexity of our and related models in Table 1. The empirical comparison of running speed is in Section 5.4. By taking advantage of both network sparsity and node attribute sparsity, our models are more efficient than the competitors, especially on large sparse networks with large sets of attributes.

4. Related work

Compared with the node-attribute models such as NMDR and niMM whose methods result in complicated inference, our Sym-NARM is much more efficient on large sparse networks, illustrated in Table 1.

The most closely related model to our Asym-NARM, also extending the BGGPF algorithm, is SSI-PF. But it uses the gamma additivity to construct the prior of node factor loadings with the sum of attribute factor loadings. Our model has several advantages over SSI-PF: (1) The derivation of Gibbs sampling of SSI-PF requires that each column of Θ is normalised (Eq. (18)). This limits the application of SSI-PF to other models such as EPM which is an unnormalised model. (2) Shown in Table 1, Asym-NARM enjoys more efficient computational complexity. (3) Shown

Table 1. The computational complexity for the compared models. N : number of nodes. K : number of latent factors. L : number of node attributes. D : the average degree (number of edges) per node ($D \ll N$ in sparse networks). D' : the average number of nodes that an attribute is active with (usually, $D' < N$). For the models that incorporate node attributes (marked with a *), the complexity with one level attributes is shown.

Model	Complexity
Models with the block matrix	
*NMDR (Kim et al., 2012)	$\mathcal{O}(N^2K + NKL)$
*niMM (Fan et al., 2016)	$\mathcal{O}(N^2K^2 + NKL)$
EPM (Zhou, 2015)	$\mathcal{O}(NK^2D)$
*Sym-NARM	$\mathcal{O}(NK^2D + D'KL)$
Models without the block matrix	
BGGPF (Zhou et al., 2012)	$\mathcal{O}(NKD)$
*SSI-PF (Hu et al., 2016a)	$\mathcal{O}(NKDL)$
*Asym-NARM	$\mathcal{O}(NKD + D'KL)$

in Section 5, our model is more effective especially when a node has multiple attributes.

There are also models that extend PF and collective matrix factorisation (Singh & Gordon, 2008) to jointly factorise relational networks and document-word matrices such as (Gopalan et al., 2014a; Zhang & Wang, 2015; Acharya et al., 2015). Our NARM models incorporate general node attributes (not only texts) as the priors of the factor loading matrix in a supervised manner, rather than jointly modelling the side information in an unsupervised manner.

Another related area is supervised topic models such as (Mcauliffe & Blei, 2008; Ramage et al., 2009; Lim & Buntine, 2016). The Dirichlet Multinomial Regression (DMR) model (Mimno & McCallum, 2012) is the most related one to ours. It models document attributes on the priors of the topic proportions with the logistic-normal transform. For comparison, we propose DMR-MMSB, extending MMSB with the DMR technique to incorporate side information on the mixed-membership distribution of each node.

5. Experiments

In this section we evaluate Sym-NARM and Asym-NARM with a set of the link prediction tasks on 10 real-world relational datasets with different sizes and various kinds of node attributes. We compare our models with the state-of-the-art relational models, demonstrating that our models outperform the competitors on those datasets in terms of link prediction performance and per-iteration running time. We report the average area under the curve of both the receiver operating characteristic (AUC-ROC) and precision recall (AUC-PR) for quantitatively analysing the models. Moreover, we perform qualitative analysis by comparing the link probabilities estimated by the compared models.

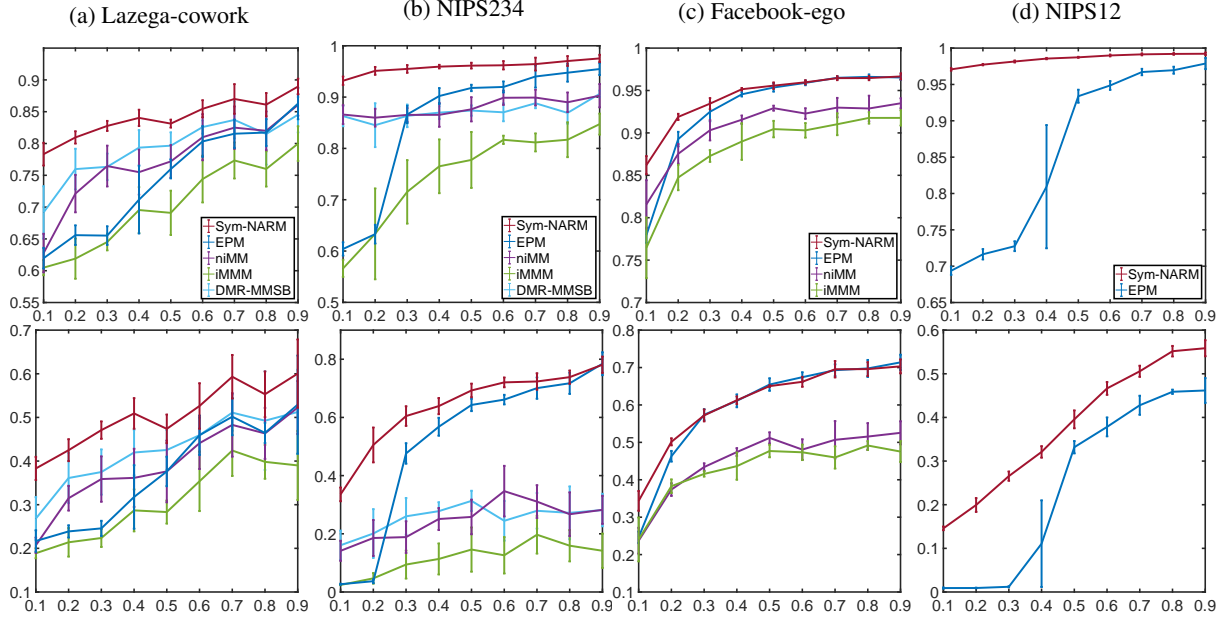


Figure 3. The AUC-ROC (the first row) and AUC-PR (the second row) scores on the undirected networks. The values on the horizontal axis are the proportions of the training data and each of the error bars is the standard deviation over the five random splits for one proportion. DMR-MMSB achieves its best performance at $K = 5$ and 10 on Lazega-cowork and NIPS234 respectively.

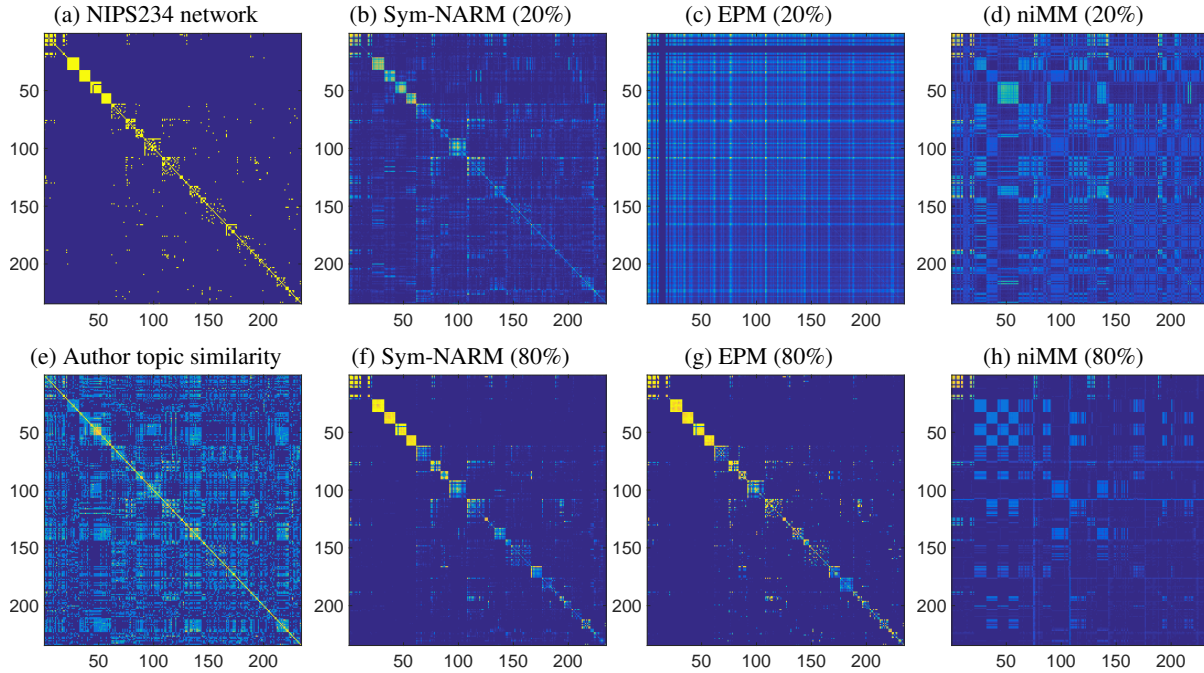


Figure 4. The link probability estimations in NIPS234. Similar to (Zhou, 2015), the nodes are reordered to make a node with a larger index belong to the same or a smaller-size community, where the disjoint community assignments are obtained by analysing the results of Sym-NARM. (a) The original NIPS234 network. (e) The topic similarity of the authors, obtained by the pairwise cosine distances of the topic proportions, with a brighter colour representing a closer distance. (b)-(d) and (f)-(h) Estimated link probabilities with 20% and 80% training data respectively for each compared model.

5.1. Link Prediction on Undirected Networks

For the link prediction task on undirected network data, we compared our **Sym-NARM** with two models that do

not consider node attributes, **EPM** (Zhou, 2015), a state-of-the-art relational model, and **iMM** (Koutsourelakis & Eliassi-Rad, 2008), a non-parametric version of MMSB,

and two node attribute models, **niMM** (Fan et al., 2016), a non-parametric relational model which has been demonstrated to outperform NMDR (Kim et al., 2012), and **DMR-MMSB**, our extension to MMSB using the Dirichlet Multinomial Regression (Mimno & McCallum, 2012). Sym-NAMR was implemented in MATLAB on top of the EPM code and we used the code released by the original authors for EPM and niMM. iMMM was implemented by Fan et al. (2016) as a variant of niMM.

The description of the four datasets used is given below:

- **Lazega-cowork:** This dataset (Lazega, 2001) contains 378 links of the co-work relationship among 71 attorneys. Each attorney is associated with attributes such as gender, office location, and age. After discretisation and binarisation, we derived a 71×18 binary node attribute matrix with 497 non-zero entries.
- **NIPS234:** This is a co-author network of the 234 authors with 598 links extracted from NIPS 1-17 conferences (Zhou, 2015). We merged all the papers written by the same author as a document, and then trained a LDA model with 100 topics. The 5 most frequent topics were used as the attributes, which gives us a 234×100 attribute matrix with 1170 non-zero entries.
- **Facebook-ego:** The original dataset (McAuley & Leskovec, 2012) was collected from survey participants of Facebook users. Out of the 10 circles (i.e., friend lists), we used the first circle that contains 347 users with 2519 links. Each user is associated with 227 binary attributes, encoding side information such as age, gender, and education. We got a 347×227 binary node attribute matrix with 3318 non-zero entries.
- **NIPS12:** NIPS12 was collected from NIPS papers in vols 0-12. It is a median-size co-author network with 2037 authors and 3134 links. Similar to NIPS234, we used the 5 most frequent topics as the attributes for each author. We got a 2037×100 binary node attribute matrix with 10185 non-zero entries.

5.1.1. EXPERIMENTAL SETTINGS

For each dataset, we varied the training data from 10% to 90% and used the remaining in testing. For each proportion, to generate five random splits, we used the code in the EPM package (Zhou, 2015) which splits a network in terms of its nodes. The reported AUC-ROC/PR scores were averaged over the five splits. We used the default hyper-parameter settings enclosed in the released code for EPM, niMM and iMMM. For our Sym-NARM, we set $\mu_0 = 1$ and all the other hyper-parameters the same as those in EPM. Note that the models in comparison except DMR-MMSB are non-parametric models. For Sym-NARM and EPM, we set the truncation level large enough for each dataset: $K_{max} = 50, 100, 256$ for Lazega-

cowork, Facebook-ego and NIPS234, NIPS12 respectively. For DMR-MMSB, we varied K in $\{5, 10, 25, 50\}$ and reported the best one. Following (Zhou, 2015), we used 3000 MCMC iterations and computed AUC-ROC/PR with the average probability over the last 1500. The performance of iMMM and niMM on NIPS12 and DMR-MMSB on Facebook-ego and NIPS12 are not reported as the datasets are too large for them given our computational resources.

5.1.2. RESULTS

The AUC-ROC/PR scores are reported in Figure 3. Overall, our Sym-NARM model performs significantly better than niMM, iMMM, and DMR-MMSB on all the datasets, and EPM on 3 datasets (except Facebook-ego with large training proportions). It is interesting that the performance of EPM on Facebook-ego gradually approaches ours when more than 30% training data were used. Note that Facebook-ego is much denser than the others, which means the network information itself could be rich enough for EPM to reconstruct the network and the node attributes contribute less. However in general, when relational data are highly incomplete (with less training data), our model is able to achieve improved link prediction performance.

To illustrate how side information helps, we qualitatively compared our model with EPM and niMM by estimating the link probabilities on NIPS234, shown in Figure 4. With 20% training data, EPM does not give a meaningful reconstruction of the original network, but it starts to with more data presented. The similarity of the authors' topics in Figure 4e matches the original network, demonstrating the usefulness of the topics, but with some error. Using the topics as the authors' attributes, our Sym-NARM achieves reasonably good reconstruction of the network with only 20% training data, further improving with 80% training data. Although niMM uses the same node attributes, its performance is not as good and is even outperformed by EPM with 80% training data.

5.2. Link Prediction on Directed Networks

Here we compared our **Asym-NARM** (implemented in MATLAB on top of the BGGPF code) with two models that do not consider node attributes, **BGGPF** (Zhou et al., 2012) and **iMMM**, and three node-attribute models, **niMM**, **SSI-PF** (Hu et al., 2016a) and **DMR-MMSB**. We used the following four datasets:

- **Lazega-advice:** This dataset is a directed network with 892 links of the advice relation among the attorneys. The node attributes are the same as in Lazega-cowork.
- **CiteSeer:** This dataset² contains a citation network with

²<http://lings.umiacs.umd.edu/projects/projects/lbc/index.html>

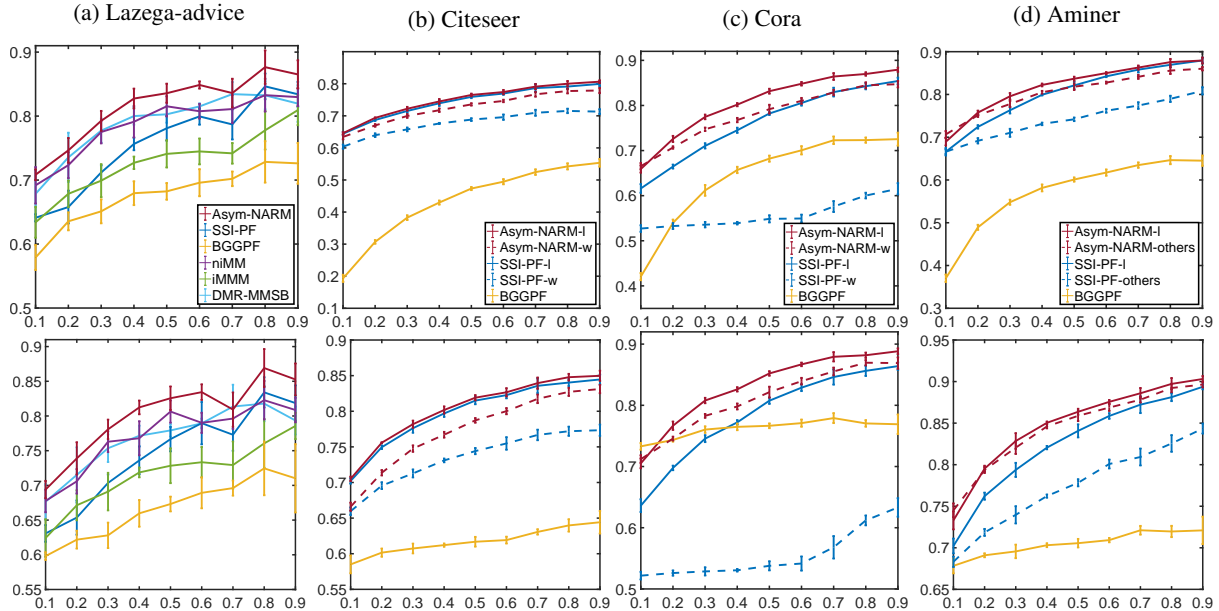


Figure 5. The AUC-ROC (the first row) and AUC-PR (the second row) scores on the directed networks. The models with “-l” and “-w” use the labels and the words as attributes respectively. The models with “-others” in Aminer use the extra attributes. DMR-MMSB achieves its best performance at $K = 10$ on Lazega-advice.

4591 links of 3312 papers, labelled with one of 6 categories. For each paper, we used both the category label and the presence/absence of 500 most frequent words as two separate attribute sets. We got a 3312×500 word attribute matrix with 65674 non-zero entries.

- **Cora:** This dataset² contains a citation network with 5429 links of 2708 papers in machine learning, labelled with one of 7 categories. Similar to Citeseer, we used both the category label and the 500 most frequent words as two separate attribute sets. We got a 2708×500 word attribute matrix with 39268 non-zero entries.
- **Aminer:** The Aminer dataset (Tang et al., 2009) contains a citation network with 2555 papers labelled with 10 categories and 5967 links. We further collected information of each paper via the Aminer’s API, including the authors’ names (2597 unique authors), abstract, venue, year, and number of citations. For the abstract, we extract the 5 most frequent topics for each paper in a similar way to NIPS234. In total, we prepared two sets of attributes: the labels and the others formed with the combination of all collected information.

5.2.1. EXPERIMENTAL SETTINGS

For fair comparison, we generated training/testing data with the code in the SSI-PF package, which splits a network in terms of its links. We used the default hyper-parameter settings of BGGPF, SSI-PF, and niMM, provided by the original authors. K_{max} was set to 50 on Lazega-advice and 200 (same as (Hu et al., 2016a)) on all the other three datasets. For our Asym-NARM, we set $\mu_0 = 1$ and the

other hyper-parameters the same as those used in (Zhou et al., 2012; Hu et al., 2016a). Following the suggestion of Hu et al. (2016a), we used 1500 MCMC iterations in total and the last 500 samples to compute the AUC-ROC/PR scores. Since Citeseer, Cora, and Aminer are already too large for niMM, iMMM, and DMR-MMSB to produce results in reasonable time given our computational resources, we reported their performance only on Lazega-advice.

5.2.2. RESULTS

Shown in Figure 5a, Asym-NARM gains better results in terms of AUC-ROC/PR on Lazega-advice in most of the training proportions. Overall, the node-attribute models perform better than the models that do not consider node attributes, showing the usefulness of node attributes. On the other three datasets, we used different sets of attributes to study how different attributes influence the performance of Asym-NARM and SSI-PF.

In general, Asym-NARM performs better than SSI-PF regardless of which set of attributes is used. The performance of SSI-PF approaches ours in Citeseer with the labels as attributes (indicated by “-l”). But the gap between SSI-PF and our model becomes larger when the words are used as attributes (indicated by “-w”). In Cora, SSI-PF with the words does not perform as well as its non-node-attribute counterpart, BGGPF, indicating it may not be as robust as our model with large sets of attributes. To investigate this, we varied the number of the most frequent words from 10 to 500 for Asym-NARM and SSI-PF on Citeseer and Cora. With more words, the AUC-ROC/PR score of SSI-PF de-

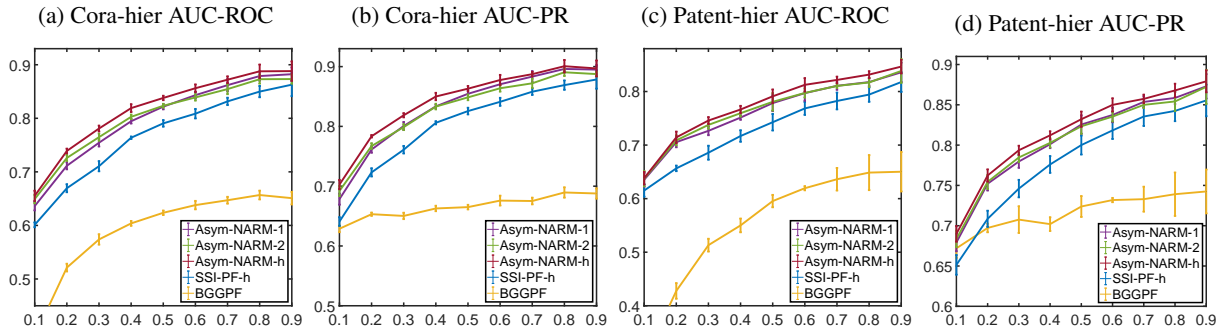


Figure 6. The AUC-ROC and AUC-PR scores on the networks with hierarchical attributes. The models with the first level attributes only, the second level attributes only, and the hierarchical attributes are marked with “-1”, “-2”, and “-h” respectively.

grades increasingly. We further checked the prior of the node factor loadings in SSI-PF (the variable that incorporates node attributes and corresponds to $g_{i,k}$ in our model) and found that the coefficient of variation of each node’s prior drops dramatically, indicating with more words, SSI-PF is failing to use the supervised information in the words.

5.3. Link Prediction with Hierarchical Node Attributes

Here we used two datasets with hierarchical node attributes: (1) **Cora-hier**: a citation network with 1712 papers and 6308 links extracted from the original Cora dataset³. The papers are labelled with one of 63 sub-areas (first level) and each sub-area belongs to one of 10 primary areas (second level), such as “machine learning in artificial intelligence” and “memory management in operating systems”; (2) **Patent-hier**: a citation network with 1461 patents and 2141 links from the National Bureau of Economic Research where the hierarchical International Patent Classification (IPC) code of a patent is used as attributes.

The AUC-ROC/PR scores in Figure 6 show that our Asym-NARM with hierarchical attributes outperforms the others, which demonstrates leveraging hierarchical side information is beneficial to link prediction. Although SSI-PF also models the hierarchical attributes, its performance in these two datasets is not comparable with our model’s.

5.4. Running Time

In this section, we compare the running time of the models for directed networks (all implemented in MATLAB and running on a desktop with 3.40 GHz CPU and 16GB RAM). Using 80% data for training, the running time for Asym-NARM, SSI-PF, and niMM on Aminer with different sets of node attributes is reported in Table 2. Note DMR-MMSB did not complete with “Authors” and “All” due to our computational resources. Asym-NARM is about 10 times faster than SSI-PF with all the attributes and about

³<https://people.cs.umass.edu/~mccallum/data.html>

Table 2. The running time (seconds per iteration) of the compared models on Aminer. AT: the topics extracted from the abstracts. All: the combination of all the attributes we have.

Attr	Non-zeros & attr size	Asym-NARM	SSI-PF	niMM	DMR-MMSB $K = 50$
Label	2660 2555*10	0.26	0.48	134.11	89.12
AT	12775 2555*100	0.29	0.87	135.22	126.44
Authors	5647 2555*2597	0.33	2.99	136.41	-
All	31273 2555*3058	0.51	5.21	136.14	-

2 times faster with the labels. Thus Asym-NARM is more efficient, especially with large sets of attributes, supporting the complexity analysis in Table 1.

6. Conclusion

As a summary of the experiments, Asym/Sym-NARM achieved better link prediction performance with faster inference. While EPM, a non-node-attribute model, performed well on nearly complete networks, it degraded with less training data. niMM and DMR-MMSB, extensions to MMSB with the logistic-normal transform, had similar results to Sym-NARM but scaled inefficiently. SSI-PF’s performance and scalability were not as good as Asym-NARM in the presented cases with flat and hierarchical attributes and it was less effective with larger numbers of attributes.

Thus NARM is a comparatively simple yet effective and efficient way of incorporating node attributes, including hierarchical attributes, for relational models with Poisson likelihood. This leads to improved link prediction and matrix completion for less complete relational data of both directed and undirected networks. With the efficient inference, our models can be used to model large sparse relational networks with node attributes.

NARM can easily be extended to multi-relational networks such as (Hu et al., 2016b) and topic models with document and word attributes, which is left for our future work.

References

- Acharya, A., Teffer, D., Henderson, J., Tyler, M., Zhou, M., and Ghosh, J. Gamma process Poisson factorization for joint modeling of network and documents. In *Machine Learning and Knowledge Discovery in Databases, European Conference, Part I*, pp. 283–299. Springer, 2015.
- Airoldi, E.M., Blei, D.M., Fienberg, S.E., and Xing, E.P. Mixed membership stochastic blockmodels. *Journal of Machine Learning Research*, 9(Sep):1981–2014, 2008.
- Buntine, W. and Hutter, M. A Bayesian view of the Poisson-Dirichlet process. *arXiv preprint arXiv:1007.0296v2 [math.ST]*, 2012.
- Chen, C., Du, L., and Buntine, W. Sampling table configurations for the hierarchical Poisson-Dirichlet process. In *Machine Learning and Knowledge Discovery in Databases. European Conference, Part I*, pp. 296–311. Springer, 2011.
- Du, L., Buntine, W., and Jin, H. A segmented topic model based on the two-parameter Poisson-Dirichlet process. *Machine Learning*, 81:5–19, 2010.
- Dunson, D.B. and Herring, A.H. Bayesian latent variable models for mixed discrete outcomes. *Biostatistics*, 6(1): 11–25, 2005.
- Fan, X., Xu, D., Yi, R., Cao, L., and Song, Y. Learning nonparametric relational models by conjugately incorporating node information in a network. *IEEE transactions on cybernetics*, 2016.
- Gopalan, P., Charlin, L., and Blei, D. Content-based recommendations with Poisson factorization. In *Advances in Neural Information Processing Systems*, pp. 3176–3184, 2014a.
- Gopalan, P., Ruiz, F.J., Ranganath, R., and Blei, D.M. Bayesian nonparametric Poisson factorization for recommendation systems. In *17th International Conference on Artificial Intelligence and Statistics*, pp. 275–283, 2014b.
- Gopalan, P., Hofman, J.M., and Blei, D.M. Scalable recommendation with hierarchical Poisson factorization. In *31st Conference on Uncertainty in Artificial Intelligence*, pp. 326–335, 2015.
- Hu, C., Rai, P., and Carin, L. Non-negative matrix factorization for discrete data with hierarchical side-information. In *19th International Conference on Artificial Intelligence and Statistics*, pp. 1124–1132, 2016a.
- Hu, C., Rai, P., and Carin, L. Topic-based embeddings for learning from large knowledge graphs. In *19th International Conference on Artificial Intelligence and Statistics*, pp. 1133–1141, 2016b.
- Huopaniemi, I., Suviataival, T., Nikkilä, J., Orešič, M., and Kaski, S. Multivariate multi-way analysis of multi-source data. *Bioinformatics*, 26(12):i391–i398, 2010.
- Kemp, C., Tenenbaum, J.B., Griffiths, T.L., Yamada, T., and Ueda, N. Learning systems of concepts with an infinite relational model. In *21st National Conference on Artificial Intelligence*, pp. 381–388. AAAI, 2006.
- Kim, D.I., Hughes, M., and Sudderth, E. The nonparametric metadata dependent relational model. In *29th International Conference on Machine Learning*, pp. 1559–1566, 2012.
- Koutsourelakis, P.-S. and Eliassi-Rad, T. Finding mixed-memberships in social networks. In *AAAI Spring Symposium: Social Information Processing*, pp. 48–53, 2008.
- Lazega, E. *The collegial phenomenon: The social mechanisms of cooperation among peers in a corporate law partnership*. Oxford University Press on Demand, 2001.
- Lim, K. and Buntine, W. Bibliographic analysis on research publications using authors, categorical labels and the citation network. *Machine Learning*, 103(2):185–213, 2016.
- Lin, J., Sugiyama, K., Kan, M.-Y., and Chua, T.-S. Addressing cold-start in app recommendation: Latent user models constructed from Twitter followers. In *36th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 283–292, 2013.
- McAuley, J.J. and Leskovec, J. Learning to discover social circles in ego networks. In *Advances in Neural Information Processing Systems*, pp. 548–56, 2012.
- Mcauliffe, J.D. and Blei, D.M. Supervised topic models. In *Advances in Neural Information Processing Systems*, pp. 121–128, 2008.
- Miller, K., Jordan, M.I., and Griffiths, T.L. Nonparametric latent feature models for link prediction. In *Advances in Neural Information Processing Systems*, pp. 1276–1284, 2009.
- Mimno, D. and McCallum, A. Topic models conditioned on arbitrary features with Dirichlet-multinomial regression. In *24th Conference on Uncertainty in Artificial Intelligence*, pp. 411–418, 2012.
- Mnih, A. and Salakhutdinov, R. Probabilistic matrix factorization. In *Advances in Neural Information Processing Systems*, pp. 1257–1264, 2008.
- Nickel, M., Murphy, K., Tresp, V., and Gabrilovich, E. A Review of Relational Machine Learning for Knowledge Graphs. *Proceedings of the IEEE*, 104(1):11–33, 2016.

- Nowicki, K. and Snijders, T.A.B. Estimation and prediction for stochastic blockstructures. *Journal of the American Statistical Association*, 96(455):1077–1087, 2001.
- Ramage, D., Hall, D., Nallapati, R., and Manning, C.D. Labeled LDA: A supervised topic model for credit attribution in multi-labeled corpora. In *2009 Conference on Empirical Methods in Natural Language Processing: Volume 1*, pp. 248–256. ACL, 2009.
- Sedhain, S., Sanner, S., Braziunas, D., Xie, L., and Christensen, J. Social collaborative filtering for cold-start recommendations. In *8th ACM Conference on Recommender Systems*, pp. 345–348, 2014.
- Singh, A.P. and Gordon, G.J. Relational learning via collective matrix factorization. In *14th ACM SIGKDD International Conference on Knowledge discovery and data mining*, pp. 650–658. ACM, 2008.
- Tang, J., Sun, J., Wang, C., and Yang, Z. Social influence analysis in large-scale networks. In *15th ACM SIGKDD International Conference on Knowledge discovery and data mining*, pp. 807–816. ACM, 2009.
- Teh, Y.W., Jordan, M.I., Beal, M.J., and Blei, D.M. Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 101(476):1566–1581, 2012.
- Wang, Y.J. and Wong, G.Y. Stochastic blockmodels for directed graphs. *Journal of the American Statistical Association*, 82(397):8–19, 1987.
- Xiang, R., Neville, J., and Rogati, M. Modeling relationship strength in online social networks. In *19th International Conference on World Wide Web*, pp. 981–990. ACM, 2010.
- Zhang, W. and Wang, J. A collective Bayesian Poisson factorization model for cold-start local event recommendation. In *21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 1455–1464, 2015.
- Zhou, M. Infinite edge partition models for overlapping community detection and link prediction. In *18th International Conference on Artificial Intelligence and Statistics*, pp. 1135–1143, 2015.
- Zhou, M. and Carin, L. Negative binomial process count and mixture modeling. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(2):307–320, 2015.
- Zhou, M., Hannah, L., Dunson, D.B., and Carin, L. Beta-negative binomial process and Poisson factor analysis. In *15th International Conference on Artificial Intelligence and Statistics*, pp. 1462–1471, 2012.