Multi-Task Metric Learning on Network Data

Chen Fang
Department of Computer Science
Dartmouth College
Hanover, NH, 03755, U.S.A.
chenfang@cs.dartmouth.edu

Daniel N. Rockmore
Department of Computer Science
Department of Mathematics
Dartmouth College
Hanover, NH, 03755
rockmore@cs.dartmouth.edu

ABSTRACT

Multi-task learning (MTL) has been shown to improve prediction performance in a number of different contexts by learning models jointly on multiple different, but related tasks. Network data, which are a priori data with a rich relational structure, provide an important context for applying MTL. In particular, the explicit relational structure implies that network data is not i.i.d. data. Network data also often comes with significant metadata (i.e., attributes) associated with each entity (node). Moreover, due to the diversity and variation in network data (e.g., multi-relational links or multi-category entities), various tasks can be performed and often a rich correlation exists between them. Learning algorithms should exploit all of these additional sources of information for better performance. In this work we take a metric-learning point of view for the MTL problem in the network context. Our approach builds on structure preserving metric learning (SPML) [3]. In particular SPML learns a Mahalanobis distance metric for node attributes using network structure as supervision, so that the learned distance function encodes the structure and can be used to predict link patterns from attributes. In the fundamental paper [3] SPML is described for single-task learning on single network. Herein, we propose a multi-task version of SPML, abbreviated as MT-SPML, which is able to learn across multiple related tasks on multiple networks via shared intermediate parametrization. MT-SPML learns a specific metric for each task and a common metric for all tasks. The task correlation is carried through the common metric and the individual metrics encode task specific information. When combined together, they are structure-preserving with respect to individual tasks. MT-SPML works on general networks, thus is suitable for a wide variety of problems. In experiments, we challenge MT-SPML with two common real-word applications: citation prediction for Wikipedia articles and social circle prediction in Google+. Our results show that MT-SPML achieves significant improvement over other competing methods.

Categories and Subject Descriptors

 $\mathrm{H.2.8}$ [DATABASE MANAGEMENT]: Database Applications— $Data\ mining$

General Terms

Algorithm, performance, theory

Keywords

Multi-task learning, metric learning, social network, link prediction

1. INTRODUCTION

Multi-task learning (MTL) [8, 4, 24, 2, 7] considers the problem of learning models jointly and simultaneously over multiple, different but related tasks. Compared to single-task learning (STL), which learns a model for each task independently using only task specific data, MTL leverages all available data and shares knowledge among tasks, thereby resulting in better model generalization and prediction performance. The underlying principle of MTL is that highly correlated tasks can benefit from each other via joint training, but additional care should be taken to respect the distinctiveness of each task, i.e., it is usually inappropriate to pool all available data and learn a single model for all tasks.

Despite the popularity and value of MTL, most MTL methods are developed for tasks on i.i.d. data. Standard examples include phoneme recognition [18] and image recognition [22]. Explicitly correlated data, often represented in the form of a network, provide a rich source of new applications contexts wherein the explicit relatedness of the data might be leveraged to improve performance on similarly correlated tasks. That is, although each task bears its own distinctiveness, relatedness cannot be ignored and should be exploited for good! The following two scenarios, provide two important examples where it is beneficial to exploit the correlation between tasks. These scenarios are in fact the settings for the experiments using real-world data that we present in Section 4.

Scenario 1: Article citation prediction

Articles tend to cite each other, especially those in the same subject area. The citation prediction problem has been studied extensively [10, 13, 21, 1, 11]. People either build a predictive model for a unified network [13] (i.e., a citation network that contains papers of all subject areas,) or build predictive models for each area independently [3]. Since article content and citation pattern varies across different ar-

eas, the former methodology ignores the difference between areas. However, some areas, while labeled as different are still related, in the sense of both their content and citation pattern. Thus the latter methodology fails to exploit the correlation among subject areas. For example, computer science and electrical engineering articles may be classified or tagged as different areas, but in many cases they may still have much in common, or at least have significant similarity or overlap. In this case, to build predictive models for citations, a learning algorithm that is capable of utilizing these overlaps and explicit commonalities has advantages over traditional methods.

Scenario 2: Social circle prediction

Members of online social networks tend to categorize their links to followers/followees. For example, many social networking platforms enable coarse-scale categorizations such as "family members," or "friends and colleagues." Finer gradations allow for categorizations such as colleagues at particular companies or classmates at specific schools. A person's social circle, studied in [15], is the ego network of a social network user (or "ego") in which all links belong to the same category. I.e., the induced subgraph of a given entity containing only links of a given type. Given a friend or stranger, the goal of social circle prediction is to assign him/her to appropriate social circles. Because some social circles are related to each other (e.g., family members and childhood friends may share some common informative features such as geological proximity), advantages may very well accrue if the relatedness of the entities was used for the various predictions, instead of building predictive assignment model for each social circle independently.

As these scenarios suggest, there should be advantages to developing methods that can leverage the correlations among tasks on network data. In what follows we show that MTL is a natural direction to pursue and that it does in fact provide some significant improvements.

Different from i.i.d. data, network data not only has attributes (metadata) associated with each entity (node), but rich structural information, mainly encoded in the links. Both attributes and structure should be exploited in learning. Structure preserving metric learning (SPML), originally developed for single-task learning [3] is such a method. It learns a Mahalanobis distance metric for node attributes by using network structure as supervision, so that the learned distance function encodes the structure and can be used to predict link patterns from attributes. Inspired by the use of SPML in the single task context, we propose its multi-task version, MT-SPML, which learns Mahalanobis distance metrics jointly over all tasks. More precisely, it learns a common metric for all tasks and one metric for each individual task. The common metric construction follows the methodology of shared intermediate parameterization [8, 16], which allows sharing knowledge between tasks. While a single task specific metric captures task specific information, when combined, they work together to preserve the connectivity structure of the corresponding network, thus are useful for link prediction from attributes. We further show that as in the case of SPML, MT-SPML can be optimized with efficient online methods similar to OASIS [5] and PEGASOS [19] via stochastic gradient descent. Finally, MT-SPML is designed for general networks, and in experiments we apply MT-SPML to two common, but different real-world prediction problems (citation prediction and social circle prediction) with promising results.

2. RELATED WORK

There is a large body of work on MTL for i.i.d. data. Yu et al. [24] applied hierarchical Bayesian modeling to nonparametric Gaussian processes, and the resulting method was used for text categorization. Evgeniou et al. [8] extended Support Vector Machines (SVMs) to MTL via parameter sharing, and the method was applied to learn predictive models for exam scores of student at different schools. Following the same intuition as [8], Parameswaran et al. [16] proposed the multi-task version of large margin nearest neighbor metric learning [23], which was tested on speech recognition. In [22], Wang et al. applied MTL to help face recognition and image retrieval. Very recently, Seltzer et al. [18] showed how multi-task deep neural network can further help phoneme recognition.

Researchers also have been studying the problem of learning across multiple graph data for various purposes. Zhou et al. [25] improved document recommendation by finding an embedding for multiple graphs via matrix factorization. In [20], Tang et al. attempted to do clustering jointly over different graphs. Prakash et al. [6] developed an algorithm to jointly do clustering and classification on networks. In the area of relational learning, tensor decomposition-based methods are usually applied [14] for problems on multirelational data.

Of greatest relevance for our work is [17], wherein Qi et al. carefully developed a mechanism to sample across networks to predict missing links in a target network. Our paper differs from it in (at least) the following ways. First and foremost, we aim at improving prediction performance of all networks, while [17] targets at a specific network and uses sources networks to help link prediction on it. Second, MT-SPML essentially learns a joint embedding of both attribute features and network topological structure, while [17] tries to linearly combine attribute features and local structure information, e.g. the number of shared neighbors between a pair of nodes. This is another important difference. In this paper, we are looking at improving the performance of predicting links only from attributes. Hence, given an incoming node at test time, we do not have any prior knowledge of its connectivity information, which is not an uncommon scenario in practice, e.g. nascent network or sparse network. In these scenarios, [17] is more likely to waste its modeling power on local structure features, which is often unavailable. The lack of initial structure information also makes our problem somewhat more difficult than the traditional link prediction problem, which has a snapshot of current network, which is usually not sparse, and predicts future links among already observed nodes. Nevertheless, the metric learned by our approach can help the traditional link prediction as well if attributes are available. Last but not least, we aim to naturally marry MTL and network/graph/relational data, to take advantage of MTL's ability of handling relatedness and heterogeneity. The proposed MT-SPML is a general method, which can handle different types of correlations and variations among tasks (e.g., the marginal distribution

of node attributes differs from task to task, and the semantics or types of links can also vary depending on specific task). Thus, our approach can be applied to general network data, like article citation, social networks and email networks.

3. OUR APPROACH

In this section, we will first cover the technical details of Structure Preserving Metric Learning (SPML). Then, both derivation and sketched proof of MT-SPML are provided.

3.1 Notations and preliminaries

Given a network, we model it as a graph $\mathbf{G} = (\mathbf{X}, \mathbf{A})$, where $\mathbf{X} \in \mathbb{R}^{d \times n}$ represents the node attributes and $\mathbf{A} \in \mathbb{R}^{n \times n}$ is the binary adjacency matrix, whose entry \mathbf{A}_{ij} indicates the linkage information between node i and node j. A Mahalanobis distance is parameterized by a positive semidefinite (PSD) matrix $\mathbf{M} \in \mathbb{R}^{d \times d}$, where $\mathbf{M} \succeq 0$. The corresponding distance function is define as $d_{\mathbf{M}}(x_i, x_j) = (x_i - x_j)^{\top} \mathbf{M}(x_i - x_j)$. This is equivalent to the existence of a linear transformation $\mathbf{L} \in \mathbb{R}^{d \times d}$ on the feature space such that $\mathbf{M} = \mathbf{L}^{\top} \mathbf{L}$. Given a metric \mathbf{M} , to predict the structure pattern of \mathbf{X} , we adopt the simple k-nearest neighbor algorithm, which is denoted as \mathcal{C} , meaning each node is connected with its top-k nearest neighbors under the defined metric. Mathematically, we denote it as $\mathcal{C}(\mathbf{X}, \mathbf{M}) = \mathbf{A}$, and we say \mathbf{M} is structure preserving or that it preserves \mathbf{A} . The Frobenius norm is demoted as $\|\cdot\|_F$.

We denote a set of networks as $\mathcal{G} = \{\mathbf{G}_1, \mathbf{G}_2, \dots, \mathbf{G}_Q\}$. Each individual network \mathbf{G}_q has its own \mathbf{X}_q and \mathbf{A}_q . We use q to index each network and (i,j) for its element. Thus, for notational simplicity \mathbf{A}_{qij} stands for element (i,j) in \mathbf{A}_q . Similarly, X_{qi} represents the feature of node i in \mathbf{X}_q . In algorithms, we will use a superscript to index over iteration, e.g., \mathbf{M}^k refers to the kth iteration of \mathbf{M} under the relevant iterative process.

3.2 SPML

The goal of SPML is to learn M from a network G = (X, A), such that M preserves A. This problem has a semidefinite max margin learning formulation, which is as follows:

$$\min_{\mathbf{M} \succeq 0} \frac{\lambda}{2} ||\mathbf{M}||_F^2 + \xi \tag{1}$$

subject to the following constraints:

$$\forall_{i,j}, d_{\mathbf{M}}(x_i, x_j) \ge (1 - \mathbf{A}_{ij}) \max_{l} (\mathbf{A}_{il} d_{\mathbf{M}}(x_i, x_l)) + 1 - \xi$$
 (2)

In (1) the Frobenius norm is a regularizer on \mathbf{M} and λ is the corresponding weight parameter. The key idea to achieve structure preserving is the set of linear constraints in (2). This essentially enforces that from node i, the distances to all disconnected nodes must be larger than the distance to the furthest connected node. Thus, when the constraints in (2) are all satisfied, $\mathcal{C}(\mathbf{X}, \mathbf{M})$ will exactly produce \mathbf{A} . Furthermore, to allow for violation (with penalty), the slack variable ξ is introduced to both (1) and (2).

With the many constraints in (2), optimizing (1) becomes unfeasible when the network has few hundreds of nodes. But

a rewriting of the problem as follows allows us to use stochastic subgradient descent (see Algorithm 1):

$$f(\mathbf{M}) = \frac{\lambda}{2} ||\mathbf{M}||_F^2 + \frac{1}{|S|} \sum_{(i,j,l) \in S} \max(\Delta_{\mathbf{M}}(x_i, x_j, x_l) + 1, 0)$$

where $\Delta_{\mathbf{M}}(x_i, x_j, x_l) = d_{\mathbf{M}}(x_i, x_l) - d_{\mathbf{M}}(x_i, x_j)$. S is defined as $S = \{(i, j, l) | \mathbf{A}_{i,l} = 1 \land \mathbf{A}_{i,j} = 0\}$, so the triplet (i, j, l) means that there is a link between node i and node l, but not between i and j. The subgradient of (3) can be calculated as

$$\nabla f = \lambda \mathbf{M} + \frac{1}{|S|} \sum_{(i,j,l) \in S_{+}^{+}} \left((x_i - x_l) (x_i - x_l)^{\top} - (x_i - x_j) (x_i - x_j)^{\top} \right)$$
(4)

where S_+ is the set of triplets whose hinge losses are positive. At every iteration t of Algorithm 1, B triplets are randomly sampled and the corresponding stochastic subgradient is calculated with regard to the current metric \mathbf{M}^t and these triplets. Since Algorithm 1 is a variant of PEGASOS [19], its complexity does not depend on training set size n, but on feature dimensionality d. For the number of iterations T needed to reach convergence, proved by [3, 19], it depends on parameter λ and optimization error, which measures how close the final objective value is to the global optimal objective value. Notice that after updating \mathbf{M} , there is an optional step, in which the current \mathbf{M} is projected to its PSD cone. Experiments in [3] show that delaying this operation to the end of the algorithm works well in practice and further reduces computational complexity.

Algorithm 1 Stochastic subgradient descent for SPML

```
Input: G = (X, A), \lambda, T, B
Output: M \succ 0
 1: \mathbf{M}^0 \leftarrow \mathbf{I}^{d \times d}
 2: for t = 1, 2, \dots, T do

3: \eta^t \leftarrow \frac{1}{\lambda \times t}

4: s \leftarrow \emptyset
  5:
             for b = 1, 2, ..., B do
  6:
                   Random sample (i, j, l) from S
  7:
                   s \leftarrow s \cup (i, j, l)
  8:
             \mathbf{M}^t \leftarrow \mathbf{M}^{t-1} - \eta^t \nabla f(\mathbf{M}^{t-1}, s)
  9:
             \mathbf{M}^t \leftarrow [\mathbf{M}^t]_+ (Optional: Project to PSD cone)
10:
11: end for
12: return \mathbf{M}^T
```

3.3 MT-SPML

In this section, we show how MT-SPML extends SPML to multi-task setting.

The input is a set of networks $\mathcal{G} = \{\mathbf{G}_1, \mathbf{G}_2, \dots, \mathbf{G}_Q\}$. Each network is $\mathbf{G}_q = (\mathbf{X}_q, \mathbf{A}_q)$. Our approach is a general method, thus it works for both problems with or without nodes overlapping. Note that, the nodes of all networks have common feature spaces. MT-SPML treats each network as a task. It follows the idea of *shared intermediate parametrization* as in [16] to enable knowledge transfer between tasks. The goal is to learn jointly over \mathcal{G} a task specific metric \mathbf{M}_q for each task and a common metric \mathbf{M}_0 , through

which knowledge transfers among tasks, so that the combined metric $(\mathbf{M}_0 + \mathbf{M}_q)$ respects the structure of \mathbf{G}_q , for all $\mathbf{G}_q \in \mathcal{G}$. Thus, the distance between two nodes in \mathbf{G}_q is defined as

$$d_q(x_i, x_j) = (x_i - x_j)^{\top} (\mathbf{M}_0 + \mathbf{M}_q)(x_i - x_j).$$
 (5)

MT-SPML is formulated as a regularized learning problem as follows:

$$\min_{\mathbf{M}_0, \mathbf{M}_1, \dots, \mathbf{M}_Q} \frac{\gamma_0}{2} ||\mathbf{M}_0 - \mathbf{I}||_F^2 + \sum_{q=1}^Q \frac{\gamma_q}{2} ||\mathbf{M}_q||_F^2 + \sum_{q=1}^Q \xi_q$$
 (6)

subject to the following constraints:

$$\forall q, i, j, :$$
 (7)

$$d_q(x_{qi}, x_{qj}) \ge (1 - \mathbf{A}_{qij}) \max_l(\mathbf{A}_{qil}d_q(x_{qi}, x_{ql})) + 1 - \xi_q.$$

In order to solve it, we rewrite it as following by incorporating the constraints:

$$f(\mathbf{M}_{0}, \mathbf{M}_{1}, \dots, \mathbf{M}_{Q}) = \frac{\gamma_{0}}{2} ||\mathbf{M}_{0} - \mathbf{I}||_{F}^{2} + \sum_{q=1}^{Q} \frac{\gamma_{q}}{2} ||\mathbf{M}_{q}||_{F}^{2}$$
$$+ \sum_{q=1}^{Q} \frac{1}{|S_{q}|} \sum_{(i,j,l) \in S_{q}} \max(\Delta_{q}(x_{qi}, x_{qj}, x_{ql}) + 1, 0)$$
(8)

Algorithm 2 Stochastic subgradient descent for MT-SPML

```
Input: \mathcal{G} = \{\mathbf{G}_1, \mathbf{G}_2, \dots, \mathbf{G}_Q\}, where \mathbf{G}_q = (\mathbf{X}_q, \mathbf{A}_q), \gamma_0, \gamma_1, \dots, \gamma_Q, T, B
Output: \mathbf{M}_0, \mathbf{M}_1, \dots, \mathbf{M}_Q \succeq 0
  1: for q=0,1,\ldots,Q do 2: \mathbf{M}_q^0 \leftarrow \mathbf{I}^{d \times d}
   3: end for
  \begin{array}{ll} 4: \ \mathbf{for} \ t=1,2,\ldots,T \ \mathbf{do} \\ 5: \ \mathbf{for} \ q=1,2,\ldots,Q \ \mathbf{do} \\ 6: \ \eta_q^t \leftarrow \frac{1}{\lambda \times t} \\ 7: \ s_q \leftarrow \emptyset \end{array}
   8:
                                  for b = 1, 2, ..., B do
                                            Random sample (i, j, l) from S_q
  9:
                                 s_q \leftarrow s_q \cup (i,j,l)
end for
\mathbf{M}_q^t \leftarrow \mathbf{M}_q^{t-1} - \eta_q^t \bigtriangledown_{\mathbf{M}_q} f(\mathbf{M}_q^{t-1},s_q)
\mathbf{M}_q^t \leftarrow [\mathbf{M}_q^t]_+ \text{ (Optional: Project to PSD cone)}
10:
11:
12:
13:
14:
                       \mathbf{M}_0^t \leftarrow \mathbf{M}_0^{t-1} - \eta_0^t \bigtriangledown_{\mathbf{M}_0} f(\mathbf{M}_0^{t-1}, \{s_1, s_2, \dots, s_Q\})
\mathbf{M}_0^t \leftarrow [\mathbf{M}_0^t]_+ \text{ (Optional: Project to PSD cone)}
15:
16:
17: end for
18: return \mathbf{M}_0^T, \mathbf{M}_1^T, \dots, \mathbf{M}_Q^T \succeq 0
```

where $\Delta_q(x_{qi}, x_{qj}, x_{ql}) = d_q(x_{qi}, x_{ql}) - d_q(x_{qi}, x_{qj})$. Although (8) has more unknown variables than (3), with respect to each unknown, it is in the same form as (3). Therefore, (8) can be solved with the same stochastic subgradient descent method using partial subgradient. The partial subgradients

of (8) with respect to \mathbf{M}_0 and \mathbf{M}_q are

$$\nabla_{\mathbf{M}_{0}} f = \gamma_{0} (\mathbf{M}_{0} - \mathbf{I})$$

$$+ \sum_{q=1}^{Q} \frac{1}{|S_{q}|} \sum_{(i,j,l) \in S_{q+}} \left((x_{qi} - x_{ql}) (x_{qi} - x_{ql})^{\top} - (x_{qi} - x_{qj}) (x_{qi} - x_{qj})^{\top} \right)$$
(9)

and

$$\nabla_{\mathbf{M}_{q}} f = \gamma_{q} \mathbf{M}_{q} + \frac{1}{|S_{q}|} \sum_{(i,j,l) \in S_{q+}} \left((x_{qi} - x_{ql}) (x_{qi} - x_{ql})^{\top} - (x_{qi} - x_{qj}) (x_{qi} - x_{qj})^{\top} \right). \quad (10)$$

With all necessary information, the optimization algorithm is outlined in Algorithm 2. Algorithm 2 runs for T iterations. Within each iteration, it does two things: (1) Randomly sample B triplets for each task, so as to calculate the partial subgradient and update the corresponding unknown; (2) Calculate the partial subgradient of the common metric \mathbf{M}_0 and update it using the $Q \times B$ triplets already sampled. Optionally, the metric matrices can be projected to their PSD cones. The analysis of Algorithm 1 still holds for Algorithm 2, thus it scales up with regard to feature dimensionality, optimization error and the parameters γ_q , but not the training set size.

4. EXPERIMENTS

In this section, we present experimental results on real-world data. We apply MT-SPML to the two scenarios mentioned in the Introduction: article citation prediction and social circle prediction. We show that in both cases, MT-SPML can significantly improve prediction performance. ¹

4.1 Citation prediction on Wikipedia

The data is obtained from [3]. The articles from the following three areas are crawled from Wikipedia: search engine, graph theory and philosophy. The citations between articles within each area are also crawled. The goal is, given an article, to predict the referencing of other articles within its area solely from its content. Therefore, at test time, no reference information of the test article is made available at all. The challenge of this problem is the fact that: (1) there is little node overlap between networks (i.e., an article belongs to only one area), thus the marginal distribution of node attributes $P(\mathbf{X}_q)$ may vary drastically from area to area, which poses difficulty for knowledge transfer; (2) the conditional probability of structure on attributes $P(\mathbf{A}|\mathbf{X})$ may also vary, because some words are informative and indicative for some areas, but not for others. The statistics of this dataset is detailed in Table 1. Bag-of-words is used to capture article content and the dimensionality is 6695. The high dimensionality reduces the need to learn full matrices. Therefore, we choose to learn diagonal metric matrices. This further reduces computational complexity. We split the dataset 80%/20% as training and testing respectively, then fix the testing part and vary the size of training set by sampling from the training part. We end up sampling

¹Code and data are available for download at http://www.cs.dartmouth.edu/~chenfang/temp/MT_SPML/demo_code.ta

20%, 40%, 60%, 80%, and finally 100% of the training part. Model selection is carried out on the sampled training set via 5-fold cross-validation. At test stage, for every test example our algorithm suggests other articles for citation according to their distances to test article. We build the receiver operator characteristic (ROC) curve for every test article, and use the average area under the curve (AUC) of the entire test set as performance measurement. We compare our results with two families of methods:

SVM methods We apply SVM-based methods in order to show the importance of modeling network structure. Since SVM-based methods do not model network structure, we need to construct features to encode this piece of information. The training examples are constructed by taking the pairwise difference of the attributes between two nodes. The training labels are binary, with 1 representing the existence of a link between a pair of nodes and 0 the absence. For a given edge, we measure its distance/length using the output of the classification score, which represents the confidence of having a link. Although the classification score is inversely proportional to the notion of distance, a simple conversion can make the two variables proportional to each other. Thus ROC and AUC can be calculated. The following specific methods are included:

- ST-SVM: This is the normal single-task SVM. An SVM is trained for each network independently. It does not explore the correlation between tasks. The model is trained and tested with LIBLINEAR [9].
- U-SVM: We train one SVM for all networks by pooling all data together. We use the capital letter "U" to represent the naive strategy of data pooling. This is essentially ignoring the fact that training examples are from different tasks and treating it as a simple single task learning problem. The model is also trained and tested using LIBLINEAR [9].
- MT-SVM: This is the multi-task SVM in [8]. Similar to our model, it jointly learns a common decision boundary for all and a specific boundary for each task. At test time, the common and task specific decision boundary together form the final classification model for each task. This method exploits task correlations via intermediate parameter sharing, but does not use network structure at the model level. We used the software from [12] for training and testing.

SPML methods: We apply three methods that are based on SPML. Compared to SVM-based methods, these methods explicitly model the network structure information. Therefore, the feature used here is simply the node attributes and links become linear constraints. Given an edge, its distance is just the Mahalanobis distance defined by learned metrics. The following methods are included:

• ST-SPML: This is the single-task SPML [3]. A metric is learned for each network independently. It models network structure but not task correlations.

Areas	# of nodes	# of edges	# of features
Search Engine	269	332	6695
Graph Theory	223	917	6695
Philosophy	303	921	6695

Table 1: Statistics of Wikipedia article data

- U-SPML: "U" means data pooling. Training examples from all tasks are pooled together and the learning procedure is simply ST-SPML. This is a naive way of sharing knowledge between tasks, but it does not respect the differences between and distinctiveness of tasks. Thus we expect inferior results, particularly for less related tasks.
- MT-SPML: This is our method. By comparing it to other methods, we can demonstrate the fact that MT-SPML not only models the structure of all networks nicely, but also exploits relatedness while respecting the distinctiveness of tasks.

Finally, we also compare to the direct use of the original feature vector, i.e., using Euclidean distance between feature vectors as the distance. While we are aware of the existence of other link prediction methods, such as Adamic-Adar [13]. As we have already mentioned, Adamic-Adar [13] only predicts potential links for nodes that are already present in the network and thus rely heavily on a snapshot of dense network structure. Thus, it is not suitable for our experiments, where no initial links for the test node are provided. Moreover, CNLP [17] targets at improving link prediction on one specific network by using other networks. Thus, its goal is fundamentally different from ours. Our methods encode observed structures in the learned metric and can be used for both unobserved test nodes and sparse graphs.

All the results are reported in Fig.1. The first thing we see is that SVM-based methods perform the worst when there are fewer training examples while the SPML family achieves good results in all settings, due to its ability to model structure information. We also find that among the SPML methods, MT-SPML consistently outperforms the others, which implies that MT-SPML is better at exploiting task correlations. Interestingly, we notice that the least amount of improvement from MT-SPML is found for philosophy articles. This observation seems to be aligned with the intuition that search engine-related and graph theory-related papers probably have more in common with each other than with philosophy papers.

We also show the convergence behavior of MT-SPML by plotting the value of $|S_{q+}|$, number of violated constraints among those randomly sampled ones, for every iteration for each task. The fewer the number of violated constraints, the better the new metric respects the network structure. In experiments we set B, the time of random sampling, to be 10. In order to make a clearer demonstration, in Fig.2 we set B to be 100. As shown by Fig.2, the numbers of violated constraints of all tasks drop quickly within the first 1000 iterations and stabilizes after 4000 iterations.

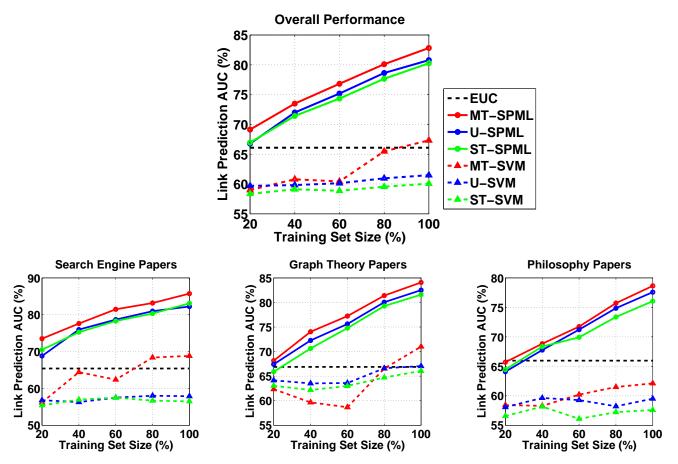


Figure 1: Link prediction performance on Wikipedia article data. Training set size is varied. Smaller figures in the lower half separate out the individual performance for each area. The bigger figure on the top is the average AUC performances over all three areas.

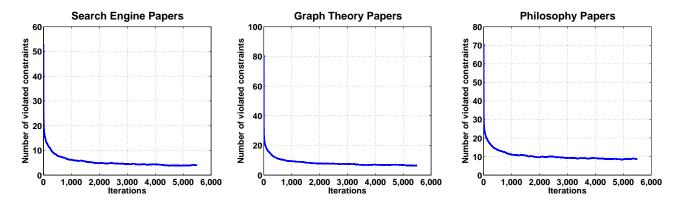


Figure 2: Number of violated constraints within first 5500 iterations. 100 constraints are sampled at every iteration.

4.2 Social circle prediction on Google+

Every member of an online social network (e.g., Google+) is the ego of his/her (sub-)network and tends - or may be forced - to categorize his/her relationships (e.g. family members, college friends or childhood friends). For each class of relationships, there is a sub-network associate with it, the social circle, which is directly formalized in the online structures of Google+ (see [15]). In this section, given a social network user (the ego) and his/her friends, we want to predict his/her social circles, namely the type of relationships between ego and ego's friends based on profile information. We are only interested in the ego network, meaning that we do not predict the links between friends. A similar topic is studied by McAuley et al [15], where the setup is very different from ours. They assume the observation of an entire ego network, including attributes and structure, but not any social circle labels, and the goal is to assign social circle labels to links in an unsupervised manner. Our problem uses a supervised learning setting, where we observe only parts of the network and the corresponding social circle labels. For the prediction of each social circle, we treat it as link prediction. However, as mentioned in our introduction (Section 1), the correlation between social circles should be exploited. Thus, we treat the prediction of each social circle as a task, and MT-SPML is applied to learn metrics jointly over the underlying ego networks of all social circles. Note that, as reported in [15], social circles largely overlap with each other, which implies strong correlations and MTL is thus likely to achieve a more significant performance gain. We obtain data from [15], which was from Google+ users and information is anonymous. We randomly pick one user and his/her social circles for our experiment. The entire ego network has 4402 nodes and 5 social circles. The profile of all nodes is also preserved. There are 6 classes of feature types, including gender, institution, job title, last name, place, and university. We build a bag-of-words feature for all feature types and concatenate them all, resulting in a feature vector of 2969 dimensions.

In this experiment, we adopt a different procedure. We start with using ST-SPML to learn a metric for each social circle independently. Then, to show the advantage of doing MTL jointly over multiple tasks, we run MT-SPML on various numbers of social circles. To avoid the exponential number of social circle combinations, we index them from 1 to 5. We begin by running on $\{1,2\}$ and add one more social circle at a time in order, resulting in the following four combinations: $\{1,2\}, \{1,2,3\}, \{1,2,3,4\}, \{1,2,3,4,5\}$ which we will continue to use in the following experiments. In this way, we can see the behavior of the algorithms as more relevant tasks joining. In Fig. 3, we compare ST-SPML to MT-SPML on the four combinations of social circles. Note that, because of the inferior performance of SVM based methods on Wikipedia article data, we entirely omit them in this experiment. Clearly, as shown in Fig. 3, all social circles benefit from MTL and the improvement is significant, except for Social Circle 2, whose performance gain is slight. We speculate that Social Circle 2 is not closely related to other circles (e.g., in terms of the number of overlapping nodes). We will discuss the case of Social Circle 2 later.

Now we compare MT-SPML to U-SPML, which simply pools all data together and estimates a model for all tasks. Both

Social Circles	1	2	3	4	5
1	0	1.1%	81.9%	89.6%	84.1%
2	1.1%	0	0.9%	1.1%	1.1%
3	81.9%	0.9%	0	73.5%	68.9%
4	89.6%	1.1%	73.5%	0	93.7%
5	84.1%	1.1%	68.9%	93.7%	0

Table 2: Statistics of overlapping nodes between social circles. Overlapping ratios are presented.

MT-SPML and U-SPML are applied to the four combination settings of different social circles. As shown by Fig. 4, MT-SPML consistently and significantly outperforms U-SPML at all locations.

Now we would like to further investigate Social Circle 2. We first show some statistics in Table 2, where we show the percentage of node overlapping between each pair of social circles. The overlap is defined as the intersection of nodes over the union. As we can see, some circles are largely overlapped (e.g., $\{1,3\}$ have 81.9% nodes in common), while Social Circle 2 barely overlaps with the others. Although overlapping is not the only quantitative measurement of correlations between social circles, a substantial set of common nodes suggests that there are some shared semantics between two relationships. The statistics of Table 2 supports our earlier speculation as to why Social Circle 2 does not benefit from joint learning as much as the others.

Furthermore, we would like to again show the advantage of MT-SPML by showing the results of a pair of tasks that are less correlated to each other. We choose Social Circles $\{1,2\}$, since they have only 1.1% nodes in common. In Fig. 5, MT-SPML is jointly learned on {1,2}, U-SPML is learned via data pooling, and ST-SPML is trained on 1 and 2 independently. The prediction performances of two tasks are reported in the two groups of bars respectively. As shown in Fig. 5, MT-SPML still gets 2%-5% performance improvement over ST-SPML (bars with circles on top). However, the strategy of simple data pooling used by U-SPML (bars with down pointing triangle) reduces the performance (produces results worse than ST-SPML). This observation suggests that on difficult cases where tasks are less relevant, MTL is still able to utilize useful correlations, while respecting the boundaries between tasks.

5. CONCLUSIONS

In this paper, we proposed MT-SPML, a large margin-based multi-task learning method for network data. It operates on networks with node attributes. It learns a task specific distance metric for every task and a common distance metric for all. By combining a task specific metric with the common distance, the final metric preserves the structure of the corresponding network, thus it can be used to predict link patterns on sparse nascent networks or for incoming nodes at test time. We applied MT-SPML to two common real-world problems, article citation prediction and social circle prediction. Better results were achieved (as compared to reasonable baselines) and detailed analysis was provided. The importance of our work lies in the fact that network data has large variation and diversity, thus many related

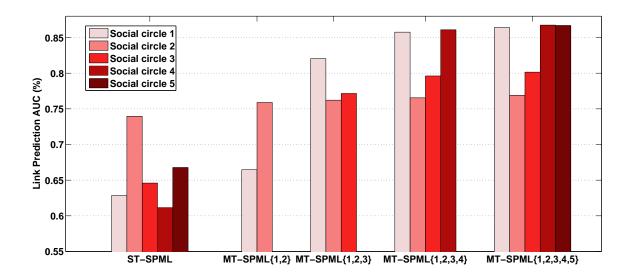


Figure 3: Link prediction performance on Google+ data. Social circles are color coded. The comparison is between ST-SPML and MT-SPML. The first group contains the prediction performance of ST-SPML on all social circles, while the other groups show the performance of MT-SPML that learned and tested on multiple combinations of social circles, for example, MT-SPML{1,2,3} means learning and testing on Social Circle 1, 2 and 3.

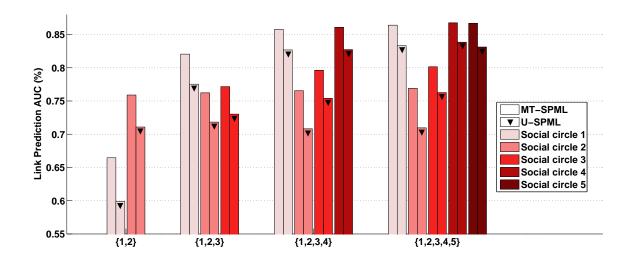


Figure 4: Link prediction performance on Google+ data. The comparison is between MT-SPML and U-SPML. Social circles are color coded. Different methods for the same task are compared side by side. U-SMPL is indicated by a downward pointing triangle. Each group is trained and tested on a set of social circles. For example, $\{1,2,3\}$ means learning and testing on Social Circle 1, 2 and 3.

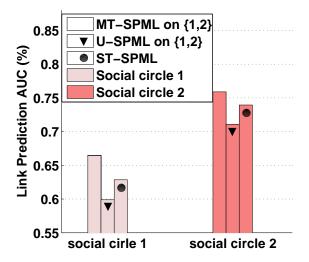


Figure 5: U-SPML negatively impacts performance when training on {1,2}, two less relevant tasks. MT-SPML is able to improve performance compared to ST-SPML by exploiting useful correlations, while U-SPML gets inferior results.

tasks can be performed, and we are able to better exploit the useful correlations. Moreover, since MT-SPML is a general method and can be optimized via stochastic gradient descent with good convergence behavior, it is suitable for general (and large) network data and can be widely applied to real-world problems. All the code and data used in experiments are available for download at the link given in the paper.

6. ACKNOWLEDGMENTS

We are grateful to our funding sources. The authors were supported by AFOSR Award FA9550-11-1-0166 and the Neukom Institute for Computational Science.

7. REFERENCES

- S. F. Adafre and M. de Rijke. Discovering missing links in wikipedia. In *Proceedings of the 3rd* International Workshop on Link Discovery, LinkKDD '05, 2005.
- [2] A. Agarwal, H. Daume III, and S. Gerber. Learning multiple tasks using manifold regularization. In NIPS. 2010.
- [3] B. H. Blake Shaw and T. Jebara. Learning a distance metric from a network. In NIPS, 2011.
- [4] R. Caruana. Multitask learning. In *Machine Learning*, pages 41–75, 1997.
- [5] G. Chechik, V. Sharma, U. Shalit, and S. Bengio. Large scale online learning of image similarity through ranking. *Journal of Machine Learning Research*, 11:1109–1135, 2010.
- [6] P. M. Comar, P.-N. Tan, and A. K. Jain. Multi task learning on multiple related networks. In CIKM, 2010.
- [7] H. Daumé, III. Bayesian multitask learning with latent hierarchies. In UAI, 2009.

- [8] T. Evgeniou and M. Pontil. Regularized multi-task learning. In KDD, 2004.
- [9] R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.-J. Lin. Liblinear: A library for large linear classification. *Journal of Machine Learning Research*, 9:1871–1874, 2008.
- [10] M. A. Hasan, V. Chaoji, S. Salem, and M. Zaki. Link prediction using supervised learning. In *In Proc. of* SDM 06 workshop on Link Analysis, Counterterrorism and Security, 2006.
- [11] H. Kashima and N. Abe. A parameterized probabilistic model of network evolution for supervised link prediction. In *ICDM*, 2006.
- [12] L. Liang and V. Cherkassky. Connection between sym+ and multi-task learning. In *IJCNN*, 2008.
- [13] D. Liben-Nowell and J. Kleinberg. The link prediction problem for social networks. In CIKM, 2003.
- [14] B. London, T. Rekatsinas, B. Huang, and L. Getoor. Multi-relational learning using weighted tensor decomposition with modular loss. *CoRR*, abs/1303.1733, 2013.
- [15] J. J. McAuley and J. Leskovec. Learning to discover social circles in ego networks. In NIPS, 2012.
- [16] S. Parameswaran and K. Weinberger. Large margin multi-task metric learning. In J. Lafferty, C. K. I. Williams, J. Shawe-Taylor, R. Zemel, and A. Culotta, editors, NIPS. 2010.
- [17] G.-J. Qi, C. C. Aggarwal, and T. S. Huang. Link prediction across networks by biased cross-network sampling. In *ICDE*, 2013.
- [18] M. Seltzer and J. Droppo. Multi-task learning in deep neural networks for improved phoneme recognition. In Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on, 2013.
- [19] S. Shalev-Shwartz, Y. Singer, N. Srebro, and A. Cotter. Pegasos: primal estimated sub-gradient solver for svm. *Math. Program.*, 127(1):3–30, 2011.
- [20] W. Tang, Z. Lu, and I. S. Dhillon. Clustering with multiple graphs. In *ICDM*, 2009.
- [21] B. Taskar, M. fai Wong, P. Abbeel, and D. Koller. Link prediction in relational data. In NIPS, 2003.
- [22] X. Wang, C. Zhang, and Z. Zhang. Boosted multi-task learning for face verification with applications to web image and video search. In CVPR, 2009.
- [23] K. Q. Weinberger and L. K. Saul. Distance metric learning for large margin nearest neighbor classification. *Journal of Machine Learning Research*, 10, 2009.
- [24] K. Yu, V. Tresp, and A. Schwaighofer. Learning gaussian processes from multiple tasks. In *ICML*, 2005.
- [25] D. Zhou, S. Zhu, K. Yu, X. Song, B. L. Tseng, H. Zha, and C. L. Giles. Learning multiple graphs for document recommendations. In WWW, 2008.