

Network Discovery and Recommendation via Joint Network and Topic Modeling

Ayan Acharya
Dept. of ECE, UT Austin
aacharya@utexas.edu

Dean Teffer
ARL, UT Austin
dean.teffer@arlut.utexas.edu

Mingyuan Zhou
Dept. of IROM, UT Austin
mzhou@utexas.edu

Joydeep Ghosh
Dept. of ECE, UT Austin
ghosh@ece.utexas.edu

ABSTRACT

Many complex social and biological interactions can be represented as graphs. Often such interaction networks also come with count-valued side-information that influence these interactions. Similarly, in applications like recommender systems, the count-valued rating matrix, capturing the interaction between users and items, may also be associated with auxiliary information in the form of a social network of users, which can conveniently be represented as a binary matrix. Jointly modeling both the primary as well as side information can result in better predictive performance. This paper introduces one such model, Joint Gamma Process Poisson Factorization (J-GPPF), which jointly models a binary symmetric interaction matrix \mathbf{B} and a count matrix \mathbf{Y} . The model discovers the ideal number of latent factors for modeling both \mathbf{B} and \mathbf{Y} from the data itself and outperforms strong baselines that do not use the side-information at all. We derive closed form updates for Gibbs sampling and predict the missing values in both \mathbf{B} and \mathbf{Y} in a mathematically consistent way.

Keywords

Network modeling, Poisson factorization, Gamma Process

1. INTRODUCTION

Social networks and other relational datasets can conveniently be represented using a binary symmetric adjacency matrix $\mathbf{B} \in \{0, 1\}^{N \times N}$. 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, all of which can be represented as count-valued side information 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. In this particular setting, \mathbf{B} can be considered as “side information” and used for modeling \mathbf{Y} better. Incorporating such side information can result in bet-

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

Copyright 2015 ACM for this paper by its authors. Copying permitted for private and academic purposes. ...\$15.00.

ter modeling of both \mathbf{B} and \mathbf{Y} when either of these matrices lacks sufficient information.

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:

- We present a fast and effective model that uses both \mathbf{B} and \mathbf{Y} to help discover better network structures and cater better recommendation.
- We perform nonparametric Bayesian modeling for discovering latent structures in both \mathbf{B} and \mathbf{Y} and predict on missing entries in both \mathbf{B} and \mathbf{Y} .
- Our model scales with the number of non-zero entries $S_{\mathbf{B}}$ in the network and the number of non-zero entries $S_{\mathbf{Y}}$ in the count matrix as $O(S_{\mathbf{B}}K_{\mathbf{B}} + S_{\mathbf{Y}}K_{\mathbf{Y}})$, where $K_{\mathbf{B}}$ is the number of network groups and $K_{\mathbf{Y}}$ is the number of latent factors in the count matrix, present at a high value while performing nonparametric modeling.

The remainder of the paper is organized as follows. We present background material in Section 2. J-GPPF and its inference algorithm are explained in Section 3. The related works are presented in Section 4. Experimental results are reported in Section 5, followed by conclusions and scope of future work in Section 6.

2. BACKGROUND

2.1 Negative Binomial Distribution

The definition of the negative binomial (NB) distribution is given in [1]. Here we introduce a few relevant lemmas which we use to derive the Gibbs sampling updates in Section 3.

Lemma 2.1 ([31]). *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:*

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

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 2.2. *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 2.3. If $x_i \sim \text{Pois}(m_i \lambda)$, $\lambda \sim \text{Gam}(r, 1/c)$, then $x = \sum_i x_i \sim \text{NB}(r, p)$, where $p = (\sum_i m_i)/(c + \sum_i m_i)$.

Lemma 2.4. If $x_i \sim \text{Pois}(m_i \lambda)$, $\lambda \sim \text{Gam}(r, 1/c)$, then

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

Lemma 2.5. If $r_i \sim \text{Gam}(a_i, 1/b) \forall i$, $b \sim \text{Gam}(c, 1/d)$, then we have:

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

The proofs of Lemmas 2.3, 2.4 and 2.5 follow from the definitions of Gamma, Poisson and Negative Binomial distributions.

Lemma 2.6. If $x_i \sim \text{Pois}(m_i r_2)$, $r_2 \sim \text{Gam}(r_1, 1/d)$, $r_1 \sim \text{Gam}(a, 1/b)$, then $(r_1 | -) \sim \text{Gam}(a + \ell, 1/(b - \log(1 - p)))$ where $(\ell | x, r_1) \sim \text{CRT}(\sum_i x_i, r_1)$, $p = \sum_i m_i / (d + \sum_i m_i)$. The proof and illustration can be found in Section 3.3 of [1].

2.2 Gamma Process

For non-parametric modeling, we employ Gamma processes, a detailed description of which is provided in [1; 11; 28].

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 and hence can be considered as a special case of the document-author framework, which we pursue to describe the model. 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., $G_{\mathbf{B}}$).

We employ two separate Gamma Processes. The first one models the latent factors in the network. A draw from this Gamma Process $G_{\mathbf{B}} \sim \text{GP}(c_{\mathbf{B}}, H_{\mathbf{B}})$ 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 $\phi_{k_{\mathbf{B}}} \sim \prod_{n=1}^N \text{Gam}(a_{\mathbf{B}}, 1/\sigma_n)$ and $\rho_{k_{\mathbf{B}}} = G_{\mathbf{B}}(\phi_{k_{\mathbf{B}}})$ is the associated weight. The second Gamma Process models the latent groups of side information. A draw from this gamma process $G_{\mathbf{Y}} \sim \text{GP}(c_{\mathbf{Y}}, H_{\mathbf{Y}})$ 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 $\beta_{k_{\mathbf{Y}}} \sim \prod_{w=1}^V \text{Gam}(\xi_{\mathbf{Y}}, 1/\zeta_w)$ and $r_{k_{\mathbf{Y}}} = G_{\mathbf{Y}}(\beta_{k_{\mathbf{Y}}})$ is the associated weight. 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 $H_{\mathbf{Y}}$. The $(n, m)^{\text{th}}$ entry in the matrix \mathbf{B} is assumed to be derived from a latent count as:

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

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; 30]. 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 2.2, 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, providing several potential advantages. First, it is more interpretable because $\rho_{k_{\mathbf{B}}}$ and $\phi_{k_{\mathbf{B}}}$ 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 $S_{\mathbf{B}}$ in the observed binary adjacency matrix \mathbf{B} .

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

$$y_{dw} \sim \text{Pois}(\zeta_{dw}), \zeta_{dw} = \left(\sum_{k_{\mathbf{Y}}} \zeta_{\mathbf{Y}d w k_{\mathbf{Y}}} + \sum_{k_{\mathbf{B}}} \zeta_{\mathbf{B}d w k_{\mathbf{B}}} \right),$$

where $\zeta_{\mathbf{Y}d w k_{\mathbf{Y}}} = r_{k_{\mathbf{Y}}} \theta_{d k_{\mathbf{Y}}} \beta_{w k_{\mathbf{Y}}}$, $Z_{nd} \in \{0, 1\}$ and $Z_{nd} = 1$ if and only if author n is one of the authors of paper d and $\zeta_{\mathbf{B}d w k_{\mathbf{B}}} = \epsilon \rho_{k_{\mathbf{B}}} (\sum_n Z_{nd} \phi_{n k_{\mathbf{B}}}) \psi_{w k_{\mathbf{B}}}$. 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 \mathbf{Y} is factorized in such a way that $y_{dw} \sim \text{Pois}(\sum_{k_{\mathbf{Y}}} r_{k_{\mathbf{Y}}} \theta_{d k_{\mathbf{Y}}} \beta_{w k_{\mathbf{Y}}})$. 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 $r_{k_{\mathbf{Y}}}$ signifies the weight of the $k_{\mathbf{Y}}^{\text{th}}$ factor. The number of latent factors, possibly smaller than both D and V , would be inferred from the data.

In the proposed joint model, \mathbf{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. In the user-rating framework, the entries in \mathbf{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. $\rho_{k_{\mathbf{B}}}$ signifies the interaction strength corresponding to the $k_{\mathbf{B}}^{\text{th}}$ 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 \mathbf{B} and \mathbf{Y} might have different levels of sparsity and the range of integers in \mathbf{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. To complete the generative process, we put Gamma

$$\begin{aligned}
(\phi_{nk_B}|-) &\sim \text{Gam} \left(a_B + \sum_{m=1}^{(n-1)} x_{nmk_B} + \sum_{m=(n+1)}^N x_{nmk_B} + y_{n.k_B}, \left(\sigma_n + \rho_{k_B} \left(\sum_{\substack{m=1 \\ m \neq n}}^N \phi_{mk_B} + \epsilon \sum_{d,w} Z_{nd} \psi_{wk_B} \right) \right)^{-1} \right), \\
(\psi_{wk_B}|-) &\sim \text{Gam} \left(\xi_B + y_{..wk_B}, \frac{1}{\zeta_w + \epsilon \rho_{k_B} \sum_{d,n} Z_{nd} \phi_{nk_B}} \right), (\theta_{dk_Y}|-) \sim \text{Gam} \left(a_Y + y_{d.k_Y}, \frac{1}{\zeta_d + r_{k_Y} \beta_{.k_Y}} \right), (\zeta_d|-) \sim \text{Gam} \left(\alpha_Y + K_Y a_Y, \frac{1}{\epsilon_Y + \theta_d} \right), \\
(\rho_{k_B}|-) &\sim \text{Gam} \left(\frac{\gamma_B}{K_B} + \sum_{\substack{(n,m) \\ n < m}} x_{nmk_B} + y_{..k_B}, \left(c_B + \sum_{\substack{(n,m) \\ n < m}} \phi_{nk_B} \phi_{mk_B} + \epsilon \sum_{n,d,w} Z_{nd} \phi_{nk_B} \psi_{wk_B} \right)^{-1} \right), \\
(r_{k_Y}|-) &\sim \text{Gam} \left(\frac{\gamma_Y}{K_Y} + y_{..k_Y}, \frac{1}{c_Y + \theta_{.k_Y} \beta_{.k_Y}} \right), (\beta_{wk_Y}|-) \sim \text{Gam} \left(\xi_Y + y_{.wk_Y}, \frac{1}{\eta_w + r_{k_Y} \theta_{.k_Y}} \right), \\
(c_B|-) &\sim \text{Gam} \left(g_B + \gamma_B, \frac{1}{h_B + \sum_{k_B} \rho_{k_B}} \right), (c_Y|-) \sim \text{Gam} \left(g_Y + \gamma_Y, \frac{1}{h_Y + \sum_{k_Y} r_{k_Y}} \right), (\epsilon|-) \sim \text{Gam} \left(f_0 + \sum_{k=1}^{K_B} y_{..k}, \left(g_0 + \sum_{k=1}^{K_B} \rho_{k_B} \sum_{n=1}^N |Z_n| \phi_{nk_B} \right)^{-1} \right), \\
(\zeta_w|-) &\sim \text{Gam} \left(a_0 + K_B \xi_B, \frac{1}{b_0 + \psi_w} \right), (\eta_w|-) \sim \text{Gam} \left(c_0 + K_Y \xi_Y, \frac{1}{d_0 + \beta_w} \right),
\end{aligned}$$

Table 1: Gibbs sampling updates in J-GPPF

$$\begin{aligned}
(\phi_{nk_B}|-) &\sim \text{Gam} \left(a_B + \sum_{\substack{m=1 \\ (n,m) \notin \mathcal{M}_B}}^{(n-1)} x_{nmk_B} + \sum_{\substack{m=(n+1) \\ (n,m) \notin \mathcal{M}_B}}^N x_{nmk_B} + \sum_{\substack{d,w \\ (d,w) \notin \mathcal{M}_Y}} y_{dnwk_B}, \left(\sigma_n + \rho_{k_B} \left(\sum_{\substack{m=1 \\ (n,m) \notin \mathcal{M}_B}}^N \phi_{mk_B} + \epsilon \sum_{\substack{d,w \\ (d,w) \notin \mathcal{M}_Y}} Z_{nd} \psi_{wk_B} \right) \right)^{-1} \right), \\
(\psi_{wk_B}|-) &\sim \text{Gam} \left(\xi_B + \sum_{\substack{d \\ (d,w) \notin \mathcal{M}_Y}} y_{dwk_B}, \left(\eta_w + \epsilon \rho_{k_B} \sum_{\substack{n,d \\ (d,w) \notin \mathcal{M}_Y}} Z_{nd} \phi_{nk_B} \right)^{-1} \right), \\
(\rho_{k_B}|-) &\sim \text{Gam} \left(\frac{\gamma_B}{K_B} + \sum_{\substack{n=1, n < m \\ (n,m) \notin \mathcal{M}_B}}^{(N-1)} x_{nmk_B} + \sum_{\substack{n,d,w \\ (d,w) \notin \mathcal{M}_Y}} y_{dnwk_B}, \left(c_B + \sum_{\substack{n=1, n < m \\ (n,m) \notin \mathcal{M}_B}}^{(N-1)} \phi_{nk_B} \phi_{mk_B} + \epsilon \sum_{\substack{n,d,w \\ (d,w) \notin \mathcal{M}_Y}} Z_{nd} \phi_{nk_B} \psi_{wk_B} \right)^{-1} \right).
\end{aligned}$$

Table 2: Sampling of ϕ_{nk_B} , ψ_{wk_B} , ρ_{k_B} in J-GPPF with missing entries

priors over c_B , c_Y , σ_n , ζ_d and ϵ as:

$$\begin{aligned}
c_B &\sim \text{Gam}(g_B, 1/h_B), c_Y \sim \text{Gam}(g_Y, 1/h_Y), \\
\epsilon &\sim \text{Gam}(g_0, 1/f_0), \\
\sigma_n &\sim \text{Gam}(\alpha_B, 1/\epsilon_B), \zeta_d \sim \text{Gam}(\alpha_Y, 1/\epsilon_Y).
\end{aligned}$$

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. 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 $\rho_{k_B} \sim \text{Gam}(\gamma_B/K_B, 1/c_B)$ and $r_{k_Y} \sim \text{Gam}(\gamma_Y/K_Y, 1/c_Y)$. Such approximation approaches the original infinite model as both K_B and K_Y approach infinity.

Sampling of $(x_{nmk_B})_{k_B=1}^{K_B}$: First, the total latent count corresponding to the non-zero entries can be derived as:

$$(x_{nm}|-) \sim b_{nm} \text{Pois}_+ \left(\sum_{k_B=1}^{K_B} \lambda_{nmk_B} \right). \quad (5)$$

After which, following Lemma 2.2 one can derive:

$$((x_{nmk_B})_{k_B=1}^{K_B}|-) \sim \text{Mult} \left(x_{nm}, \left(\frac{\lambda_{nmk_B}}{\sum_{k_B=1}^{K_B} \lambda_{nmk_B}} \right)_{k_B=1}^{K_B} \right). \quad (6)$$

Sampling of $(y_{dwk_B})_k$: Again, following Lemma 2.2, we have:

$$\begin{aligned}
&\left((y_{dwk_Y})_{k_Y=1}^{K_Y}, (y_{dnwk_B})_{k_B=1, n \in \mathcal{Z}_d}^{K_B} \right) \sim \quad (7) \\
&\text{Mult} \left(y_{dw}, \frac{\{\zeta_{dwk_Y}\}_{k_Y}, \{\zeta_{dnwk_B}\}_{n \in \mathcal{Z}_d, k_B}}{\sum_{k_Y} \zeta_{dwk_Y} + \sum_{n \in \mathcal{Z}_d} \sum_{k_B} \zeta_{dnwk_B}} \right).
\end{aligned}$$

Sampling of ϕ_{nk_B} , ψ_{wk_B} , ρ_{k_B} , θ_{dk_Y} , β_{wk_Y} , r_{k_Y} and ϵ : Sampling of these parameters follow from Lemma 2.4 and are given in Table 1. 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 σ_n , ζ_d , ϵ , ζ_w , η_w , c_B and c_Y : Sampling of these parameters follow from Lemma 2.5 and are given in Table 1.

Sampling of γ_B : Using Lemma 2.2, one can show that $x_{..k_B} \sim \text{Pois}(\rho_{k_B})$. Integrating ρ_{k_B} and using Lemma 2.4, 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 2.4, 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 2.6 sample:

$$(\gamma_B|-) \sim \text{Gam} \left(e_B + \sum_{k_B} l_{k_B}, (f_B - q_B)^{-1} \right), \quad (8)$$

where $q_B = \sum_{k_B} \log(c_B / (c_B + \sum_n \phi_{nk_B} \phi_{k_B}^{-n})) / K_B$.

Sampling of γ_Y : Using Lemma 2.2, one can show that $y_{..(K_B+k_Y)} \sim$

$$\begin{aligned}
(\theta_{dk_Y} | -) &\sim \text{Gam} \left(a_Y + \sum_{w:(d,w) \notin \mathcal{M}_Y} y_{dwk_Y}, \left(\zeta_d + r_{k_Y} \sum_{w:(d,w) \notin \mathcal{M}_Y} \beta_{wk_Y} \right)^{-1} \right), \\
(\beta_{wk_Y} | -) &\sim \text{Gam} \left(\xi_Y + \sum_{d:(d,w) \notin \mathcal{M}_Y} y_{dwk_Y}, \left(\eta_w + r_{k_Y} \sum_{d:n:(d,w) \notin \mathcal{M}_Y} Z_{dn} \phi_{nk_Y} \right)^{-1} \right), \\
(r_{k_Y} | -) &\sim \text{Gam} \left(\frac{\gamma_Y}{K_Y} + \sum_{d,w:(d,w) \notin \mathcal{M}_Y} y_{dwk_Y}, \left(c_Y + \sum_{d,w:(d,w) \notin \mathcal{M}_Y} \theta_{dk_Y} \beta_{wk_Y} \right)^{-1} \right).
\end{aligned}$$

Table 3: Sampling of $\theta_{dk_Y}, \beta_{wk_Y}, r_{k_Y}$ in J-GPPF with missing entries

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

Table 4: Generative Process of N-GPPF

$$\begin{aligned}
y_{dw} &\sim \text{Pois} \left(\sum_{k_Y=1}^{\infty} r_{k_Y} \theta_{dk_Y} \beta_{wk_Y} \right), \\
\theta_{k_Y} &\sim \prod_{d=1}^D \text{Gam}(a_Y, 1/\zeta_d), \beta_{k_Y} \sim \prod_{w=1}^V \text{Gam}(\xi_Y, 1/\zeta_w), \\
\zeta_d &\sim \text{Gam}(\alpha, 1/\varepsilon), r_{k_Y} \sim \text{Gam}(\gamma_Y/K_Y, 1/c_Y), \\
\gamma_{Y_Y} &\sim \text{Gam}(e_Y, 1/f_Y), c_Y \sim \text{Gam}(g_Y, 1/h_Y).
\end{aligned}$$

Table 5: Generative Process of C-GPPF

$\text{Pois}(r_{k_Y})$ and after integrating r_{k_Y} and using Lemma 2.4, 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 2.6 sample:

$$(\gamma_Y | -) \sim \text{Gam} \left(e_Y + \sum_{k_Y} m_{k_Y}, (f_Y - q_Y)^{-1} \right), \quad (9)$$

where $q_Y = \sum_{k_Y} \log(c_Y/(c_Y + \theta_{.k_Y}))/K_Y$.

3.2 Gibbs Sampling for J-GPPF with Missing Entries

Parameters whose update get affected in presence of missing entries are $\rho_{k_B}, \phi_{nk_B}, \psi_{wk_B}, r_{k_Y}, \theta_{dk_Y}, \beta_{wk_Y}$. Sampling of these parameters follow from Lemma 2.4 and are given in Table 2 and 3. Here \mathcal{M}_B and \mathcal{M}_Y denote the set of missing entries in B and Y respectively.

3.3 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 B is modeled without the auxiliary matrix Y . The update equations of variables corresponding to N-GPPF can be obtained with $Z = \mathbf{0}$. Another special case of J-GPPF appears when only the count matrix Y is modeled without using the contribution from the network matrix B . The generative model of N-GPPF and C-GPPF are given in Table 4 and 5 respectively.

3.4 Computation Complexity

The Gibbs sampling updates of J-GPPF can be calculated in $O(K_B S_B + (K_B + K_Y) S_Y + N K_B + D K_Y + V(K_B + K_Y))$ time, where S_B is the number of non-zero entries in B and S_Y is the number of non-zero entries in Y . It is obvious that for large matrices the computation is primarily of the order of $K_B S_B + (K_B + K_Y) S_Y$. Such complexity is a huge saving when compared to other methods like MMSB [2], that only models B and incurs computation cost of $O(N^2 K_B)$; and standard matrix factorization approaches [25] that work with the matrix Y and incur $O(DV K_Y)$ computation cost. Interestingly, the inference in [13] incurs cost $O(K_Y^2 D + K_Y V + K_Y S_Y)$ with K_Y signifying the

termination point of stick breaking construction in their model. C-GPPF incurs computation cost $O(DK_Y + K_Y S_Y + VK_Y)$, an apparent improvement over that of [13].

4. RELATED WORK

The Infinite Relational Model (IRM [14]) 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 but the computational complexity of the underlying inference mechanism is $O(N^2)$. Such quadratic computation complexity is also a problem with many other existing latent variable network models, such as the latent feature relational model [21] and its max margin version [34], and the infinite latent attribute model [23]. Some of the existing approaches handle sparsity in real-world networks by using some auxiliary information [16; 20; 29]. Recommender system and text mining researchers, in contrast, tend to take an orthogonal approach. In recommender systems [9; 17], Y may represent a user-by-item rating matrix and the objective in this setting is to predict the missing entries in Y , and the social network matrix B plays a secondary role in providing auxiliary information to facilitate this task [17]. Similarly, in the text mining community, many existing models [4; 19; 22; 27] use the network information or other forms of side information to improve the discovery of "topics" from the document-by-word matrix Y . The matrix B can represent, for example, the interaction network of authors participating in writing the documents. The Relational Topic Model [10] discovers links between documents based on their topic distributions, obtained through unsupervised exploration. The Author-Topic framework (AT [24]) and the Author-Recipient-Topic model (ART [19]) jointly model documents along with the authors of the documents. Block-LDA [4], on the other hand, provides a generative model for the links between authors and recipients in addition to documents. 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)

[33], 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, such as non-negative matrix factorization [8; 15], gamma-Poisson model [7; 26], LDA [6], and gamma-NB processes [32; 33], although constructed with different motivations and for distinct problems, can all be viewed as PFA with different prior distributions imposed on Φ and Θ . 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 [17]). 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. The collaborative topic Poisson factorization (CTPF) framework proposed in [12] 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. Although both models make use of Poisson factorization to infer low-rank matrices in order to recommend items to users, J-GPPF is a fundamentally different model. As recently summarized in [5], there are many useful approaches for employing social side information to improve recommendation. For example, in both [3; 18], the authors define a dependence for item recommendations based on each node’s neighbors or community. J-GPPF is a joint model of both the social network and the item ratings, and solves for the latent space factorization of each, and missing elements of each, simultaneously. In addition to this innovation, model inference scales with the number of observed elements in each matrix and the number of latent groups, not the full dimension of each matrix.

5. EXPERIMENTAL RESULTS

5.1 NIPS Authorship Network

This dataset contains a list of all papers and authors from NIPS 1988 to 2003. We took the 234 authors who had published with the most other people and looked at their co-authorship information. After standard pre-processing and removing words that appear less than 50 times in the over-all corpus corresponding to these users, the number of users in the graph who write at least one document, is 225 and the total number of unique words is 1354. The total number of documents is 1165.

5.2 GoodReads Data

Using the Goodreads API, we collect a base set of users with recent activity on the website. For each user in the base set, the user’s friends as well as friends of friends on the site are collected (two hops in the graph). This process is repeated over a 24-hour time period, with a new base set constructed each time (*i.e.* friends are not polled recursively). By running for a full day, multiple time zones are covered and the reviews are collected for all identified users, with a maximum of 200 reviews per user. Each review consists of a book ID and a rating from 0 to 5. Similar dataset has also been used in [9]. After standard pre-processing and removing words that appear less than 10 times in the over-all corpus, the number of users in the graph is 84 and the total number of unique words is 189.

5.3 Experimental Setup and Results

In all the experiments, we initialize ϵ to 2 and let the sampler decide what value works best for joint modeling. We use $K_{\mathbf{B}} = K_{\mathbf{Y}} = 50$ and initialize all the hyper-parameters to 1. In the first set

of experiments, for each dataset, we hold out data from \mathbf{B} only and ran 20 different experiments and display the mean AUC and one standard error. In this setup, we consider N-GPPF, the infinite relational model (IRM) of [14] and the Mixed Membership Stochastic Block Model (MMSB) [2] as the baseline algorithms. Fig. 1 and 2 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 \mathbf{B} and \mathbf{Y} .

In the second set of experiments, we hold out data from \mathbf{Y} only and run 20 different experiments and display the precision@top-20 for J-GPPF. This evaluation is structured along the lines of the work in [13]. We calculate the intersection of the top 20 predicted set of words (arranged in the decreasing order of counts) and the top 20 words in a document and divide the number by 20 to get the precision for each document. We then calculate mean average precision (MAP) by taking the average of the precision over all the documents. C-GPPF and the hierarchical Poisson matrix factorization (HPMF) [13] are considered as the baselines, both of which model only \mathbf{Y} . Fig. 3 and 4 show that \mathbf{B} helps in boosting the predictive performance in J-GPPF over a wide range of fractions of the data that is held out from \mathbf{Y} .

6. 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 \mathbf{B} , \mathbf{Y} and in latent group membership. We derived an efficient MCMC inference method, and compared our approach to several popular baselines that either work on \mathbf{B} or \mathbf{Y} .

Acknowledgement

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

References

- [1] A. Acharya, J. Ghosh, and M. Zhou. Nonparametric Bayesian Factor Analysis for Dynamic Count Matrices. In *Proc. of AISTATS*, pages 1–9, 2015.
- [2] E. M. Airoldi, D. M. Blei, S. E. Fienberg, and E. P. Xing. Mixed membership stochastic blockmodels. *JMLR*, 9:1981–2014, jun 2008.
- [3] J. Aranda, I. Givoni, J. Handcock, and D. Tarlow. An online social network-based recommendation system. *Toronto, Ontario, Canada*, 2007.
- [4] R. Balasubramanian and W. W. Cohen. Block-LDA: Jointly modeling entity-annotated text and entity-entity links. In *Proc. of SDM*, pages 450–461, 2011.
- [5] D. Bernardes, M. Diaby, R. Fournier, F. FogelmanSoulié, and E. Vennet. A social formalism and survey for recommender systems. *SIGKDD Explor. Newsl.*, 16(2):20–37, may 2015.
- [6] D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent Dirichlet Allocation. *JMLR*, 3:993–1022, 2003.
- [7] J. Canny. Gap: a factor model for discrete data. In *SIGIR*, 2004.
- [8] A. T. Cemgil. Bayesian inference for nonnegative matrix factorisation models. *Intell. Neuroscience*, 2009:4:1–4:17, Jan. 2009.

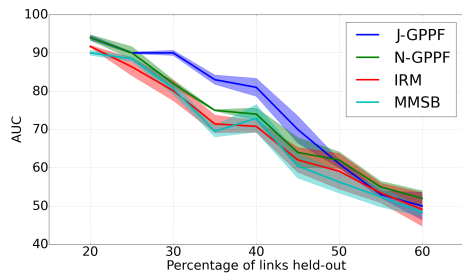


Figure 1: AUC NIPS

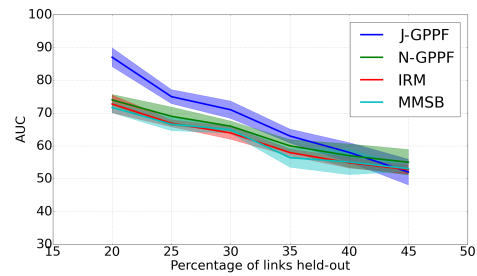


Figure 2: AUC GoodReads

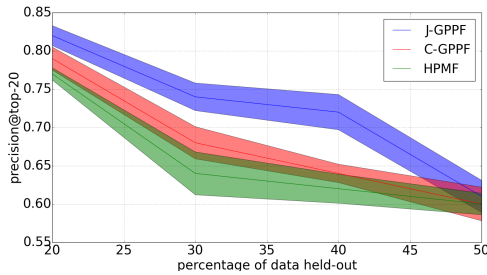


Figure 3: MAP NIPS

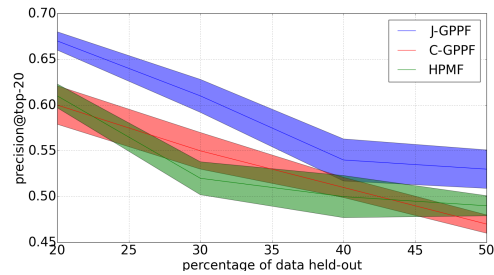


Figure 4: MAP GoodReads

- [9] A. Chaney, P. Gopalan, and D. Blei. Poisson trust factorization for incorporating social networks into personalized item recommendation. In *NIPS Workshop: What Difference Does Personalization Make?*, 2013.
- [10] J. Chang and D. Blei. Relational topic models for document networks. In *Proc. of AISTATS*, 2009.
- [11] T. S. Ferguson. A Bayesian analysis of some nonparametric problems. *Ann. Statist.*, 1973.
- [12] P. Gopalan, L. Charlin, and D. Blei. Content-based recommendations with poisson factorization. In *Proc. of NIPS*, pages 3176–3184. 2014.
- [13] P. Gopalan, F. Ruiz, R. Ranganath, and D. Blei. Bayesian nonparametric poisson factorization for recommendation systems. In *Proc. of AISTATS*, 2014.
- [14] C. Kemp, J. Tenenbaum, T. Griffiths, T. Yamada, and N. Ueda. Learning systems of concepts with an infinite relational model. In *Proc. of AAI*, pages 381–388, 2006.
- [15] D. D. Lee and H. S. Seung. Algorithms for non-negative matrix factorization. In *NIPS*, 2001.
- [16] J. Leskovec and J. Julian. Learning to discover social circles in ego networks. In *Proc. of NIPS*, pages 539–547. 2012.
- [17] H. Ma, H. Yang, M. R. Lyu, and I. King. Sorec: Social recommendation using probabilistic matrix factorization. In *Proc. of CIKM*, pages 931–940, 2008.
- [18] H. Ma, D. Zhou, C. Liu, M. Lyu, and I. King. Recommender systems with social regularization. In *Proc. of WSDM*, pages 287–296, 2011.
- [19] A. McCallum, X. Wang, and A. Corrada-Emmanuel. Topic and role discovery in social networks with experiments on enron and academic email. *J. Artif. Int. Res.*, 30(1):249–272, Oct. 2007.
- [20] A. Menon and C. Elkan. Link prediction via matrix factorization. In *Machine Learning and Knowledge Discovery in Databases*, volume 6912 of *Lecture Notes in Computer Science*, pages 437–452. Springer Berlin / Heidelberg, 2011.
- [21] K. T. Miller, T. L. Griffiths, and M. I. Jordan. Nonparametric latent feature models for link prediction. In *Proc. of NIPS*, pages 1276–1284, 2009.
- [22] R. Nallapati, A. Ahmed, E. Xing, and W. Cohen. Joint latent topic models for text and citations. In *Proc. of KDD*, pages 542–550, 2008.
- [23] K. Palla, Z. Ghahramani, and D. A. Knowles. An infinite latent attribute model for network data. In *Proc. of ICML*, pages 1607–1614, 2012.
- [24] M. Rosen-Zvi, T. Griffiths, M. Steyvers, and P. Smyth. The author-topic model for authors and documents. In *Proc. of UAI*, pages 487–494, 2004.
- [25] R. Salakhutdinov and A. Mnih. Probabilistic matrix factorization. In *Proc. of NIPS*, 2007.
- [26] M. K. Titsias. The infinite gamma-Poisson feature model. In *Proc. of NIPS*, 2008.
- [27] Z. Wen and C. Lin. Towards finding valuable topics. In *Proc. of SDM*, pages 720–731, 2010.
- [28] R. L. Wolpert, M. A. Clyde, and C. Tu. Stochastic expansions using continuous dictionaries: Lévy Adaptive Regression Kernels. *Annals of Statistics*, 2011.
- [29] T. Yoshida. Toward finding hidden communities based on user profile. *J. Intell. Inf. Syst.*, 40(2):189–209, Apr. 2013.
- [30] M. Zhou. Infinite edge partition models for overlapping community detection and link prediction. In *Proc. of AISTATS*, pages 1135–1143. 2015.
- [31] M. Zhou and L. Carin. Augment-and-conquer negative binomial processes. In *Proc. of NIPS*, 2012.
- [32] M. Zhou and L. Carin. Negative binomial process count and mixture modeling. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 2015.
- [33] M. Zhou, L. Hannah, D. Dunson, and L. Carin. Beta-negative binomial process and Poisson factor analysis. In *Proc. of AISTATS*, pages 1462–1471, 2012.
- [34] J. Zhu. Max-margin nonparametric latent feature models for link prediction. In *Proc. of ICML*, 2012.