

A Probabilistic Framework for Relational Clustering

Bo Long
Computer Science Dept.
SUNY Binghamton
Binghamton, NY 13902
blong1@binghamton.edu

Zhongfei (Mark) Zhang
Computer Science Dept.
SUNY Binghamton
Binghamton, NY 13902
zzhang@binghamton.edu

Philip S. Yu
IBM Watson Research Center
19 Skyline Drive
Hawthorne, NY 10532
psyu@us.ibm.com

ABSTRACT

Relational clustering has attracted more and more attention due to its phenomenal impact in various important applications which involve multi-type interrelated data objects, such as Web mining, search marketing, bioinformatics, citation analysis, and epidemiology. In this paper, we propose a probabilistic model for relational clustering, which also provides a principal framework to unify various important clustering tasks including traditional attributes-based clustering, semi-supervised clustering, co-clustering and graph clustering. The proposed model seeks to identify cluster structures for each type of data objects and interaction patterns between different types of objects. Under this model, we propose parametric hard and soft relational clustering algorithms under a large number of exponential family distributions. The algorithms are applicable to relational data of various structures and at the same time unifies a number of state-of-the-art clustering algorithms: co-clustering algorithms, the k-partite graph clustering, and semi-supervised clustering based on hidden Markov random fields.

Categories and Subject Descriptions: E.4 [Coding and Information Theory]: Data compaction and compression; H.3.3 [Information search and Retrieval]: Clustering; I.5.3 [Pattern Recognition]: Clustering.

General Terms: Algorithms.

Keywords: Clustering, Relational data, Relational clustering, Semi-supervised clustering, EM-algorithm, Bregman divergences, Exponential families.

1. INTRODUCTION

Most clustering approaches in the literature focus on "flat" data in which each data object is represented as a fixed-length attribute vector [38]. However, many real-world data sets are much richer in structure, involving objects of multiple types that are related to each other, such as documents and words in a text corpus, Web pages, search queries and Web users in a Web search system, and shops, customers, suppliers, shareholders and advertisement media in a marketing system.

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.

KDD'07, August 12–15, 2007, San Jose, California, USA.

Copyright 2007 ACM 978-1-59593-609-7/07/0008 ...\$5.00.

In general, relational data contain three types of information, attributes for individual objects, homogeneous relations between objects of the same type, heterogeneous relations between objects of different types. For example, for a scientific publication relational data set of papers and authors, the personal information such as affiliation for authors are attributes; the citation relations among papers are homogeneous relations; the authorship relations between papers and authors are heterogeneous relations. Such data violate the classic IID assumption in machine learning and statistics and present huge challenges to traditional clustering approaches. An intuitive solution is that we transform relational data into flat data and then cluster each type of objects independently. However, this may not work well due to the following reasons.

First, the transformation causes the loss of relation and structure information [14]. Second, traditional clustering approaches are unable to tackle influence propagation in clustering relational data, i.e., the hidden patterns of different types of objects could affect each other both directly and indirectly (pass along relation chains). Third, in some data mining applications, users are not only interested in the hidden structure for each type of objects, but also interaction patterns involving multi-types of objects. For example, in document clustering, in addition to document clusters and word clusters, the relationship between document clusters and word clusters is also useful information. It is difficult to discover such interaction patterns by clustering each type of objects individually.

Moreover, a number of important clustering problems, which have been of intensive interest in the literature, can be viewed as special cases of relational clustering. For example, graph clustering (partitioning) [7, 42, 13, 6, 20, 28] can be viewed as clustering on singly-type relational data consisting of only homogeneous relations (represented as a graph affinity matrix); co-clustering [12, 2] which arises in important applications such as document clustering and micro-array data clustering, can be formulated as clustering on bi-type relational data consisting of only heterogeneous relations. Recently, semi-supervised clustering [46, 4] has attracted significant attention, which is a special type of clustering using both labeled and unlabeled data. In section 5, we show that semi-supervised clustering can be formulated as clustering on singly-type relational data consisting of attributes and homogeneous relations.

Therefore, relational data present not only huge challenges to traditional unsupervised clustering approaches, but also great need for theoretical unification of various clustering tasks. In this paper, we propose a probabilistic model for relational clustering, which also provides a principal framework to unify various important clustering tasks includ-

ing traditional attributes-based clustering, semi-supervised clustering, co-clustering and graph clustering. The proposed model seeks to identify cluster structures for each type of data objects and interaction patterns between different types of objects. It is applicable to relational data of various structures. Under this model, we propose parametric hard and soft relational clustering algorithms under a large number of exponential family distributions. The algorithms are applicable to various relational data from various applications and at the same time unify a number of state-of-the-art clustering algorithms: co-clustering algorithms, the k-partite graph clustering, Bregman k-means, and semi-supervised clustering based on hidden Markov random fields.

2. RELATED WORK

Clustering on a special case of relational data, bi-type relational data consisting of only heterogeneous relations, such as the word-document data, is called co-clustering or bi-clustering. Several previous efforts related to co-clustering are model based [22, 23]. Spectral graph partitioning has also been applied to bi-type relational data [11, 25]. These algorithms formulate the data matrix as a bipartite graph and seek to find the optimal normalized cut for the graph. Due to the nature of a bipartite graph, these algorithms have the restriction that the clusters from different types of objects must have one-to-one associations. Information-theory based co-clustering has also attracted attention in the literature. [12] proposes a co-clustering algorithm to maximize the mutual information between the clustered random variables subject to the constraints on the number of row and column clusters. A more generalized co-clustering framework is presented by [2] wherein any Bregman divergence can be used in the objective function. Recently, co-clustering has been addressed based on matrix factorization. [35] proposes an EM-like algorithm based on multiplicative updating rules.

Graph clustering (partitioning) clusters homogeneous data objects based on pairwise similarities, which can be viewed as homogeneous relations. Graph partitioning has been studied for decades and a number of different approaches, such as spectral approaches [7, 42, 13] and multilevel approaches [6, 20, 28], have been proposed. Some efforts [17, 43, 21, 21, 1] based on stochastic block modeling also focus on homogeneous relations.

Compared with co-clustering and homogeneous-relation-based clustering, clustering on general relational data, which may consist of more than two types of data objects with various structures, has not been well studied in the literature. Several noticeable efforts are discussed as follows. [45, 19] extend the probabilistic relational model to the clustering scenario by introducing latent variables into the model; these models focus on using attribute information for clustering. [18] formulates star-structured relational data as a star-structured m -partite graph and develops an algorithm based on semi-definite programming to partition the graph. [34] formulates multi-type relational data as K-partite graphs and proposes a family of algorithms to identify the hidden structures of a k-partite graph by constructing a relation summary network to approximate the original k-partite graph under a broad range of distortion measures. The above graph-based algorithms do not consider attribute information.

Some efforts on relational clustering are based on inductive logic programming [37, 24, 31]. Base on the idea of mutual reinforcement clustering, [51] proposes a framework

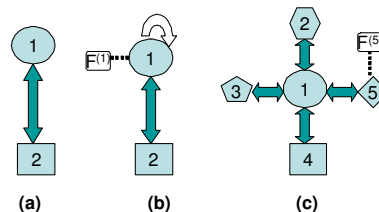


Figure 1: Examples of the structures of relational data.

for clustering heterogeneous Web objects and [47] presents an approach to improve the cluster quality of interrelated data objects through an iterative reinforcement clustering process. There are no sound objective function and theoretical proof on the effectiveness and correctness (convergence) of the mutual reinforcement clustering. Some efforts [26, 50, 49, 5] in the literature focus on how to measure the similarities or choosing cross-relational attributes.

To summarize, the research on relational data clustering has attracted substantial attention, especially in the special cases of relational data. However, there is still limited and preliminary work on general relational data clustering.

3. MODEL FORMULATION

With different compositions of three types of information, attributes, homogeneous relations and heterogeneous relations, relational data could have very different structures. Figure 1 shows three examples of the structures of relational data. Figure 1(a) refers to a simple bi-type of relational data with only heterogeneous relations such as word-document data. Figure 1(b) represents a bi-type data with all types of information, such as actor-movie data, in which actors (type 1) have attributes such as gender; actors are related to each other by collaboration in movies (homogeneous relations); actors are related to movies (type 2) by taking roles in movies (heterogeneous relations). Figure 1(c) represents the data consisting of companies, customers, suppliers, shareholders and advertisement media, in which customers (type 5) have attributes.

In this paper, we represent a relational data set as a set of matrices. Assume that a relational data set has m different types of data objects, $\mathcal{X}^{(1)} = \{x_i^{(1)}\}_{i=1}^{n_1}, \dots, \mathcal{X}^{(m)} = \{x_i^{(m)}\}_{i=1}^{n_m}$, where n_j denotes the number of objects of the j th type and $x_p^{(j)}$ denotes the name of the p th object of the j th type. We represent the observations of the relational data as three sets of matrices, attribute matrices $\{\mathbf{F}^{(j)} \in \mathbb{R}^{d_j \times n_j}\}_{j=1}^m$, where d_j denotes the dimension of attributes for the j th type objects and $\mathbf{F}_p^{(j)}$ denotes the attribute vector for object $x_p^{(j)}$; homogeneous relation matrices $\{\mathbf{S}^{(j)} \in \mathbb{R}^{n_j \times n_j}\}_{j=1}^m$, where $\mathbf{S}_{pq}^{(j)}$ denotes the relation between $x_p^{(j)}$ and $x_q^{(j)}$; heterogeneous relation matrices $\{\mathbf{R}^{(ij)} \in \mathbb{R}^{n_i \times n_j}\}_{i,j=1}^m$, where $\mathbf{R}_{pq}^{(ij)}$ denotes the relation between $x_p^{(i)}$ and $x_q^{(j)}$. The above representation is a general formulation. In real applications, not every type of objects has attributes, homogeneous relations and heterogeneous relations. For example, the relational data set in Figure 1(a) is represented by only one heterogeneous matrix $\mathbf{R}^{(12)}$, and the one in Figure 1(b) is represented by three matrices, $\mathbf{F}^{(1)}$, $\mathbf{S}^{(1)}$ and $\mathbf{R}^{(12)}$. Moreover, for a specific clustering task, we may not use all available attributes and relations after feature or relation selection pre-processing.

Mixed membership models, which assume that each object has mixed membership denoting its association with

classes, have been widely used in the applications involving soft classification [16], such as matching words and pictures [39], race genetic structures [39, 48], and classifying scientific publications [15].

In this paper, we propose a relational mixed membership model to cluster relational data (we refer to the model as *mixed membership relational clustering* or MMRC throughout the rest of the paper).

Assume that each type of objects $\mathcal{X}^{(j)}$ has k_j latent classes. We represent the membership vectors for all the objects in $\mathcal{X}^{(j)}$ as a membership matrix $\Lambda^{(j)} \in [0, 1]^{k_j \times n_j}$ such that the sum of elements of each column $\Lambda_{\cdot p}^{(j)}$ is 1 and $\Lambda_{\cdot p}^{(j)}$ denotes the membership vector for object $x_p^{(j)}$, i.e., $\Lambda_{gp}^{(j)}$ denotes the probability that object $x_p^{(j)}$ associates with the g th latent class. We also write the parameters of distributions to generate attributes, homogeneous relations and heterogeneous relations in matrix forms. Let $\Theta^{(j)} \in \mathbb{R}^{d_j \times k_j}$ denote the distribution parameter matrix for generating attributes $\mathbf{F}^{(j)}$ such that $\Theta_{\cdot g}^{(j)}$ denotes the parameter vector associated with the g th latent class. Similarly, $\Gamma^{(j)} \in \mathbb{R}^{k_j \times k_j}$ denotes the parameter matrix for generating homogeneous relations $\mathbf{S}^{(j)}$; $\Upsilon^{(ij)} \in \mathbb{R}^{k_i \times k_j}$ denotes the parameter matrix for generating heterogeneous relations $\mathbf{R}^{(ij)}$. In summary, the parameters of MMRC model are

$$\Omega = \{\{\Lambda^{(j)}\}_{j=1}^m, \{\Theta^{(j)}\}_{j=1}^m, \{\Gamma^{(j)}\}_{j=1}^m, \{\Upsilon^{(ij)}\}_{i,j=1}^m\}.$$

In general, the meanings of the parameters, Θ , Λ , and Υ , depend on the specific distribution assumptions. However, in Section 4.1, we show that for a large number of exponential family distributions, these parameters can be formulated as expectations with intuitive interpretations.

Next, we introduce the latent variables into the model. For each object x_p^j , a latent cluster indicator vector is generated based on its membership parameter $\Lambda_{\cdot p}^{(j)}$, which is denoted as $\mathbf{C}_p^{(j)}$, i.e., $\mathbf{C}^{(j)} \in \{0, 1\}^{k_j \times n_j}$ is a latent indicator matrix for all the j th type objects in $\mathcal{X}^{(j)}$.

Finally, we present the generative process of observations, $\{\mathbf{F}^{(j)}\}_{j=1}^m$, $\{\mathbf{S}^{(j)}\}_{j=1}^m$, and $\{\mathbf{R}^{(ij)}\}_{i,j=1}^m$ as follows:

1. For each object $x_p^{(j)}$
 - Sample $\mathbf{C}_p^{(j)} \sim \text{Multinomial}(\Lambda_{\cdot p}^{(j)}, 1)$.
2. For each object $x_p^{(j)}$
 - Sample $\mathbf{F}_p^{(j)} \sim \text{Pr}(\mathbf{F}_p^{(j)} | \Theta^{(j)} \mathbf{C}_p^{(j)})$.
3. For each pair of objects $x_p^{(j)}$ and $x_q^{(j)}$
 - Sample $\mathbf{S}_{pq}^{(j)} \sim \text{Pr}(\mathbf{S}_{pq}^{(j)} | (\mathbf{C}_p^{(j)})^T \Gamma^{(j)} \mathbf{C}_q^{(j)})$.
4. For each pair of objects $x_p^{(i)}$ and $x_q^{(j)}$
 - Sample $\mathbf{R}_{pq}^{(ij)} \sim \text{Pr}(\mathbf{R}_{pq}^{(ij)} | (\mathbf{C}_p^{(i)})^T \Upsilon^{(ij)} \mathbf{C}_q^{(j)})$.

In the above generative process, a latent indicator vector for each object is generated based on multinomial distribution with the membership vector as parameters. Observations are generated independently conditioning on latent indicator variables. The parameters of condition distributions are formulated as products of the parameter matrices and latent indicators, i.e., $\text{Pr}(\mathbf{F}_p^{(j)} | \mathbf{C}_p^{(j)}, \Theta^{(j)}) = \text{Pr}(\mathbf{F}_p^{(j)} | \Theta^{(j)} \mathbf{C}_p^{(j)})$, $\text{Pr}(\mathbf{S}_{pq}^{(j)} | \mathbf{C}_p^{(j)}, \mathbf{C}_q^{(j)}, \Gamma^{(j)}) = \text{Pr}(\mathbf{S}_{pq}^{(j)} | (\mathbf{C}_p^{(j)})^T \Gamma^{(j)} \mathbf{C}_q^{(j)})$, and $\text{Pr}(\mathbf{R}_{pq}^{(ij)} | \mathbf{C}_p^{(i)}, \mathbf{C}_q^{(j)}, \Upsilon^{(ij)}) = \text{Pr}(\mathbf{R}_{pq}^{(ij)} | (\mathbf{C}_p^{(i)})^T \Upsilon^{(ij)} \mathbf{C}_q^{(j)})$. Under this formulation, an observation is sampled from the

distributions of its associated latent classes. For example, if $\mathbf{C}_p^{(i)}$ indicates that $x_p^{(i)}$ is with the g th latent class and $\mathbf{C}_q^{(j)}$ indicates that $x_q^{(j)}$ is with the h th latent class, then $(\mathbf{C}_p^{(i)})^T \Upsilon^{(ij)} \mathbf{C}_q^{(j)} = \Upsilon_{gh}^{(ij)}$. Hence, we have $\text{Pr}(\mathbf{R}_{pq}^{(ij)} | \Upsilon_{gh}^{(ij)})$ implying that the relation between $x_p^{(i)}$ and $x_q^{(j)}$ is sampled by using the parameter $\Upsilon_{gh}^{(ij)}$.

With matrix representation, the joint probability distribution over the observations and the latent variables can be formulated as follows,

$$\begin{aligned} \text{Pr}(\Psi | \Omega) &= \prod_{j=1}^m \text{Pr}(\mathbf{C}^{(j)} | \Lambda^{(j)}) \prod_{j=1}^m \text{Pr}(\mathbf{F}^{(j)} | \Theta^{(j)} \mathbf{C}^{(j)}) \\ &\quad \prod_{j=1}^m \text{Pr}(\mathbf{S}^{(j)} | (\mathbf{C}^{(j)})^T \Gamma^{(j)} \mathbf{C}^{(j)}) \quad (1) \\ &\quad \prod_{i=1}^m \prod_{j=1}^m \text{Pr}(\mathbf{R}^{(ij)} | (\mathbf{C}^{(i)})^T \Upsilon^{(ij)} \mathbf{C}^{(j)}) \end{aligned}$$

where $\Psi = \{\{\mathbf{C}^{(j)}\}_{j=1}^m, \{\mathbf{F}^{(j)}\}_{j=1}^m, \{\mathbf{S}^{(j)}\}_{j=1}^m, \{\mathbf{R}^{(ij)}\}_{i,j=1}^m\}$, $\text{Pr}(\mathbf{C}^{(j)} | \Lambda^{(j)}) = \prod_{p=1}^{n_j} \text{Multinomial}(\Lambda_{\cdot p}^{(j)}, 1)$, $\text{Pr}(\mathbf{F}^{(j)} | \Theta^{(j)} \mathbf{C}^{(j)}) = \prod_{p=1}^{n_j} \text{Pr}(\mathbf{F}_p^{(j)} | \Theta^{(j)} \mathbf{C}_p^{(j)})$, $\text{Pr}(\mathbf{S}^{(j)} | (\mathbf{C}^{(j)})^T \Gamma^{(j)} \mathbf{C}^{(j)}) = \prod_{p,q=1}^{n_j} \text{Pr}(\mathbf{S}_{pq}^{(j)} | (\mathbf{C}_p^{(j)})^T \Gamma^{(j)} \mathbf{C}_q^{(j)})$, and similarly for $\mathbf{R}^{(ij)}$.

4. ALGORITHM DERIVATION

In this section, based on the MMRC model we derive parametric soft and hard relational clustering algorithms under a large number of exponential family distributions.

4.1 MMRC with Exponential Families

To avoid clutter, instead of general relational data, we use relational data similar to the one in Figure 1(b), which is a representative relational data set containing all three types of information for relational data, attributes, homogeneous relations and heterogeneous relations. However, the derivation and algorithms are applicable to general relational data.

For the relational data set in Figure 1(b), we have two types of objects, one attribute matrix \mathbf{F} , one homogeneous relation matrix \mathbf{S} and one heterogeneous relation matrix \mathbf{R} . Based on Eq.(1), we have the following likelihood function,

$$\begin{aligned} \mathcal{L}(\Omega | \Psi) &= \text{Pr}(\mathbf{C}^{(1)} | \Lambda^{(1)}) \text{Pr}(\mathbf{C}^{(2)} | \Lambda^{(2)}) \text{Pr}(\mathbf{F} | \Theta \mathbf{C}^{(1)}) \\ &\quad \text{Pr}(\mathbf{S} | (\mathbf{C}^{(1)})^T \Gamma \mathbf{C}^{(1)}) \text{Pr}(\mathbf{R} | (\mathbf{C}^{(1)})^T \Upsilon \mathbf{C}^{(2)}) \quad (2) \end{aligned}$$

Our goal is to maximize the likelihood function in Eq. (2) to estimate unknown parameters.

For the likelihood function in Eq.(2), the specific forms of condition distributions for attributes and relations depend on specific applications. Presumably, for a specific likelihood function, we need to derive a specific algorithm. However, a large number of useful distributions, such as normal distribution, Poisson distribution, and Bernoulli distributions, belong to exponential families and the distribution functions of exponential families can be formulated as a general form. This nice property facilitates us to derive a general EM algorithm for the MMRC model.

It is shown in the literature [3, 9] that there exists bijection between exponential families and Bregman divergences [40]. For example, the normal distribution, Bernoulli distribution, multinomial distribution and exponential distribution correspond to Euclidean distance, logistic loss, KL-divergence

and Itakura-Satio distance, respectively. Based on the bijection, an exponential family density $Pr(\mathbf{x})$ can always be formulated as the following expression with a Bregman divergence D_ϕ ,

$$Pr(\mathbf{x}) = \exp(-D_\phi(\mathbf{x}, \mu))f_\phi(\mathbf{x}), \quad (3)$$

where $f_\phi(\mathbf{x})$ is a uniquely determined function for each exponential probability density, and μ is the expectation parameter. Therefore, for the MMRC model under exponential family distributions, we have the following,

$$Pr(\mathbf{F}|\Theta\mathbf{C}^{(1)}) = \exp(-D_{\phi_1}(\mathbf{F}, \Theta\mathbf{C}^{(1)}))f_{\phi_1}(\mathbf{F}) \quad (4)$$

$$Pr(\mathbf{S}|\mathbf{C}^{(1)T}\Gamma\mathbf{C}^{(1)}) = \exp(-D_{\phi_2}(\mathbf{S}, (\mathbf{C}^{(1)})^T\Gamma\mathbf{C}^{(1)}))f_{\phi_2}(\mathbf{S}) \quad (5)$$

$$Pr(\mathbf{R}|\mathbf{C}^{(1)T}\Upsilon\mathbf{C}^{(2)}) = \exp(-D_{\phi_3}(\mathbf{R}, (\mathbf{C}^{(1)})^T\Upsilon\mathbf{C}^{(2)}))f_{\phi_3}(\mathbf{R}) \quad (6)$$

In the above equations, a Bregman divergence of two matrices is defined as the sum of the Bregman divergence of each pair of elements from the two matrices. Another advantage of the above formulation is that under this formulation, the parameters, Θ , Λ , and Υ , are expectations of intuitive interpretations. Θ consists of center vectors of attributes; Γ provides an intuitive summary of cluster structure within the same type objects, since $\Gamma_{gh}^{(1)}$ implies expectation relations between the g th cluster and the h th cluster of type 1 objects; similarly, Υ provides an intuitive summary for cluster structures between the different type objects. In the above formulation, we use different Bregman divergences, D_{ϕ_1} , D_{ϕ_2} , and D_{ϕ_3} , for the attributes, homogeneous relations and heterogeneous relations, since they could have different distributions in real applications. For example, suppose we have $\Theta^{(1)} = \begin{bmatrix} 1.1 & 2.3 \\ 1.5 & 2.5 \end{bmatrix}$ for normal distribution, $\Gamma^{(1)} = \begin{bmatrix} 0.9 & 0.1 \\ 0.1 & 0.7 \end{bmatrix}$ for Bernoulli distribution, and $\Upsilon^{(12)} = \begin{bmatrix} 1 & 3 \\ 3 & 1 \end{bmatrix}$ for Poisson distribution, then the cluster structures of the data are very intuitive. First, the center attribute vectors for the two clusters of type 1 are $\begin{bmatrix} 1.1 \\ 1.5 \end{bmatrix}$ and $\begin{bmatrix} 2.3 \\ 2.5 \end{bmatrix}$; second, by $\Gamma^{(1)}$ we know that the type 1 nodes from different clusters are barely related and cluster 1 is denser than cluster 2; third, by $\Upsilon^{(12)}$ we know that cluster 1 of type 1 nodes are related to cluster 2 of type 2 nodes more strongly than to cluster 1 of type 2, and so on so forth.

Since the distributions of $\mathbf{C}^{(1)}$ and $\mathbf{C}^{(2)}$ are modeled as multinomial distributions, we have the following

$$Pr(\mathbf{C}^{(1)}|\Lambda^{(1)}) = \prod_{p=1}^{n_1} \prod_{g=1}^{k_1} (\Lambda_{gp}^{(1)})^{C_{gp}^{(1)}} \quad (7)$$

$$Pr(\mathbf{C}^{(2)}|\Lambda^{(2)}) = \prod_{q=1}^{n_2} \prod_{h=1}^{k_2} (\Lambda_{hq}^{(2)})^{C_{hq}^{(2)}} \quad (8)$$

Substituting Eqs. (4), (5), (6), (7), and (8) into Eq. (2) and taking some algebraic manipulations, we obtain the following log-likelihood function for MMRC under exponential families,

$$\begin{aligned} \log(\mathcal{L}(\Omega|\Psi)) &= \sum_{p=1}^{n_1} \sum_{g=1}^{k_1} C_{gp}^{(1)} \log \Lambda_{gp}^{(1)} + \sum_{q=1}^{n_2} \sum_{h=1}^{k_2} C_{hq}^{(2)} \log \Lambda_{hq}^{(2)} \\ &\quad - D_{\phi_1}(\mathbf{F}, \Theta\mathbf{C}^{(1)}) - D_{\phi_2}(\mathbf{S}, (\mathbf{C}^{(1)})^T\Gamma\mathbf{C}^{(1)}) \\ &\quad - D_{\phi_3}(\mathbf{R}, (\mathbf{C}^{(1)})^T\Upsilon\mathbf{C}^{(2)}) + \tau \end{aligned} \quad (9)$$

where $\tau = \log f_{\phi_1}(\mathbf{F}) + \log f_{\phi_2}(\mathbf{S}) + \log f_{\phi_3}(\mathbf{R})$, which is a constant in the log-likelihood function.

Expectation Maximization (EM) is a general approach to find the maximum-likelihood estimate of the parameters when the model has latent variables. EM does maximum likelihood estimation by iteratively maximizing the expectation of the complete (log-)likelihood, which is the following under the MMRC model,

$$\mathbf{Q}(\Omega, \tilde{\Omega}) = \mathbf{E}[\log(\mathcal{L}(\Omega|\Psi))|\mathbf{C}^{(1)}, \mathbf{C}^{(2)}, \tilde{\Omega}], \quad (10)$$

where $\tilde{\Omega}$ denotes the current estimation of the parameters and Ω is the new parameters that we optimize to increase \mathbf{Q} . Two steps, E-step and M-step, are alternatively performed to maximize the objective function in Eq. (10).

4.2 Monte Carlo E-step

In the E-step, based on Bayes's rule, the posterior probability of the latent variables,

$$\begin{aligned} Pr(\mathbf{C}^{(1)}, \mathbf{C}^{(2)}|\mathbf{F}, \mathbf{S}, \mathbf{R}, \tilde{\Omega}) &= \\ \frac{Pr(\mathbf{C}^{(1)}, \mathbf{C}^{(2)}, \mathbf{F}, \mathbf{S}, \mathbf{R}|\tilde{\Omega})}{\sum_{\mathbf{C}^{(1)}, \mathbf{C}^{(2)}} Pr(\mathbf{C}^{(1)}, \mathbf{C}^{(2)}, \mathbf{F}, \mathbf{S}, \mathbf{R}|\tilde{\Omega})}, \end{aligned} \quad (11)$$

is updated using the current estimation of the parameters. However, conditioning on observations, the latent variables are not independent, i.e., there exist dependencies between the posterior probabilities of $\mathbf{C}^{(1)}$ and $\mathbf{C}^{(2)}$, and between those of $\mathbf{C}_{.p}^{(1)}$ and $\mathbf{C}_{.q}^{(1)}$. Hence, directly computing the posterior based on Eq. (11) is prohibitively expensive.

There exist several techniques for computing intractable posterior, such as Monte Carlo approaches, belief propagation, and variational methods. We follow a Monte Carlo approach, Gibbs sampler, which is a method of constructing a Markov chain whose stationary distribution is the distribution to be estimated.

It is easy to compute the posterior of a latent indicator vector while fixing all other latent indicator vectors, i.e.,

$$\begin{aligned} Pr(\mathbf{C}_{.-p}^{(1)}|\mathbf{C}_{.-p}^{(1)}, \mathbf{C}^{(2)}, \mathbf{F}, \mathbf{S}, \mathbf{R}, \tilde{\Omega}) &= \\ \frac{Pr(\mathbf{C}^{(1)}, \mathbf{C}^{(2)}, \mathbf{F}, \mathbf{S}, \mathbf{R}|\tilde{\Omega})}{\sum_{\mathbf{C}_{.-p}^{(1)}} Pr(\mathbf{C}^{(1)}, \mathbf{C}^{(2)}, \mathbf{F}, \mathbf{S}, \mathbf{R}|\tilde{\Omega})}, \end{aligned} \quad (12)$$

where $\mathbf{C}_{.-p}^{(1)}$ denotes all the latent indicator vectors except for $\mathbf{C}_{.p}^{(1)}$. Therefore, we present the following Markov chain to estimate the posterior in Eq. (11).

- Sample $\mathbf{C}_{.1}^{(1)}$
from distribution $Pr(\mathbf{C}_{.1}^{(1)}|\mathbf{C}_{.-1}^{(1)}, \mathbf{C}^{(2)}, \mathbf{F}, \mathbf{S}, \mathbf{R}, \tilde{\Omega})$;
-
- Sample $\mathbf{C}_{.n_1}^{(1)}$
from distribution $Pr(\mathbf{C}_{.n_1}^{(1)}|\mathbf{C}_{.-n_1}^{(1)}, \mathbf{C}^{(2)}, \mathbf{F}, \mathbf{S}, \mathbf{R}, \tilde{\Omega})$;
- Sample $\mathbf{C}_{.1}^{(2)}$
from distribution $Pr(\mathbf{C}_{.1}^{(2)}|\mathbf{C}_{.-1}^{(2)}, \mathbf{C}^{(1)}, \mathbf{F}, \mathbf{S}, \mathbf{R}, \tilde{\Omega})$;
-
- Sample $\mathbf{C}_{.n_2}^{(2)}$
from distribution $Pr(\mathbf{C}_{.n_2}^{(2)}|\mathbf{C}_{.-n_2}^{(2)}, \mathbf{C}^{(1)}, \mathbf{F}, \mathbf{S}, \mathbf{R}, \tilde{\Omega})$;

Note that at each sampling step in the above procedure, we use the latent indicator variables sampled from previous

steps. The above procedure iterates until the stop criterion is satisfied. It can be shown that the above procedure is a Markov chain converging to $Pr(\mathbf{C}^{(1)}, \mathbf{C}^{(2)} | \mathbf{F}, \mathbf{S}, \mathbf{R}, \tilde{\Omega})$. Assume that we keep l samples for estimation; then the posterior can be obtained simply by the empirical joint distribution of $\mathbf{C}^{(1)}$ and $\mathbf{C}^{(2)}$ in the l samples.

4.3 M-step

After the E-step, we have the posterior probability of latent variables to evaluate the expectation of the complete log-likelihood,

$$\mathbf{Q}(\Omega, \tilde{\Omega}) = \sum_{\mathbf{C}^{(1)}, \mathbf{C}^{(2)}} \log(\mathcal{L}(\Omega | \Psi)) Pr(\mathbf{C}^{(1)}, \mathbf{C}^{(2)} | \mathbf{F}, \mathbf{S}, \mathbf{R}, \tilde{\Omega}). \quad (13)$$

In the M-step, we optimize the unknown parameters by

$$\Omega^* = \arg \max_{\Omega} \mathbf{Q}(\Omega, \tilde{\Omega}). \quad (14)$$

First, we derive the update rules for membership parameters $\Lambda^{(1)}$ and $\Lambda^{(2)}$. To derive the expression for each $\Lambda_{hp}^{(1)}$, we introduce the Lagrange multiplier α with the constraint $\sum_{g=1}^{k_1} \Lambda_{gp}^{(1)} = 1$, and solve the following equation,

$$\frac{\partial}{\partial \Lambda_{hp}^{(1)}} \{ \mathbf{Q}(\Omega, \tilde{\Omega}) + \alpha (\sum_{g=1}^{k_1} \Lambda_{gp}^{(1)} - 1) \} = 0. \quad (15)$$

Substituting Eqs. (9) and (13) into Eq. (15), after some algebraic manipulations, we have

$$Pr(\mathbf{C}_{hp}^{(1)} = 1 | \mathbf{F}, \mathbf{S}, \mathbf{R}, \tilde{\Omega}) - \alpha \Lambda_{hp}^{(1)} = 0. \quad (16)$$

Summing both sides over h , we obtain $\alpha = 1$ resulting in the following update rule,

$$\Lambda_{hp}^{(1)} = Pr(\mathbf{C}_{hp}^{(1)} = 1 | \mathbf{F}, \mathbf{S}, \mathbf{R}, \tilde{\Omega}), \quad (17)$$

i.e., $\Lambda_{hp}^{(1)}$ is updated as the posterior probability that the p th object is associated with the h th cluster. Similarly, we have the following update rule for $\Lambda_{hp}^{(2)}$

$$\Lambda_{hp}^{(2)} = Pr(\mathbf{C}_{hp}^{(2)} = 1 | \mathbf{F}, \mathbf{S}, \mathbf{R}, \tilde{\Omega}). \quad (18)$$

Second, we derive the update rule for Θ . Based on Eqs. (9) and (13), optimizing Θ is equivalent to the following optimization,

$$\arg \min_{\Theta} \sum_{\mathbf{C}^{(1)}, \mathbf{C}^{(2)}} D_{\phi_1}(\mathbf{F}, \Theta \mathbf{C}^{(1)}) Pr(\mathbf{C}^{(1)}, \mathbf{C}^{(2)} | \mathbf{F}, \mathbf{S}, \mathbf{R}, \tilde{\Omega}). \quad (19)$$

We reformulated the above expression as,

$$\arg \min_{\Theta} \sum_{\mathbf{C}^{(1)}} \sum_{g=1}^{k_1} \sum_{p: \mathbf{C}_{gp}^{(1)}=1} D_{\phi_1}(\mathbf{F}_{\cdot p}, \Theta_{\cdot g}) Pr(\mathbf{C}_{gp}^{(1)} = 1 | \mathbf{F}, \mathbf{S}, \mathbf{R}, \tilde{\Omega}). \quad (20)$$

To solve the above optimization, we make use of an important property of Bregman divergence presented in the following theorem.

THEOREM 1. *Let X be a random variable taking values in $\mathcal{X} = \{\mathbf{x}_i\}_{i=1}^n \subset S \subseteq \mathbb{R}^d$ following v . Given a Bregman divergence $D_{\phi} : S \times \text{int}(S) \mapsto [0, \infty)$, the problem*

$$\min_{s \in S} E_v[D_{\phi}(X, s)] \quad (21)$$

has a unique minimizer given by $s^ = E_v[X]$.*

The proof of Theorem 1 is omitted (please refer [3, 40]). Theorem 1 states that the Bregman representative of a random variable is always the expectation of the variable. Based on Theorem 1 and the objective function in (20), we update $\Theta_{\cdot g}$ as follows,

$$\Theta_{\cdot g} = \frac{\sum_{p=1}^{n_1} \mathbf{F}_{\cdot p} Pr(\mathbf{C}_{gp}^{(1)} = 1 | \mathbf{F}, \mathbf{S}, \mathbf{R}, \tilde{\Omega})}{\sum_{p=1}^{n_1} Pr(\mathbf{C}_{gp}^{(1)} = 1 | \mathbf{F}, \mathbf{S}, \mathbf{R}, \tilde{\Omega})}. \quad (22)$$

Third, we derive the update rule for Γ . Based on Eqs. (9) and (13), we formulate optimizing Γ as the following optimization,

$$\arg \min_{\Gamma} \sum_{\mathbf{C}^{(1)}} \sum_{g=1}^{k_1} \sum_{h=1}^{k_1} \sum_{p: \mathbf{C}_{gp}^{(1)}=1, q: \mathbf{C}_{hq}^{(1)}=1} D_{\phi_2}(\mathbf{S}_{pq}, \Gamma_{gh}) \tilde{p}, \quad (23)$$

where \tilde{p} denotes $Pr(\mathbf{C}_{gp}^{(1)} = 1, \mathbf{C}_{hq}^{(1)} = 1 | \mathbf{F}, \mathbf{S}, \mathbf{R}, \tilde{\Omega})$ and $1 \leq p, q \leq n_1$. Based on Theorem 1, we update each Γ_{gh} as follows,

$$\Gamma_{gh} = \frac{\sum_{p, q=1}^{n_1} \mathbf{S}_{pq} Pr(\mathbf{C}_{gp}^{(1)} = 1, \mathbf{C}_{hq}^{(1)} = 1 | \mathbf{F}, \mathbf{S}, \mathbf{R}, \tilde{\Omega})}{\sum_{p, q=1}^{n_1} Pr(\mathbf{C}_{gp}^{(1)} = 1, \mathbf{C}_{hq}^{(1)} = 1 | \mathbf{F}, \mathbf{S}, \mathbf{R}, \tilde{\Omega})}. \quad (24)$$

Fourth, we derive the update rule for Υ . Based on Eqs. (9) and (13), we formulate optimizing Υ as the following optimization,

$$\arg \min_{\Upsilon} \sum_{\mathbf{C}^{(1)}, \mathbf{C}^{(2)}} \sum_{g=1}^{k_1} \sum_{h=1}^{k_2} \sum_{p: \mathbf{C}_{gp}^{(1)}=1, q: \mathbf{C}_{hq}^{(2)}=1} D_{\phi_3}(\mathbf{R}_{pq}, \Upsilon_{gh}) \tilde{p}, \quad (25)$$

where \tilde{p} denotes $Pr(\mathbf{C}_{gp}^{(1)} = 1, \mathbf{C}_{hq}^{(2)} = 1 | \mathbf{F}, \mathbf{S}, \mathbf{R}, \tilde{\Omega})$, $1 \leq p \leq n_1$ and $1 \leq q \leq n_2$. Based on Theorem 1, we update each Υ_{gh} as follows,

$$\Upsilon_{gh} = \frac{\sum_{p=1}^{n_1} \sum_{q=1}^{n_2} \mathbf{R}_{pq} Pr(\mathbf{C}_{gp}^{(1)} = 1, \mathbf{C}_{hq}^{(2)} = 1 | \mathbf{F}, \mathbf{S}, \mathbf{R}, \tilde{\Omega})}{\sum_{p=1}^{n_1} \sum_{q=1}^{n_2} Pr(\mathbf{C}_{gp}^{(1)} = 1, \mathbf{C}_{hq}^{(2)} = 1 | \mathbf{F}, \mathbf{S}, \mathbf{R}, \tilde{\Omega})}. \quad (26)$$

Combining the E-step and M-step, we have a general relational clustering algorithm, Exponential Family MMRC (EF-MMRC) algorithm, which is summarized in Algorithm 1. Since it is straightforward to apply our algorithm derivation to a relational data set of any structure, Algorithm 1 is proposed based on the input of a general relational data set. Despite that the input relational data could have various structures, EF-MMRC works simply as follows: in the E-step, EF-MMRC iteratively updates the posterior probabilities that an object is associated with the clusters (the Markov chain in Section 4.2); in the M-step, based on the current cluster association (posterior probabilities), the cluster representatives of attributes and relations are updated as the weighted mean of the observations no matter which exponential distributions are assumed.

Therefore, with the simplicity of the traditional centroid-based clustering algorithms, EF-MMRC is capable of making use of all attribute information and homogeneous and heterogeneous relation information to learn hidden structures from various relational data. Since EF-MMRC simultaneously clusters multi-type interrelated objects, the cluster structures of different types of objects may interact with each other directly or indirectly during the clustering process to automatically deal with the influence propagation. Besides the local cluster structures for each type of objects,

Algorithm 1 Exponential Family MMRC Algorithm

Input: A relational data set $\{\{\mathbf{F}^{(j)}\}_{j=1}^m, \{\mathbf{S}^{(j)}\}_{j=1}^m, \{\mathbf{R}^{(ij)}\}_{i,j=1}^m\}$, a set of exponential family distributions (Bregman divergences) assumed for the data set.

Output: Membership Matrices $\{\Lambda^{(j)}\}_{j=1}^m$, attribute expectation matrices $\{\Theta^{(j)}\}_{j=1}^m$, homogeneous relation expectation matrices $\{\Gamma^{(j)}\}_{j=1}^m$, and heterogeneous relation expectation matrices $\{\Upsilon^{(ij)}\}_{i,j=1}^m$.

Method:

```
1: Initialize the parameters as  $\tilde{\Omega} = \{\{\tilde{\Lambda}^{(j)}\}_{j=1}^m, \{\tilde{\Theta}^{(j)}\}_{j=1}^m, \{\tilde{\Gamma}^{(j)}\}_{j=1}^m, \{\tilde{\Upsilon}^{(ij)}\}_{i,j=1}^m\}$ .
2: repeat
3:   {E-step}
4:   Compute the posterior  $Pr(\{\mathbf{C}^{(j)}\}_{j=1}^m | \{\mathbf{F}^{(j)}\}_{j=1}^m, \{\mathbf{S}^{(j)}\}_{j=1}^m, \{\mathbf{R}^{(ij)}\}_{i,j=1}^m, \tilde{\Omega})$  using the Gibbs sampler.
5:   {M-step}
6:   for  $j = 1$  to  $m$  do
7:     Compute  $\Lambda^{(j)}$  using update rule (17).
8:     Compute  $\Theta^{(j)}$  using update rule (22).
9:     Compute  $\Gamma^{(j)}$  using update rule (24).
10:    for  $i = 1$  to  $m$  do
11:      Compute  $\Upsilon^{(ij)}$  using update rule (26).
12:    end for
13:  end for
14:   $\tilde{\Omega} = \Omega$ 
15: until convergence
```

the output of EF-MMRC also provides the summary of the global hidden structure for the data, i.e., based on Γ and Υ , we know how the clusters of the same type and different types are related to each other. Furthermore, relational data from different applications may have different probabilistic distributions on the attributes and relations; it is easy for EF-MMRC to adapt to this situation by simply using different Bregman divergences corresponding to different exponential family distributions.

If we assume $O(m)$ types of heterogeneous relations among m types of objects, which is typical in real applications, and let $n = \Theta(n_i)$ and $k = \Theta(k_i)$, the computational complexity of EF-MMRC can be shown to be $O(tmn^2k)$ for t iterations. If we apply the k-means algorithm to each type of nodes individually by transforming the relations into attributes for each type of nodes, the total computational complexity is also $O(tmn^2k)$.

4.4 Hard MMRC Algorithm

Due to its simplicity, scalability, and broad applicability, k-means algorithm has become one of the most popular clustering algorithms. Hence, it is desirable to extend k-means to relational data. Some efforts [47, 2, 12, 33] in the literature work in this direction. However, these approaches apply to only some special and simple cases of relational data, such as bi-type heterogeneous relational data.

As traditional k-means can be formulated as a hard version of Gaussian mixture model EM algorithm [29], we propose the hard version of MMRC algorithm as a general relational k-means algorithm (from now on, we call Algorithm 1 as soft EF-MMRC), which applies to various relational data.

To derive the hard version MMRC algorithm, we omit soft membership parameters $\Lambda^{(j)}$ in the MMRC model ($\mathbf{C}^{(j)}$ in

the model provides the hard membership for each object). Next, we change the computation of the posterior probabilities in the E-step to reassignment procedure, i.e., in the E-step, based on the estimation of the current parameters, we re-assign cluster labels, $\{\mathbf{C}^{(j)}\}_{j=1}^m$, to maximize the objective function in (9). In particular, for each object, while fixing the cluster assignments of all other objects, we assign it to each cluster to find the optimal cluster assignment maximizing the objective function in (9), which is equivalent to minimizing the Bregman distances between the observations and the corresponding expectation parameters. After all objects are assigned, the re-assignment process is repeated until no object changes its cluster assignment between two successive iterations.

In the M-step, we estimate the parameters based on the cluster assignments from the E-step. A simple way to derive the update rules is to follow the derivation in Section 4.3 but replace the posterior probabilities by its hard versions. For example, after the E-step, if the object $x_p^{(j)}$ is assigned to the g th cluster, i.e., $\mathbf{C}_{gp}^{(j)} = 1$, then the posterior $Pr(\mathbf{C}_{gp}^{(1)} = 1 | \mathbf{F}, \mathbf{S}, \mathbf{R}, \tilde{\Omega}) = 1$ and $Pr(\mathbf{C}_{hp}^{(1)} = 1 | \mathbf{F}, \mathbf{S}, \mathbf{R}, \tilde{\Omega}) = 0$ for $h \neq g$.

Using the hard versions of the posterior probabilities, we derive the following update rule for $\Theta^{(j)}$,

$$\Theta_{\cdot g}^{(j)} = \frac{\sum_{p: \mathbf{C}_{gp}^{(j)}=1} \mathbf{F}_{\cdot p}^{(j)}}{\sum_{p=1}^{n_j} \mathbf{C}_{gp}^{(j)}}. \quad (27)$$

In the above update rule, since $\sum_{p=1}^{n_1} \mathbf{C}_{gp}^{(j)}$ is the size of the g th cluster, $\Theta_{\cdot g}^{(j)}$ is actually updated as the mean of the attribute vectors of the objects assigned to the g th cluster.

Similarly, we have the following update rule for $\Gamma^{(j)}$

$$\Gamma_{gh}^{(j)} = \frac{\sum_{p: \mathbf{C}_{gp}^{(j)}=1, q: \mathbf{C}_{hq}^{(j)}=1} \mathbf{S}_{pq}^{(j)}}{\sum_{p=1}^{n_j} \mathbf{C}_{gp}^{(j)} \sum_{q=1}^{n_j} \mathbf{C}_{hq}^{(j)}}, \quad (28)$$

i.e., $\Gamma_{gh}^{(j)}$ is updated as the mean of the relations between the objects of the j th type from the g th cluster and from the h th cluster.

Each heterogeneous relation expectation parameter $\Upsilon_{gh}^{(ij)}$ is updated as the mean of the objects of the i th type from the g th cluster and of the j th type from the h th cluster,

$$\Upsilon_{gh}^{(ij)} = \frac{\sum_{p: \mathbf{C}_{gp}^{(i)}=1, q: \mathbf{C}_{hq}^{(j)}=1} \mathbf{R}_{pq}^{(ij)}}{\sum_{p=1}^{n_i} \mathbf{C}_{gp}^{(i)} \sum_{q=1}^{n_j} \mathbf{C}_{hq}^{(j)}}. \quad (29)$$

The hard version of EF-MMRC algorithm is summarized in Algorithm 2. It works simply as the classic k-means. However, it is applicable to various relational data under various Bregman distance functions corresponding to various assumptions of probability distributions. Based on the EM framework, its convergence is guaranteed. When applied to some special cases of relational data, it provides simple and new algorithms for some important data mining problems. For example, when applied to the data of one homogeneous relation matrix representing a graph affinity matrix, it provides a simple and new graph partitioning algorithm.

Based on Algorithms 1 and 2, there is another version of EF-MMRC, i.e., we may combine soft and hard EF-MMRC together to have mixed EF-MMRC. For example, we first run hard EF-MMRC several times as initialization, then run soft EF-MMRC.

Algorithm 2 Hard MMRC Algorithm

Input: A relational data set $\{\{\mathbf{F}^{(j)}\}_{j=1}^m, \{\mathbf{S}^{(j)}\}_{j=1}^m, \{\mathbf{R}^{(ij)}\}_{i,j=1}^m\}$, a set of exponential family distributions (Bregman divergences) assumed for the data set.

Output: Cluster indicator matrices $\{\mathbf{C}^{(j)}\}_{j=1}^m$, attribute expectation matrices $\{\Theta^{(j)}\}_{j=1}^m$, homogeneous relation expectation matrices $\{\Gamma^{(j)}\}_{j=1}^m$, and heterogeneous relation expectation matrices $\{\Upsilon^{(ij)}\}_{i,j=1}^m$.

Method:

- 1: Initialize the parameters as $\tilde{\Omega} = \{\{\tilde{\Lambda}^{(j)}\}_{j=1}^m, \{\tilde{\Theta}^{(j)}\}_{j=1}^m, \{\tilde{\Gamma}^{(j)}\}_{j=1}^m, \{\tilde{\Upsilon}^{(ij)}\}_{i,j=1}^m\}$.
- 2: **repeat**
- 3: {E-step}
- 4: Based on the current parameters, reassign cluster labels for each objects, i.e., update $\{\mathbf{C}^{(j)}\}_{j=1}^m$, to maximize the objective function in Eq. (9).
- 5: {M-step}
- 6: **for** $j = 1$ to m **do**
- 7: Compute $\Theta^{(j)}$ using update rule (27).
- 8: Compute $\Gamma^{(j)}$ using update rule (28).
- 9: **for** $i = 1$ to m **do**
- 10: Compute $\Upsilon^{(ij)}$ using update rule (29).
- 11: **end for**
- 12: **end for**
- 13: $\Omega = \tilde{\Omega}$
- 14: **until** convergence

5. A UNIFIED VIEW TO CLUSTERING

In this section we discuss the connections between existing clustering approaches and the MMRF model and EF-MMRF algorithms. By considering them as special cases or variations of the MMRF model, we show that MMRF provides a unified view to the existing clustering approaches from various important data mining applications.

5.1 Semi-supervised Clustering

Recently, semi-supervised clustering has become a topic of significant interest [4, 46], which seeks to cluster a set of data points with a set of pairwise constraints.

Semi-supervised clustering can be formulated as a special case of relational clustering, clustering on the single-type relational data set consisting of attributes \mathbf{F} and homogeneous relations \mathbf{S} . For semi-supervised clustering, S_{pq} denotes the pairwise constraint on the p th object and the q th object.

[4] provides a general model for semi-supervised clustering based on Hidden Markov Random Fields (HMRFs). We show that it can be formulated as a special case of MMRC model. As in [4], we define the homogeneous relation matrix S as follows,

$$\mathbf{S}_{pq} = \begin{cases} f_M(x_p, x_q) & \text{if } (x_p, x_q) \in \mathcal{M} \\ f_C(x_p, x_q) & \text{if } (x_p, x_q) \in \mathcal{C} \\ 0 & \text{otherwise} \end{cases}$$

where \mathcal{M} denotes a set of must-link constraints; \mathcal{C} denotes a set of cannot-link constraints; $f_M(x_p, x_q)$ is a function that penalizes the violation of must-link constraint; $f_C(x_p, x_q)$ is a penalty function for cannot-links. If we assume Gibbs distribution [41] for \mathbf{S} ,

$$Pr(S) = \frac{1}{z_1} \exp\left(-\sum_{p,q} \mathbf{S}_{pq}\right). \quad (30)$$

where z_1 is the normalization constant. Since [4] focuses on

only hard clustering, we omit the soft member parameters in the MMRC model to consider hard clustering. Based on Eq.(30) and Eq.(4), the likelihood function of hard semi-supervised clustering under MMRC model is

$$L(\Theta|\mathbf{F}) = \frac{1}{z} \exp\left(-\sum_{p,q} \mathbf{S}_{pq}\right) \exp(-D_\phi(F, \Lambda \mathbf{C})) \quad (31)$$

Since \mathbf{C} is an indicator matrix, Eq. (31) can be formulated as

$$L(\Theta|\mathbf{F}) = \frac{1}{z} \exp\left(-\sum_{p,q} \mathbf{S}_{pq}\right) \exp\left(-\sum_{g=1}^k \sum_{p: \mathbf{C}_{gp}=1} D_\phi(F_p, \Lambda \cdot g)\right) \quad (32)$$

The above likelihood function is equivalent to the objective function of semi-supervised clustering based on HMRFs [4]. Furthermore, when applied to optimizing the objective function in Eq.(32), hard MMRC provides a family of semi-supervised clustering algorithms similar to HMRF-KMeans in [4]; on the other hand, soft EF-MMRC provides new and soft version semi-supervised clustering algorithms.

5.2 Co-clustering

Co-clustering or bi-clustering arise in many important applications, such as document clustering, micro-array data clustering. A number of approaches [12, 8, 33, 2] have been proposed for co-clustering. These efforts can be generalized as solving the following matrix approximation problem [34],

$$\arg \min_{\mathbf{C}, \Upsilon} \mathfrak{D}(\mathbf{R}, (\mathbf{C}^{(1)})^T \Upsilon \mathbf{C}^{(2)}) \quad (33)$$

where $R \in \mathbb{R}^{n_1 \times n_2}$ is the data matrix, $\mathbf{C}^{(1)} \in \{0, 1\}^{k_1 \times n_1}$ and $\mathbf{C}^{(2)} \in \{0, 1\}^{k_2 \times n_2}$ are indicator matrices, $\Upsilon \in \mathbb{R}^{k_1 \times k_2}$ is the relation representative matrix, and \mathfrak{D} is a distance function. For example, [12] uses KL-divergences as the distance function; [8, 33] use Euclidean distances.

Co-clustering is equivalent to clustering on relational data of one heterogeneous relation matrix \mathbf{R} . Based on Eq.(9), by omitting the soft membership parameters, maximizing log-likelihood function of hard clustering on a heterogeneous relation matrix under the MMRC model is equivalent to the minimization in (33). The algorithms proposed in [12, 8, 33, 2] can be viewed as special cases of hard EF-MMRC. At the same time, soft EF-MMRC provides another family of new algorithms for co-clustering.

Our previous work [34] proposes the relation summary network model for clustering k-partite graphs, which can be shown to be equivalent on clustering on relational data of multiple heterogeneous relation matrices. The proposed algorithms in [34] can also be viewed as special cases of the hard EF-MMRC algorithm.

5.3 Graph Clustering

Graph clustering (partitioning) is an important problem in many domains, such as circuit partitioning, VLSI design, task scheduling. Existing graph partitioning approaches are mainly based on edge cut objectives, such as Kernighan-Lin objective [30], normalized cut [42], ratio cut [7], ratio association [42], and min-max cut [13].

Graph clustering is equivalent to clustering on single-type relational data of one homogeneous relation matrix \mathbf{S} . The log-likelihood function of the hard clustering under MMRC model is $-D_\phi(\mathbf{S}, (\mathbf{C})^T \Gamma \mathbf{C})$. We propose the following theorem to show that the edge cut objectives are mathematically equivalent to a special case of the MMRC model. Since most

graph partitioning objective functions use weighted indicator matrix such that $\mathbf{C}\mathbf{C}^T = \mathbf{I}_k$, where \mathbf{I}_k is an identity matrix, we follow this formulation in the following theorem.

THEOREM 2. *With restricting Γ to be the form of $r\mathbf{I}_k$ for $r > 0$, maximizing the log-likelihood of hard MMRC clustering on \mathbf{S} under normal distribution, i.e.,*

$$\max_{\mathbf{C} \in \{0,1\}^{k \times n}, \mathbf{C}\mathbf{C}^T = \mathbf{I}_k} -\|\mathbf{S} - (\mathbf{C})^T(r\mathbf{I}_k)\mathbf{C}\|^2, \quad (34)$$

is equivalent to the trace maximization

$$\max \text{tr}(\mathbf{C}\mathbf{S}\mathbf{C}^T), \quad (35)$$

where tr denotes the trace of a matrix.

PROOF. Let L denote the objective function in Eq. (34).

$$\begin{aligned} L &= -\|\mathbf{S} - r\mathbf{C}^T\mathbf{C}\|^2 \\ &= -\text{tr}((\mathbf{S} - r\mathbf{C}^T\mathbf{C})^T(\mathbf{S} - r\mathbf{C}^T\mathbf{C})) \\ &= -\text{tr}(\mathbf{S}^T\mathbf{S}) + 2r\text{tr}(\mathbf{C}^T\mathbf{C}\mathbf{S}) - r^2\text{tr}(\mathbf{C}^T\mathbf{C}\mathbf{C}^T\mathbf{C}) \\ &= -\text{tr}(\mathbf{S}^T\mathbf{S}) + 2r\text{tr}(\mathbf{C}\mathbf{S}\mathbf{C}^T) - r^2k \end{aligned}$$

The above deduction uses the property of trace $\text{tr}(\mathbf{X}\mathbf{Y}) = \text{tr}(\mathbf{Y}\mathbf{X})$. Since $\text{tr}(\mathbf{S}^T\mathbf{S})$, r and k are constants, the maximization of L is equivalent to the maximization of $\text{tr}(\mathbf{C}\mathbf{S}\mathbf{C}^T)$. The proof is completed. \square

Since it is shown in the literature [10] that the edge cut objectives can be formulated as the trace maximization, Theorem 2 states that edge-cut based graph clustering is equivalent to MMRC model under normal distribution with the diagonal constraint on the parameter matrix Γ . This connection provides not only a new understanding for graph partitioning but also a family of new algorithms (soft and hard MMRC algorithms) for graph clustering.

Finally, we point out that MMRC model does not exclude traditional attribute-based clustering. When applied to an attribute data matrix under Euclidean distances, hard MMRC algorithm is actually reduced to the classic k-means; soft MMRC algorithm is very close to the traditional mixture model EM clustering except that it does not involve mixing proportions in the computation.

In summary, MMRC model provides a principal framework to unify various important clustering tasks including traditional attributes-based clustering, semi-supervised clustering, co-clustering and graph clustering; soft and hard EF-MMRC algorithms unify a number of state-of-the-art clustering algorithms and at the same time provide new solutions to various clustering tasks.

6. EXPERIMENTS

This section provides empirical evidence to show the effectiveness of the MMRC model and algorithms. Since a number of state-of-the-art clustering algorithms [12, 8, 33, 2, 3, 4] can be viewed as special cases of EF-MMRC model and algorithms, the experimental results in these efforts also illustrate the effectiveness of the MMRC model and algorithms. In this paper, we apply MMRC algorithms to tasks of graph clustering, bi-clustering, tri-clustering, and clustering on a general relational data set of all three types of information. In the experiments, we use mixed version MMRC, i.e., hard MMRC initialization followed by soft MMRC. Although MMRC can adopt various distribution assumptions, due to space limit, we use MMRC under normal or Poisson distribution assumption in the experiments. However, this

Table 1: Summary of relational data for Graph Clustering.

Name	n	k	Balance	Source
tr11	414	9	0.046	TREC
tr23	204	6	0.066	TREC
NG1-20	14000	20	1.0	20-newsgroups
k1b	2340	6	0.043	WebACE

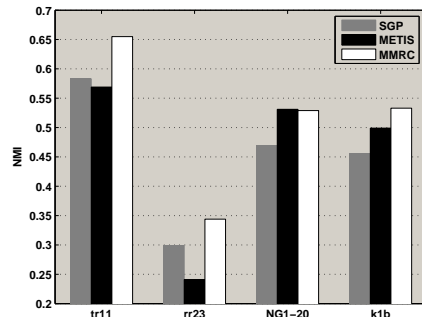


Figure 2: NMI comparison of SGP, METIS and MMRC algorithms.

does not imply that they are optimal distribution assumptions for the data. How to decide the optimal distribution assumption is beyond the scope of this paper.

For performance measure, we elect to use the Normalized Mutual Information (NMI) [44] between the resulting cluster labels and the true cluster labels, which is a standard way to measure the cluster quality. The final performance score is the average of ten runs.

6.1 Graph Clustering

In this section, we present experiments on the MMRC algorithm under normal distribution in comparison with two representative graph partitioning algorithms, the spectral graph partitioning (SGP) from [36] that is generalized to work with both normalized cut and ratio association, and the classic multilevel algorithm, METIS [28].

The graphs based on the text data have been widely used to test graph partitioning algorithms [13, 11, 25]. In this study, we use various data sets from the 20-newsgroups [32], WebACE and TREC [27], which cover data sets of different sizes, different balances and different levels of difficulties. The data are pre-processed by removing the stop words and each document is represented by a term-frequency vector using TF-IDF weights. Then we construct relational data for each text data set such that objects (documents) are related to each other with cosine similarities between the term-frequency vectors. A summary of all the data sets to construct relational data used in this paper is shown in Table 1, in which n denotes the number of objects in the relational data, k denotes the number of true clusters, and *balance* denotes the size ratio of the smallest clusters to the largest clusters.

For the number of clusters k , we simply use the number of the true clusters. Note that how to choose the optimal number of clusters is a nontrivial model selection problem and beyond the scope of this paper.

Figure 2 shows the NMI comparison of the three algorithms. We observe that although there is no single winner on all the graphs, overall the MMRC algorithm performs better than SGP and METIS. Especially on the difficult data set tr23, MMRC increases performance about 30%. Hence, MMRC under normal distribution provides a new graph partitioning algorithm which is viable and competi-

Data set	Taxonomy structure
<i>TT-TM1</i>	{rec.sport.baseball, rec.sport.hockey}, {talk.politics.guns, talk.politics.mideast, talk.politics.misc}
<i>TT-TM2</i>	{comp.graphics, comp.os.ms-windows.misc}, {rec.autos, rec.motorcycles}, {sci.crypt, sci.electronics}

Table 3: Taxonomy structures of two data sets for constructing tri-partite relational data

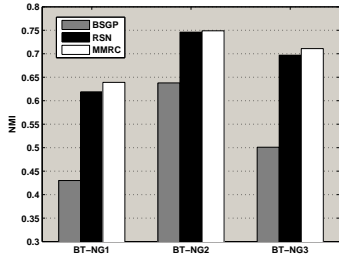


Figure 3: NMI comparison of BSGP, RSN and MMRC algorithms for bi-type data.

tive compared with the two existing state-of-the-art graph partitioning algorithms. Note that although the normal distribution is most popular, MMRC under other distribution assumptions may be more desirable in specific graph clustering applications depends on the statistical properties of the graphs.

6.2 Bi-clustering and Tri-clustering

In this section, we apply the MMRC algorithm under Poisson distribution to clustering bi-type relational data, word-document data, and tri-type relational data, word-document-category data. Two algorithms, Bi-partite Spectral Graph partitioning (BSGP) [11] and Relation Summary Network under Generalized I-divergence (RSN-GI) [34], are used as comparison in bi-clustering. For tri-clustering, Consistent Bipartite Graph Co-partitioning (CBGC) [18] and RSN-GI are used as comparison.

The bi-type relational data, word-document data, are constructed based on various subsets of the 20-News group data. We pre-process the data by selecting the top 2000 words by the mutual information. The document-word matrix is based on *tf.idf* weighting scheme and each document vector is normalized to a unit L_2 norm vector. Specific details of the data sets are listed in Table 2. For example, for the data set *BT-NG3* we randomly and evenly sample 200 documents from the corresponding newsgroups; then we formulate a bi-type relational data set of 1600 document and 2000 word.

The tri-type relational data are built based on the 20-newsgroups data for hierarchical taxonomy mining. In the field of text categorization, hierarchical taxonomy classifica-

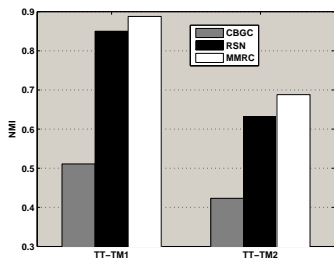


Figure 4: NMI comparison of CBGC, RSN and MMRC algorithms for tri-type data.

cluster 23 of actors
Viggo Mortensen, Sean Bean, Miranda Otto, Ian Holm, Christopher Lee, Cate Blanchett, Ian McKellen, Liv Tyler, David Wenham, Brad Dourif, John Rhys-Davies, Elijah Wood, Bernard Hill, Sean Astin, Andy Serkis, Dominic Monaghan, Karl Urban, Orlando Bloom, Billy Boyd, John Noble, Sala Baker
cluster 118 of movies
The Lord of the Rings: The Fellowship of the Ring (2001) The Lord of the Rings: The Two Towers (2002) The Lord of the Rings: The Return of the King (2003)

Table 4: Two clusters from actor-movie data

tion is widely used to obtain a better trade-off between effectiveness and efficiency than flat taxonomy classification. To take advantage of hierarchical classification, one must mine a hierarchical taxonomy from the data set. We see that words, documents, and categories formulate a sandwich structure tri-type relational data set, in which documents are central type nodes. The links between documents and categories are constructed such that if a document belongs to k categories, the weights of links between this document and these k category nodes are $1/k$ (please refer [18] for details). The true taxonomy structures for two data sets, *TP-TM1* and *TP-TM2*, are documented in Table 3.

Figure 3 and Figure 4 show the NMI comparison of the three algorithms on bi-type and tri-type relational data, respectively. We observe that the MMRC algorithm performs significantly better than BSGP and CBGC. MMRC performs slightly better than RSN on some data sets. Since RSN is a special case of hard MMRC, this shows that mixed MMRC improves hard MMRC’s performance on the data sets. Therefore, compared with the existing stated-of-the-art algorithms, the MMRC algorithm performs more effectively on these bi-clustering or tri-clustering tasks and on the other hand, it is flexible for different types of multi-clustering tasks which may be more complicated than tri-type clustering.

6.3 A Case Study on Actor-movie Data

We also run the MMRC algorithm on the actor-movie relational data based on IMDB movie data set for a case study. In the data, actors are related to each other by collaboration (homogeneous relations); actors are related to movies by taking roles in movies (heterogeneous relations); movies have attributes such as release time and rating (note that there is no links between movies). Hence the data have all the three types of information. We formulate a data set of 20000 actors and 4000 movies. We run experiments with $k = 200$. Although there is no ground truth for the data’s cluster structure, we observe that most resulting clusters that are actors or movies of the similar style such as action, or tight groups from specific movie serials. For example, Table 4 shows cluster 23 of actors and cluster 118 of movies; the parameter $\Upsilon_{23,118}$ shows that these two clusters are strongly related to each other. In fact, the actor cluster contains the actors in the movie series “The Lord of the Rings”. Note that if we only have one type of actor objects, we only get the actor clusters, but with two types of nodes, although there is no links between the movies, we also get the related movie clusters to explain how the actors are related.

7. CONCLUSIONS

In this paper, we propose a probabilistic model for relational clustering, which provides a principal framework to unify various important clustering tasks including tradi-

Dataset Name	Newsgroups Included	# Documents per Group	Total # Documents
<i>BT-NG1</i>	rec.sport.baseball, rec.sport.hockey	200	400
<i>BT-NG2</i>	comp.os.ms-windows.misc, comp.windows.x, rec.motorcycles, sci.crypt, sci.space	200	1000
<i>BT-NG3</i>	comp.os.ms-windows.misc, comp.windows.x, misc.forsale, rec.motorcycles, rec.motorcycles.sci.crypt, sci.space, talk.politics.mideast, talk.religion.misc	200	1600

Table 2: Subsets of Newsgroup Data for bi-type relational data

tional attributes-based clustering, semi-supervised clustering, co-clustering and graph clustering. Under this model, we propose parametric hard and soft relational clustering algorithms under a large number of exponential family distributions. The algorithms are applicable to relational data of various structures and at the same time unify a number of state-of-the-art clustering algorithms. The theoretic analysis and experimental evaluation show the effectiveness and great potential of the proposed model and algorithms.

8. REFERENCES

- [1] E. Airoldi, D. Blei, E. Xing, and S. Fienberg. Mixed membership stochastic block models for relational data with application to protein-protein interactions. In *ENAR-2006*.
- [2] A. Banerjee, I. S. Dhillon, J. Ghosh, S. Merugu, and D. S. Modha. A generalized maximum entropy approach to bregman co-clustering and matrix approximation. In *KDD*, pages 509–514, 2004.
- [3] A. Banerjee, S. Merugu, I. S. Dhillon, and J. Ghosh. Clustering with bregman divergences. *J. Mach. Learn. Res.*, 6:1705–1749, 2005.
- [4] S. Basu, M. Bilenko, and R. J. Mooney. A probabilistic framework for semi-supervised clustering. In *KDD04*, pages 59–68, 2004.
- [5] I. Bhattacharya and L. Getor. Entity resolution in graph data. Technical Report CS-TR-4758, University of Maryland, 2005.
- [6] T. N. Bui and C. Jones. A heuristic for reducing fill-in in sparse matrix factorization. In *PPSC*, pages 445–452, 1993.
- [7] P. K. Chan, M. D. F. Schlag, and J. Y. Zien. Spectral k-way ratio-cut partitioning and clustering. In *DAC '93*.
- [8] H. Cho, I. Dhillon, Y. Guan, and S. Sra. Minimum sum squared residue co-clustering of gene expression data. In *SDM*, 2004.
- [9] M. Collins, S. Dasgupta, and R. Reina. A generalization of principal component analysis to the exponential family. In *NIPS'01*, 2001.
- [10] I. Dhillon, Y. Guan, and B. Kulis. A unified view of kernel k-means, spectral clustering and graph cuts. Technical Report TR-04-25, University of Texas at Austin, 2004.
- [11] I. S. Dhillon. Co-clustering documents and words using bipartite spectral graph partitioning. In *KDD'01*.
- [12] I. S. Dhillon, S. Mallela, and D. S. Modha. Information-theoretic co-clustering. In *KDD'03*, pages 89–98.
- [13] C. H. Q. Ding, X. He, H. Zha, M. Gu, and H. D. Simon. A min-max cut algorithm for graph partitioning and data clustering. In *Proceedings of ICDM 2001*, pages 107–114, 2001.
- [14] S. Dzeroski and N. Lavrac, editors. *Relational Data Mining*. Springer, 2001.
- [15] E. Erosheva, S. Fienberg, and J. Lafferty. Mixed membership models of scientific publications. In *NAS*.
- [16] E. Erosheva and S. E. Fienberg. Bayesian mixed membership models for soft clustering and classification. *Classification-The Ubiquitous Challenge*, pages 11–26, 2005.
- [17] S. E. Fienberg, M. M. Meyer, and S. Wasserman. Statistical analysis of multiple sociometric relations. *Journal of American Statistical Association*, 80:51–87, 1985.
- [18] B. Gao, T.-Y. Liu, X. Zheng, Q.-S. Cheng, and W.-Y. Ma. Consistent bipartite graph co-partitioning for star-structured high-order heterogeneous data co-clustering. In *KDD '05*, pages 41–50, 2005.
- [19] L. Getor. An introduction to probabilistic graphical models for relational data. *Data Engineering Bulletin*, 29, 2006.
- [20] B. Hendrickson and R. Leland. A multilevel algorithm for partitioning graphs. In *Supercomputing '95*, page 28, 1995.
- [21] P. Hoff, A. Raftery, and M. Handcock. Latent space approaches to social network analysis. *Journal of American Statistical Association*, 97:1090–1098, 2002.
- [22] T. Hofmann. Probabilistic latent semantic analysis. In *Proc. of Uncertainty in Artificial Intelligence, UAI'99*, Stockholm, 1999.
- [23] T. Hofmann and J. Puzicha. Latent class models for collaborative filtering. In *IJCAI'99*, Stockholm, 1999.
- [24] L. B. Holder and D. J. Cook. Graph-based relational learning: current and future directions. *SIGKDD Explor. Newsl.*, 5(1):90–93, 2003.
- [25] M. X. H. Zha, C. Ding and H. Simon. Bi-partite graph partitioning and data clustering. In *ACM CIKM'01*, 2001.
- [26] G. Jeh and J. Widom. Simrank: A measure of structural-context similarity. In *KDD-2002*, 2002.
- [27] G. Karypis. A clustering toolkit, 2002.
- [28] G. Karypis and V. Kumar. A fast and high quality multilevel scheme for partitioning irregular graphs. *SIAM J. Sci. Comput.*, 20(1):359–392, 1998.
- [29] M. Kearns, Y. Mansour, and A. Ng. An information-theoretic analysis of hard and soft assignment methods for clustering. In *UAI'97*, pages 282–293, 2004.
- [30] B. Kernighan and S. Lin. An efficient heuristic procedure for partitioning graphs. *The Bell System Technical Journal*, 49(2):291–307, 1970.
- [31] M. Kirsten and S. Wrobel. Relational distance-based clustering. In *Proc. Fachgruppentreffen Maschinelles Lernen (FGML-98)*, pages 119 – 124, 1998.
- [32] K. Lang. News weeder: Learning to filter netnews. In *ICML*, 1995.
- [33] T. Li. A general model for clustering binary data. In *KDD'05*, 2005.
- [34] B. Long, X. Wu, Z. M. Zhang, and P. S. Yu. Unsupervised learning on k-partite graphs. In *KDD-2006*, 2006.
- [35] B. Long, Z. Zhang, and P. Yu. Co-clustering by block value decomposition. In *KDD'05*, 2005.
- [36] A. Ng, M. Jordan, and Y. Weiss. On spectral clustering: Analysis and an algorithm. In *Advances in Neural Information Processing Systems 14*, 2001.
- [37] L. D. Raedt and H. Blockeel. Using logical decision trees for clustering. In *Proceedings of the 7th International Workshop on Inductive Logic Programming*, 1997.
- [38] R.O.Duda, P.E.Hart, and D.G.Stork. *Pattern Classification*. John Wiley & Sons, New York, 2000.
- [39] N. Rosenberg, J. Pritchard, J. Weber, and H. Cann. Genetic structure of human population. *Science*, 298, 2002.
- [40] J. S.D.Pietra, V.D.Pietera. Duality and auxiliary functions for bregman distances. Technical Report CMU-CS-01-109, Carnegie Mellon University, 2001.
- [41] S.Geman and D.Geman. Stochastic relaxation, gibbs distribution, and the bayesian restoration of images. *Pattern Analysis and Machine Intelligence*, 6:721–742, 1984.
- [42] J. Shi and J. Malik. Normalized cuts and image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(8):888–905, 2000.
- [43] T. Snijders. Markov chain monte carlo estimation of exponential random graph models. *Journal of Social Structure*, 2002.
- [44] A. Strehl and J. Ghosh. Cluster ensembles – a knowledge reuse framework for combining partitionings. In *AAAI 2002*, pages 93–98, 2002.
- [45] B. Taskar, E. Segal, and D. Koller. Probabilistic classification and clustering in relational data. In *Proceeding of IJCAI-01*, 2001.
- [46] K. Wagstaff, C. Cardie, S. Rogers, and S. Schroedl. Constrained k-means clustering with background knowledge. In *ICML-2001*, pages 577–584, 2001.
- [47] J. Wang, H. Zeng, Z. Chen, H. Lu, L. Tao, and W.-Y. Ma. Recom: reinforcement clustering of multi-type interrelated data objects. In *SIGIR '03*, pages 274–281, 2003.
- [48] E. King, A. Ng, M. Jorda, and S. Russel. Distance metric learning with applications to clustering with side information. In *NIPS'03*, volume 16, 2003.
- [49] X. Yin, J. Han, and P. Yu. Cross-relational clustering with user's guidance. In *KDD-2005*, 2005.
- [50] X. Yin, J. Han, and P. Yu. Linkclus: Efficient clustering via heterogeneous semantic links. In *VLDB-2006*, 2006.
- [51] H.-J. Zeng, Z. Chen, and W.-Y. Ma. A unified framework for clustering heterogeneous web objects. In *WISE '02*, pages 161–172, 2002.