

Gamma Process Poisson Factorization for Joint Modeling of Network and Documents

Ayan Acharya¹, Dean Teffer², Jette Henderson², Marcus Tyler², Mingyuan Zhou³, and Joydeep Ghosh¹

¹Department of ECE, University of Texas at Austin, USA,
{aacharya,jghosh}@utexas.edu,

²Applied Research Laboratories, University of Texas at Austin, USA,
{dean.teffer, jhende, mtyler}@arlut.utexas.edu,

³Department of IROM, University of Texas at Austin, USA,
mzhou@utexas.edu

Abstract. Developing models to discover, analyze, and predict clusters within networked entities is an area of active and diverse research. However, many of the existing approaches do not take into consideration pertinent auxiliary information. This paper introduces Joint Gamma Process Poisson Factorization (J-GPPF) to jointly model network and side-information. J-GPPF naturally fits sparse networks, accommodates separately-clustered side information in a principled way, and effectively addresses the computational challenges of analyzing large networks. Evaluated with hold-out link prediction performance on sparse networks (both synthetic and real-world) with side information, J-GPPF is shown to clearly outperform algorithms that only model the network adjacency matrix.

Keywords: Network modeling, Poisson factorization, Gamma Process

1 Introduction

Social networks and other relational datasets often involve a large number of nodes N with sparse connections between them. If the relationship is symmetric, it can be represented compactly using a binary symmetric adjacency matrix $\mathbf{B} \in \{0, 1\}^{N \times N}$, where $b_{ij} = b_{ji} = 1$ if and only if nodes i and j are linked. Often, the nodes in such datasets are also associated with “side information,” such as documents read or written, movies rated, or messages sent by these nodes. Integer-valued side information are commonly observed and can be naturally represented by a count matrix $\mathbf{Y} \in \mathbb{Z}^{D \times V}$, where $\mathbb{Z} = \{0, 1, \dots\}$. For example, \mathbf{B} may represent a coauthor network and \mathbf{Y} may correspond to a document-by-word count matrix representing the documents written by all these authors. In another example, \mathbf{B} may represent a user-by-user social network and \mathbf{Y} may represent a user-by-item rating matrix that adds nuance and support to the network data. Incorporating such side information can result in better community identification

and superior link prediction performance as compared to modeling only the network adjacency matrix \mathbf{B} , especially for sparse networks.

Many of the popular network models [18, 2, 28, 25, 13] are demonstrated to work well for small size networks. However, small networks are often dense, while larger real-world networks tend to be much sparser and hence challenge existing modeling approaches. Incorporating auxiliary information associated with the nodes has the potential to address such challenges, as it may help better identify latent communities and predict missing links. A model that takes advantage of such side information has the potential to outperform network-only models. However, the side information may not necessarily suggest the same community structure as the existing links. Thus a network latent factor model that allows separate factors for side information and network interactions, but at the same time is equipped with a mechanism to capture dependencies between the two types of factors, is desirable.

This paper proposes **Joint Gamma Process Poisson Factorization (J-GPPF)** to jointly factorize \mathbf{B} and \mathbf{Y} in a nonparametric Bayesian manner. The paper makes the following contributions: 1) we present a fast and effective model that uses side information to help discover latent network structures, 2) we perform nonparametric Bayesian modeling for discovering latent structures in both \mathbf{B} and \mathbf{Y} , and 3) our model scales with the number of non-zero entries in the network $S_{\mathbf{B}}$ as $O(S_{\mathbf{B}}K_{\mathbf{B}})$, where $K_{\mathbf{B}}$ is the number of network groups inferred from the data.

The remainder of the paper is organized as follows. We present background material and related work in Section 2. J-GPPF and its inference algorithm are explained in Section 3. Experimental results are reported in Section 4, followed by conclusions in Section 5.

2 Background and Related Work

This section presents the related literature and the background materials that are useful for understanding the framework described in Section 3.

2.1 Negative Binomial Distribution

The negative binomial (NB) distribution $m \sim \text{NB}(r, p)$, with probability mass function (PMF) $\Pr(M = m) = \frac{\Gamma(m+r)}{m!\Gamma(r)} p^m (1-p)^r$ for $m \in \mathbb{Z}$, can be augmented into a gamma-Poisson construction as $m \sim \text{Pois}(\lambda)$, $\lambda \sim \text{Gamma}(r, p/(1-p))$, where the gamma distribution is parameterized by its shape r and scale $p/(1-p)$. It can also be augmented under a compound Poisson representation as $m = \sum_{t=1}^l u_t$, $u_t \stackrel{iid}{\sim} \text{Log}(p)$, $l \sim \text{Pois}(-r \ln(1-p))$, where $u \sim \text{Log}(p)$ is the logarithmic distribution [17]. Consequently, we have the following Lemma.

Lemma 1 ([41]). *If $m \sim \text{NB}(r, p)$ is represented under its compound Poisson representation, then the conditional posterior of l given m and r has PMF:*

$$\Pr(l = j | m, r) = \frac{\Gamma(r)}{\Gamma(m+r)} |s(m, j)| r^j, \quad j = 0, 1, \dots, m, \quad (1)$$

where $|s(m, j)|$ are unsigned Stirling numbers of the first kind. We denote this conditional posterior as $(l|m, r) \sim \text{CRT}(m, r)$, a Chinese restaurant table (CRT) count random variable, which can be generated via $l = \sum_{n=1}^m z_n$, $z_n \sim \text{Bernoulli}(r/(n-1+r))$.

Lemma 2. Let $X = \sum_{k=1}^K x_k$, $x_k \sim \text{Pois}(\zeta_k) \forall k$, and $\zeta = \sum_{k=1}^K \zeta_k$. If $(y_1, \dots, y_K|X) \sim \text{Mult}(X, \zeta_1/\zeta, \dots, \zeta_K/\zeta)$ and $X \sim \text{Pois}(\zeta)$, then the following holds:

$$P(X, x_1, \dots, x_K) = P(X, y_1, \dots, y_K). \quad (2)$$

Lemma 3. If $x_i \sim \text{Pois}(m_i \lambda)$, $\lambda \sim \text{Gamma}(r, 1/c)$, then $x = \sum_i x_i \sim \text{NB}(r, p)$, where $p = (\sum_i m_i)/(c + \sum_i m_i)$.

Lemma 4. If $x_i \sim \text{Pois}(m_i \lambda)$, $\lambda \sim \text{Gamma}(r, 1/c)$, then

$$(\lambda|\{x_i\}, r, c) \sim \text{Gamma}\left(r + \sum_i x_i, 1/(c + \sum_i m_i)\right). \quad (3)$$

Lemma 5. If $r_i \sim \text{Gamma}(a_i, 1/b) \forall i$, $b \sim \text{Gamma}(c, 1/d)$, then we have:

$$(b|\{r_i, a_i\}, c, d) \sim \text{Gamma}\left(\sum_i a_i + c, 1/(\sum_i r_i + d)\right). \quad (4)$$

The proofs of Lemmas 3, 4 and 5 follow from the definitions of Gamma, Poisson and Negative Binomial distributions.

Lemma 6. If $x_i \sim \text{Pois}(m_i r_2)$, $r_2 \sim \text{Gamma}(r_1, 1/d)$, $r_1 \sim \text{Gamma}(a, 1/b)$, then $(r_1| -) \sim \text{Gamma}(a + \ell, 1/(b - \log(1 - p)))$ where $(\ell|x, r_1) \sim \text{CRT}(\sum_i x_i, r_1)$ and $p = \sum_i m_i/(d + \sum_i m_i)$. The proof and illustration can be found in Section 3.3 of [1].

2.2 Gamma Process

The Gamma Process [12, 36] $G \sim \Gamma\text{P}(c, H)$ is a completely random measure defined on the product space $\mathbb{R}_+ \times \Omega$, with concentration parameter c and a finite and continuous base measure H over a complete separable metric space Ω , such that $G(A_i) \sim \text{Gamma}(H(A_i), 1/c)$ are independent gamma random variables for disjoint partition $\{A_i\}_i$ of Ω . The Lévy measure of the Gamma Process can be expressed as $\nu(dr d\omega) = r^{-1} e^{-cr} dr H(d\omega)$. Since the Poisson intensity $\nu^+ = \nu(\mathbb{R}_+ \times \Omega) = \infty$ and the value of $\int_{\mathbb{R}_+ \times \Omega} r \nu(dr d\omega)$ is finite, a draw from the Gamma Process consists of countably infinite atoms, which can be expressed as follows:

$$G = \sum_{k=1}^{\infty} r_k \delta_{\omega_k}, (r_k, \omega_k) \stackrel{iid}{\sim} \pi(dr d\omega), \pi(dr d\omega) \nu^+ \equiv \nu(dr d\omega). \quad (5)$$

A gamma process based model has an inherent shrinkage mechanism, as in the prior the number of atoms with weights greater than $\tau \in \mathbb{R}_+$ follows a Poisson distribution with parameter $H(\Omega) \int_{\tau}^{\infty} r^{-1} \exp(-cr) dr$, the value of which decreases as τ increases.

2.3 Network Modeling, Topic Modeling and Count Matrix Factorization

The Infinite Relational Model (IRM [18]) allows for multiple types of relations between entities in a network and an infinite number of clusters, but restricts these entities to belong to only one cluster. The Mixed Membership Stochastic Blockmodel (MMSB [2]) assumes that each node in the network can exhibit a mixture of communities. Though the MMSB has been applied successfully to discover complex network structure in a variety of applications, the computational complexity of the underlying inference mechanism is in the order of N^2 , which limits its use to small networks. Computation complexity is also a problem with many other existing latent variable network models, such as the latent feature relational model [25] and its max margin version [44], and the infinite latent attribute model [28]. The Assortative Mixed-Membership Stochastic Blockmodel (a-MMSB [13]) bypasses the quadratic complexity of the MMSB by making certain assumptions about the network structure that might not be true in general. The hierarchical Dirichlet process relational model [19] allows mixed membership with an unbounded number of latent communities; however, it is built on the a-MMSB whose assumptions could be restrictive.

Some of the existing approaches handle sparsity in real-world networks by using some auxiliary information [21, 39, 24]. For example, in a protein-protein interaction network, the features describing the biological properties of each protein can be used [24]. In an extremely sparse social network, information about each user’s profile can be used to better recommend friends [21]. Recommender system and text mining researchers, in contrast, tend to take an orthogonal approach. In recommender systems [22, 10], \mathbf{Y} may represent a user-by-item rating matrix and the objective in this setting is to predict the missing entries in \mathbf{Y} , and the social network matrix \mathbf{B} plays a secondary role in providing auxiliary information to facilitate this task [22]. Similarly, in the text mining community, many existing models [23, 26, 35, 3] use the network information or other forms of side information to improve the discovery of “topics” from the document-by-word matrix \mathbf{Y} . The matrix \mathbf{B} can represent, for example, the interaction network of authors participating in writing the documents. The Relational Topic Model [11] discovers links between documents based on their topic distributions, obtained through unsupervised exploration. The Author-Topic framework (AT [30]) and the Author-Recipient-Topic model (ART [23]) jointly model documents along with the authors of the documents. Block-LDA [3], on the other hand, provides a generative model for the links between authors and recipients in addition to documents. The Group-Topic model [34] addresses the task of modeling events pertaining to pairs of entities with textual attributes that annotate the event. J-GPPF differs from these existing approaches in mathematical formulation, including more effective modeling of both sparsity and the dependence between network interactions and side information.

A large number of discrete latent variable models for count matrix factorization can be united under Poisson factor analysis (PFA) [43], which factorizes a count matrix $\mathbf{Y} \in \mathbb{Z}^{D \times V}$ under the Poisson likelihood as $\mathbf{Y} \sim \text{Pois}(\Phi\Theta)$,

where $\Phi \in \mathbb{R}_+^{D \times K}$ is the factor loading matrix or dictionary, $\Theta \in \mathbb{R}_+^{K \times V}$ is the factor score matrix. A wide variety of algorithms, although constructed with different motivations and for distinct problems, can all be viewed as PFA with different prior distributions imposed on Φ and Θ . For example, non-negative matrix factorization [20, 9], with the objective to minimize the Kullback-Leibler divergence between \mathbf{N} and its factorization $\Phi\Theta$, is essentially PFA solved with maximum likelihood estimation. LDA [6] is equivalent to PFA, in terms of both block Gibbs sampling and variational inference, if Dirichlet distribution priors are imposed on both $\phi_k \in \mathbb{R}_+^D$, the columns of Φ , and $\theta_k \in \mathbb{R}_+^V$, the columns of Θ . The gamma-Poisson model [8, 32] is PFA with gamma priors on Φ and Θ . A family of negative binomial (NB) processes, such as the beta-NB [43, 7] and gamma-NB processes [41, 42], impose different gamma priors on $\{\theta_{vk}\}$, the marginalization of which leads to differently parameterized NB distributions to explain the latent counts. Both the beta-NB and gamma-NB process PFAs are nonparametric Bayesian models that allow K to grow without limits [16].

J-GPPF models both \mathbf{Y} and \mathbf{B} using Poisson factorization. As discussed in [1], Poisson factorization has several practical advantages over other factorization methods that use Gaussian assumptions (*e.g.* in [22]). First, zero-valued observations could be efficiently processed during inference, so the model can readily accommodate large, sparse datasets. Second, Poisson factorization is a natural representation of count data. Additionally, the model allows mixed membership across an unbounded number of latent communities using the gamma Process as a prior. The authors in [4] also use Poisson factorization to model a binary interaction matrix. However, this is a parametric model and a KL-divergence based objective is optimized w.r.t. the latent factors without using any prior information. To model the binary observations of the network matrix \mathbf{B} , J-GPPF additionally uses a novel Poisson-Bernoulli (PoBe) link, discussed in detail in Section 3, that transforms the count values from the Poisson factorization to binary values. Similar transformation has also been used in the BigCLAM model [37] which builds on the works of [4]. This model was later extended to include non-network information in the form of binary attributes [38]. Neither BigCLAM nor its extension allows non-parametric modeling or imposing prior structure on the latent factors, thereby limiting the flexibility of the models and making the obtained solutions more sensitive to initialization. The collaborative topic Poisson factorization (CTPF) framework proposed in [15] solves a different problem where the objective is to recommend articles to users of similar interest. CTPF is a parametric model and variational approximation is adopted to solve the inference.

3 Joint Gamma Process Poisson Factorization (J-GPPF)

Let there be a network of N users encoded in an $N \times N$ binary matrix \mathbf{B} . The users in the network participate in writing D documents summarized in a $D \times V$ count matrix \mathbf{Y} , where V is the size of the vocabulary. Additionally, a binary matrix \mathbf{Z} of dimension $D \times N$ can also be maintained, where the unity

entries in each column indicate the set of documents in which the corresponding user contributes. In applications where \mathbf{B} represents a user-by-user social network and \mathbf{Y} represents a user-by-item rating matrix, \mathbf{Z} turns out to be an N -dimensional identity matrix. However, in the following model description we consider the more general document-author framework. Also, to make the notations more explicit, the variables associated with the side information have \mathbf{Y} as a subscript (e.g., $G_{\mathbf{Y}}$) and those associated with the network make similar use of the subscript \mathbf{B} (e.g., $G_{\mathbf{B}}$). Also, if \mathbf{Y} represents a matrix of dimension $D \times V$, then $y_{.w}$ represents the sum over all the rows for the entries in the w^{th} column, and $y_{d.}$ represents the sum over all the columns for the entries in the d^{th} row.

Before providing an explicit description of the model, we introduce two separate Gamma Processes. The first one models the latent factors in the network and also contributes to generate the count matrix. A draw from this Gamma Process $G_{\mathbf{B}} \sim \Gamma P(c_{\mathbf{B}}, H_{\mathbf{B}})$ is expressed as: $G_{\mathbf{B}} = \sum_{k_{\mathbf{B}}=1}^{\infty} \rho_{k_{\mathbf{B}}} \delta_{\phi_{k_{\mathbf{B}}}}$, where $\phi_{k_{\mathbf{B}}} \in \Omega_{\mathbf{B}}$ is an atom drawn from an N -dimensional base distribution as $\phi_{k_{\mathbf{B}}} \sim \prod_{n=1}^N \text{Gamma}(a_{\mathbf{B}}, 1/\sigma_n)$, $\rho_{k_{\mathbf{B}}} = G_{\mathbf{B}}(\phi_{k_{\mathbf{B}}})$ is the associated weight, and $H_{\mathbf{B}}$ is the corresponding base measure. The second Gamma Process models the latent groups of side information. A draw from this gamma process $G_{\mathbf{Y}} \sim \Gamma P(c_{\mathbf{Y}}, H_{\mathbf{Y}})$ is expressed as: $G_{\mathbf{Y}} = \sum_{k_{\mathbf{Y}}=1}^{\infty} r_{k_{\mathbf{Y}}} \delta_{\beta_{k_{\mathbf{Y}}}}$, where $\beta_{k_{\mathbf{Y}}} \in \Omega_{\mathbf{Y}}$ is an atom drawn from a V -dimensional base distribution as $\beta_{k_{\mathbf{Y}}} \sim \text{Dir}(\xi_{\mathbf{Y}})$, $r_{k_{\mathbf{Y}}} = G_{\mathbf{Y}}(\beta_{k_{\mathbf{Y}}})$ is the associated weight, and $H_{\mathbf{Y}}$ is the corresponding base measure. Also, $\gamma_{\mathbf{B}} = H_{\mathbf{B}}(\Omega_{\mathbf{B}})$ is defined as the mass parameter corresponding to the base measure $H_{\mathbf{B}}$ and $\gamma_{\mathbf{Y}} = H_{\mathbf{Y}}(\Omega_{\mathbf{Y}})$ is defined as the mass parameter corresponding to the base measure $H_{\mathbf{Y}}$. In the following paragraphs, we explain how these Gamma processes, with the atoms and their associated weights, are used for modeling both \mathbf{B} and \mathbf{Y} .

The $(n, m)^{\text{th}}$ entry in the matrix \mathbf{B} is assumed to be derived from a latent count as:

$$b_{nm} = \mathbb{I}_{\{x_{nm} \geq 1\}}, \quad x_{nm} \sim \text{Pois}(\lambda_{nm}), \quad \lambda_{nm} = \sum_{k_{\mathbf{B}}} \lambda_{nmk_{\mathbf{B}}}, \quad (6)$$

where $\lambda_{nmk_{\mathbf{B}}} = \rho_{k_{\mathbf{B}}} \phi_{nk_{\mathbf{B}}} \phi_{mk_{\mathbf{B}}}$. This is called as the Poisson-Bernoulli (PoBe) link in [1, 40]. The distribution of b_{nm} given λ_{nm} is named as the Poisson-Bernoulli distribution, with the PMF: $f(b_{nm}|\lambda_{nm}) = e^{-\lambda_{nm}(1-b_{nm})}(1-e^{-\lambda_{nm}})^{b_{nm}}$. One may consider $\lambda_{nmk_{\mathbf{B}}}$ as the strength of mutual latent community membership between nodes n and m in the network for latent community $k_{\mathbf{B}}$, and λ_{nm} as the interaction strength aggregating all possible community membership. Using Lemma 2, one may augment the above representation as $x_{nm} = \sum_{k_{\mathbf{B}}} x_{nmk_{\mathbf{B}}}$, $x_{nmk_{\mathbf{B}}} \sim \text{Pois}(\lambda_{nmk_{\mathbf{B}}})$. Thus each interaction pattern contributes a count and the total latent count aggregates the countably infinite interaction patterns.

Unlike the usual approach that links the binary observations to latent Gaussian random variables with a logistic or probit function, the above approach links the binary observations to Poisson random variables. Thus, this approach transforms the problem of modeling binary network interaction into a count modeling

problem, providing several potential advantages. First, it is more interpretable because ρ_{k_B} and ϕ_{k_B} are non-negative and the aggregation of different interaction patterns increases the probability of establishing a link between two nodes. Second, the computational benefit is significant since the computational complexity is approximately linear in the number of non-zeros S_B in the observed binary adjacency matrix B . This benefit is especially pertinent in many real-world datasets where S_B is significantly smaller than N^2 .

To model the matrix Y , its $(d, w)^{\text{th}}$ entry y_{dw} is generated as:

$$y_{dw} \sim \text{Pois}(\zeta_{dw}), \zeta_{dw} = \left(\sum_{k_Y} \zeta_{Y dw k_Y} + \sum_{k_B} \zeta_{B dw k_B} \right),$$

$$\zeta_{Y dw k_Y} = r_{k_Y} \theta_{dk_Y} \beta_{wk_Y}, \zeta_{B dw k_B} = \epsilon \rho_{k_B} \left(\sum_n Z_{nd} \phi_{nk_B} \right) \psi_{wk_B},$$

where $Z_{nd} \in \{0, 1\}$ and $Z_{nd} = 1$ if and only if author n is one of the authors of paper d . One can consider ζ_{dw} as the affinity of document d for word w . This affinity is influenced by two different components, one of which comes from the network modeling. Without the contribution from network modeling, the joint model reduces to a gamma process Poisson matrix factorization model, in which

the matrix Y is factorized in such a way that $y_{dw} \sim \text{Pois} \left(\sum_{k_Y} r_{k_Y} \theta_{dk_Y} \beta_{wk_Y} \right)$.

Here, $\Theta \in \mathbb{R}_+^{D \times \infty}$ is the factor score matrix, $\beta \in \mathbb{R}_+^{V \times \infty}$ is the factor loading matrix (or dictionary) and r_{k_Y} signifies the weight of the k_Y^{th} factor. The number of latent factors, possibly smaller than both D and V , would be inferred from the data.

In the proposed joint model, Y is also determined by the users participating in writing the d^{th} document. We assume that the distribution over word counts for a document is a function of *both* its topic distribution *as well as* the characteristics of the users associated with it. In the author-document framework, the authors employ different writing styles and have expertise in different domains. For example, an author from machine learning and statistics would use words like “probability”, “classifiers”, “patterns”, “prediction” more often than an author with an economics background. Frameworks such as author-topic model [30, 23] were motivated by a related concept. In the user-rating framework, the entries in Y are also believed to be influenced by the interaction network of the users. Such influence of the authors is modeled by the interaction of the authors in the latent communities *via* the latent factors $\phi \in \mathbb{R}_+^{N \times \infty}$ and $\psi \in \mathbb{R}_+^{V \times \infty}$, which encodes the writing style of the authors belonging to different latent communities. Since an infinite number of network communities is maintained, each entry y_{dw} is assumed to come from an infinite dimensional interaction. ρ_{k_B} signifies the interaction strength corresponding to the k_B^{th} network community. The contributions of the interaction from all the authors participating in a given document are accumulated to produce the total contribution from the networks in generating y_{dw} . Since B and Y might have different levels of sparsity and the range of integers in Y can be quite large, a parameter ϵ is required to balance

the contribution of the network communities in dictating the structure of \mathbf{Y} . A low value of ϵ forces disjoint modeling of \mathbf{B} and \mathbf{Y} , while a higher value implies joint modeling of \mathbf{B} and \mathbf{Y} where information can flow both ways, from network discovery to topic discovery and vice-versa. We present a thorough discussion of the effect of ϵ in Section 4.1. To complete the generative process, we put Gamma priors over σ_n , ς_d , c_B , c_Y and ϵ as:

$$c_B \sim \text{Gamma}(g_B, 1/h_B), c_Y \sim \text{Gamma}(g_Y, 1/h_Y), \epsilon \sim \text{Gamma}(g_0, 1/f_0). \quad (7)$$

$$\sigma_n \sim \text{Gamma}(\alpha_B, 1/\varepsilon_B), \varsigma_d \sim \text{Gamma}(\alpha_Y, 1/\varepsilon_Y). \quad (8)$$

3.1 Inference *via* Gibbs Sampling

Though J-GPPF supports countably infinite number of latent communities for network modeling and infinite number of latent factors for topic modeling, in practice it is impossible to instantiate all of them. Instead of marginalizing out the underlying stochastic process [5, 27] or using slice sampling [33] for non-parametric modeling, for simplicity, we consider a finite approximation of the infinite model by truncating the number of graph communities and the latent topics to K_B and K_Y respectively, by letting $\rho_{k_B} \sim \text{Gamma}(\gamma_B/K_B, 1/c_B)$ and $r_{k_Y} \sim \text{Gamma}(\gamma_Y/K_Y, 1/c_Y)$. Such approximation approaches the original infinite model as both K_B and K_Y approach infinity. With such finite approximation, the generative process of J-GPPF is summarized in Table 1. For notational convenience, we represent the set of documents the n^{th} author contributes to as \mathcal{Z}_n and the set of authors contributing to the d^{th} document as \mathcal{Z}_d .

$b_{nm} = I_{\{x_{nm} \geq 1\}}, x_{nm} \sim \text{Pois} \left(\sum_{k_B} \rho_{k_B} \phi_{nk_B} \phi_{mk_B} \right),$ $y_{dw} \sim \text{Pois} \left(\sum_{k_Y} r_{k_Y} \theta_{dk_Y} \beta_{wk_Y} + \epsilon \sum_{k_B} \rho_{k_B} \left(\sum_n Z_{nd} \phi_{nk_B} \right) \psi_{wk_B} \right),$ $\phi_{k_B} \sim \prod_{n=1}^N \text{Gamma}(a_B, 1/\sigma_n), \psi_{k_B} \sim \text{Dir}(\xi_B),$ $\theta_{k_Y} \sim \prod_{d=1}^D \text{Gamma}(a_Y, 1/\varsigma_d), \beta_{k_Y} \sim \text{Dir}(\xi_Y), \epsilon \sim \text{Gamma}(f_0, 1/g_0),$ $\sigma_n \sim \text{Gamma}(\alpha_B, 1/\varepsilon_B), \rho_{k_B} \sim \text{Gamma}(\gamma_B/K_B, 1/c_B),$ $\gamma_B \sim \text{Gamma}(e_B, 1/f_B), c_B \sim \text{Gamma}(g_B, 1/h_B),$ $\varsigma_d \sim \text{Gamma}(\alpha_Y, 1/\varepsilon_Y), r_{k_Y} \sim \text{Gamma}(\gamma_Y/K_Y, 1/c_Y),$ $\gamma_Y \sim \text{Gamma}(e_Y, 1/f_Y), c_Y \sim \text{Gamma}(g_Y, 1/h_Y).$

Table 1. Generative Process of J-GPPF

Sampling of $(x_{nmk_B})_{k_B=1}^{K_B}$: We first sample the network links according to the following:

$$(x_{nm} | -) \sim b_{nm} \text{Pois}_+ \left(\sum_{k_B=1}^{K_B} \lambda_{nmk_B} \right). \quad (9)$$

Sampling from a truncated Poisson distribution is described in detail in [40]. Since, one can augment $x_{nm} \sim \text{Pois} \left(\sum_{k_B=1}^{K_B} \lambda_{nmk_B} \right)$ as $x_{nm} = \sum_{k_B=1}^{K_B} x_{nmk_B}$,

where $x_{nmk_B} \sim \text{Pois}(\lambda_{nmk_B})$, equivalently, one obtains the following:

$$\left((x_{nmk_B})_{k_B=1}^{K_B} \mid -\right) \sim \text{Mult} \left(x_{nm}, \left(\lambda_{nmk_B} / \sum_{k_B=1}^{K_B} \lambda_{nmk_B} \right)_{k_B=1}^{K_B} \right). \quad (10)$$

Sampling of $(y_{dwk})_k$: Since, one can augment $y_{dw} \sim \text{Pois}(\zeta_{dw})$ as $y_{dw} = \sum_{k_Y=1}^{K_Y} y_{dwk_Y} + \sum_{n \in \mathcal{Z}_d} \sum_{k_B=1}^{K_B} y_{dnwk_B}$, $y_{dwk_Y} \sim \text{Pois}(\zeta_{dwk_Y})$, $y_{dnwk_B} \sim \text{Pois}(\zeta_{dnwk_B})$, again following Lemma 2, we have:

$$\left((y_{dwk_Y})_{k_Y=1}^{K_Y}, (y_{dnwk_B})_{k_B=1, n \in \mathcal{Z}_d}^{K_B} \mid -\right) \sim \text{Mult} \left(y_{dw}, \frac{\{\zeta_{dwk_Y}\}_{k_Y}, \{\zeta_{dnwk_B}\}_{n \in \mathcal{Z}_d, k_B}}{\sum_{k_Y} \zeta_{dwk_Y} + \sum_{n \in \mathcal{Z}_d} \sum_{k_B} \zeta_{dnwk_B}} \right). \quad (11)$$

Sampling of ϕ_{nk_B} , ρ_{k_B} , θ_{dk_Y} , r_{k_Y} and ϵ : Sampling of these parameters follow from Lemma 4 and are given as follows:

$$(\phi_{nk_B} \mid -) \sim \text{Gamma} \left(a_B + x_{n \cdot k_B} + y_{n \cdot k_B}, \frac{1}{\sigma_n + \rho_{k_B}(\phi_{k_B}^{-n} + \epsilon | \mathcal{Z}_n |)} \right), \quad (12)$$

$$(\rho_{k_B} \mid -) \sim \text{Gamma} \left(\frac{\gamma_B}{K_B} + x_{\cdot \cdot k_B} + y_{\cdot \cdot k_B}, \frac{1}{c_B + \sum_n \phi_{nk_B} \phi_{k_B}^{-n} + \epsilon \sum_n |\mathcal{Z}_n| \phi_{nk_B}} \right), \quad (13)$$

$$(\theta_{dk_Y} \mid -) \sim \text{Gamma} \left(a_Y + y_{d \cdot k_Y}, \frac{1}{\varsigma_d + r_{k_Y}} \right), (r_{k_Y} \mid -) \sim \text{Gamma} \left(\frac{\gamma_Y}{K_Y} + y_{\cdot \cdot k_Y}, \frac{1}{c_Y + \theta_{\cdot k_Y}} \right), \quad (14)$$

$$(\epsilon \mid -) \sim \text{Gamma} \left(f_0 + \sum_{k=1}^{K_B} y_{\cdot \cdot k}, \frac{1}{g_0 + q_0} \right), q_0 = \sum_{k=1}^{K_B} \rho_{k_B} \sum_{n=1}^N |\mathcal{Z}_n| \phi_{nk_B}. \quad (15)$$

The sampling of parameters ϕ_{nk_B} and ρ_{k_B} exhibits how information from the count matrix \mathbf{Y} influences the discovery of the latent network structure. The latent counts from \mathbf{Y} impact the shape parameters for both the posterior gamma distribution of ϕ_{nk_B} and ρ_{k_B} , while \mathbf{Z} influences the corresponding scale parameters.

Sampling of ψ_{k_B} : Since $y_{dnwk_B} \sim \text{Pois}(\epsilon \rho_{k_B} \phi_{nk_B} \psi_{wk_B})$, using Lemma 2 we have: $(y_{\cdot \cdot wk_B})_{w=1}^V \sim \text{Mult}(y_{\cdot \cdot k_B}, (\psi_{wk_B})_{w=1}^V)$. Since the Dirichlet distribution is conjugate to the multinomial, the posterior of ψ_{k_B} also becomes a Dirichlet distribution and can be sampled as:

$$(\psi_{k_B} \mid -) \sim \text{Dir}(\boldsymbol{\xi}_{B1} + y_{\cdot 1 k_B}, \dots, \boldsymbol{\xi}_{BV} + y_{\cdot V k_B}). \quad (16)$$

Sampling of β_{k_Y} : Since $y_{dwk_Y} \sim \text{Pois}(r_{k_Y} \theta_{dk_Y} \beta_{wk_Y})$, again using Lemma 2, we have:

$$(y_{\cdot \cdot wk_Y})_{w=1}^V \sim \text{Mult}(y_{\cdot \cdot k_Y}, (\beta_{wk_Y})_{w=1}^V).$$

Using conjugacy, the posterior of β_{k_Y} can be sampled as:

$$(\beta_{k_Y} \mid -) \sim \text{Dir}(\boldsymbol{\xi}_{Y1} + y_{\cdot 1 k_Y}, \dots, \boldsymbol{\xi}_{YV} + y_{\cdot V k_Y}). \quad (17)$$

Sampling of σ_n , ς_d , c_B and c_Y : Sampling of these parameters follow from Lemma 5 and are given as:

$$(\sigma_n|-) \sim \text{Gamma}\left(\alpha_B + K_B a_B, \frac{1}{\varepsilon_B + \phi_n}\right), (\varsigma_d|-) \sim \text{Gamma}\left(\alpha_Y + K_Y a_Y, \frac{1}{\varepsilon_Y + \theta_d}\right), \quad (18)$$

$$(c_B|-) \sim \text{Gamma}\left(g_B + \gamma_B, \frac{1}{h_B + \sum_{k_B} \rho_{k_B}}\right), (c_Y|-) \sim \text{Gamma}\left(g_Y + \gamma_Y, \frac{1}{h_Y + \sum_{k_Y} r_{k_Y}}\right). \quad (19)$$

Sampling of γ_B : Using Lemma 2, one can show that $x_{..k_B} \sim \text{Pois}(\rho_{k_B})$. Integrating ρ_{k_B} and using Lemma 4, one can have $x_{..k_B} \sim \text{NB}(\gamma_B, p_B)$, where $p_B = 1/(c_B + 1)$. Similarly, $y_{..k_B} \sim \text{Pois}(\rho_{k_B})$ and after integrating ρ_{k_B} and using Lemma 4, we have $y_{..k_B} \sim \text{NB}(\gamma_B, p_B)$. We now augment $l_{k_B} \sim \text{CRT}(x_{..k_B} + y_{..k_B}, \gamma_B)$ and then following Lemma 6 sample

$$(\gamma_B|-) \sim \text{Gamma}\left(e_B + \sum_{k_B} l_{k_B}, \frac{1}{f_B - q_B}\right), q_B = \sum_{k_B} \frac{q_{k_B}}{K_B}, q_{k_B} = \log\left(\frac{c_B}{c_B + \sum_n \phi_{nk_B} \phi_{k_B}^{-n}}\right). \quad (20)$$

Sampling of γ_Y : Using Lemma 2, one can show that $y_{..(K_B + k_Y)} \sim \text{Pois}(r_{k_Y})$ and after integrating r_{k_Y} and using Lemma 4, we have $y_{..(K_B + k_Y)} \sim \text{NB}(\gamma_Y, p_Y)$, where $p_Y = 1/(c_Y + 1)$. We now augment $m_{k_Y} \sim \text{CRT}(y_{..(K_B + k_Y)}, \gamma_Y)$ and then following Lemma 6 sample

$$(\gamma_Y|-) \sim \text{Gamma}\left(e_Y + \sum_{k_Y} m_{k_Y}, \frac{1}{f_Y - q_Y}\right), q_Y = \sum_{k_Y} \frac{q_{k_Y}}{K_Y}, q_{k_Y} = \log\left(\frac{c_Y}{c_Y + \theta_{k_Y}}\right). \quad (21)$$

$b_{nm} = I_{\{x_{nm} \geq 1\}}, x_{nm} \sim \text{Pois}\left(\sum_{k_B=1}^{\infty} \lambda_{nmk_B}\right), r_{k_B} \sim \text{Gamma}(\gamma_B/K_B, 1/c_B),$ $\phi_{k_B} \sim \prod_{n=1}^N \text{Gamma}(a_B, 1/\sigma_n), \sigma_n \sim \text{Gamma}(\alpha_B, 1/\varepsilon_B),$ $\gamma_B \sim \text{Gamma}(e_B, 1/f_B), c_B \sim \text{Gamma}(g_B, 1/h_B),$

Table 2. Generative Process of N-GPPF

$y_{dw} \sim \text{Pois}\left(\sum_{k_Y=1}^{\infty} r_{k_Y} \theta_{dk_Y} \beta_{wk_Y}\right),$ $\theta_{k_Y} \sim \prod_{d=1}^D \text{Gamma}(a_Y, 1/\varsigma_d), \beta_{k_Y} \sim \text{Dir}(\xi_Y),$ $\varsigma_d \sim \text{Gamma}(\eta_Y, 1/\xi_Y), r_{k_Y} \sim \text{Gamma}(\gamma_Y/K_Y, 1/c_Y),$ $\gamma_{Y_Y} \sim \text{Gamma}(e_Y, 1/f_Y), c_Y \sim \text{Gamma}(g_Y, 1/h_Y).$
--

Table 3. Generative Process of C-GPPF

3.2 Special cases: Network Only GPPF (N-GPPF) and Corpus Only GPPF (C-GPPF)

A special case of J-GPPF appears when only the binary matrix \mathbf{B} is modeled without the auxiliary matrix \mathbf{Y} . The generative model of N-GPPF is given in Table 2. The update equations of variables corresponding to N-GPPF can be obtained with $\mathbf{Z} = \mathbf{0}$ in the corresponding equations. As mentioned in Section 2.3, N-GPPF can be considered as the gamma process infinite edge partition model (EPM) proposed in [40], which is shown to well model assortative networks but not necessarily disassortative ones. Using the techniques developed in [40] to capture community-community interactions, it is relatively straightforward to extend J-GPPF to model disassortative networks. Another special case of J-GPPF appears when only the count matrix \mathbf{Y} is modeled without using the contribution from the network matrix \mathbf{B} . The generative model of C-GPPF is given in Table 3.

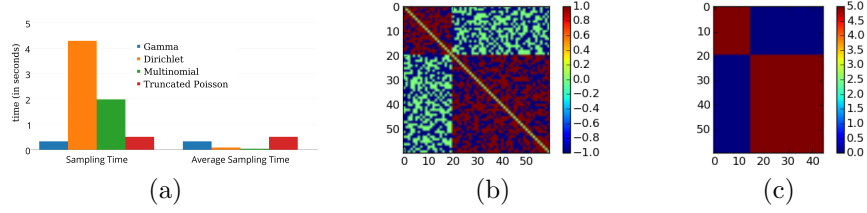


Fig. 1. (a) Time to generate a million of samples, (b) \mathbf{B} with held-out data, (c) \mathbf{Y}

3.3 Computation Complexity

The Gibbs sampling updates of J-GPPF can be calculated in $O(K_{\mathbf{B}}S_{\mathbf{B}} + (K_{\mathbf{B}} + K_{\mathbf{Y}})S_{\mathbf{Y}} + NK_{\mathbf{B}} + DK_{\mathbf{Y}} + V(K_{\mathbf{B}} + K_{\mathbf{Y}}))$ time, where $S_{\mathbf{B}}$ is the number of non-zero entries in \mathbf{B} and $S_{\mathbf{Y}}$ is the number of non-zero entries in \mathbf{Y} . It is obvious that for large matrices the computation is primarily of the order of $K_{\mathbf{B}}S_{\mathbf{B}} + (K_{\mathbf{B}} + K_{\mathbf{Y}})S_{\mathbf{Y}}$. Such complexity is a huge saving when compared to other methods like MMSB [2], that only models \mathbf{B} and incurs computation cost of $O(N^2K_{\mathbf{B}})$; and standard matrix factorization approaches [31] that work with the matrix \mathbf{Y} and incur $O(DVK_{\mathbf{Y}})$ computation cost. Interestingly, the inference in [14] incurs cost $O(K_{\mathbf{Y}}^2D + K_{\mathbf{Y}}V + K_{\mathbf{Y}}S_{\mathbf{Y}})$ with $K_{\mathbf{Y}}$ signifying the termination point of stick breaking construction in their model. C-GPPF incurs computation cost $O(DK_{\mathbf{Y}} + K_{\mathbf{Y}}S_{\mathbf{Y}} + VK_{\mathbf{Y}})$, an apparent improvement over that of [14]. However, one needs to keep in mind that [14] use variational approximation for which the updates are available in closed form solution. Our method does not use any approximation to joint distribution but uses Gibbs sampling, the computation cost of which should also be taken into account. In Fig. 1(a), we show the computation time for generating one million samples from Gamma, Dirichlet (of dimension 50), multinomial (of dimension 50) and truncated Poisson distributions using the samplers available from GNU Scientific Library (GSL) on an Intel 2127U machine with 2 GB of RAM and 1.90 GHz of processor base frequency. To highlight the average complexity of sampling from Dirichlet and multinomial distributions, we further display another plot where the computation time is divided by 50 for these samplers only. One can see that to draw one million samples, our implementation of the sampler for truncated Poisson distribution takes the longest, though the difference from the Gamma sampler in GSL is not that significant.

4 Experimental Results

4.1 Experiments with Synthetic Data

We generate a synthetic network of size 60×60 (\mathbf{B}) and a count data matrix of size 60×45 (\mathbf{Y}). Each user in the network writes exactly one document and a user and the corresponding document are indexed by the same row-index in \mathbf{B} and \mathbf{Y} respectively. To evaluate the quality of reconstruction in presence of side-information and less of network structure, we hold-out 50% of links and

equal number of non-links from \mathbf{B} . This is shown in Fig. 1(b) where the links are presented by brown, the non-links by green and the held-out data by deep blue. Clearly, the network consists of two groups. $\mathbf{Y} \in \{0, 5\}^{60 \times 45}$, shown in Fig 1(c), is also assumed to consist of the same structure as \mathbf{B} where the zeros are presented by deep blue and the non-zeros are represented by brown. The performance of N-GPPF is displayed in Fig. 2(a). Evidently, there is not much structure visible in the discovered partition of \mathbf{B} from N-GPPF and that is reflected in the poor value of AUC in Fig. 3(a). The parameter ϵ , when fixed at a given value, plays an important role in determining the quality of reconstruction for J-GPPF. As $\epsilon \rightarrow 0$, J-GPPF approaches the performance of N-GPPF on \mathbf{B} and we observe as poor a quality of reconstruction as in Fig. 2(a). When ϵ is increased and set to 1.0, J-GPPF departs from N-GPPF and performs much better in terms of both structure recovery and prediction on held-out data as shown in Fig. 2(e) and Fig. 3(b). With $\epsilon = 10.0$, perfect reconstruction and prediction are recorded as shown in Fig. 2(i) and Fig. 3(c) respectively. In this synthetic example, \mathbf{Y} is purposefully designed to reinforce the structure of \mathbf{B} when most of its links and non-links are held-out. However, in real applications, \mathbf{Y} might not contain as much of information and the Gibbs sampler needs to find a suitable value of ϵ that can carefully glean information from \mathbf{Y} .

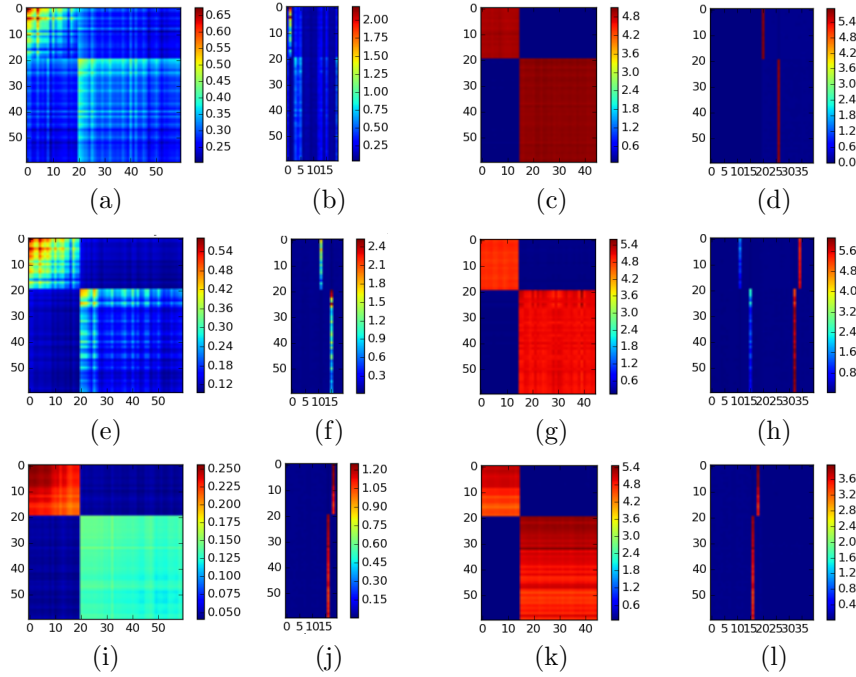


Fig. 2. Performance of J-GPPF: $\epsilon = 10^{-3}$ (top row), $\epsilon = 1$ (middle row), $\epsilon = 10$ (bottom row)

There are few more interesting observations from the experiment with synthetic data that characterize the behavior of the model and match our intuitions.

In our experiment with synthetic data $K_{\mathbf{B}} = K_{\mathbf{Y}} = 20$ is used. Fig. 2(b) demonstrates the assignment of the users in the network communities and Fig. 2(d) illustrates the assignment of the documents to the combined space of network communities and the topics (with the network communities appearing before the topics in the plot). For $\epsilon = 0.001$, we observe disjoint modeling of \mathbf{B} and \mathbf{Y} , with two latent factors modeling \mathbf{Y} and multiple latent factors modeling \mathbf{B} without any clear assignment. As we increase ϵ , we start observing joint modeling of \mathbf{B} and \mathbf{Y} . For $\epsilon = 1.0$, as Fig. 2(h) reveals, two of the network latent factors and two of the factors for count data together model \mathbf{Y} , the contribution from the network factors being expectedly small. Fig. 2(f) shows how two of the exact same latent factors model \mathbf{B} as well. Fig. 2(j) and Fig. 2(l) show how two of the latent factors corresponding to \mathbf{B} dictate the modeling of both \mathbf{B} and \mathbf{Y} when $\epsilon = 10.0$. This implies that the discovery of latent groups in \mathbf{B} is dictated mostly by the information contained in \mathbf{Y} . In all these cases, however, we observe perfect reconstruction of \mathbf{Y} as shown in Fig. 2(c), Fig. 2(g) and Fig. 2(k).

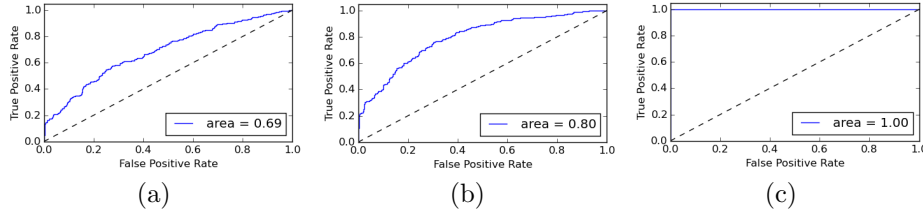


Fig. 3. (a) AUC with $\epsilon = 0.001$, (b) AUC with $\epsilon = 1.0$, (c) AUC with $\epsilon = 10.0$

4.2 Experiments with Real World Data

To evaluate the performance of J-GPPF, we consider N-GPPF, the infinite relational model (IRM) of [18] and the Mixed Membership Stochastic Block Model (MMSB) [2] as the baseline algorithms.

NIPS Authorship Network This dataset contains the papers and authors from NIPS 1988 to 2003. We took the 234 authors who published with the most other people and looked at their co-authors. After pre-processing and removing words that appear less than 50 times, the number of users in the graph is 225 and the total number of unique words is 1354. The total number of documents is 1165.

GoodReads Data Using the Goodreads API, we collected a base set of users with recent activity on the website. The friends and friends of friends of these users were collected. Up to 200 reviews were saved per user, each consisting of a book ID and a rating from 0 to 5. A similar dataset was used in [10]. After pre-processing and removing words that appear less than 10 times, the number of users in the graph is 84 and the total number of unique words is 189.

Twitter Data The Twitter dataset is a set of geo-tagged tweets collected by the authors in [29]. We extracted a subset of users located in San Francisco for our analysis. We created a graph within the subset by collecting follower information

from the Twitter API. The side information consists of tweets aggregated by user, with one document per user. After pre-processing and removing words that appear less than 25 times, the number of users in the graph is 670 and the total number of unique words is 538.

Experimental Setup and Results

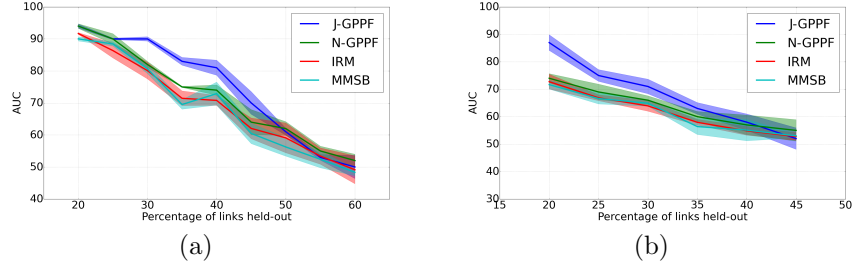


Fig. 4. (a) NIPS Data, (b) GoodReads Data

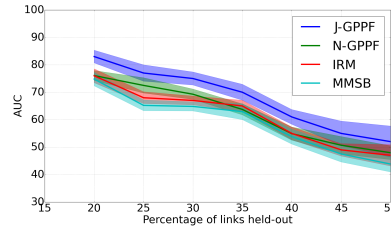


Fig. 5. Twitter Data

In all the experiments, we initialized ϵ to 2 and let the sampler decide what value works best for joint modeling. We used $K_{\mathbf{B}} = K_{\mathbf{Y}} = 50$ and initialized all the hyper-parameters to 1. For each dataset, we ran 20 different experiments and display the mean AUC and one standard error. Fig. 4 and 5 demonstrate the performances of the models in predicting the held-out data. J-GPPF clearly has advantage over other network-only models when the network is sparse enough and the auxiliary information is sufficiently strong. However, all methods fail when the sparsity increases beyond a certain point. The performance of J-GPPF also drops below the performances of network-only models in highly sparse networks, as the sampler faces additional difficulty in extracting information from both the network and the count matrix.

5 Conclusion

We propose J-GPPF that jointly factorizes the network adjacency matrix and the associated side information that can be represented as a count matrix. The model has the advantage of representing true sparsity in adjacency matrix, in latent group membership, and in the side information. We derived an efficient MCMC inference method, and compared our approach to several popular network algorithms that model the network adjacency matrix. Experimental results confirm the efficiency of the proposed approach in utilizing side information to improve the performance of network models.

Acknowledgement

This work is supported by the United States Office of Naval Research, grant No. N00014-14-1-0039.

References

1. Acharya, A., Ghosh, J., Zhou, M.: Nonparametric Bayesian Factor Analysis for Dynamic Count Matrices. In: Proc. of AISTATS (to appear) (2015)
2. Airoldi, E.M., Blei, D.M., Fienberg, S.E., Xing, E.P.: Mixed membership stochastic blockmodels. *JMLR* 9, 1981–2014 (jun 2008)
3. Balasubramanyan, R., Cohen, W.W.: Block-LDA: Jointly modeling entity-annotated text and entity-entity links. In: Proc. of SDM. pp. 450–461 (2011)
4. Ball, B., Karrer, B., Newman, M.: Efficient and principled method for detecting communities in networks. *Phys. Rev. E* 84 (Sep 2011)
5. Blackwell, D., MacQueen, J.: Ferguson distributions via Pólya urn schemes. *The Annals of Statistics* (1973)
6. Blei, D.M., Ng, A.Y., Jordan, M.I.: Latent Dirichlet Allocation. *JMLR* 3, 993–1022 (2003)
7. Broderick, T., Mackey, L., Paisley, J., Jordan, M.I.: Combinatorial clustering and the beta negative binomial process. *arXiv:1111.1802v5* (2013)
8. Canny, J.: Gap: a factor model for discrete data. In: SIGIR (2004)
9. Cemgil, A.T.: Bayesian inference for nonnegative matrix factorisation models. *Intell. Neuroscience* (2009)
10. Chaney, A., Gopalan, P., Blei, D.: Poisson trust factorization for incorporating social networks into personalized item recommendation. In: NIPS Workshop: What Difference Does Personalization Make? (2013)
11. Chang, J., Blei, D.: Relational topic models for document networks. In: Proc. of AISTATS (2009)
12. Ferguson, T.S.: A Bayesian analysis of some nonparametric problems. *Ann. Statist.* (1973)
13. Gopalan, P., Mimno, D.M., Gerrish, S., Freedman, M.J., Blei, D.M.: Scalable inference of overlapping communities. In: Proc. of NIPS. pp. 2258–2266 (2012)
14. Gopalan, P., Ruiz, F., Ranganath, R., Blei, D.: Bayesian nonparametric poisson factorization for recommendation systems. In: Proc. of AISTATS (2014)
15. Gopalan, P., Charlin, L., Blei, D.: Content-based recommendations with poisson factorization. In: Proc. of NIPS, pp. 3176–3184 (2014)
16. Hjort, N.L.: Nonparametric Bayes estimators based on beta processes in models for life history data. *Ann. Statist.* (1990)
17. Johnson, N.L., Kemp, A.W., Kotz, S.: *Univariate Discrete Distributions*. John Wiley & Sons (2005)
18. Kemp, C., Tenenbaum, J., Griffiths, T., Yamada, T., Ueda, N.: Learning systems of concepts with an infinite relational model. In: Proc. of AAAI. pp. 381–388 (2006)
19. Kim, D.I., Gopalan, P., Blei, D.M., Sudderth, E.B.: Efficient online inference for bayesian nonparametric relational models. In: Proc. of NIPS. pp. 962–970 (2013)
20. Lee, D.D., Seung, H.S.: Algorithms for non-negative matrix factorization. In: NIPS (2001)
21. Leskovec, J., Julian, J.: Learning to discover social circles in ego networks. In: Proc. of NIPS, pp. 539–547 (2012)
22. Ma, H., Yang, H., Lyu, M.R., King, I.: Sorec: Social recommendation using probabilistic matrix factorization. In: Proc. of CIKM. pp. 931–940 (2008)

23. McCallum, A., Wang, X., Corrada-Emmanuel, A.: Topic and role discovery in social networks with experiments on enron and academic email. *J. Artif. Int. Res.* 30(1), 249–272 (Oct 2007)
24. Menon, A., Elkan, C.: Link prediction via matrix factorization. In: *Machine Learning and Knowledge Discovery in Databases, Lecture Notes in Computer Science*, vol. 6912, pp. 437–452. Springer Berlin / Heidelberg (2011)
25. Miller, K.T., Griffiths, T.L., Jordan, M.I.: Nonparametric latent feature models for link prediction. In: *Proc. of NIPS*. pp. 1276–1284 (2009)
26. Nallapati, R., Ahmed, A., Xing, E., Cohen, W.: Joint latent topic models for text and citations. In: *Proc. of KDD*. pp. 542–550 (2008)
27. Neal, R.M.: Markov chain sampling methods for Dirichlet process mixture models. *Journal of computational and graphical statistics* (2000)
28. Palla, K., Ghahramani, Z., Knowles, D.A.: An infinite latent attribute model for network data. In: *Proc. of ICML*. pp. 1607–1614 (2012)
29. Roller, S., Speriosu, M., Rallapalli, S., Wing, B., Baldridge, J.: Supervised text-based geolocation using language models on an adaptive grid. In: *Proc. of EMNLP-CoNLL*. pp. 1500–1510 (2012)
30. Rosen-Zvi, M., Griffiths, T., Steyvers, M., Smyth, P.: The author-topic model for authors and documents. In: *Proc. of UAI*. pp. 487–494 (2004)
31. Salakhutdinov, R., Mnih, A.: Probabilistic matrix factorization. In: *Proc. of NIPS* (2007)
32. Titsias, M.K.: The infinite gamma-Poisson feature model. In: *Proc. of NIPS* (2008)
33. Walker, S.G.: Sampling the Dirichlet mixture model with slices. *Communications in Statistics Simulation and Computation* (2007)
34. Wang, X., Mohanty, N., McCallum, A.: Group and topic discovery from relations and their attributes. In: *Proc. of NIPS*. pp. 1449–1456 (2006)
35. Wen, Z., Lin, C.: Towards finding valuable topics. In: *Proc. of SDM*. pp. 720–731 (2010)
36. Wolpert, R.L., Clyde, M.A., Tu, C.: Stochastic expansions using continuous dictionaries: Lévy Adaptive Regression Kernels. *Annals of Statistics* (2011)
37. Yang, J., Leskovec, J.: Overlapping community detection at scale: A nonnegative matrix factorization approach. In: *Proc. of WSDM*. pp. 587–596 (2013)
38. Yang, J., McAuley, J.J., Leskovec, J.: Community detection in networks with node attributes. In: *Proc. of ICDM*. pp. 1151–1156 (2013)
39. Yoshida, T.: Toward finding hidden communities based on user profile. *J. Intell. Inf. Syst.* 40(2), 189–209 (Apr 2013)
40. Zhou, M.: Infinite edge partition models for overlapping community detection and link prediction. In: *Proc. of AISTATS* (to appear) (2015)
41. Zhou, M., Carin, L.: Augment-and-conquer negative binomial processes. In: *Proc. of NIPS* (2012)
42. Zhou, M., Carin, L.: Negative binomial process count and mixture modeling. *IEEE Trans. Pattern Analysis and Machine Intelligence* (2015)
43. Zhou, M., Hannah, L., Dunson, D., Carin, L.: Beta-negative binomial process and Poisson factor analysis. In: *Proc. of AISTATS* (2012)
44. Zhu, J.: Max-margin nonparametric latent feature models for link prediction. In: *Proc. of ICML* (2012)