

Discriminative Relational Topic Models

Ning Chen, Jun Zhu, *Member, IEEE*, Fei Xia, and Bo Zhang

Abstract—Many scientific and engineering fields involve analyzing network data. For document networks, relational topic models (RTMs) provide a probabilistic generative process to describe both the link structure and document contents, and they have shown promise on predicting network structures and discovering latent topic representations. However, existing RTMs have limitations in both the restricted model expressiveness and incapability of dealing with imbalanced network data. To expand the scope and improve the inference accuracy of RTMs, this paper presents three extensions: 1) unlike the common link likelihood with a diagonal weight matrix that allows the-same-topic interactions only, we generalize it to use a full weight matrix that captures all pairwise topic interactions and is applicable to asymmetric networks; 2) instead of doing standard Bayesian inference, we perform regularized Bayesian inference (RegBayes) with a regularization parameter to deal with the imbalanced link structure issue in common real networks and improve the discriminative ability of learned latent representations; and 3) instead of doing variational approximation with strict mean-field assumptions, we present collapsed Gibbs sampling algorithms for the generalized relational topic models by exploring data augmentation without making restricting assumptions. Under the generic RegBayes framework, we carefully investigate two popular discriminative loss functions, namely, the logistic log-loss and the max-margin hinge loss. Experimental results on several real network datasets demonstrate the significance of these extensions on improving the prediction performance, and the time efficiency can be dramatically improved with a simple fast approximation method.

Index Terms—statistical network analysis, relational topic models, data augmentation, regularized Bayesian inference



arXiv:1310.2409v1 [cs.LG] 9 Oct 2013

1 INTRODUCTION

MANY scientific and engineering fields involve analyzing large collections of data that can be well described by networks, where vertices represent entities and edges represent relationships or interactions between entities; and to name a few, such data include online social networks, communication networks, protein interaction networks, academic paper citation and coauthorship networks, etc. As the availability and scope of network data increase, statistical network analysis (SNA) has attracted a considerable amount of attention (see [17] for a comprehensive survey). Among the many tasks studied in SNA, link prediction [25], [4] is a most fundamental one that attempts to estimate the link structure of networks based on partially observed links and/or entity attributes (if exist). Link prediction could provide useful predictive models for suggesting friends to social network users or citations to scientific articles.

Many link prediction methods have been proposed, including the early work on designing good similarity measures [25] that are used to rank unobserved links and those on learning supervised classifiers with well-conceived features [19], [26]. Though specific domain knowledge can be used to design effective feature representations, feature engineering is generally a labor-intensive process. In order to expand the scope and

ease of applicability of machine learning methods, fast growing interests have been spent on learning feature representations from data [6]. Along this line, recent research on link prediction has focused on learning latent variable models, including both parametric [20], [21], [2] and nonparametric Bayesian methods [31], [41]. Though these methods could model the network structures well, little attention has been paid to account for observed attributes of the entities, such as the text contents of papers in a citation network or the contents of web pages in a hyperlinked network. One work that accounts for both text contents and network structures is the relational topic models (RTMs) [8], an extension of latent Dirichlet allocation (LDA) [7] to predicting link structures among documents as well as discovering their latent topic structures.

Though powerful, existing RTMs have some assumptions that could limit their applicability and inference accuracy. First, RTMs define a symmetric link likelihood model with a diagonal weight matrix that allows the-same-topic interactions only, and the symmetric nature could also make RTMs unsuitable for asymmetric networks. Second, by performing standard Bayesian inference under a generative modeling process, RTMs do not explicitly deal with the common imbalance issue in real networks, which normally have only a few observed links while most entity pairs do not have links, and the learned topic representations could be weak at predicting link structures. Finally, RTMs and other variants [27] apply variational methods to estimate model parameters with mean-field assumptions [24], which are normally too restrictive to be realistic in practice.

To address the above limitations, this paper

• N. Chen[†], J. Zhu[†], F. Xia[‡] and B. Zhang[†] are with the Department of Computer Science and Technology, National Lab of Information Science and Technology, State Key Lab of Intelligent Technology and Systems, Tsinghua University, Beijing, 100084 China.
E-mail: [†]{ningchen, dcszj, dcszb}@mail.tsinghua.edu.cn,
[‡]xia.fe09@gmail.com.

presents discriminative relational topic models, which consist of three extensions to improving RTMs:

- 1) we relax the symmetric assumption and define the generalized relational topic models (gRTMs) with a full weight matrix that allows all pairwise topic interactions and is more suitable for asymmetric networks;
- 2) we perform regularized Bayesian inference (Reg-Bayes) [43] that introduces a regularization parameter to deal with the imbalance problem in common real networks;
- 3) we present a collapsed Gibbs sampling algorithm for gRTMs by exploring the classical ideas of data augmentation [11], [40], [14].

Our methods are quite generic, in the sense that we can use various loss functions to learn discriminative latent representations. In this paper, we particularly focus on two types of popular loss functions, namely, logistic log-loss and max-margin hinge loss. For the max-margin loss, the resulting max-margin RTMs are themselves new contributions to the field of statistical network analysis.

For posterior inference, we present efficient Markov Chain Monte Carlo (MCMC) methods for both types of loss functions by introducing auxiliary variables. Specifically, for the logistic log-loss, we introduce a set of Polya-Gamma random variables [34], one per training link, to derive an exact mixture representation of the logistic link likelihood; while for the max-margin hinge loss, we introduce a set of generalized inverse Gaussian variables [12] to derive a mixture representation of the corresponding unnormalized pseudo-likelihood. Then, we integrate out the intermediate Dirichlet variables and derive the local conditional distributions for collapsed Gibbs sampling analytically. These “augment-and-collapse” algorithms are simple and efficient. More importantly, they do not make any restricting assumptions on the desired posterior distribution. Experimental results on several real networks demonstrate that these extensions are important and can significantly improve the performance.

The rest paper is structured as follows. Section 2 summarizes the related work. Section 3 presents the generalized RTMs with both the log-loss and hinge loss. Section 4 presents the “augment-and-collapse” Gibbs sampling algorithms for both types of loss functions. Section 5 presents experimental results. Finally, Section 6 concludes with future directions discussed.

2 RELATED WORK

Probabilistic latent variable models, e.g., latent Dirichlet allocation (LDA) [7], have been widely developed for modeling link relationships between documents, as they share nice properties on dealing with missing attributes as well as discovering representative latent structures. For instance, RTMs [8] capture both

text contents and network relations for document link prediction; Topic-Link LDA [27] performs topic modeling and author community discovery in one unified framework; Link-PLSA-LDA [32] combines probabilistic latent semantic analysis (PLSA) [23] and LDA into a single framework to explicitly model the topical relationship between documents; Others include Pairwise-Link-LDA [33], Copycat and Citation Influence models [13], Latent Topic Hypertext Models (LTHM) [1], Block-LDA models [5], etc. One shared goal of the aforementioned models is link prediction. For static networks, our focus in this paper, this problem is usually formulated as inferring the missing links given the other observed ones. However, very few work explicitly imposes discriminative training, and many models suffer from the common imbalance issue in sparse networks (e.g., the number of unobserved links is much larger than that of the observed ones). In this paper, we build our approaches by exploring the nice framework of regularized Bayesian inference (RegBayes) [44], under which one could easily introduce posterior regularization and do discriminative training in a cost sensitive manner.

Another under-addressed problem in most probabilistic topic models for link prediction [8], [27] is the intractability of posterior inference due to the non-conjugacy between the prior and link likelihood (e.g., logistic likelihood). Existing approaches using variational inference with mean field assumption are often too restrictive in practice. Recently, [34] and [35] show that by making use of the ideas of data augmentation, the intractable likelihood (either a logistic likelihood or the one induced from a hinge loss) could be expressed as a marginal of a higher-dimensional distribution with augmented variables that leads to a scale mixture of Gaussian components. These strategies have been successfully explored to develop efficient Gibbs samplers for supervised topic models [45], [42]. This paper further explores data augmentation techniques to do collapsed Gibbs sampling for our discriminative relational topic models. Please note that our methods could also be applied to many of the aforementioned relational latent variable models. Finally, this paper is a systematical generalization of the conference paper [9].

3 GENERALIZED RTMS

We consider document networks with binary link structures. Let $\mathcal{D} = \{(\mathbf{w}_i, \mathbf{w}_j, y_{ij})\}_{(i,j) \in \mathcal{I}}$ be a labeled training set, where $\mathbf{w}_i = \{w_{in}\}_{n=1}^{N_i}$ denote the words within document i and the response variable y_{ij} takes values from the binary output space $\mathcal{Y} = \{0, 1\}$. A relational topic model (RTM) consists of two parts — an LDA model [7] for describing the words $\mathbf{W} = \{\mathbf{w}_i\}_{i=1}^D$ and a classifier for considering link structures $\mathbf{y} = \{y_{ij}\}_{(i,j) \in \mathcal{I}}$. Let K be the number of topics and each topic Φ_k is a multinomial distribution over a V -word vocabulary. For Bayesian RTMs, the topics are

TABLE 1

Learned diagonal weight matrix of 10-topic RTM and representative words corresponding with topics.

36.6	36.9	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1	learning, bound, PAC, hypothesis, algorithm
17.9	0.0	-74.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	2	numerical, solutions, extensions, approach, remark
-0.9	0.0	0.0	34.8	0.0	0.0	0.0	0.0	0.0	0.0	3	mixtures, experts, EM, Bayesian, probabilistic
-19.6	0.0	0.0	0.0	44.1	0.0	0.0	0.0	0.0	0.0	4	features, selection, case-based, networks, model
-38.4	0.0	0.0	0.0	0.0	42.5	0.0	0.0	0.0	0.0	5	planning, learning, acting, reinforcement, dynamic
-57.1	0.0	0.0	0.0	0.0	0.0	41.1	0.0	0.0	0.0	6	genetic, algorithm, evolving, evolutionary, learning
	0.0	0.0	0.0	0.0	0.0	0.0	-61.1	0.0	0.0	7	plateau, feature, performance, sparse, networks
	0.0	0.0	0.0	0.0	0.0	0.0	0.0	29.3	0.0	8	modulo, schedule, parallelism, control, processor
	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	41.3	9	neural, cortical, networks, learning, feedforward
	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	10	markov, models, monte, carlo, Gibbs, sampler

TABLE 2

Learned weight matrix of 10-topic gRTM and representative words corresponding with topics.

28.3	20.0	-4.5	-6.2	3.0	-8.1	-9.0	-8.6	-10.8	-7.4	-10.8	1	genetic, evolving, algorithm, coding, programming
21.4	-5.6	21.6	-3.4	1.8	-3.6	-0.5	-10.5	-10.6	-5.3	-9.2	2	logic, grammars, FOIL, EBG, knowledge, clauses
14.6	-5.9	-4.8	21.0	5.9	-7.4	-8.1	-9.4	-9.9	-8.8	-7.5	3	reinforcement, learning, planning, act, exploration
7.7	2.9	3.0	5.3	16.7	3.6	6.9	13.0	5.3	1.3	3.7	4	mixtures, EM, Bayesian, networks, learning, genetic
0.9	-8.1	-5.0	-6.6	2.4	23.1	-8.3	-3.4	-7.6	-7.2	-8.8	5	images, visual, scenes, mixtures, networks, learning
-5.9	-8.9	-1.9	-8.4	7.4	-8.5	29.6	-11.3	-10.0	-6.4	-8.2	6	decision-tree, rules, induction, learning, features
	-8.5	-9.7	-9.8	14.1	-4.4	-10.5	22.3	-6.9	-10.6	-11.2	7	wake-sleep, learning, networks, cortical, inhibition
	-10.5	-11.7	-9.4	5.8	-6.9	-10.3	-7.1	24.7	-8.0	-7.8	8	monte, carlo, hastings, markov, chain, sampler
	-7.7	-5.7	-7.7	-0.3	-6.9	-5.1	-11.5	-7.1	28.3	-7.6	9	case-based, reasoning, CBR, event-based, cases
	-11.2	-7.0	-8.7	2.7	-9.1	-7.9	-12.1	-8.8	-8.3	30.0	10	markov, learning, bayesian, networks, distributions

samples drawn from a prior, e.g., $\Phi_k \sim \text{Dir}(\beta)$, a Dirichlet distribution. The generating process can be described as

- 1) For each document $i = 1, 2, \dots, D$:
 - a) draw a topic mixing proportion $\theta_i \sim \text{Dir}(\alpha)$
 - b) for each word $n = 1, 2, \dots, N_i$:
 - i) draw a topic assignment $z_{in} \sim \text{Mult}(\theta_i)$
 - ii) draw the observed word $w_{in} \sim \text{Mult}(\Phi_{z_{in}})$
- 2) For each pair of documents $(i, j) \in \mathcal{L}$:
 - a) draw a link indicator $y_{ij} \sim p(\cdot | \mathbf{z}_i, \mathbf{z}_j, \boldsymbol{\eta})$, where $\mathbf{z}_i = \{z_{in}\}_{n=1}^{N_i}$.

We have used $\text{Mult}(\cdot)$ to denote a multinomial distribution; and used $\Phi_{z_{in}}$ to denote the topic selected by the non-zero entry of z_{in} , a K -dimensional binary vector with only one entry equaling to 1.

Previous work has defined the link likelihood as

$$p(y_{ij} = 1 | \mathbf{z}_i, \mathbf{z}_j, \boldsymbol{\eta}) = \sigma(\boldsymbol{\eta}^\top (\bar{\mathbf{z}}_i \circ \bar{\mathbf{z}}_j)), \quad (1)$$

where $\bar{\mathbf{z}}_i = \frac{1}{N_i} \sum_{n=1}^{N_i} \mathbf{z}_{in}$ is the average topic assignments of document i ; σ is the sigmoid function; and \circ denotes elementwise product. In [8], other choices of σ such as the exponential function and the cumulative distribution function of the normal distribution were also used, as long as it is a monotonically increasing function with respect to the weighted inner product between $\bar{\mathbf{z}}_i$ and $\bar{\mathbf{z}}_j$. Here, we focus on the commonly used logistic likelihood model [31], [27], as no one has

shown consistently superior performance than others.

3.1 The Full RTM Model

Since $\boldsymbol{\eta}^\top (\bar{\mathbf{z}}_i \circ \bar{\mathbf{z}}_j) = \bar{\mathbf{z}}_i^\top \text{diag}(\boldsymbol{\eta}) \bar{\mathbf{z}}_j$, the standard RTM learns a diagonal weight matrix which only captures the-same-topic interactions (i.e., there is a non-zero contribution to the link likelihood only when documents i and j have the same topic). One example of the fitted diagonal matrix on the Cora citation network [8] is shown in Table 1, where each row corresponds to a topic and we show the representative words for the topic at the right hand side. Due to the positiveness of the latent features (i.e., $\bar{\mathbf{z}}_i$) and the competition between the diagonal entries, some of η_k will have positive values while some are negative. The negative interactions may conflict our intuitions of understanding a citation network, where we would expect that papers with the same topics tend to have citation links. Furthermore, by using a diagonal weight matrix, the model is symmetric, i.e., the probability of a link from document i to j is the same as the probability of a link from j to i . The symmetry property does not hold for many networks, e.g., citation networks.

To make RTMs more expressive and applicable to asymmetric networks, the first simple extension is to define the link likelihood as

$$p(y_{ij} = 1 | \mathbf{z}_i, \mathbf{z}_j, U) = \sigma(\bar{\mathbf{z}}_i^\top U \bar{\mathbf{z}}_j), \quad (2)$$

using a full $K \times K$ weight matrix U . Using the algorithm to be presented, an example of the learned U matrix on the same Cora citation network is shown in Table 2. We can see that by allowing all pairwise topic interactions, all the diagonal entries are positive, while most off-diagonal entries are negative. This is consistent with our intuition that documents with the same topics tend to have citation links, while documents with different topics are less likely to have citation links. We also note that there are some documents with generic topics (e.g., topic 4) that have positive link interactions with almost all others.

3.2 Regularized Bayesian Inference

Given \mathcal{D} , we let $\mathbf{Z} = \{\mathbf{z}_i\}_{i=1}^D$ and $\Theta = \{\theta_i\}_{i=1}^D$ denote all the topic assignments and mixing proportions respectively. To fit RTM models, maximum likelihood estimation (MLE) has been used with an EM algorithm [8]. We consider Bayesian inference [21], [31] to get the posterior distribution

$$p(\Theta, \mathbf{Z}, \Phi, U | \mathcal{D}) \propto p_0(\Theta, \mathbf{Z}, \Phi, U) p(\mathcal{D} | \mathbf{Z}, \Phi, U),$$

where $p(\mathcal{D} | \mathbf{Z}, \Phi, U) = p(\mathbf{W} | \mathbf{Z}, \Phi) p(\mathbf{y} | \mathbf{Z}, U)$ is the likelihood of the observed data and $p_0(\Theta, \mathbf{Z}, \Phi, U) = p_0(U) [\prod_i p(\theta_i | \alpha) \prod_n p(z_{in} | \theta_i)] \prod_k p(\Phi_k | \beta)$ is the prior distribution defined by the model. One common issue with this estimation is that real networks are highly imbalanced—the number of positive links is much smaller than the number of negative links. For example, less than 0.1% document pairs in the Cora network have positive links.

To deal with this imbalance issue, we propose to do regularized Bayesian inference (RegBayes) [43] which offers an extra freedom to handle the imbalance issue in a cost-sensitive manner. Specifically, we define a Gibbs classifier for binary links as follows.

- 1) **A Latent Predictor:** If the weight matrix U and topic assignments \mathbf{Z} are given, we build a classifier using the likelihood (2) and the *latent* prediction rule is

$$\hat{y}_{ij} | \mathbf{z}_i, \mathbf{z}_j, U = \mathbb{I}(\bar{\mathbf{z}}_i^\top U \bar{\mathbf{z}}_j > 0), \quad (3)$$

where $\mathbb{I}(\cdot)$ is an indicator function that equals to 1 if predicate holds otherwise 0. Then, the training error of this latent prediction rule is

$$\text{Err}(U, \mathbf{Z}) = \sum_{(i,j) \in \mathcal{I}} \mathbb{I}(y_{ij} \neq \hat{y}_{ij} | \mathbf{z}_i, \mathbf{z}_j, U).$$

Since directly optimizing the training error is hard, a convex surrogate loss is commonly used in machine learning. Here, we consider two popular examples, namely, the logistic log-loss and the hinge loss

$$\begin{aligned} \mathcal{R}_1(U, \mathbf{Z}) &= - \sum_{(i,j) \in \mathcal{I}} \log p(y_{ij} | \mathbf{z}_i, \mathbf{z}_j, U), \\ \mathcal{R}_2(U, \mathbf{Z}) &= \sum_{(i,j) \in \mathcal{I}} \max(0, \ell - \tilde{y}_{ij} \mathbf{z}_i^\top U \mathbf{z}_j), \end{aligned}$$

where $\ell (\geq 1)$ is a cost parameter that penalizes a wrong prediction and $\tilde{y}_{ij} = 2y_{ij} - 1$ is a transformation of the 0/1 binary links to be $-1/+1$ for notation convenience.

- 2) **Expected Loss:** Since both U and \mathbf{Z} are hidden variables, we infer a posterior distribution $q(U, \mathbf{Z})$ that has the minimal expected loss

$$\mathcal{R}_1(q(U, \mathbf{Z})) = \mathbb{E}_q[\mathcal{R}_1(U, \mathbf{Z})] \quad (4)$$

$$\mathcal{R}_2(q(U, \mathbf{Z})) = \mathbb{E}_q[\mathcal{R}_2(U, \mathbf{Z})]. \quad (5)$$

Remark 1: Note that both loss functions $\mathcal{R}_1(U, \mathbf{Z})$ and $\mathcal{R}_2(U, \mathbf{Z})$ are convex over the parameters U when the latent topics \mathbf{Z} are fixed. The hinge loss is an upper bound of the training error, while the log-loss is not. Many comparisons have been done in the context of classification [36]. Our results will provide a careful comparison of these two loss functions in the context of relational topic models.

Remark 2: Both $\mathcal{R}_1(q(U, \mathbf{Z}))$ and $\mathcal{R}_2(q(U, \mathbf{Z}))$ are good surrogate loss for the expected link prediction error

$$\text{Err}(q(U, \mathbf{Z})) = \mathbb{E}_q[\text{Err}(U, \mathbf{Z})],$$

of a Gibbs classifier that randomly draws a model U from the posterior distribution q and makes predictions [28][16]. The expected hinge loss $\mathcal{R}_2(q(U, \mathbf{Z}))$ is also an upper bound of $\text{Err}(q(U, \mathbf{Z}))$.

With the above Gibbs classifiers, we define the generalized relational topic models (gRTM) as solving the regularized Bayesian inference problem

$$\min_{q(U, \Theta, \mathbf{Z}, \Phi) \in \mathcal{P}} \mathcal{L}(q(U, \Theta, \mathbf{Z}, \Phi)) + c\mathcal{R}(q(U, \mathbf{Z})) \quad (6)$$

where $\mathcal{L}(q) = \text{KL}(q(U, \Theta, \mathbf{Z}, \Phi) || p_0(U, \Theta, \mathbf{Z}, \Phi)) - \mathbb{E}_q[\log p(\mathbf{W} | \mathbf{Z}, \Phi)]$ is an information theoretical objective; c is a positive regularization parameter controlling the influence from link structures; and \mathcal{P} is the space of normalized distributions. In fact, minimizing the single term of $\mathcal{L}(q)$ results in the posterior distribution of the vanilla LDA without considering link information. For the second term, we have used \mathcal{R} to denote a generic loss function, which can be either the log-loss \mathcal{R}_1 or the hinge-loss \mathcal{R}_2 in this paper. Note that the Gibbs classifiers and the LDA likelihood are coupled by sharing the latent topic assignments \mathbf{Z} , and the strong coupling makes it possible to learn a posterior distribution that can describe the observed words well and make accurate predictions.

To better understand the above formulation, we define the un-normalized pseudo-likelihood¹ for links:

$$\psi_1(y_{ij} | \mathbf{z}_i, \mathbf{z}_j, U) = p^c(y_{ij} | \mathbf{z}_i, \mathbf{z}_j, U) = \frac{e^{cy_{ij}\omega_{ij}}}{(1 + e^{\omega_{ij}})^c}, \quad (7)$$

$$\psi_2(y_{ij} | \mathbf{z}_i, \mathbf{z}_j, U) = \exp(-2c \max(0, 1 - y_{ij}\omega_{ij})), \quad (8)$$

1. Pseudo-likelihood has been used as an approximate maximum likelihood estimation procedure [39]. Here, we use it to denote an unnormalized likelihood of empirical data.

where $\omega_{ij} = \bar{\mathbf{z}}_i^\top U \bar{\mathbf{z}}_j$ is the discriminant function value. The pseudo-likelihood ψ_1 is un-normalized if $c \neq 1$. Then, the inference problem (6) can be written as

$$\min_{q(U, \Theta, \mathbf{Z}, \Phi) \in \mathcal{P}} \mathcal{L}(q(U, \Theta, \mathbf{Z}, \Phi)) - \mathbb{E}_q [\log \psi(\mathbf{y}|\mathbf{Z}, U)] \quad (9)$$

where $\psi(\mathbf{y}|\mathbf{Z}, U) = \prod_{(i,j) \in \mathcal{I}} \psi_1(y_{ij}|\mathbf{z}_i, \mathbf{z}_j, U)$ if using log-loss and $\psi(\mathbf{y}|\mathbf{Z}, U) = \prod_{(i,j) \in \mathcal{I}} \psi_2(y_{ij}|\mathbf{z}_i, \mathbf{z}_j, U)$ if using hinge loss.

We can show that the optimum solution of problem (6) or the equivalent problem (9) is the posterior distribution with link information

$$q(U, \Theta, \mathbf{Z}, \Phi) = \frac{p_0(U, \Theta, \mathbf{Z}, \Phi) p(\mathbf{W}|\mathbf{Z}, \Phi) \psi(\mathbf{y}|\mathbf{Z}, U)}{\phi(\mathbf{y}, \mathbf{W})}$$

where $\phi(\mathbf{y}, \mathbf{W})$ is the normalization constant to make q as a normalized distribution.

Therefore, by solving problem (6) or (9) we are in fact doing Bayesian inference with a generalized pseudo-likelihood, which is a powered version of the likelihood (2) in the case of using the log-loss. The flexibility of using regularization parameters can play a significant role in dealing with imbalanced network data as we shall see in the experiments. For example, we can use a larger c value for the sparse positive links, while using a smaller c for the dense negative links. This simple strategy has been shown effective in learning classifiers [3] and link prediction models [41] with highly imbalanced data. Finally, for the logistic log-loss an ad hoc generative story can be described as in RTMs, where c can be understood as the pseudo-count of a link.

4 AUGMENT AND COLLAPSE SAMPLING

For gRTMs with either the log-loss or the hinge loss, exact posterior inference is intractable due to the non-conjugacy between the prior and pseudo-likelihood. Previous inference methods for the standard RTMs use variational techniques with mean-field assumptions. For example, a variational EM algorithm was developed in [8] with the factorization assumption that $q(U, \Theta, \mathbf{Z}, \Phi) = q(U) [\prod_i q(\theta_i) \prod_n q(z_{in})] \prod_k q(\Phi_k)$ which can be too restrictive to be realistic in practice. In this section, we present simple and efficient Gibbs sampling algorithms without any restricting assumptions on q . Our ‘‘augment-and-collapse’’ sampling algorithm relies on a data augmentation reformulation of the RegBayes problem (9).

Before a full exposition of the algorithms, we summarize the high-level ideas. For the pseudo-likelihood $\psi(\mathbf{y}|\mathbf{Z}, U)$, it is not easy to derive a sampling algorithm directly. Instead, we develop our algorithms by introducing auxiliary variables, which lead to a scale mixture of Gaussian components and analytic conditional distributions for Bayesian inference without an accept/reject ratio. Below, we present the algorithms for the log-loss and hinge loss in turn.

4.1 Sampling Algorithm for the Log-Loss

For the case with the log-loss, our algorithm represents an extension of Polson et al.’s approach [34] to deal with the highly non-trivial Bayesian latent variable models for relational data analysis.

4.1.1 Formulation with Data Augmentation

Let us first introduce the Polya-Gamma variables [34].

Definition 3: A random variable X has a Polya-Gamma distribution, denoted by $X \sim \mathcal{PG}(a, b)$, if

$$X = \frac{1}{2\pi^2} \sum_{m=1}^{\infty} \frac{g_m}{(m-1/2)^2 + b^2/(4\pi^2)},$$

where $(a > 0, b \in \mathcal{R})$ are parameters and each $g_m \sim \mathcal{G}(a, 1)$ is an independent Gamma random variable.

Then, using the ideas of data augmentation [34], we have the following results

Lemma 4: The pseudo-likelihood can be expressed as

$$\psi_1(y_{ij}|\mathbf{z}_i, \mathbf{z}_j, U) = \frac{1}{2^c} e^{\kappa_{ij}\omega_{ij}} \int_0^\infty e^{-\frac{\lambda_{ij}\omega_{ij}^2}{2}} p(\lambda_{ij}|c, 0) d\lambda_{ij},$$

where $\kappa_{ij} = c(y_{ij} - 1/2)$ and λ_{ij} is a Polya-Gamma variable with parameters $a = c$ and $b = 0$.

Lemma 4 indicates that the posterior distribution of the generalized Bayesian logistic relational topic models, i.e., $q(U, \Theta, \mathbf{Z}, \Phi)$, can be expressed as the marginal of a higher dimensional distribution that includes the augmented variables λ . The complete posterior distribution is

$$q(U, \lambda, \Theta, \mathbf{Z}, \Phi) = \frac{p_0(U, \Theta, \mathbf{Z}, \Phi) p(\mathbf{W}|\mathbf{Z}, \Phi) \psi(\mathbf{y}, \lambda|\mathbf{Z}, U)}{\phi(\mathbf{y}, \mathbf{W})}$$

where $\psi(\mathbf{y}, \lambda|\mathbf{Z}, U) = \prod_{(i,j) \in \mathcal{I}} \exp(\kappa_{ij}\omega_{ij} - \frac{\lambda_{ij}\omega_{ij}^2}{2}) p(\lambda_{ij}|c, 0)$ is the joint pseudo-distribution² of \mathbf{y} and λ .

4.1.2 Inference with Collapsed Gibbs Sampling

Although we can do Gibbs sampling to infer the complete posterior $q(U, \lambda, \Theta, \mathbf{Z}, \Phi)$ and thus $q(U, \Theta, \mathbf{Z}, \Phi)$ by ignoring λ , the mixing rate would be slow due to the large sample space. An effective way to reduce the sample space and improve mixing rates is to integrate out the intermediate Dirichlet variables (Θ, Φ) and build a Markov chain whose equilibrium distribution is the collapsed distribution $q(U, \lambda, \mathbf{Z})$. Such a collapsed Gibbs sampling procedure has been successfully used in LDA [18]. For gRTMs, the collapsed posterior distribution is

$$\begin{aligned} q(U, \lambda, \mathbf{Z}) &\propto p_0(U) p(\mathbf{W}, \mathbf{Z}|\alpha, \beta) \psi(\mathbf{y}, \lambda|\mathbf{Z}, U) \\ &= p_0(U) \prod_{k=1}^K \frac{\delta(\mathbf{C}_k + \beta)}{\delta(\beta)} \prod_{i=1}^D \frac{\delta(\mathbf{C}_i + \alpha)}{\delta(\alpha)} \\ &\quad \times \prod_{(i,j) \in \mathcal{I}} \exp\left(\kappa_{ij}\omega_{ij} - \frac{\lambda_{ij}\omega_{ij}^2}{2}\right) p(\lambda_{ij}|c, 0), \end{aligned}$$

2. Not normalized appropriately.

where $\delta(\mathbf{x}) = \frac{\prod_{i=1}^{\dim(\mathbf{x})} \Gamma(x_i)}{\Gamma(\sum_{i=1}^{\dim(\mathbf{x})} x_i)}$, C_k^t is the number of times the term t being assigned to topic k over the whole corpus and $\mathbf{C}_k = \{C_k^t\}_{t=1}^V$; C_i^k is the number of times that terms are associated with topic k within the i -th document and $\mathbf{C}_i = \{C_i^k\}_{k=1}^K$. Then, the conditional distributions used in collapsed Gibbs sampling are as follows.

For U : for notation clarity, we define $\bar{\mathbf{z}}_{ij} = \text{vec}(\bar{\mathbf{z}}_i \bar{\mathbf{z}}_j^\top)$ and $\boldsymbol{\eta} = \text{vec}(U)$, where $\text{vec}(A)$ is a vector concatenating the row vectors of matrix A . Then, we have the discriminant function value $\omega_{ij} = \boldsymbol{\eta}^\top \bar{\mathbf{z}}_{ij}$. For the commonly used isotropic Gaussian prior $p_0(U) = \prod_{kk'} \mathcal{N}(U_{kk'}; 0, \nu^2)$, i.e., $p_0(\boldsymbol{\eta}) = \prod_m \mathcal{N}(\eta_m; 0, \nu^2)$, we have

$$\begin{aligned} q(\boldsymbol{\eta}|\mathbf{Z}, \boldsymbol{\lambda}) &\propto p_0(\boldsymbol{\eta}) \prod_{(i,j) \in \mathcal{I}} \exp\left(\kappa_{ij} \boldsymbol{\eta}^\top \bar{\mathbf{z}}_{ij} - \frac{\lambda_{ij} (\boldsymbol{\eta}^\top \bar{\mathbf{z}}_{ij})^2}{2}\right) \\ &= \mathcal{N}(\boldsymbol{\eta}; \boldsymbol{\mu}, \boldsymbol{\Sigma}), \end{aligned} \quad (10)$$

where $\boldsymbol{\Sigma} = \left(\frac{1}{\nu^2} I + \sum_{(i,j) \in \mathcal{I}} \lambda_{ij} \bar{\mathbf{z}}_{ij} \bar{\mathbf{z}}_{ij}^\top\right)^{-1}$ and $\boldsymbol{\mu} = \boldsymbol{\Sigma} \left(\sum_{(i,j) \in \mathcal{I}} \kappa_{ij} \bar{\mathbf{z}}_{ij}\right)$. Therefore, we can easily draw a sample from a K^2 -dimensional multivariate Gaussian distribution. The inverse can be robustly done using Cholesky decomposition. Since K is normally not large, the inversion is relatively efficient, especially when the number of documents is large. We will provide empirical analysis in the experiment section. Note that for large K this step can be a practical limitation. But fortunately, there are good parallel algorithms for Cholesky decomposition [15], which can be used for applications with large K values.

For \mathbf{Z} : the conditional distribution of \mathbf{Z} is

$$q(\mathbf{Z}|U, \boldsymbol{\lambda}) \propto \prod_{k=1}^K \frac{\delta(\mathbf{C}_k + \boldsymbol{\beta})}{\delta(\boldsymbol{\beta})} \prod_{i=1}^D \frac{\delta(\mathbf{C}_i + \boldsymbol{\alpha})}{\delta(\boldsymbol{\alpha})} \prod_{(i,j) \in \mathcal{I}} \psi_1(y_{ij}|\boldsymbol{\lambda}, \mathbf{Z})$$

where $\psi_1(y_{ij}|\boldsymbol{\lambda}, \mathbf{Z}) = \exp(\kappa_{ij} \omega_{ij} - \frac{\lambda_{ij} \omega_{ij}^2}{2})$. By canceling common factors, we can derive the local conditional of one variable z_{in} given others \mathbf{Z}_{-n} as:

$$\begin{aligned} q(z_{in}^k = 1 | \mathbf{Z}_{-n}, U, \boldsymbol{\lambda}, w_{in} = t) &\propto \frac{(C_{k,-n}^t + \beta_t)(C_{i,-n}^k + \alpha_k)}{\sum_t C_{k,-n}^t + \sum_{t=1}^V \beta_t} \\ &\times \prod_{j \in \mathcal{N}_i^+} \psi_1(y_{ij}|\boldsymbol{\lambda}, \mathbf{Z}_{-n}, z_{in}^k = 1) \\ &\times \prod_{j \in \mathcal{N}_i^-} \psi_1(y_{ji}|\boldsymbol{\lambda}, \mathbf{Z}_{-n}, z_{in}^k = 1), \end{aligned} \quad (11)$$

where $C_{:, -n}$ indicates that term n is excluded from the corresponding document or topic; and $\mathcal{N}_i^+ = \{j : (i, j) \in \mathcal{I}\}$ and $\mathcal{N}_i^- = \{j : (j, i) \in \mathcal{I}\}$ denote the neighbors of document i in the training network. For symmetric networks, $\mathcal{N}_i^+ = \mathcal{N}_i^-$, only one part is sufficient. We can see that the first term is from the LDA model for observed word counts and the second term is from the link structures \mathbf{y} .

Algorithm 1 Collapsed Gibbs Sampling Algorithm for Generalized RTMs with Logistic Log-loss

- 1: **Initialization:** set $\boldsymbol{\lambda} = 1$ and randomly draw z_{dn} from a uniform distribution.
 - 2: **for** $m = 1$ **to** M **do**
 - 3: draw the classifier from the normal distribution (10)
 - 4: **for** $i = 1$ **to** D **do**
 - 5: **for** each word n in document i **do**
 - 6: draw the topic using distribution (11)
 - 7: **end for**
 - 8: **end for**
 - 9: **for** $(i, j) \in \mathcal{I}$ **do**
 - 10: draw λ_{ij} from distribution (12).
 - 11: **end for**
 - 12: **end for**
-

For $\boldsymbol{\lambda}$: the conditional distribution of the augmented variables $\boldsymbol{\lambda}$ is a Polya-Gamma distribution

$$\begin{aligned} q(\lambda_{ij}|\mathbf{Z}, U) &\propto \exp\left(-\frac{\lambda_{ij} \omega_{ij}^2}{2}\right) p(\lambda_{ij}|c, 0) \\ &= \mathcal{PG}(\lambda_{ij}; c, \omega_{ij}). \end{aligned} \quad (12)$$

The equality is achieved by using the construction definition of the general $\mathcal{PG}(a, b)$ class through an exponential tilting of the $\mathcal{PG}(a, 0)$ density [34]. To draw samples from the Polya-Gamma distribution, a naive implementation using the infinite sum-of-Gamma representation is not efficient and it also involves a potentially inaccurate step of truncating the infinite sum. Here we adopt the method proposed in [34], which draws the samples from the closely related exponentially tilted Jacobi distribution.

With the above conditional distributions, we can construct a Markov chain which iteratively draws samples of $\boldsymbol{\eta}$ (i.e., U) using Eq. (10), \mathbf{Z} using Eq. (11) and $\boldsymbol{\lambda}$ using Eq. (12) as shown in Alg. 1, with an initial condition. In our experiments, we initially set $\boldsymbol{\lambda} = 1$ and randomly draw \mathbf{Z} from a uniform distribution. In training, we run the Markov chain for M iterations (i.e., the so called burn-in stage). Then, we draw a sample \hat{U} as the final classifier to make predictions on testing data. As we shall see in practice, the Markov chain converges to stable prediction performance with a few burn-in iterations.

4.2 Sampling Algorithm for the Hinge Loss

Now, we present an ‘‘augment-and-collapse’’ Gibbs sampling algorithm for the gRTMs with the hinge loss. The algorithm represents an extension of the recent techniques [42] to relational data analysis.

4.2.1 Formula with Data Augmentation

As we do not have a closed-form of the expected margin loss, it is hard to deal with the expected hinge

loss in Eq. (5). Here, we develop a collapsed Gibbs sampling method based on a data augmentation formulation of the expected margin loss to infer the posterior distribution

$$q(U, \Theta, \mathbf{Z}, \Phi) = \frac{p_0(U, \Theta, \mathbf{Z}, \Phi)p(\mathbf{W}|\mathbf{Z}, \Phi)\psi(\mathbf{y}|\mathbf{Z}, U)}{\phi(\mathbf{y}, \mathbf{W})},$$

where $\phi(\mathbf{y}, \mathbf{W})$ is the normalization constant and $\psi(\mathbf{y}|\mathbf{Z}, U) = \prod_{(i,j) \in \mathcal{I}} \psi_2(y_{ij}|\mathbf{z}_i, \mathbf{z}_j, U)$ in this case. Specifically, we have the following data augmentation representation of the pseudo-likelihood:

$$\begin{aligned} \psi_2(y_{ij}|\mathbf{z}_i, \mathbf{z}_j, U) \\ = \int_0^\infty \frac{1}{\sqrt{2\pi\lambda_{ij}}} \exp\left\{-\frac{(\lambda_{ij} + c\zeta_{ij})^2}{2\lambda_{ij}}\right\} d\lambda_{ij}, \end{aligned} \quad (13)$$

where $\zeta_{ij} = \ell - y_{ij}\omega_{ij}$. Eq. (13) can be derived following [35], and it indicates that the posterior distribution $q(U, \Theta, \mathbf{Z}, \Phi)$ can be expressed as the marginal of a higher dimensional posterior distribution that includes the augmented variables λ :

$$q(U, \lambda, \Theta, \mathbf{Z}, \Phi) = \frac{p_0(U, \Theta, \mathbf{Z}, \Phi)p(\mathbf{W}|\mathbf{Z}, \Phi)\psi(\mathbf{y}, \lambda|\mathbf{Z}, U)}{\phi(\mathbf{y}, \mathbf{W})},$$

where the unnormalized distribution of \mathbf{y} and λ is

$$\psi(\mathbf{y}, \lambda|\mathbf{Z}, U) = \prod_{(i,j) \in \mathcal{I}} \frac{1}{\sqrt{2\pi\lambda_{ij}}} \exp\left(-\frac{(\lambda_{ij} + c\zeta_{ij})^2}{2\lambda_{ij}}\right).$$

4.2.2 Inference with Collapsed Gibbs Sampling

Similar as in the log-loss case, although we can sample the complete distribution $q(U, \lambda, \Theta, \mathbf{Z}, \Phi)$, the mixing rate would be slow due to the high dimensional sample space. Thus, we reduce the sample space and improve mixing rate by integrating out the intermediate Dirichlet variables (Θ, Φ) and building a Markov chain whose equilibrium distribution is the resulting marginal distribution $q(U, \lambda, \mathbf{Z})$. Specifically, the collapsed posterior distribution is

$$\begin{aligned} q(U, \lambda, \mathbf{Z}) &\propto p_0(U)p(\mathbf{W}, \mathbf{Z}|\alpha, \beta) \prod_{i,j} \phi(y_{ij}, \lambda_{ij}|\mathbf{z}_i, \mathbf{z}_j, U) \\ &= p_0(U) \prod_{i=1}^D \frac{\delta(\mathbf{C}_i + \alpha)}{\delta(\alpha)} \times \prod_{k=1}^K \frac{\delta(\mathbf{C}_k + \beta)}{\delta(\beta)} \\ &\quad \times \prod_{(i,j) \in \mathcal{I}} \frac{1}{\sqrt{2\pi\lambda_{ij}}} \exp\left\{-\frac{(\lambda_{ij} + c\zeta_{ij})^2}{2\lambda_{ij}}\right\}. \end{aligned}$$

Then we could get the conditional distribution using the collapsed Gibbs sampling as following:

For U : we use the similar notations, $\boldsymbol{\eta} = \text{vec}(U)$ and $\bar{\mathbf{z}}_{ij} = \text{vec}(\bar{\mathbf{z}}_i \bar{\mathbf{z}}_j^\top)$. For the commonly used isotropic Gaussian prior, $p_0(U) = \prod_{k,k'} \mathcal{N}(U_{k,k'}; 0, \nu^2)$, the posterior distribution of $q(U|\mathbf{Z}, \lambda)$ or $q(\boldsymbol{\eta}|\mathbf{Z}, \lambda)$ is still a Gaussian distribution:

$$\begin{aligned} q(\boldsymbol{\eta}|\mathbf{Z}, \lambda) &\propto p_0(U) \prod_{(i,j) \in \mathcal{I}} \exp\left(-\frac{(\lambda_{ij} + c\zeta_{ij})^2}{2\lambda_{ij}}\right) \\ &= \mathcal{N}(\boldsymbol{\eta}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \end{aligned} \quad (14)$$

where $\boldsymbol{\Sigma} = \left(\frac{1}{\sigma^2}I + c^2 \sum_{i,j} \frac{\bar{\mathbf{z}}_{ij} \bar{\mathbf{z}}_{ij}^\top}{\lambda_{ij}}\right)^{-1}$ and $\boldsymbol{\mu} = \boldsymbol{\Sigma}(c \sum_{i,j} y_{ij} \frac{(\lambda_{ij} + c\ell)}{\lambda_{ij}} \bar{\mathbf{z}}_{ij})$.

For \mathbf{Z} : the conditional posterior distribution of \mathbf{Z} is

$$q(\mathbf{Z}|U, \lambda) \propto \prod_{i=1}^D \frac{\delta(\mathbf{C}_i + \alpha)}{\delta(\alpha)} \prod_{k=1}^K \frac{\delta(\mathbf{C}_k + \beta)}{\delta(\beta)} \prod_{(i,j) \in \mathcal{I}} \psi_2(y_{ij}|\lambda, \mathbf{Z}),$$

where $\psi_2(y_{ij}|\lambda, \mathbf{Z}) = \exp(-\frac{(\lambda_{ij} + c\zeta_{ij})^2}{2\lambda_{ij}})$. By canceling common factors, we can derive the local conditional of one variable z_{in} given others \mathbf{Z}_{-} as:

$$\begin{aligned} q(z_{in}^k = 1|\mathbf{Z}_{-}, U, \lambda, w_{in} = t) \\ \propto \frac{(C_{k,-n}^t + \beta_t)(C_{i,-n}^k + \alpha_k)}{\sum_t C_{k,-n}^t + \sum_{t=1}^V \beta_t} \\ \times \prod_{j \in \mathcal{N}_i^+} \psi_2(y_{ij}|\lambda, \mathbf{Z}_{-}, z_{in}^k = 1) \\ \times \prod_{j \in \mathcal{N}_i^-} \psi_2(y_{ji}|\lambda, \mathbf{Z}_{-}, z_{in}^k = 1). \end{aligned} \quad (15)$$

Again, we can see that the first term is from the LDA model for observed word counts and the second term is from the link structures \mathbf{y} .

For λ : due to the independence structure among the augmented variables when \mathbf{Z} and U are given, we can derive the conditional posterior distribution of each augmented variable λ_{ij} as:

$$\begin{aligned} q(\lambda_{ij}|\mathbf{Z}, U) &\propto \frac{1}{\sqrt{2\pi\lambda_{ij}}} \exp\left(-\frac{(\lambda_{ij} + c\zeta_{ij})^2}{2\lambda_{ij}}\right) \\ &= \mathcal{GIG}\left(\lambda_{ij}; \frac{1}{2}, 1, c^2\zeta_{ij}^2\right) \end{aligned} \quad (16)$$

where $\mathcal{GIG}(x; p, a, b) = C(p, a, b)x^{p-1} \exp(-\frac{1}{2}(\frac{b}{x} + ax))$ is a generalized inverse Gaussian distribution [12] and $C(p, a, b)$ is a normalization constant. Therefore, we can derive that λ_{ij}^{-1} follows an inverse Gaussian distribution

$$p(\lambda_{ij}^{-1}|\mathbf{Z}, U) = \mathcal{IG}\left(\lambda_{ij}^{-1}; \frac{1}{c|\zeta_{ij}|}, 1\right),$$

where $\mathcal{IG}(x; a, b) = \sqrt{\frac{b}{2\pi x^3}} \exp(-\frac{b(x-a)^2}{2a^2x})$ for $a, b > 0$.

With the above conditional distributions, we can construct a Markov chain which iteratively draws samples of the weights $\boldsymbol{\eta}$ (i.e., U) using Eq. (14), the topic assignments \mathbf{Z} using Eq. (15) and the augmented variables λ using Eq. (16), with an initial condition which is the same as in the case of the logistic log-loss. To sample from an inverse Gaussian distribution, we apply the efficient transformation method with multiple roots [30].

Remark 5: We note that the Gibbs sampling algorithms for both the hinge loss and logistic loss have a similar structure. But they have different distributions for the augmented variables. As we shall see in experiments, drawing samples from the different distributions for λ will have different efficiency.

4.3 Prediction

Since gRTMs account for both text contents and network structures, we can make predictions for each of them conditioned on the other [8]. For link prediction, given a test document \mathbf{w} , we need to infer its topic assignments \mathbf{z} in order to apply the classifier (3). This can be done with a collapsed Gibbs sampling method, where the conditional distribution is

$$p(z_n^k = 1 | \mathbf{z}_{-n}) \propto \hat{\phi}_{kw_n} (C_{-n}^k + \alpha_k);$$

C_{-n}^k is the times that the terms in this document \mathbf{w} are assigned to topic k with the n -th term excluded; and $\hat{\Phi}$ is a point estimate of the topics, with $\hat{\phi}_{kt} \propto C_{kt}^t + \beta_t$. To initialize, we randomly set each word to a topic, and then run the Gibbs sampler until some stopping criterion is met, e.g., the relative change of likelihood is less than a threshold (e.g., $1e-4$ in our experiments).

For word prediction, we need to infer the distribution

$$p(w_n | \mathbf{y}, \mathcal{D}, \hat{\Phi}, \hat{U}) = \sum_k \hat{\phi}_{kw_n} p(z_n^k = 1 | \mathbf{y}, \mathcal{D}, \hat{U}).$$

This can be done by drawing a few samples of z_n and compute the empirical mean of $\hat{\phi}_{kw_n}$ using the sampled z_n . The number of samples is determined by running a Gibbs sampler until some stopping criterion is met, e.g., the relative change of likelihood is less than $1e-4$ in our experiments.

5 EXPERIMENTS

Now, we present both quantitative and qualitative results on several real network datasets to demonstrate the efficacy of the generalized discriminative relational topic models. We also present extensive sensitivity analysis with respect to various parameters.

5.1 Data sets and Setups

We present experiments on three public datasets of document networks³:

- 1) The *Cora* data [29] consists of abstracts of 2,708 computer science research papers, with links between documents that cite each other. In total, the Cora citation network has 5,429 positive links, and the dictionary consists of 1,433 words.
- 2) The *WebKB* data [10] contains 877 webpages from the computer science departments of different universities, with links between webpages that are hyper-linked. In total, the WebKB network has 1,608 positive links and the dictionary has 1,703 words.
- 3) The *Citeseer* data [37] consists of 3,312 scientific publications with 4,732 positive links, and the dictionary contains 3,703 unique words.

Since many baseline methods have been outperformed by RTMs on the same datasets [8], we focus

on evaluating the effects of the various extensions in the discriminative gRTMs with log-loss (denoted by Gibbs-gRTM) and hinge loss (denoted by Gibbs-gMMRTM) by comparing with various special cases:

- 1) **Var-RTM**: the standard RTMs (i.e., $c = 1$) with a diagonal logistic likelihood and a variational EM algorithm with mean-field assumptions [8];
- 2) **Gibbs-RTM**: the Gibbs-RTM model with a diagonal weight matrix and a Gibbs sampling algorithm for the logistic link likelihood;
- 3) **Gibbs-gRTM**: the Gibbs-gRTM model with a full weight matrix and a Gibbs sampling algorithm for the logistic link likelihood;
- 4) **Approx-gRTM**: the Gibbs-gRTM model with fast approximation on sampling \mathbf{Z} , by computing the link likelihood term in Eq. (10) for once and caching it for sampling all the word topics in each document;
- 5) **Gibbs-MMRTM**: the Gibbs-MMRTM model with a diagonal weight matrix and a Gibbs sampling algorithm for the hinge loss;
- 6) **Gibbs-gMMRTM**: the Gibbs-gMMRTM model with a full weight matrix and a Gibbs sampling algorithm for the hinge loss;
- 7) **Approx-gMMRTM**: the Gibbs-gMMRTM model with fast approximation on sampling \mathbf{Z} , which is similar to Approx-gRTM.

For Var-RTM, we follow the setup [8] and use positive links only as training data; to deal with the one-class problem, a regularization penalty was used, which in effect injects some number of pseudo-observations (each with a fixed uniform topic distribution). For the other proposed models, including Gibbs-gRTM, Gibbs-RTM, Approx-gRTM, Gibbs-gMMRTM, Gibbs-MMRTM, and Approx-gMMRTM, we instead draw some unobserved links as negative examples. Though subsampling normally results in imbalanced datasets, the regularization parameter c in our discriminative gRTMs can effectively address it, as we shall see. Here, we fix c at 1 for negative examples, while we tune it for positive examples. All the training and testing time are fairly calculated on a desktop computer with four 3.10GHz processors and 4G RAM.

5.2 Quantitative Results

We first report the overall results of *link rank*, *word rank* and *AUC* (area under the ROC curve) to measure the prediction performance, following the setups in [8]. Link rank is defined as the average rank of the observed links from the held-out test documents to the training documents, and word rank is defined as the average rank of the words in testing documents given their links to the training documents. Therefore, lower link rank and word rank are better, and higher AUC value is better. The test documents are completely new that are not observed during training. In the training phase all the words along with their links of the test documents are removed.

3. <http://www.cs.umd.edu/projects/linqs/projects/lbc/index.html>

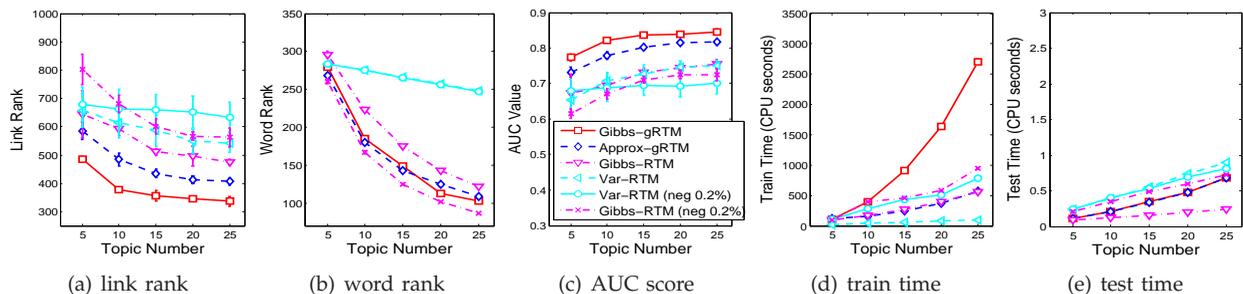


Fig. 1. Results of various models with different numbers of topics on the Cora citation dataset.

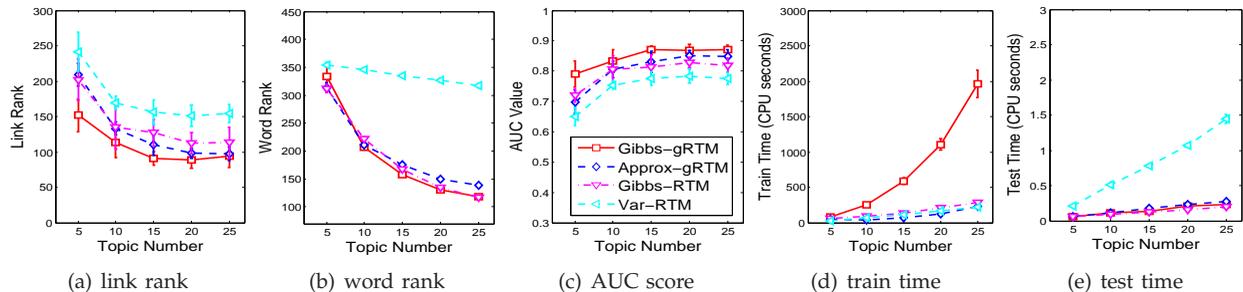


Fig. 2. Results of various models with different numbers of topics on the WebKB dataset.

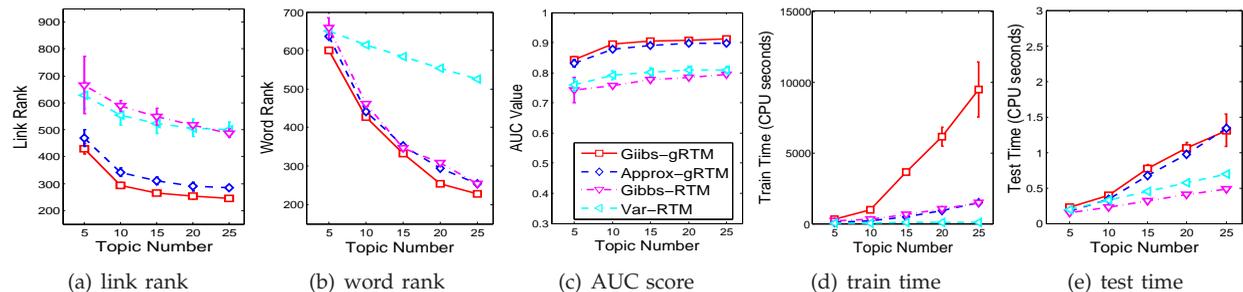


Fig. 3. Results of various models with different numbers of topics on the Citeseer dataset.

5.2.1 Results with the Log-loss

Fig. 1, Fig. 2 and Fig. 3 show the 5-fold average results and standard deviations of various models on all the three datasets with varying numbers of topic. For the RTM models using collapsed Gibbs sampling, we randomly draw 1% of the unobserved links as negative training examples, which lead to imbalanced training sets. We can see that the generalized Gibbs-gRTM can effectively deal with the imbalance and achieve significantly better results on link rank and AUC scores than all other competitors. For word rank, all the RTM models using Gibbs sampling perform better than the RTMs using variational EM methods when the number of topics is larger than 5.

The outstanding performance of Gibbs-gRTM is due to many possible factors. For example, the superior performance of Gibbs-gRTM over the diagonal Gibbs-RTM demonstrates that it is important to consider all pairwise topic interactions to fit real network data; and the superior performance of Gibbs-RTM over Var-RTM shows the benefits of using the regularization parameter c in the regularized Bayesian framework

TABLE 3
Split of training time on Cora dataset.

	Sample Z	Sample λ	Sample U
K=10	331.2 (73.55%)	55.3 (12.29%)	67.8 (14.16%)
K=15	746.8 (76.54%)	55.0 (5.64%)	173.9 (17.82%)
K=20	1300.3 (74.16%)	55.4 (3.16%)	397.7 (22.68%)

and a collapsed Gibbs sampling algorithm without restricting mean-field assumptions⁴.

To single out the influence of the proposed Gibbs sampling algorithm, we also present the results of Var-RTM and Gibbs-RTM with $c = 1$, both of which randomly sample 0.2% unobserved links⁵ as negative examples on the Cora dataset. We can see that by using Gibbs sampling without restricting mean-field

4. Gibbs-RTM doesn't outperform Var-RTM on Citeseer because they use different strategies of drawing negative samples. If we use the same strategy (e.g., randomly drawing 1% negative samples), Gibbs-RTM significantly outperforms Var-RTM.

5. Var-RTM performs much worse if using 1% negative links, while Gibbs-RTM could obtain similar performance (see Fig. 13) due to its effectiveness in dealing with imbalance.

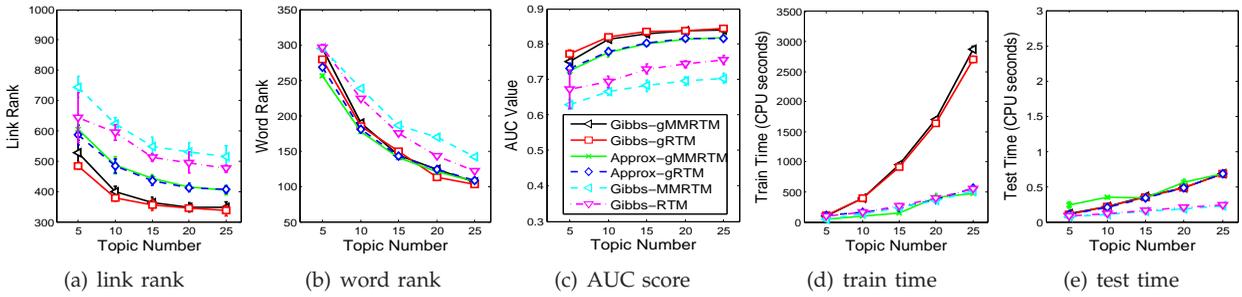


Fig. 4. Results of various models with different numbers of topics on the Cora dataset.

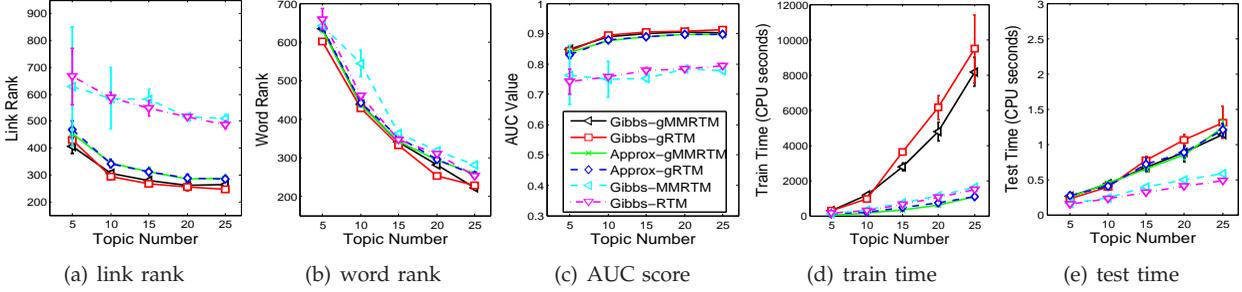


Fig. 5. Results of various models with different numbers of topics on the Citeseer dataset.

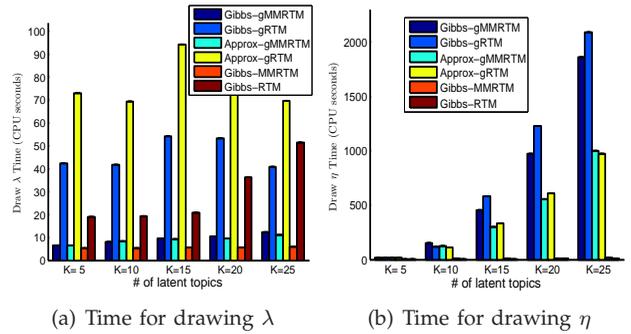
assumptions, Gibbs-RTM (neg 0.2%) outperforms Var-RTM (neg 0.2%) that makes mean-field assumptions when the number of topics is larger than 10. We defer more careful analysis of other factors in the next section, including c and the subsampling ratio.

We also note that the cost we pay for the outstanding performance of Gibbs-gRTM is on training time, which is much longer than that of Var-RTM because Gibbs-gRTM has K^2 latent features in the logistic likelihood and more training link pairs, while Var-RTM has K latent features and only uses the sparse positive links as training examples. Fortunately, we can apply a simple approximate method in sampling \mathbf{Z} as in Approx-gRTM to significantly improve the training efficiency, while the prediction performance is not sacrificed much. In fact, Approx-gRTM is still significantly better than Var-RTM in all cases, and it has comparable link prediction performance with Gibbs-gRTM on the WebKB dataset, when K is large. Table 3 further shows the training time spent on each sub-step of the Gibbs sampling algorithm of Gibbs-gRTM. We can see that the step of sampling \mathbf{Z} takes most of the time ($> 70\%$); and the steps of sampling \mathbf{Z} and $\boldsymbol{\eta}$ take more time as K increases, while the step of sampling $\boldsymbol{\lambda}$ takes almost a constant time when K changes.

5.2.2 Results with the Hinge Loss

Fig. 4 and Fig. 5 show the 5-fold average results with standard deviations of the discriminative RTMs with hinge loss, comparing with the RTMs with log-loss on Cora and Citeseer datasets⁶. We can see that

6. The result on WebKB dataset is similar, but omitted for saving space. Please refer to Fig. 15 in Appendix.

Fig. 6. Time complexity of drawing $\boldsymbol{\lambda}$ and $\boldsymbol{\eta}$ on the Citeseer dataset.

the discriminative RTM models with hinge loss (i.e., Gibbs-gMMRTM and Gibbs-MMRTM) obtain comparable predictive results (e.g., link rank and AUC scores) with the RTMs using log-loss (i.e., Gibbs-gRTM and Gibbs-RTM). And owing to the use of a full weight matrix, Gibbs-gMMRTM obtains superior performance over the diagonal Gibbs-MMRTM. These results verify the fact that the max-margin RTMs can be used as a competing alternative approach for statistical network link prediction. For word rank, all the RTM models using Gibbs sampling perform similarly.

As shown in Fig. 6, one superiority of the max margin Gibbs-gMMRTM is that the time cost of drawing $\boldsymbol{\lambda}$ is cheaper than that in Gibbs-gRTM with log-loss. Specifically, the time of drawing $\boldsymbol{\lambda}$ in Gibbs-gRTM is about 10 times longer than Gibbs-gMMRTM (Fig. 6(a)). This is because sampling from a Poly-gamma distribution in Gibbs-gRTM needs a few steps

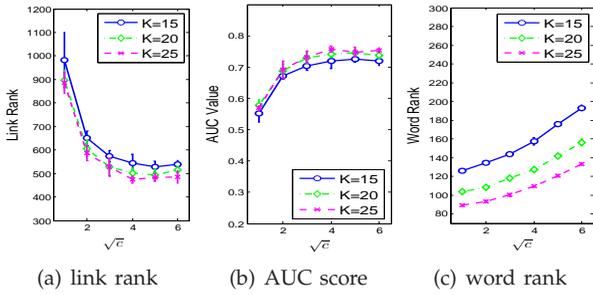


Fig. 7. Performance of Gibbs-RTM with different c values on the Cora dataset.

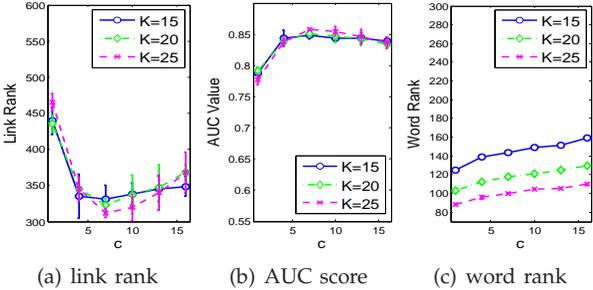


Fig. 8. Performance of Gibbs-gRTM with different c values on the Cora dataset.

of iteration for convergence, which takes more time than the constant time sampler of an inverse Gaussian distribution [30] in Gibbs-gMMRTM. We also observe that the time costs for drawing η (Fig. 6(b)) in Gibbs-gRTM and Gibbs-gMMRTM are comparable⁷. As most of the time is spent on drawing \mathbf{Z} and η , the total training time of the RTMs with the two types of losses are similar (gMMRTM is slightly faster on Citeseer). Fortunately, we can also develop Approx-gMMRTM by using a simple approximate method in sampling \mathbf{Z} to greatly improve the time efficiency (Fig. 4 and Fig. 5), and the prediction performance is still very compelling, especially on the Citeseer dataset.

5.3 Sensitivity Analysis

To provide more insights about the behaviors of our discriminative RTMs, we present a careful analysis of various factors.

5.3.1 Hyper-parameter c

Fig. 7 and 9 show the prediction performance of the diagonal Gibbs-RTM and Gibbs-MMRTM with different c values on both Cora and Citeseer datasets⁸, and Fig. 8 and 10 show the results of the generalized Gibbs-gRTM and Gibbs-gMMRTM. For Gibbs-RTM and Gibbs-MMRTM, we can see that the link rank decreases and AUC scores increase when c becomes larger and the prediction performance is stable in a

⁷ Sampling \mathbf{Z} also takes comparable time. Omitted for saving space.

⁸ We have similar observations on the WebKB dataset, again omitted for saving space.

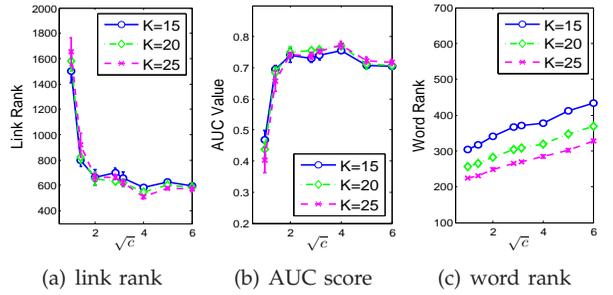


Fig. 9. Performance of Gibbs-MMRTM with different c values on the Citeseer dataset.

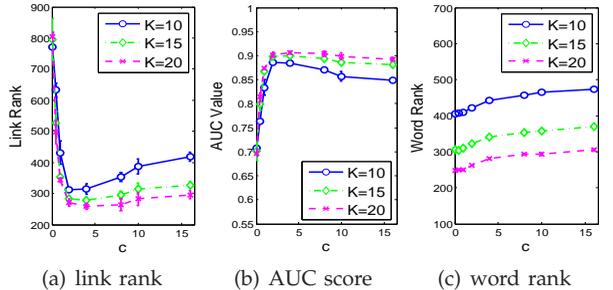


Fig. 10. Performance of Gibbs-gMMRTM with different c values on the Citeseer dataset.

wide range (e.g., $2 \leq \sqrt{c} \leq 6$). But the RTM model (i.e., $c = 1$) using Gibbs sampling doesn't perform well due to its ineffectiveness in dealing with imbalanced network data. In Fig. 8 and 10, we can observe that when $2 \leq c \leq 10$, the link rank and AUC scores of Gibbs-gRTM achieve the local optimum, which performs much better than the performance of Gibbs-gRTM when $c = 1$. In general, we can see that both Gibbs-gRTM and Gibbs-gMMRTM need a smaller c to get the best performance. This is because by allowing all pairwise topic interactions, Gibbs-gRTM and Gibbs-gMMRTM are much more expressive than Gibbs-RTM and Gibbs-MMRTM with a diagonal weight matrix; and thus easier to over-fit when c gets large.

For all the proposed models, the word rank increases slowly with the growth of c . This is because a larger c value makes the model more concentrated on fitting link structures and thus the fitness of observed words sacrifices a bit. But if we compare with the variational RTM (i.e., Var-RTM) as shown in Fig. 1 and Fig. 3, the word ranks of all the four proposed RTMs using Gibbs sampling are much lower for all the c values we have tested. This suggests the advantages of the collapsed Gibbs sampling algorithms.

5.3.2 Burn-In Steps

Fig. 11 and Fig. 12 show the sensitivity of Gibbs-gRTM and Gibbs-gMMRTM with respect to the number of burn-in iterations, respectively. We can see that the link rank and AUC scores converge fast to stable optimum points with about 300 iterations. The training time grows almost linearly with respect to the

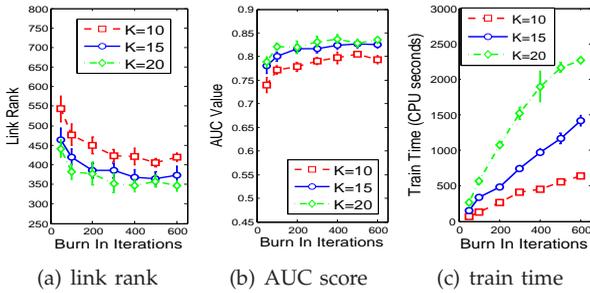


Fig. 11. Performance of Gibbs-gRTM with different burn-in steps on the Cora dataset.

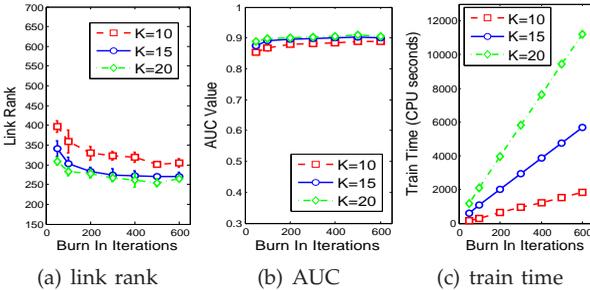


Fig. 12. Performance of Gibbs-gMMRTM with different Burn-In steps on the Citeseer dataset.

number of burn-in iterations. We have similar observations for the diagonal Gibbs-RTM, Gibbs-MMRTM and Approximate RTMs with fast approximation. In the previous experiments, we have set the burn-in steps at 400 for Cora and Citeseer, which is sufficiently large.

5.3.3 Subsample ratio

Fig. 13 shows the influence of the subsample ratio on the performance of Gibbs-gRTM on the Cora data. In total, less than 0.1% links are positive on the Cora networks. We can see that by introducing the regularization parameter c , Gibbs-gRTM can effectively fit various imbalanced network data and the different subsample ratios have a weak influence on the performance of Gibbs-gRTM. Since a larger subsample ratio leads to a bigger training set, the training time increases as expected. We have similar observations on Gibbs-gMMRTM and other models.

5.3.4 Dirichlet prior α

Fig. 14 shows the sensitivity of the generalized Gibbs-gRTM model and diagonal Gibbs-RTM on the Cora dataset with different α values. We can see that the results are quite stable in a wide range of α (i.e., $1 \leq \alpha \leq 10$) for three different topic numbers. We have similar observations for Gibbs-gMMRTM. In the previous experiments, we set $\alpha = 5$ for both Gibbs-gRTM and Gibbs-gMMRTM.

5.4 Link Suggestion

As in [8], Gibbs-gRTM could perform the task of suggesting links for a new document (i.e., test data)

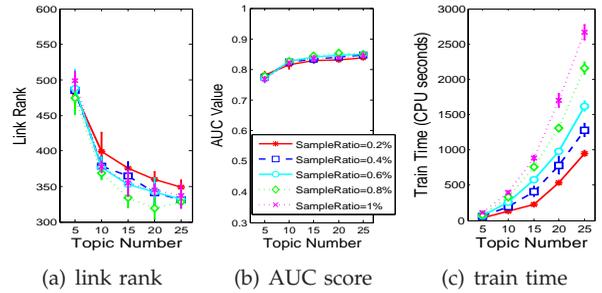


Fig. 13. Performance of Gibbs-gRTM with different numbers of negative training links on the Cora dataset.

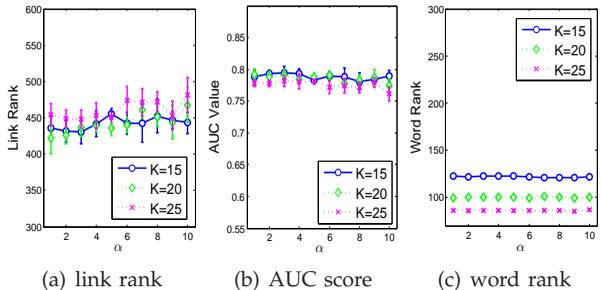


Fig. 14. Performance of Gibbs-gRTM ($c = 1$) with different α values on the Cora dataset.

based on its text contents. Table 4 shows the example suggested citations for two query documents: 1) “Competitive environments evolve better solutions for complex tasks” and 2) “Planning by Incremental Dynamic Programming” in Cora data using Gibbs-gRTM and Var-RTM. The query documents are not observed during training, and suggestion results are ranked by the values of link prediction likelihood between the training documents and the given query. We can see that Gibbs-gRTM outperforms Var-RTM in terms of identifying more ground-truth links. For query 1, Gibbs-gRTM finds 4 truly linked documents (5 in total) in the top-8 suggested results, while Var-RTM finds 3. For query 2, Gibbs-gRTM finds 2 while Var-RTM does not find any. In general, Gibbs-gRTM outperforms Var-RTM on the link suggestion task across the whole corpus. We also observe that the suggested documents which are not truly linked to the query document are also very related to it semantically.

6 CONCLUSIONS AND DISCUSSIONS

We have presented discriminative relational topic models (gRTMs and gMMRTMs) which consider all pairwise topic interactions and are suitable for asymmetric networks. We perform regularized Bayesian inference that introduces a regularization parameter to control the imbalance issue in common real networks and gives a freedom to incorporate two popular loss functions (i.e., logistic log-loss and hinge loss). We also presented a simple “augment-and-collapse” sampling algorithm for the proposed discriminative RTMs

TABLE 4

Top 8 link predictions made by Gibbs-gRTM and Var-RTM on the Cora dataset. (Papers with titles in bold have ground-truth links with the query document.)

Query: Competitive environments evolve better solutions for complex tasks	
Coevolving High Level Representations Strongly typed genetic programming in evolving cooperation strategies Genetic Algorithms in Search, Optimization and Machine Learning Improving tactical plans with genetic algorithms Some studies in machine learning using the game of Checkers Issues in evolutionary robotics: From Animals to Animats Strongly Typed Genetic Programming Evaluating and improving steady state evolutionary algorithms on constraint satisfaction problems	Gibbs-gRTM
Coevolving High Level Representations A survey of Evolutionary Strategies Genetic Algorithms in Search, Optimization and Machine Learning Strongly typed genetic programming in evolving cooperation strategies Solving combinatorial problems using evolutionary algorithms A promising genetic algorithm approach to job-shop scheduling, rescheduling, and open-shop scheduling problems Evolutionary Module Acquisition An Empirical Investigation of Multi-Parent Recombination Operators in Evolution Strategies	Var-RTM
Query: Planning by Incremental Dynamic Programming	
Learning to predict by the methods of temporal differences Neuronlike adaptive elements that can solve difficult learning control problems Learning to Act using Real- Time Dynamic Programming A new learning algorithm for blind signal separation Planning with closed-loop macro actions Some studies in machine learning using the game of Checkers Transfer of Learning by Composing Solutions of Elemental Sequential Tasks Introduction to the Theory of Neural Computation	Gibbs-gRTM
Causation, action, and counterfactuals Learning Policies for Partially Observable Environments Asynchronous modified policy iteration with single-sided updates Hidden Markov models in computational biology: Applications to protein modeling Exploiting structure in policy construction Planning and acting in partially observable stochastic domains A qualitative Markov assumption and its implications for belief change Dynamic Programming and Markov Processes	Var-RTM

without restricting assumptions on the posterior distribution. Experiments on real network data demonstrate significant improvements on prediction tasks. The time efficiency can be significantly improved with a simple approximation method.

For future work, we are interested in making the sampling algorithm scalable to large networks by using distributed architectures [38] or doing online inference [22]. Moreover, developing nonparametric RTMs to avoid model selection problems (i.e., automatically resolve the number of latent topics in RTMs) is an interesting direction. Finally, our current focus is on static networks, and it is interesting to extend the models to deal with dynamic networks, where incorporating time varying dependencies is a challenging problem to address.

ACKNOWLEDGMENTS

This work is supported by National Key Project for Basic Research of China (Grant Nos. 2013CB329403, 2012CB316301), and Tsinghua Self-innovation Project (Grant Nos: 20121088071, 20111081111), and China Postdoctoral Science Foundation Grant (Grant Nos: 2013T60117, 2012M520281).

REFERENCES

- [1] M. Rosen-zvi A. Gruber and Y. Weiss. Latent Topic Models for Hypertext. In *Proceedings of Uncertainty in Artificial Intelligence*, 2008.
- [2] E. Airoldi, D. M. Blei, S. E. Fienberg, and E. P. Xing. Mixed Membership Stochastic Blockmodels. In *Advances in Neural Information Processing Systems*, 2008.
- [3] R. Akbani, S. Kwek, and N. Japkowicz. Applying Support Vector Machines to Imbalanced Datasets. In *European Conference on Machine Learning*, 2004.
- [4] L. Backstrom and J. Leskovec. Supervised Random Walks: Predicting and Recommending Links in Social Networks. In *International Conference on Web Search and Data Mining*, 2011.
- [5] R. Balasubramanyan and W. Cohen. Block-LDA: Jointly Modeling Entity-Annotated Text and Entity-entity Links. In *Proceeding of the SIAM International Conference on Data Mining*, 2011.
- [6] Y. Bengio, A. Courville, and P. Vincent. Representation Learning: A Review and New Perspectives. *arXiv:1206.5538v2*, 2012.
- [7] D. Blei, A. Ng, and M. I. Jordan. Latent Dirichlet Allocation. *Journal of Machine Learning Research*, (3):993–1022, 2003.
- [8] J. Chang and D. Blei. Relational Topic Models for Document Networks. In *International Conference on Artificial Intelligence and Statistics*, 2009.
- [9] N. Chen, J. Zhu, F. Xia, and B. Zhang. Generalized Relational Topic Models with Data Augmentation. In *International Joint Conference on Artificial Intelligence*, 2013.
- [10] M. Craven, D. Distasco, D. Freitag, and A. McCallum. Learning to Extract Symbolic Knowledge from the World Wide Web. In *AAAI Conference on Artificial Intelligence*, 1998.
- [11] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum Likelihood Estimation from Incomplete Data via the EM Algorithm. *Journal of the Royal Statistical Society, Ser. B*, (39):1–38, 1977.

- [12] L. Devroye. *Non-uniform random variate generation*. Springer-Verlag, 1986.
- [13] L. Dietz, S. Bickel, and T. Scheffer. Unsupervised Prediction of Citation Influences. In *Proceedings of the 24th Annual International Conference on Machine Learning*, 2007.
- [14] D. Van Dyk and X. Meng. The Art of Data Augmentation. *Journal of Computational and Graphical Statistics*, 10(1):1–50, 2001.
- [15] A. George, M. Heath, and J. Liu. Parallel Cholesky Factorization on a Shared-memory Multiprocessor. *Linear Algebra and Its Applications*, 77:165–187, 1986.
- [16] P. Germain, A. Lacasse, F. Laviolette, and M. Marchand. PAC-Bayesian Learning of Linear Classifiers. In *International Conference on Machine Learning*, pages 353–360, 2009.
- [17] A. Goldenberg, A. X. Zheng, S. E. Fienberg, and E. M. Airoldi. A Survey of Statistical Network Models. *Foundations and Trends in Machine Learning*, 2(2):129–233, 2010.
- [18] T. L. Griffiths and M. Steyvers. Finding Scientific Topics. *Proceedings of the National Academy of Sciences*, 2004.
- [19] M. A. Hasan, V. Chaoji, S. Salem, and M. Zaki. Link prediction using supervised learning. In *SIAM Workshop on Link Analysis, Counterterrorism and Security*, 2006.
- [20] P. D. Hoff, A. E. Raftery, and M. S. Handcock. Latent Space Approaches to Social Network Analysis. *Journal of American Statistical Association*, 97(460), 2002.
- [21] P.D. Hoff. Modeling Homophily and Stochastic Equivalence in Symmetric Relational Data. In *Advances in Neural Information Processing Systems*, 2007.
- [22] M. Hoffman, D. Blei, and F. Bach. Online Learning for Latent Dirichlet Allocation. In *Advances in Neural Information Processing Systems*, 2010.
- [23] T. Hofmann. Probabilistic Latent Semantic Analysis. In *Proceedings of Uncertainty in Artificial Intelligence*, 1999.
- [24] M. Jordan, Z. Ghahramani, T. Jaakkola, and L. Saul. *An introduction to variational methods for graphical models*. MIT Press, Cambridge, MA, 1999.
- [25] D. Liben-Nowell and J.M. Kleinberg. The Link Prediction Problem for Social Networks. In *ACM Conference of Information and Knowledge Management*, 2003.
- [26] R. N. Lichtenwalter, J. T. Lussier, and N. V. Chawla. New Perspectives and Methods in Link Prediction. In *Proceedings of the 19th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2010.
- [27] Y. Liu, A. Niculescu-Mizil, and W. Gryc. Topic-Link LDA: Joint Models of Topic and Author Community. In *International Conference on Machine Learning*, 2009.
- [28] D. McAllester. PAC-Bayesian stochastic model selection. *Machine Learning*, 51:5–21, 2003.
- [29] A. McCallum, K. Nigam, J. Rennie, and K. Seymore. Automating the Construction of Internet Portals with Machine Learning. *Information Retrieval*, 2000.
- [30] J.R. Michael, W.R. Schucany, and R.W. Haas. Generating Random Variates Using Transformations with Multiple Roots. *The American Statistician*, 30(2):88–90, 1976.
- [31] K. Miller, T. Griffiths, and M. Jordan. Nonparametric Latent Feature Models for Link Prediction. In *Advances in Neural Information Processing Systems*, 2009.
- [32] R. Nallapati and W. Cohen. Link-PLSA-LDA: A New Unsupervised Model for Topics and Influence in Blogs. In *Proceedings of International Conference on Weblogs and Social Media*, 2008.
- [33] R. M. Nallapati, A. Ahmed, E. P. Xing, and W. Cohen. Joint Latent Topic Models for Text and Citations. In *Proceeding of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2008.
- [34] N. G. Polson, J. G. Scott, and J. Windle. Bayesian Inference for Logistic Models using Polya-Gamma Latent Variables. *arXiv:1205.0310v1*, 2012.
- [35] N. G. Polson and S. L. Scott. Data Augmentation for Support Vector Machines. *Bayesian Analysis*, 6(1):1–24, 2011.
- [36] L. Rosasco, E. De Vito, A. Caponnetto, M. Piana, and A. Verri. Are Loss Functions All the Same? *Neural Computation*, (16):1063–1076, 2004.
- [37] P. Sen, G. Namata, M. Bilgic, and L. Getoor. Collective classification in network data. *AI Magazine*, 29(3):93–106, 2008.
- [38] A. Smola and S. Narayanamurthy. An Architecture for Parallel Topic Models. *International Conference on Very Large Data Bases*, 2010.
- [39] D. Strauss and M. Ikeda. Pseudolikelihood Estimation for Social Networks. *Journal of American Statistical Association*, 85(409):204–212, 1990.
- [40] M. A. Tanner and W. H. Wong. The Calculation of Posterior Distributions by Data Augmentation. *Journal of the American Statistical Association*, 82(398):528–540, 1987.
- [41] J. Zhu. Max-Margin Nonparametric Latent Feature Models for Link Prediction. In *International Conference on Machine Learning*, 2012.
- [42] J. Zhu, N. Chen, H. Perkins, and B. Zhang. Gibbs Max-margin Topic Models with Fast Sampling Algorithms. In *International Conference on Machine Learning*, 2013.
- [43] J. Zhu, N. Chen, and E.P. Xing. Infinite Latent SVM for Classification and Multi-task Learning. In *Advances in Neural Information Processing Systems*, 2011.
- [44] J. Zhu, N. Chen, and E.P. Xing. Bayesian Inference with Posterior Regularization and applications to Infinite Latent SVMs. *arXiv:1210.1766v2*, 2013.
- [45] J. Zhu, X. Zheng, and B. Zhang. Bayesian Logistic Supervised Topic Models with Data Augmentation. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics*, 2013.



Ning Chen received her BS from China Northwestern Polytechnical University, and PhD degree in the Department of Computer Science and Technology at Tsinghua University, China, where she is currently a post-doc fellow. She was a visiting researcher in the Machine Learning Department of Carnegie Mellon University. Her research interests are primarily in machine learning, especially probabilistic graphical models, Bayesian Nonparametrics with applications on data mining and computer vision.



Jun Zhu received his BS, MS and PhD degrees all from the Department of Computer Science and Technology in Tsinghua University, China, where he is currently an associate professor. He was a project scientist and postdoctoral fellow in the Machine Learning Department, Carnegie Mellon University. His research interests are primarily on developing statistical machine learning methods to understand scientific and engineering data arising from various fields. He

is a member of the IEEE.



Fei Xia received his BS from School of Software, Tsinghua University, China. He is currently working toward his MS degree in the Language Technologies Institute, School of Computer Science, Carnegie Mellon University, USA. His research interests are primarily on machine learning especially on probabilistic graphical models, Bayesian nonparametrics and data mining problems such as social networks.



Bo Zhang graduated from the Department of Automatic Control, Tsinghua University, Beijing, China, in 1958. Currently, he is a Professor in the Department of Computer Science and Technology, Tsinghua University and a Fellow of Chinese Academy of Sciences, Beijing, China. His main interests are artificial intelligence, pattern recognition, neural networks, and intelligent control. He has published over 150 papers and four monographs in these fields.

APPENDIX

In this section, we present additional experimental results.

A.1 Prediction Performance on WebKB Dataset

Fig. 15 shows the 5-fold average results with standard deviations of the discriminative RTMs (with both the log-loss and hinge loss) on the WebKB dataset. We have similar observations as shown in Section 5.2.2 on the other two datasets. Discriminative RTMs with the hinge loss (i.e., Gibbs-gMMRTM and Gibbs-MMRTM) obtain comparable predictive results with the RTMs using the log-loss (i.e., Gibbs-gRTM and Gibbs-RTM). And generalized gRTMs achieve superior performance over the diagonal RTMs, especially when the topic numbers are less than 25. We also develop Approx-gMMRTM and Approx-gRTM by using a simple approximation method in sampling \mathbf{Z} (see Section 5.1) to greatly improve the time efficiency without sacrificing much prediction performance.

A.2 Topic Discovery

Table. 5 shows 7 example topics discovered by the 10-topic Gibbs-gRTM on the Cora dataset. For each topic, we show the 6 top-ranked document titles that yield higher values of Θ . In order to qualitatively illustrate the semantic meaning of each topic among the documents from 7 categories⁹, in the left part of Table 5, we show the average probability of each category distributed on the particular topic. Note that category labels are not considered in all the models in this paper, we use it here just to visualize the discovered semantic meanings of the proposed Gibbs-gRTM. We can observe that most of the discovered topics are representative for documents from one or several categories. For example, topics T1 and T2 tend to represent documents about “Genetic Algorithms” and “Rule Learning”, respectively. Similarly, topics T3 and T6 are good at representing documents about “Reinforcement Learning” and “Theory”, respectively.

9. The seven categories are *Case Based*, *Genetic Algorithms*, *Neural Networks*, *Probabilistic Methods*, *Reinforcement Learning*, *Rule Learning* and *Theory*.

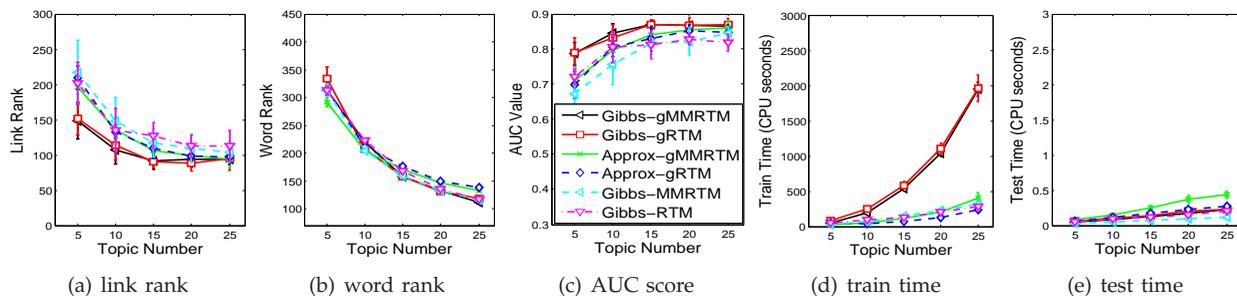


Fig. 15. Results of various models with different numbers of topics on the WebKB dataset.

TABLE 5

Example topics discovered by a 10-topic Gibbs-gRTM on the Cora dataset. For each topic, we show 6 top-ranked documents as well as the average probabilities of that topic on representing documents from 7 categories.

Topic	Top-6 Document Titles
<p>T1: Genetic Algorithms</p>	<ol style="list-style-type: none"> 1. Stage scheduling: A tech. to reduce the register requirements of a modulo schedule. 2. Optimum modulo schedules for minimum register requirements. 3. Duplication of coding segments in genetic programming. 4. Genetic programming and redundancy. 5. A cooperative coevolutionary approach to function optimization. 6. Evolving graphs and networks with edge encoding: Preliminary report.
<p>T2: Rule Learning</p>	<ol style="list-style-type: none"> 1. Inductive Constraint Logic. 2. The difficulties of learning logic programs with cut. 3. Learning se-mantic grammars with constructive inductive logic programming. 4. Learning Singly Recursive Relations from Small Datasets. 5. Least generalizations and greatest specializations of sets of clauses. 6. Learning logical definitions from relations.
<p>T3: reinforcement learning</p>	<ol style="list-style-type: none"> 1. Integ. Architect. for Learning, Planning & Reacting by Approx. Dynamic Program. 2. Multiagent reinforcement learning: Theoretical framework and an algorithm. 3. Learning to Act using Real- Time Dynamic Programming. 4. Learning to predict by the methods of temporal differences. 5. Robot shaping: Developing autonomous agents though learning. 6. Planning and acting in partially observable stochastic domains.
<p>T6: Theory</p>	<ol style="list-style-type: none"> 1. Learning with Many Irrelevant Features. 2. Learning decision lists using homogeneous rules. 3. An empirical comparison of selection measures for decision-tree induction. 4. Learning active classifiers. 5. Using Decision Trees to Improve Case-based Learning. 6. Utilizing prior concepts for learning.
<p>T7: Neural Networks</p>	<ol style="list-style-type: none"> 1. Learning factorial codes by predictability minimization. 2. The wake-sleep algorithm for unsupervised neural networks. 3. Learning to control fast-weight memories: An alternative to recurrent nets. 4. An improvement over LBG inspired from neural networks. 5. A distributed feature map model of the lexicon. 6. Self-organizing process based on lateral inhibition and synaptic resource redistribution.
<p>T8: Probabilistic Methods</p>	<ol style="list-style-type: none"> 1. Density estimation by wavelet thresholding. 2. On Bayesian analysis of mixtures with an unknown number of components. 3. Markov chain Monte Carlo methods based on "slicing" the density function. 4. Markov chain Monte Carlo convergence diagnostics: A comparative review. 5. Bootstrap C-Interv. for Smooth Splines & Comparison to Bayesian C-Interv. 6. Rates of convergence of the Hastings and Metropolis algorithms.
<p>T9: Case Based</p>	<ol style="list-style-type: none"> 1. Case Retrieval Nets: Basic ideas and extensions. 2. Case-based reasoning: Foundat. issues, methodological variat., & sys. approaches. 3. Adapter: an integrated diagnostic system combining case-based and abduct. reasoning. 4. An event-based abductive model of update. 5. Applying Case Retrieval Nets to diagnostic tasks in technical domains. 6. Introspective Reasoning using Meta-Explanations for Multistrategy Learning.