

Multi-typed Objects Multi-view Multi-instance Multi-label Learning

Yuanlin Yang^{1,2}, Guoxian Yu^{1,2,3,4,*}, Jun Wang³, Carlotta Domeniconi⁵, Xiangliang Zhang⁴

¹College of Computer and Information Sciences, Southwest University, Chongqing, China

²School of Software, Shandong University, Jinan, China

³Joint SDU-NTU Centre for Artificial Intelligence Research, Shandong University, Jinan, China

⁴CEMSE, King Abdullah University of Science and Technology, Thuwal, SA

⁵Department of Computer Science, George Mason University, VA, USA

Email: ylyang@swu.edu.cn; {gxyu, kingjun}@sdu.edu.cn, carlotta@cs.gmu.edu; xiangliang.zhang@kaust.edu.sa

Abstract—Multi-typed objects Multi-view Multi-instance Multi-label Learning (M4L) deals with interconnected multi-typed objects (or bags) that are made of diverse instances, represented with heterogeneous feature views and annotated with a set of non-exclusive but semantically related labels. M4L is more general and powerful than the typical Multi-view Multi-instance Multi-label Learning (M3L), which only accommodates single-typed bags and lacks the power to jointly model the naturally interconnected *multi-typed* objects in the physical world. To combat with this novel and challenging learning task, we develop a joint matrix factorization based solution (M4L-JMF). Particularly, M4L-JMF firstly encodes the diverse attributes and multiple inter(intra)-associations among multi-typed bags into respective data matrices, and then jointly factorizes these matrices into low-rank ones to explore the composite latent representation of each bag and its instances (if any). In addition, it incorporates a dispatch and aggregation term to distribute the labels of bags to individual instances and reversely aggregate the labels of instances to their affiliated bags in a coherent manner. Experimental results on benchmark datasets show that M4L-JMF achieves significantly better results than simple adaptations of existing M3L solutions on this novel problem.

Index Terms—Multi-typed Objects, Multi-instance Learning, Multi-view Learning, Multi-label Learning, Joint Matrix Factorization

I. INTRODUCTION

With the prosperity of Internet of Things, objects are often represented by multiple heterogeneous feature views. For example, an image is numerically encoded by its texture, shape and color features. This image can also be simultaneously tagged with several related semantic labels (i.e., sun, sea, water, bird). To learn from such multi-modal multi-label data, various multi-view multi-label learning approaches have been introduced [1]. However, a real-world object may contain variable number of inconsistent instances (sub-objects). For example, a web page includes multiple images and content paragraphs, each of which can be viewed as an instance of the image/text view. To model such complex objects, Multi-view Multi-instance Multi-label Learning (M3L) has

been invented, it aims to leverage the relationships between instances, their hosting objects (bags), semantic labels and between heterogeneous feature views to predict the labels of objects and those of individual instances [2]–[6].

Existing M3L algorithms focus on *single-typed* objects. Given the diverse interconnections between objects of multiple types, the labels of a complex object are not only determined by its own attributes, but also by its connections with objects of other types. M3L lacks the capability to simultaneously model *multi-typed* objects. One typical solution is to project the objects of other types toward the target-type of objects to form the composite features, and then learn on the composite features (networks) [7], [8]. Unfortunately, such projection may override the intrinsic structure information among multi-typed objects [7], [9]. Matrix factorization based solutions have been introduced to model interconnected multi-typed objects, and these solutions can respect the intrinsic structure of these objects and integrate multiple feature views of objects [10]–[12]. However, these solutions ignore the general case that one object is composed with multiple instances, which convey important context information for labeling the complex object [5], [13].

For example, as instantiated in Figure 1, a social network includes user objects (encoded with social connections, personal profiles), image objects, text objects, and various instances (i.e., paragraphs and image parts) affiliated with these multi-typed objects. To effectively learn from multi-view multi-typed complex objects with interconnections, we term a new learning paradigm *Multi-typed objects Multi-view Multi-instance Multi-label Learning (M4L)* and introduce a joint matrix factorization based solution to solve this challenging task. M4L takes the objects (bags), instances and labels as nodes, and constructs a heterogeneous network with diverse edge types to encode the inter- and intra-relations between multi-typed objects, associations between instances and their hosting bags, relations between bags and labels, and correlations between labels. Next, it jointly factorizes the block association data matrices of the heterogeneous network and the attribute data matrix of nodes of respective types into low-rank matrices to explore the latent representations of bags, instances and

*Corresponding author: gxyu@sdu.edu.cn (Guoxian Yu). This work is supported by NSFC (No. 61872300, 62031003 and 62072380).

semantic labels, and then exploits the latent representations to predict the relations between bags (instances) and labels. In addition, to account for the bag-instance associations, we further introduce a dispatch and aggregation term to push the bag-level labels to individual instances and aggregate the instance-level labels to their hosting bags in a coherent way.

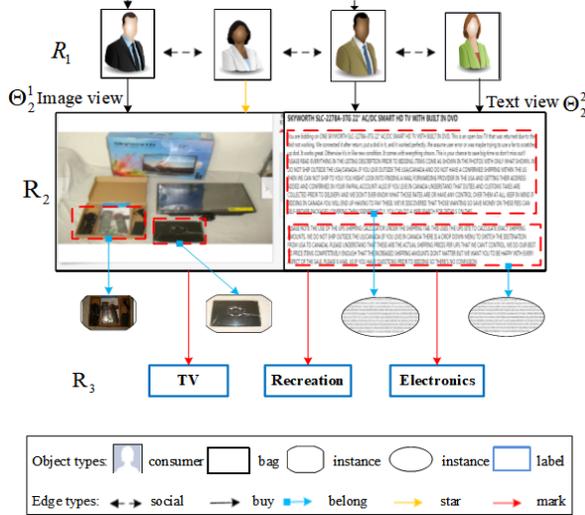


Fig. 1. An illustration of the M4 (Multi-typed objects Multi-view Multi-instance Multi-label) data with two types of objects (customers R_1 and goods R_2). The goods are encoded by the image view Θ_2^1 and text view Θ_2^2 , and the goods (bags) in each view is further made of diverse instances (i.e., remote control panel, set-top box and prices), multi-type objects and their affiliated instances can be simultaneously tagged with several semantic labels R_3 (i.e., ‘recreation’, ‘TV’, ‘electronics’). M4L aims to fuse multi-type objects to annotate the objects (instances) with the semantic labels.

The main contributions of this work are:

- We study a *new* learning paradigm Multi-typed objects Multi-view Multi-instance Multi-label Learning (M4L), which is a universal framework learning on naturally interconnected multi-typed complex objects. This learning scenario is more general than canonically studied M3L, which can only work with single-typed objects.
- We introduce a joint matrix factorization based solution M4L-JMF, which can leverage multiple inter(intra)-relations between bags of different types, associations between bags and instances, and the label correlations to annotate the complex bag and its instances.
- Experimental results on benchmark datasets show that M4L-JMF outperforms competitive and related methods.

II. RELATED WORKS

Our work is closely related with M3L and its degenerated versions (Multi-view Multi-label Learning [14], Multi-view Multi-instance Learning [13], and Multi-instance Multi-label Learning [15]), and data fusion by matrix factorization [9]. A comprehensive overview of the progress in these fast-progress areas is out of scope of this paper. Compared with these degenerated versions, M3L is less explored, due to the multiplicity and difficulty of learning from M3 data. To be

self-inclusive, we give a brief review of M3L solutions. To the authors knowledge, M3LDA [2] is the first M3L algorithm, it learns a visual-label part from the visual view and a text-label part from the text view, and forces these two parts having consistent labels. M3DN [3] separately applies a deep network for each view, and requires the bag-level predictions from different views being consistent within the same bag. M3DNS [6] extends M3DN by additionally making use of unlabeled instances and label correlations. M3Lcmf [5] utilizes a heterogeneous network to capture different types of relations between bags, instances and labels, and then collaboratively factorizes the relational data matrices of the network into low-rank ones to explore the latent relationships between bags, instances, and labels. WSM3L [16] studies M3L in a more general setting with unpaired view data and missing labels by multi-modal dictionary learning. However, these M3L methods *can only* consider single-typed bags, while in practice, these bags are also connected with objects of other types (as shown in Figure 1), which indirectly reflect the property (i.e., labels) of the target bags.

Given the huge demand of integrating multi-modal data, data fusion techniques have been extensively studied and applied in various domains [7], [12]. Compared with other data fusion solutions, such as feature view concatenation [17] and classifier ensemble [18], matrix factorization based solutions can respect the intrinsic structure among multi-typed objects without projecting toward the target objects, and train a single model to simultaneously fuse multiple information sources with less information loss [9]. To name a few, Wang *et al.* [19] applied symmetric nonnegative matrix tri-factorization (SN-MTF) to simultaneously cluster multi-typed objects. However, SNMTF has an overwhelming computation load, because it performs matrix factorization on a big matrix, whose block matrices encoded inter(intra)-relations between multi-typed objects. Data fusion by matrix factorization (DFMF) [9] collaboratively factorized these block matrices with much smaller sizes into low-rank ones and then reconstructed the target relational matrix to predict the relations between multi-typed objects. [10], [11] further considered the different relevance of multiple inter(intra)-relational block matrices toward the target prediction task and selectively fused these block data sources. All these prior studies of data fusion based on matrix factorization simply assume that each object is made of a single instance. As such, they *cannot* model the complex objects composed with diverse instances, as studied in M3L methods. M4L is different from the heterogeneous information network based data fusion [20]. The latter focuses on the heterogeneity of objects and does not consider the composition (sub-objects) of complex objects, as the consumer-product relationship shown in Figure 1. It simply takes consumers and commodities as nodes. Therefore, M4L considers a more sophisticated and practical learning scenario.

III. THE PROPOSED METHOD

A. Problem Statement

Suppose there are m types of directly or indirectly related objects (including semantic labels and sub-objects, a.k.a. instances), which are encoded by a set of inter-relational data matrices $\mathbf{R}_{ij} \in \mathbb{R}^{n_i \times n_j}$, $i, j \in \{1, 2, \dots, m\}$. One matrix \mathbf{R}_{ij} encodes the relations between n_i objects of the i -th type and n_j objects of the j -th type, and thus can be asymmetric. There are also a set of intra-association data matrices $\Theta_i^{(t)} \in \mathbb{R}^{n_i \times n_i}$, $t \in \{1, 2, \dots, t_i\}$, where t_i is the number of intra-relational data views for the i -th type of objects. Among these multi-typed objects (bags), some types of objects are further made of sub-objects (instances). Without loss of generality, we assume the i -th type of objects are instances of the b -th type, and $\mathbf{R}_{bi}(j, k) = 1$ if the k -th object of the i -th type is a member instance of the j -th object of the b -th type (e.g., the remote control panel of the 4-th type is an instance of the image object of the 2nd type.) Suppose each entity of the m -th type corresponds to a semantic label, and the label correlations are encoded by the intra-relational data matrix $\Theta_m \in \mathbb{R}^{n_m \times n_m}$. When our target object type is b , the aim of M4L is to predict the inter-relational matrix \mathbf{R}_{bm} for n_b bags and n_m labels, and/or the inter-relational matrix \mathbf{R}_{im} for n_i instances and n_m labels. The prediction is made by learning a mapping function $f(\mathcal{R}, \Theta) \in \{0, 1\}^{n_m}$ to relate the objects to n_m distinct labels. Here, \mathcal{R} collectively stores all the inter-relational data matrices \mathbf{R}_{ij} , and Θ collectively stores all the intra-relational data matrices Θ_i .

B. Joint Matrix Factorization

To complete the relational data matrix \mathbf{R}_{bm} (or \mathbf{R}_{im}) for bag (or instance)-label association prediction, we can take the target bags as anchors and then project objects of other types toward these anchors to form a composite bag-bag (instance-instance) intra-relational data matrix, and then use the known labels of bags to predict the labels of other bags (or instances) of the same type. In fact, this projection idea has been extensively used to integrate interconnected multi-type objects, and worked with multiple kernel (view) learning, classifier ensemble based data fusion solutions [1], [18], [21]. However, such projection may override the intrinsic structures among objects and cause information loss, and further compromise the performance [12].

Matrix factorization based data fusion techniques have been recently studied. They can integrate interconnected multi-typed objects of different types without conducting projection, while respect the intrinsic structures among objects [9]. This basic framework can be formulated as follows:

$$\min_{\mathbf{G} \geq 0} \Omega(\mathbf{G}, \mathbf{S}) = \sum_{\mathbf{R}_{ij} \in \mathcal{R}} \|\mathbf{R}_{ij} - \mathbf{G}_i \mathbf{S}_{ij} \mathbf{G}_j^T\|_F^2 + \sum_{t=1}^{\tau} \text{tr}(\mathbf{G}^T \Theta^{(t)} \mathbf{G}) \quad (1)$$

The minimization of the above objective function aims to reconstruct the incomplete \mathbf{R}_{bm} (or \mathbf{R}_{im}) to complete the associations between bags (instances) and labels, and thus

achieve the prediction by integrating interconnected objects of diverse types, without projecting these objects onto the target bags (instances). Here, $\mathbf{G} = \text{diag}(\mathbf{G}_1, \mathbf{G}_2, \dots, \mathbf{G}_m)$, $\mathbf{G}_i \in \mathbb{R}^{n_i \times k_i}$ is the low-rank representation of objects of the i -th type. \mathbf{S} is made of $\mathbf{S}_{ij} \in \mathbb{R}^{k_i \times k_j}$ ($k_i \ll n_i, k_j \ll n_j$), which can be viewed as a compressed data matrix that encodes latent inter-relations between objects of the i -th type and those of the j -th type. Intra-association data matrices $\Theta^{(t)} = \text{diag}(\Theta_1^{(t)}, \Theta_2^{(t)}, \dots, \Theta_m^{(t)})$ ($t \in \{1, 2, \dots, \max_i t_i\}$), where the i -th block matrix along the main diagonal of $\Theta^{(t)}$ is zero if $t > t_i$, and $\tau = \max_i t_i$, $\|\cdot\|_F^2$ is the Frobenius norm, $\text{tr}(\cdot)$ is the matrix trace operator. Entries in intra-association data matrices are positive for dissimilar objects, and negative for similar ones. The positive entries can be viewed as *cannot-link* constraints [22], which force pairs of dissimilar objects being far away from each other in the low-rank representation space. These intra(inter)-association data matrices jointly guide the learning of mutually consistent low-rank matrix \mathbf{G}_i , since \mathbf{G}_i is not only learnt by data matrices (i.e., \mathbf{R}_{ij}) directly related with the i -th type of objects, but also by data matrices (i.e., \mathbf{R}_{jk} , $j \neq i, k \neq i$) indirectly related with this object type, and by the intra-relational data matrices $\Theta_i^{(t)}$.

A real-world object (bag) may be further made of several different sub-objects (instances), and the labels of this bag are determined by the labels of its instances [15]. In many practical domains (i.e., medical image analysis and biology), the precise labels of instances are more important and interesting than those of bags, which carry more specific knowledge about the regions (i.e., local patches and functional sites) of the bag (i.e., image and molecule) [23]. Unfortunately, the labels of instances are typically unknown and while the labels of bags can be more easily collected. However, (1) overlooks the important *bag-instance* associations, and thus it can only model the degenerated fusion of multi-typed multi-view multi-label objects. Another disadvantage is that (1) does not differentiate the relevance of different relational data sources.

C. Unified Objective Function

To leverage the bag-instance association between two types of objects, we introduce a dispatch and aggregation term to push the labels of bags to their affiliated instances, and reversely aggregate the labels of instances to their hosting bags at the same time. For this purpose, we extend (1) as follows:

$$\min_{\mathbf{G} \geq 0} \Omega(\mathbf{G}, \mathbf{S}) = \sum_{i,j=1}^m \|\mathbf{R}_{ij} - \mathbf{G}_i \mathbf{S}_{ij} \mathbf{G}_j^T\|_F^2 + \sum_{p=1}^m \sum_{t=1}^{\tau} \text{tr}(\mathbf{G}_p^T \Theta_p^{(t)} \mathbf{G}_p) + \|\mathbf{R}_{bm} - \mathbf{R}_{bi} \mathbf{G}_i \mathbf{S}_{im} \mathbf{G}_m^T\|_F^2 \quad (2)$$

$\|\mathbf{R}_{bm} - \mathbf{R}_{bi} \mathbf{G}_i \mathbf{S}_{im} \mathbf{G}_m^T\|_F^2$ is added to achieve the aggregation of the labels of instances (predicted by $\mathbf{G}_i \mathbf{S}_{im} \mathbf{G}_m^T$) to their hosting bags via the bag-instance association matrix \mathbf{R}_{bi} . Given the labels of bags are typically available, this term can also dispatch the labels of bags (stored in \mathbf{R}_{bm}) to individual instances. In this way, (2) not only can integrate multi-typed objects, but also account for the complex objects that are made

of multiple sub-objects (instances) to predict the labels of bags and those of instances in a coherent way.

Multiple inter(intra)-relational data matrices contain complementary information of objects of different types, but they may also include some noisy or irrelevant data matrices. Although the low-rank matrix factorization can reduce the inner noises of individual data matrices to some extent [24], [25], it is still necessary to selectively fuse these relational data matrices with different relevance toward the target task. To concrete this, we advocate to set adaptive weights to intra(inter)-relational data matrices and thus to explicitly remove noisy data matrices as follows:

$$\begin{aligned} \min_{\mathbf{G} \geq 0} \Omega(\mathbf{G}, \mathbf{S}, \mathbf{W}^r, \mathbf{W}^h) &= \sum_{i,j=1}^m \mathbf{W}_{ij}^r \|\mathbf{R}_{ij} - \mathbf{G}_i \mathbf{S}_{ij} \mathbf{G}_j^T\|_F^2 \\ &+ \sum_{p=1}^m \sum_{t=1}^{\tau} \mathbf{W}_{pt}^h \text{tr}(\mathbf{G}_p^T \Theta_p^{(t)} \mathbf{G}_p) + \|\mathbf{R}_{bm} - \mathbf{R}_{bi} \mathbf{G}_i \mathbf{S}_{im} \mathbf{G}_m^T\|_F^2 \quad (3) \\ &+ \alpha \|\text{vec}(\mathbf{W}^r)\|_F^2 + \beta \|\text{vec}(\mathbf{W}^h)\|_F^2 \\ \text{s.t. } \mathbf{W}^r \geq 0, \mathbf{W}^h \geq 0, \sum \text{vec}(\mathbf{W}_i^r) &= 1, \sum \text{vec}(\mathbf{W}_i^h) = 1 \end{aligned}$$

where $\mathbf{W}^r \in \mathbb{R}^{m \times m}$ and $\mathbf{W}^h \in \mathbb{R}^{m \times \tau}$ are the weight matrices. \mathbf{W}^r stores the weights assigned to $|\mathcal{R}|$ inter-relational matrices and \mathbf{W}_{pt}^h encodes the weight of the t -th intra-relational matrix of the p -th object type. $\text{vec}(\mathbf{W}_i^r)$ is the vectorisation operator that stacks the i -th row of \mathbf{W}^r . $\mathbf{W}_{ij}^r = 0$ if $\mathbf{R}_{ij} \notin \mathcal{R}$. For $\Theta_p^{(t)}$, if $t \geq \max_i t_i$, $\mathbf{W}_{pt}^h = 0$. α and β are the regularization weights for these two weight matrices. They work alike the ridge regression to avoid the trivial solution that selects only one inter-relational data matrix and only one intra-relational data matrix. (3) not only can explore the contribution of different intra-relational data matrices, but also selectively fuse inter-relational matrices by assigning weights to them.

Suppose T is the maximum number of iterations, the time complexity of our model is $O(T(|\mathcal{R}| + \tau m + m\tau n_m))$. m is number of object types, n_m represents the maximum number of objects of type m and τ is the maximum number of views. The objective function of our M4L is non-convex in \mathbf{G} , \mathbf{S} , \mathbf{W}^r and \mathbf{W}^h altogether. We can use the idea of auxiliary functions frequently used in the convergence proof of approximate matrix factorization algorithms to alternatively optimize \mathbf{G} and \mathbf{S} in (3) [26], [27].

IV. EXPERIMENT

A. Experimental Setup

We used two publicly available datasets (Isoform and LncRNA) for the experiments, The two datasets have multi-typed objects and are specifically introduced below.

The Isoform dataset is collected from functional biology domain for predicting the functions of isoforms (instances), which are alternatively spliced from genes (bags) [23]. This dataset is a natural testbed of M4L. Unfortunately, limited by the wet-lab techniques, it just has the bag-level labels (a.k.a. functional annotations of genes), so we use the bag-level predictions aggregated from instance-level for a surrogate evaluation. This evaluation protocol is canonically used in isoform function prediction [23], [28]. We randomly selected

795 genes with 6,457 instances (isoforms), 495 miRNAs and 704 Gene Ontology terms (labels) as the dataset for experiments. The genes are represented with 2 feature views, and the isoforms are also represented with 2 feature views (Adrenal Gland and Esophagus Muscularis Mucosa). The more detailed information of these views can be found in [23]. The LncRNA dataset is also collected from biology domain, it contains six types of objects (LncRNAs(240), miRNAs(495), Genes(15527), Drugs(8283), and Diseases(412)), it was widely used in predicting the association between LncRNAs and diseases [10]. The genes are represented with 6 feature views, and the drugs with 10 feature views. The detailed information of multiple feature views of this dataset can be found in [10].

We perform experiments on the above benchmark datasets to quantitatively study the performance of the proposed M4L-JMF, and compare it against six representative and related approaches (M3Lcmf [5], M2IL [13] and ICM2L [29]), and data fusion solutions (SNMTF [19], DFMF [9] and MFLDA [10])). The first three compared methods are M3L solutions or the degenerated versions, while the other three methods addressed the fusion of interconnected multi-typed objects via matrix factorization, without consideration of the bag-instance associations. The input parameters of these comparison methods are specified (or optimized) according to the recommendations of the authors in their code or papers. The sensitivity of these parameters will be studied later. For each compared method, we run 5-fold cross validation for 10 independent rounds, and report the average results.

To quantitatively evaluate the M4L model's performance, three widely used multi-label evaluation metrics (AUROC, AUPRC and AvgF1) are adopted. AUROC firstly computes the value of the area under the ROC curve for each label and then takes the average for all labels. AUPRC firstly calculates the area under the Precision-Recall (PR) curve for each label and then takes the average values for all labels. For these three evaluation metrics, their larger values indicate the better the performance. AvgF1 needs to convert the bag-label association probabilistic matrix \mathbf{R}_{bm} or instance-label matrix \mathbf{R}_{im} into binary ones. Following the typical settings [10], [18], we adopt the top K labels corresponding to the largest entries of each row of \mathbf{R}_{bm} (\mathbf{R}_{im}) as the relevant labels of bags (instances), where K is the next integer of the average number of labels per bags/instances of the dataset.

B. Results on M4 data

To study the performance our M4L-JMF, we apply it on the Isoform dataset to predict the associations between isoforms and Gene Ontology terms, say the functional labels of isoforms. For comparison with M3L algorithms that are not designed to handle multi-type objects, we project objects of other types toward the genes at first and then form an M3 dataset composed with genes, isoforms and Gene Ontology labels, and then apply M3Lcmf, ICM2L, and M2IL on this projected M3 dataset. For matrix factorization based data fusion algorithms, we skip the bag-instance associations but fuse these types of objects to predict the association between

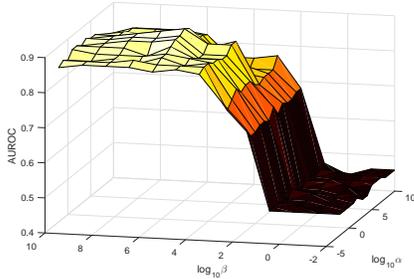


Fig. 2. AUROC of M4L-JMF under different settings of α and β on the LncRNA dataset.

labels and genes. Table I summarizes the results of M4L and compared methods.

From Table I, we have the following important observations: (i) Projection of multi-type objects along the target type causes information loss. To handle multi-typed objects, M3L methods (M3Lcmf, ICM2L and M2IL) were applied on the projection of multi-typed objects. The results in Table I show that M3L methods in general perform worse than matrix factorization based methods and M4L-JMF, evaluated on AUROC and AUPRC. This comparison corroborates the issue of fusing data by projecting data along the target type, which has been typically used in previous data fusion solutions [8], [21], [30]. Our proposed M4L-JMF can better deal with multi-type objects than the simple adaption of M3L solutions.

(ii) The bag-instance associations carry important information and can boost the performance of M4L. This observation is made by the comparison between matrix factorization(MF)-based solutions and M4L-JMF. Due to the ignorance of bag-instance associations, MF solutions perform worse than M4L-JMF on all metrics. In addition, M3Lcmf considers all the inter(intra)-relation between bags, instances and labels. Therefore, it has the highest AvgF1. However, it has lower AUROC and AUPRC than MF solutions and M4L due to its limitation on handling multi-type objects, as we discussed in the observation (i). Overall, these comparisons confirm that it is important to model the bag-instance association.

(iii) Selective fusion of data can further boost the performance. M3Lcmf and DFMF simply add up all the inter(intra)-relational data matrices, without considering the different relevance of these data sources. In contrast, our proposed M4L-JMF differentiates the relevance of these data sources. For this reason, M4L-JMF manifests much better results than them. Although MFLDA also considers the different relevance of inter-relational data sources, they ignore the intra-relational data sources. So it also gives lower results than our proposed M4L-JMF.

In summary, the results on the real M4 dataset confirm that M4L-JMF can more comprehensively model M4 data, without projecting the multi-type objects and skipping the bag-instance associations. For this advantage, it achieves better results than these compared methods.

TABLE I

RESULTS ON ISOFORM DATASET OF M4L, M3L-BASED METHODS AND MATRIX FACTORIZATION (MF)-BASED METHODS BY 5-FOLD CROSS VALIDATION. ●/○ INDICATES WHETHER M4L-JMF IS STATISTICALLY (ACCORDING TO PAIRWISE t -TEST AT 95% SIGNIFICANCE LEVEL) SUPERIOR/INFERIOR TO THE OTHER METHOD.

	Method	AvgF1	AUROC	AUPRC
M3L	M3Lcmf	0.152±0.004○	0.663±0.018●	0.154±0.013●
	ICM2L	0.074±0.001○	0.533±0.001●	0.041±0.022●
	M2IL	0.025±0.004●	0.544±0.009●	0.032±0.013●
MF	DFMF	0.051±0.001●	0.943±0.009●	0.637±0.054●
	SNMTF	0.021±0.001●	0.790±0.012●	0.015±0.002●
	MFLDA	0.029±0.002●	0.946±0.005●	0.546±0.011●
M4L	M4L-JMF	0.055±0.002	0.967±0.004	0.674±0.026

TABLE II

RESULTS OF M4L-JMF AND MATRIX FACTORIZATION BASED METHODS ON THE LNCRNA DATASET. ●/○ INDICATES WHETHER M4L-JMF IS STATISTICALLY (ACCORDING TO PAIRWISE t -TEST AT 95% SIGNIFICANCE LEVEL) SUPERIOR/INFERIOR TO THE OTHER METHOD.

Method	AvgF1	AUROC	AUPRC
DFMF	0.062±0.001●	0.872±0.007●	0.546±0.091●
SNMTF	0.023±0.001●	0.804±0.001●	0.016±0.002●
MFLDA	0.064±0.003●	0.874±0.005●	0.573±0.053●
M4L-JMF	0.067±0.002	0.895±0.004	0.616±0.026

C. Results on LncRNA dataset

We further study the performance of M4L-JMF on the LncRNA dataset (a natural test for multi-type objects fusion). For the experiments on the LncRNA dataset (without instance-label associations), we skip the bag-instance associations and compare M4L-JMF with matrix factorization based solutions only. The other pre-processes are the same as the experiments in previous subsection. The results on LncRNA dataset are given in Table II.

From Table II, we can find that, even without the important bag-instance associations, M4L-JMF still shows a good performance in fusing multi-type objects. That is because M4L-JMF can selectively fuse both the inter-relational and the intra-relational data sources, whereas these compared methods either ignore the different relevance of these data sources, or only differentiate the inter-relational ones.

Overall, these experimental results demonstrate the flexibility of M4L-JMF in diverse settings, and justify the contributions of weighting inter(intra)-relational data matrices.

D. Parameter Sensitivity Analysis

In this paper, three parameters (α , β and the low-rank size d of \mathbf{G}) in (3) should be specified for our proposed M4L-JMF. To investigate the sensitivity of first two parameters, we vary α and β in the range $\{10^{-2}, 10^{-1}, \dots, 10^{10}\}$, and report the average AUROC of M4L-JMF under different combinations of them in Fig. 2. M4L-JMF achieves the highest AUROC when $\alpha = 10^6$ and $\beta = 10^7$. The AUROC value increases when α or β rises, then it slightly decreases when $\alpha > 10^6$ or $\beta > 10^7$. A too small input value for α (or β) makes M4L-JMF only fuse one inter-relational data matrix and one intra-relational data matrix. On the other hand, a too large value for α (or β) leads M4L-JMF fusing all the relational data matrices without differentiating the relevance among them. This pattern shows

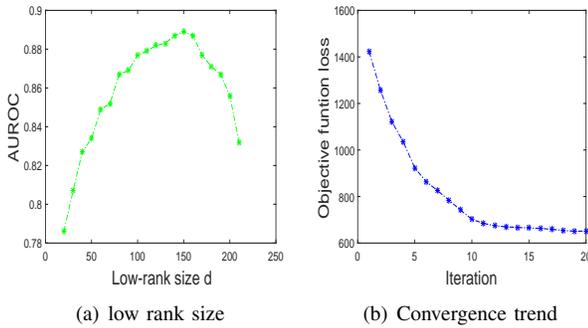


Fig. 3. AUROC vs. d (low rank size of \mathbf{G}_i) and convergence trend on the LncRNA dataset

that M4L-JMF can mine the complementary information of multiple relational data sources and account for different relevance of them. We also vary the low-rank size d ($k_i = d$ for all object types for simplicity) for the representation of objects (\mathbf{G}_i) from different ranges to study the optimal low-rank size using 5-fold cross validation, and show the AUROC under each input value of d with $\alpha = 10^6$ and $\beta = 10^7$ in Fig. 3(a). For the LncRNA dataset, we observe that a too small d can not sufficiently encode the latent feature information of multi-type objects and labels, and while a too large d may bring in some noises and thus leads to a low AUROC value. To investigate the convergence trend of M4L-JMF, we record the objective function value, i.e., the loss of (3) in each iteration on LncRNA datasets, and report the results in Fig. 3(b). We can see that the loss decreases as the iteration proceeds and comes to a convergence within 20 iterations. This trend proves that our alternative optimization procedure can quickly converge.

V. CONCLUSIONS

In this paper, we studied a novel learning paradigm (Multi-typed objects Multi-view Multi-instance Multi-label Learning) for naturally interconnected multi-type complex objects, and introduced a joint matrix factorization based approach M4L-JMF. Experimental results on real-world and benchmark datasets validated that M4L-JMF can more comprehensively fuse multi-type objects and mine complex relations between bags, instances and labels, and it achieves better results than other competitive and related methods.

REFERENCES

- [1] J. Zhao, X. Xie, X. Xu, and S. Sun, "Multi-view learning overview: Recent progress and new challenges," *Information Fusion*, vol. 38, pp. 43–54, 2017.
- [2] C.-T. Nguyen, D.-C. Zhan, and Z.-H. Zhou, "Multi-modal image annotation with multi-instance multi-label lda," in *IJCAI*, 2013, pp. 1558–1564.
- [3] Y. Yang, Y.-F. Wu, D.-C. Zhan, Z.-B. Liu, and Y. Jiang, "Complex object classification: A multi-modal multi-instance multi-label deep network with optimal transport," in *KDD*, 2018, pp. 2594–2603.
- [4] C.-T. Nguyen, X. Wang, J. Liu, and Z.-H. Zhou, "Labeling complicated objects: Multi-view multi-instance multi-label learning," in *AAAI*, 2014, pp. 2013–2019.
- [5] Y. Xing, G. Yu, C. Domeniconi, J. Wang, Z. Zhang, and M. Guo, "Multi-view multi-instance multi-label learning based on collaborative matrix factorization," in *AAAI*, 2019, pp. 5508–5515.

- [6] Y. Yang, Z.-Y. Fu, D.-C. Zhan, Z.-B. Liu, and Y. Jiang, "Semi-supervised multi-modal multi-instance multi-label deep network with optimal transport," *IEEE Transactions on Knowledge and Data Engineering*, vol. 99, no. 1, pp. 1–14, 2020.
- [7] V. Glgorijević and N. Pržulj, "Methods for biological data integration: perspectives and challenges," *Journal of the Royal Society Interface*, vol. 12, no. 112, p. 20150571, 2015.
- [8] G. Yu, G. Fu, C. Lu, Y. Ren, and J. Wang, "Brwlda: bi-random walks for predicting lncrna-disease associations," *Oncotarget*, vol. 8, no. 36, p. 60429, 2017.
- [9] M. Žitnik and B. Zupan, "Data fusion by matrix factorization," *TPAMI*, vol. 37, no. 1, pp. 41–53, 2015.
- [10] G. Fu, J. Wang, C. Domeniconi, and G. Yu, "Matrix factorization-based data fusion for the prediction of lncrna-disease associations," *Bioinformatics*, vol. 34, no. 9, pp. 1529–1537, 2018.
- [11] Y. Wang, G. Yu, J. Wang, G. Fu, M. Guo, and C. Domeniconi, "Weighted matrix factorization on multi-relational data for lncrna-disease association prediction," *Methods*, vol. 173, pp. 32–43, 2020.
- [12] M. Žitnik, F. Nguyen, B. Wang, J. Leskovec, A. Goldenberg, and M. Hoffmann, "Machine learning for integrating data in biology and medicine: Principles, practice, and opportunities," *Information Fusion*, vol. 50, pp. 71–91, 2019.
- [13] B. Li, C. Yuan, W. Xiong, W. Hu, H. Peng, X. Ding, and S. Maybank, "Multi-view multi-instance learning based on joint sparse representation and multi-view dictionary learning," *TPAMI*, vol. 39, no. 12, pp. 2554–2560, 2017.
- [14] Q. Tan, G. Yu, C. Domeniconi, J. Wang, and Z. Zhang, "Incomplete multi-view weak-label learning," in *AAAI*, 2018, pp. 4414–4421.
- [15] Z.-H. Zhou, M.-L. Zhang, S.-J. Huang, and Y.-F. Li, "Multi-instance multi-label learning," *Artificial Intelligence*, vol. 176, no. 1, pp. 2291–2320, 2012.
- [16] Y. Xing, G. Yu, J. Wang, C. Domeniconi, and X. Zhang, "Weakly-supervised multi-view multi-instance multi-label learning," in *IJCAI*, 2020, pp. 3124–3130.
- [17] Q. Zheng, J. Zhu, Z. Li, S. Pang, J. Wang, and Y. Li, "Feature concatenation multi-view subspace clustering," *Neurocomputing*, vol. 379, pp. 89–102, 2020.
- [18] G. Yu, C. Domeniconi, H. Rangwala, G. Zhang, and Z. Yu, "Transductive multi-label ensemble classification for protein function prediction," in *KDD*, 2012, pp. 1077–1085.
- [19] H. Wang, H. Huang, and C. Ding, "Simultaneous clustering of multi-type relational data via symmetric nonnegative matrix tri-factorization," in *CIKM*, 2011, pp. 279–284.
- [20] G. Yu, Y. Wang, J. Wang, C. Domeniconi, M. Guo, and X. Zhang, "Attributed heterogeneous network fusion via collaborative matrix tri-factorization," *Information Fusion*, vol. 63, pp. 153–165, 2020.
- [21] G. Yu, H. Rangwala, C. Domeniconi, G. Zhang, and Z. Zhang, "Predicting protein function using multiple kernels," in *IJCAI*, 2013, pp. 1869–1875.
- [22] M. Bilenko, S. Basu, and R. J. Mooney, "Integrating constraints and metric learning in semi-supervised clustering," in *ICML*, 2004, pp. 59–68.
- [23] G. Yu, K. Wang, C. Domeniconi, M. Guo, and J. Wang, "Isoform function prediction based on bi-random walks on a heterogeneous network," *Bioinformatics*, vol. 36, no. 1, pp. 303–310, 2020.
- [24] D. Meng and F. De La Torre, "Robust matrix factorization with unknown noise," in *ICCV*, 2013, pp. 1337–1344.
- [25] X. Chen, G. Yu, C. Domeniconi, J. Wang, Z. Li, and Z. Zhang, "Cost effective multi-label active learning via querying subexamples," in *ICDM*, 2018, pp. 905–910.
- [26] C. H. Ding, T. Li, and M. I. Jordan, "Convex and semi-nonnegative matrix factorizations," *TPAMI*, vol. 32, no. 1, pp. 45–55, 2008.
- [27] D. D. Lee and H. S. Seung, "Algorithms for non-negative matrix factorization," in *NeurIPS*, 2001, pp. 556–562.
- [28] T. Luo, W. Zhang, S. Qiu, Y. Yang, D. Yi, G. Wang, J. Ye, and J. Wang, "Functional annotation of human protein coding isoforms via non-convex multi-instance learning," in *KDD*, 2017, pp. 345–354.
- [29] Q. Tan, G. Yu, J. Wang, C. Domeniconi, and X. Zhang, "Individuality- and commonality-based multiview multilabel learning," *IEEE T. on CYB*, vol. 99, no. 1, pp. 1–13, 2020.
- [30] C. Lu, M. Yang, F. Luo, F.-X. Wu, M. Li, Y. Pan, Y. Li, and J. Wang, "Prediction of lncrna-disease associations based on inductive matrix completion," *Bioinformatics*, vol. 34, no. 19, pp. 3357–3364, 2018.