

# Autoencoder Based Sample Selection for Self-Taught Learning

Siwei Feng<sup>a</sup>, Han Yu<sup>b,c,d</sup>, Marco F. Duarte<sup>a,\*</sup>

<sup>a</sup>*Department of Electrical and Computer Engineering, University of Massachusetts,  
Amherst, MA 01003 USA*

<sup>b</sup>*School of Computer Science and Engineering, Nanyang Technological University,  
Singapore 639798*

<sup>c</sup>*Joint NTU-UBC Research Centre of Excellence in Active Living for the Elderly, Singapore  
639798*

<sup>d</sup>*Alibaba-NTU Singapore Joint Research Institute, Singapore 639798*

---

## Abstract

Self-taught learning is a technique that uses a large number of unlabeled data as source samples to improve the task performance on target samples. Compared with other transfer learning techniques, self-taught learning can be applied to a broader set of scenarios due to the loose restrictions on the source data. However, knowledge transferred from source samples that are not sufficiently related to the target domain may negatively influence the target learner, which is referred to as negative transfer. In this paper, we propose a metric for the relevance between a source sample and the target samples. To be more specific, both source and target samples are reconstructed through a single-layer autoencoder with a linear relationship between source samples and reconstructed target samples being simultaneously enforced. An  $\ell_{2,1}$ -norm sparsity constraint is imposed on the transformation matrix to identify source samples relevant to the target domain. Source domain samples that are deemed relevant are assigned pseudo-labels reflecting their relevance to target domain samples, and are combined with target samples in order to provide an expanded training set for classifier training. Local data structures are also preserved during source

---

\*Corresponding author

*Email addresses:* [siwei@umass.edu](mailto:siwei@umass.edu) (Siwei Feng), [han.yu@ntu.edu.sg](mailto:han.yu@ntu.edu.sg) (Han Yu), [mduarte@ecs.umass.edu](mailto:mduarte@ecs.umass.edu) (Marco F. Duarte)

sample selection through spectral graph analysis. Promising results in extensive experiments show the advantages of the proposed approach.

*Keywords:* Self-Taught Learning, Sample Selection, Autoencoder, Domain Mapping, Spectral Graph Analysis.

---

## 1. Introduction

Supervised learning is widely used in many machine learning tasks [1–3]. However, applications of supervised learning methods are limited in practical scenarios due to its requirements on large-scale labeled training datasets<sup>1</sup> with both training and testing data sharing the same label and feature space, which lead to high costs in collecting eligible training data [4, 5].

Several techniques have been proposed to tackle the limitations of supervised learning methods. Semi-supervised learning algorithms [6, 7] use both labeled and unlabeled data to improve performance when labeled training data are limited. However, the success of many semi-supervised learning algorithms highly depends on the validity of assumptions that the unlabeled and labeled data have the same distribution [6] or class labels [7]. Therefore, it is still difficult to gather unlabeled data that satisfy these preconditions.

In order to further loosen the restrictions on training data, many transfer learning approaches [8] have been proposed to use the knowledge obtained from auxiliary domains<sup>2</sup> to improve the performance on target domain tasks. *Self-taught learning* [9–17] is a type of transfer learning technique that employs unlabeled auxiliary data to improve the performance of a supervised learning task when labeled training data are limited. Though similar to semi-supervised learning, self-taught learning methods have fewer restrictions on unlabeled data, as they allow the label spaces and marginal probability distributions of unlabeled and labeled data to be different. In self-taught learning, unlabeled data are used

---

<sup>1</sup>“Training data” is used in the sequel to denote data used for model learning.

<sup>2</sup>“Source domain” is used in the sequel to denote auxiliary domains that are used to improve target domain task performance.

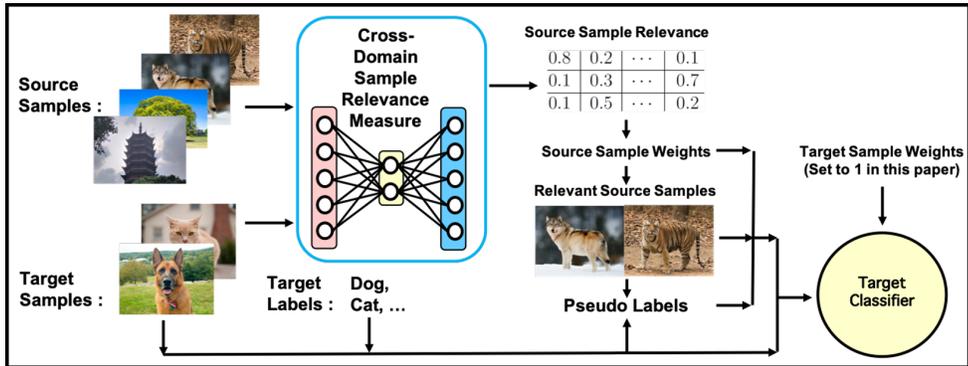


Figure 1: Block diagram for the GASTL framework.

as the source from which the knowledge learned is applied to tasks performed on labeled target data. Such a loose restriction on unlabeled data significantly simplifies learning due to the huge volume of unlabeled data we can access. However, the easily obtained unlabeled data inevitably contain samples that are only weakly related to the labeled training data, which may even harm the supervised learning performance if we treat them equally as other unlabeled samples during knowledge transfer. Target task performance degradation due to knowledge transfer from source domains that are not sufficiently related is known as *negative transfer* [18], which has been studied under several transfer learning scenarios [19, 20]. However, to the best of our knowledge, the problem of negative transfer in self-taught learning has not been studied before.

In this paper, we propose a novel algorithm for self-taught learning with unlabeled source data that are related to labeled target data to be selected with the purpose of reducing negative transfer. The algorithm leverages a linear mapping, a  $k$ -nearest neighbor ( $k$ NN) graph, and a single-layer autoencoder to obtain a metric for cross-domain relevance. We refer to this method as graph and autoencoder-based self-taught learning (GASTL). The GASTL framework includes two modules: a source sample re-weighting module and a classifier training module. In the first module, we assign each unlabeled source sample

a weight that indicates its relevance to labeled target samples.<sup>3</sup> In the second module, source samples with large weights are selected as a transfer training set to be combined with target data to train a classifier. Each selected source sample is assigned a pseudo-label from the target domain label space to be used during classifier training. The source sample weights are also used during classifier training. The trained classifier is then used to predict labels of unseen target samples. Figure 1 shows the flowchart of GASTL.

The key contributions of this paper are as follows:

- We propose a novel metric for the relevance of each source sample to the target domain in the scenario of self-taught learning based on an autoencoder and graph data regularization. To the best of our knowledge, we are the first to measure cross-domain sample relevance in order to tackle the issue of negative transfer in self-taught learning problems.
- We propose a novel classifier training scheme with both selected source samples and target samples as the training dataset with the relevance of each source sample to the target domain being considered. We are not aware of existing self-taught learning approaches that integrate cross-domain relevance into classifier training.
- We present an efficient solver for the knowledge transfer optimization problem described above that relies on an iterative scheme based on gradient descent of the proposed objective function.
- We provide multiple numerical results to demonstrate the performance improvements in terms of classification accuracy and sensitivity to parameters achieved by the proposed method compared with state-of-the-art self-taught learning methods and other relevant techniques.

---

<sup>3</sup>The setting of self-taught learning requires source samples to be unlabeled and target samples to be labeled. Therefore we do not specify the availability of label information for both source and target samples in the sequel.

The rest of this paper is organized as follows. Section 2 introduces notation. Section 3 overviews relevant techniques. The proposed framework is presented in Section 4. Experimental results are provided in Section 5. Section 6 concludes with suggestions for future work.

## 2. Notation

Vectors are denoted by bold lowercase letters while matrices are denoted by bold uppercase letters. The superscript  $T$  of a matrix denotes the transposition operation. For a matrix  $\mathbf{A}$ ,  $\mathbf{A}^{(q)}$  denotes the  $q^{\text{th}}$  column and  $\mathbf{A}_{(p)}$  denotes the  $p^{\text{th}}$  row, while  $\mathbf{A}^{(p,q)}$  denotes the entry at the  $p^{\text{th}}$  row and  $q^{\text{th}}$  column. The  $\ell_{r,p}$ -norm for a matrix  $\mathbf{W} \in \mathbb{R}^{a \times b}$  is denoted as  $\|\mathbf{W}\|_{r,p} = \left( \sum_{i=1}^a \left( \sum_{j=1}^b |\mathbf{W}^{(i,j)}|^r \right)^{p/r} \right)^{1/p}$ . The  $\ell_r$ -norm for a vector  $\mathbf{w} \in \mathbb{R}^a$  is denoted as  $\|\mathbf{w}\|_r = \left( \sum_{i=1}^a |\mathbf{w}^{(i)}|^r \right)^{1/r}$ . The trace of a matrix  $\mathbf{L} \in \mathbb{R}^{a \times a}$  is defined as  $\text{Tr}(\mathbf{L}) = \sum_{i=1}^a \mathbf{L}^{(i,i)}$ . We use  $\mathbf{1}$  and  $\mathbf{0}$  to denote an all-ones and all-zeros matrix or vector of the appropriate size, respectively. We use  $\mathbf{X} = [\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, \dots, \mathbf{X}^{(n)}] \in \mathbb{R}^{d \times n}$  to denote a sample set, where  $\mathbf{X}^{(i)} \in \mathbb{R}^d$  is the  $i^{\text{th}}$  sample in  $\mathbf{X}$  for  $i = 1, 2, \dots, n$ , and where  $d$  and  $n$  denote the sample dimensionality and the number of samples in  $\mathbf{X}$ , respectively.

In transfer learning, we use  $\mathcal{D}$  to denote a domain and  $\mathcal{T}$  to denote a task. A domain  $\mathcal{D}$  consists of a feature space  $\mathcal{X}$  and a marginal probability distribution  $P(\mathbf{X})$  over a sample set  $\mathbf{X}$ . A task  $\mathcal{T}$  consists of a label space  $\mathcal{Y}$  and an objective predictive function  $f(\mathbf{X}, \mathbf{Y})$  to predict the corresponding labels  $\mathbf{Y}$  of a sample set  $\mathbf{X}$ . We refer readers to [8] for a detailed explanation of these notations. We use  $\mathcal{D}_{\text{src}} = \{\mathcal{X}_{\text{src}}, P(\mathbf{X}_{\text{src}})\}$  and  $\mathcal{T}_{\text{src}} = \{\mathcal{Y}_{\text{src}}, f(\mathbf{X}_{\text{src}}, \mathbf{Y}_{\text{src}})\}$  to denote the source domain and task, and use  $\mathcal{D}_{\text{trg}} = \{\mathcal{X}_{\text{trg}}, P(\mathbf{X}_{\text{trg}})\}$  and  $\mathcal{T}_{\text{trg}} = \{\mathcal{Y}_{\text{trg}}, f(\mathbf{X}_{\text{trg}}, \mathbf{Y}_{\text{trg}})\}$  for the target domain and task.

### 3. Background and Related Work

#### 3.1. Self-Taught Learning

Transfer learning methods can be classified as homogeneous ( $\mathcal{X}_{\text{src}} = \mathcal{X}_{\text{trg}}$ ) or heterogeneous ( $\mathcal{X}_{\text{src}} \neq \mathcal{X}_{\text{trg}}$ ), and as transductive ( $\mathcal{T}_{\text{src}} = \mathcal{T}_{\text{trg}}$ ) or inductive ( $\mathcal{T}_{\text{src}} \neq \mathcal{T}_{\text{trg}}$ ) [8]. The self-taught learning setting assumes different label spaces between the source and target domains, and label information is assumed to be unavailable in the source domain. Therefore, the self-taught learning setting is similar to the inductive transfer learning setting when no labeled data is available in the source domain. The goal of self-taught learning is to use the unlabeled source domain data to help improve the target domain task performance. In this paper we focus on homogeneous self-taught learning.

The idea of self-taught learning was first proposed by Raina et al. [9].<sup>4</sup> In STL, a dictionary  $\mathbf{D}$  is first learned using source domain samples. After that, the sparse coefficients for target domain samples  $\mathbf{A}_{\text{trg}}$  are computed based on  $\mathbf{D}$ . Finally, a classifier is learned on  $\mathbf{A}_{\text{trg}}$  as well as target domain labels by applying a supervised learning algorithm. Robust and discriminative self-taught learning (RDSTL) [16] is an extension of STL that takes advantage of the label information of target samples during learning and makes the dictionary learning process more robust to noise and outliers by imposing an  $\ell_{2,1}$ -norm constraint on both the dictionary learning reconstruction loss and the sparse coefficient matrix. Self-taught low-rank (S-Low) coding [17] is suitable for both clustering and classification tasks in visual learning. By imposing a low-rank constraint onto the sparse coefficient matrix, S-Low coding is able to characterize the global structure information in the target domain. Self-taught clustering [10] is the first algorithm proposed to tackle unsupervised inductive transfer learning problems and aims at clustering a small amount of target unlabeled data by learning a useful feature representation with the help of large amounts of unlabeled

---

<sup>4</sup>We use abbreviation STL to denote the method of [9] in the sequel, while we use the full name “self-taught learning“ for the class of learning problems.

source domain data. Kuen et al. [11] incorporates self-taught learning into visual tracking, where local representations are learned offline on unlabeled data and transferred to the observational model of the proposed tracker. Kemker et al. [14] applies self-taught learning to hyperspectral image classification by learning models from sufficiently large quantities of unlabeled source data that are distinct from the labeled target data to extract generalizable features. The trained models are then used to extract features on the labeled target data for classification. He et al. [15] combines self-taught learning, sparse autoencoder, and radial basis functions to the field of wound infection detection to solve the problem of insufficient labeled wound infection samples. A basis vector is first learned on the easily obtained unlabeled pollutant gas samples. Then new representation of wound infection samples are learned using the basis vector under a sparsity constraint for classification. The idea of self-taught learning has also been incorporated into object localization [12, 13] to learn object detectors without or with little human supervision.

Although these self-taught learning approaches use different schemes for knowledge transfer, they all use the entire source sample set without considering their relevance to the target domain, which makes these methods potentially vulnerable to negative transfer.

### 3.2. Single-Layer Autoencoder

A single-layer autoencoder is an artificial neural network that aims to reconstruct inputs by using only a single hidden layer, and is widely used in transfer learning [21]. Given input data  $\mathbf{x} \in \mathbb{R}^d$ , an autoencoder first maps  $\mathbf{x}$  to a compressed data representation  $\mathbf{z} \in \mathbb{R}^m$  in a hidden layer, given by  $\mathbf{z} = f(\mathbf{W}_1\mathbf{x} + \mathbf{b}_1)$ , where  $\mathbf{W}_1 \in \mathbb{R}^{m \times d}$  is a weight matrix,  $\mathbf{b}_1 \in \mathbb{R}^m$  is a bias vector, and  $f(\cdot)$  is an elementary nonlinear activation function. This part is referred to as an *encoder*. The second step is to map the compressed data representation  $\mathbf{z}$  to output data  $\bar{\mathbf{x}} \in \mathbb{R}^d$ , which is  $\bar{\mathbf{x}} = g(\mathbf{W}_2\mathbf{z} + \mathbf{b}_2)$ , where  $\mathbf{W}_2 \in \mathbb{R}^{d \times m}$  and  $\mathbf{b}_2 \in \mathbb{R}^d$  are the corresponding weight matrix and bias vector, respectively. This part is referred to as a *decoder*.

The optimization problem underlying autoencoder training is to minimize the difference between the input data and the output data. To be more specific, given a set of data  $\mathbf{X} = [\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, \dots, \mathbf{X}^{(n)}]$ , the parameters  $\mathbf{W}_1$ ,  $\mathbf{W}_2$ ,  $\mathbf{b}_1$ , and  $\mathbf{b}_2$  are adapted to minimize the reconstruction error  $\sum_i \|\mathbf{X}^{(i)} - \bar{\mathbf{X}}^{(i)}\|_2^2$ , where  $\bar{\mathbf{X}}^{(i)}$  is the output of autoencoder to the input  $\mathbf{X}^{(i)}$ .

### 3.3. Instance-Based Transfer Learning

Instance-based transfer learning assumes that certain instances of data from the source domain can be reused for learning in the target domain by reweighting. Instance reweighting and sampling based on instance importance are two major techniques in this context [8]. Kernel mean matching (KMM) [22] and multiscale landmarks selection (MLS) [23] are two instance-based transfer learning methods that fit our setting of unlabeled source data since neither of them require label information from the source domain. KMM is a nonparametric method which directly produces resampling weights without requiring the estimation of biased densities or selection probabilities, or the assumption that the probabilities of the different classes are known. MLS selects landmarks that are similarly distributed in the two domains to reduce the discrepancy between the domains, where each instance from the union of source and target domain data is considered as a candidate landmark, and the candidate is kept as a landmark if its quality is sufficiently high according to some measure.

## 4. Proposed Method

In this section, we introduce our proposed GASTL approach. The basic framework of GASTL is to reconstruct both source and target samples through a single-layer autoencoder, while simultaneously enforcing a linear relationship between source samples and target samples. Both global and local data structures are preserved through a single-layer autoencoder and spectral graph analysis, respectively. We develop a metric for the relevance between each source sample and the target samples, which is used for source sample selection for

knowledge transfer. Meanwhile, a weight is assigned to each source sample reflecting its relevance to the target samples for the subsequent classifier training, during which each selected source sample is assigned a pseudo-label from the target domain label space and combined with target samples to build the classifier training sample set. Source sample weights are also considered during classifier training. Finally, the trained classifier is used to predict labels of unseen target samples.

#### 4.1. Knowledge Transfer and Relevance Measure

In this section, we present the problem formulation of our knowledge transfer scheme. We also propose a measure for relevance between each source sample and the target samples.

##### 4.1.1. Objective Function

The objective function of GASTL includes four parts: a data reconstruction term, a domain mapping term, a regularization term for sample selection, and a term based on spectral graph analysis for local data structure preservation. The details of these four terms are described below.

Many transfer learning methods perform knowledge transfer from a source domain to a target domain by finding a mapping between them [24, 25]. In particular, we can set  $h_1(\mathbf{X}_{\text{trg}}) = h_2(\mathbf{X}_{\text{src}})\mathbf{A}$ , where  $h_1(\cdot)$  and  $h_2(\cdot)$  are two transformations, while  $\mathbf{A}$  is a matrix that linearly maps transformed source samples  $h_2(\mathbf{X}_{\text{src}})$  into transformed target samples  $h_1(\mathbf{X}_{\text{trg}})$ . More specifically, the mapping is obtained from the optimization:

$$\min_{\Theta, \mathbf{A}} \mathcal{M}(\Theta, \mathbf{A}) + \lambda \mathcal{R}(\mathbf{A}),$$

where  $\Theta$  is a set of parameters used for the nonlinear mappings  $h_1$  and  $h_2$ , while  $\mathcal{M}(\Theta, \mathbf{A}) = \mathcal{L}(h_1(\mathbf{X}_{\text{trg}}), h_2(\mathbf{X}_{\text{src}})\mathbf{A})$  denotes a cost function for domain mapping, where  $\mathcal{L}(\cdot, \cdot)$  is a loss function and  $\mathcal{R}(\cdot)$  corresponds to a regularization function on  $\mathbf{A}$  to avoid overfitting.<sup>5</sup>

---

<sup>5</sup>We empirically found that regularizing  $\Theta$  did not noticeably affect the performance of

A simple way to achieve domain mapping is to assume a linear mapping between source and target data, which is  $\mathbf{X}_{\text{trg}} = \mathbf{X}_{\text{src}}\mathbf{A}$ . This requires the cost  $\mathcal{M}(\Theta, \mathbf{A}) = \mathcal{L}(\mathbf{X}_{\text{trg}}, \mathbf{X}_{\text{src}}\mathbf{A})$ . The use of a linear mapping in knowledge transfer is often computationally efficient. However, the success of this knowledge transfer scheme relies on an assumption that  $\mathbf{X}_{\text{trg}} \in \text{span}(\mathbf{X}_{\text{src}})$  [26]. Due to the ubiquitous large discrepancy between the source and target domains in self-taught learning scenarios,  $\mathbf{X}_{\text{trg}}$  is usually not in the span of  $\mathbf{X}_{\text{src}}$ , and hence a linear reconstruction scheme can hardly do well in knowledge transfer. Therefore, we aim to find a non-linear reconstruction scheme that can decrease the discrepancy between source and target domains. One possible way to do this is to find a nonlinear transformation on  $\mathbf{X}_{\text{trg}}$ , and recover the output of this transformation as a linear transformation of source samples which are relevant to the target samples. That is,  $h(\mathbf{X}_{\text{trg}}) = \mathbf{X}_{\text{src}}\mathbf{A}$ , where  $h(\cdot)$  is a nonlinear transformation. Furthermore, due to the possible large diversity of source samples compared with the target samples, we can assume that the feature space shared by both source and target domains are separated into several clusters: the source samples lie near a union of many clusters, while the target samples concentrate near a single cluster. Intuitively, negative transfer can be alleviated by using source samples that are close to target samples for knowledge transfer.

As mentioned in Section 3.2, a single-layer autoencoder aims at minimizing the reconstruction error between output and input data. We use  $\mathbf{X} = [\mathbf{X}_{\text{src}} \ \mathbf{X}_{\text{trg}}]$  as the input to a single-layer autoencoder by optimizing a reconstruction error-driven loss function:

$$\mathcal{L}(\Theta) = \frac{1}{2n} \|\mathbf{X} - h(\mathbf{X}; \Theta)\|_F^2, \quad (1)$$

where  $n = n_{\text{src}} + n_{\text{trg}}$ ,  $\Theta = [\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}_1, \mathbf{b}_2]$ , and  $h(\mathbf{X}; \Theta) = g(\mathbf{W}_2 \cdot f(\mathbf{W}_1\mathbf{X} + \mathbf{b}_1) + \mathbf{b}_2)$ .<sup>6</sup> We use the sigmoid function as the activation function:  $f(z) = g(z) = 1/(1 + \exp(-z))$ . In (1), both source and target samples share the

---

knowledge transfer. Therefore, we do not pursue such regularization in this paper.

<sup>6</sup>We often drop the dependence on  $\Theta$  for readability, i.e. we use  $h(\mathbf{X})$  to denote  $h(\mathbf{X}; \Theta)$  when no ambiguity is caused.

same parameters to train an autoencoder, which makes the reconstructed source and target samples lie in the same submanifold under the learned parameters  $\Theta$ . Meanwhile, we use the following minimization problem for the purpose of domain mapping:

$$\mathcal{C}(\Theta, \mathbf{A}) = \frac{1}{2n_{\text{trg}}} \|\mathbf{X}_{\text{src}}\mathbf{A} - \text{h}(\mathbf{X}_{\text{trg}}; \Theta)\|_F^2. \quad (2)$$

That is, we enforce the target samples in the autoencoder output to be reconstructed by a linear combination of the source samples. While it is feasible to separate the optimization of (1) and (2), we observed that a joint framework is able to provide better knowledge transfer performance. Due to the nonlinear nature of transformation featured by a single-layer autoencoder, the distribution gap can be ameliorated through minimizing  $\mathcal{C}(\Theta, \mathbf{A})$  with respect to  $\Theta$  and  $\mathbf{A}$ . Therefore, we define the mapping cost

$$\mathcal{M}(\Theta, \mathbf{A}) = \mathcal{L}(\Theta) + \mu\mathcal{C}(\Theta, \mathbf{A}),$$

where  $\mu$  is a balance parameter, to obtain a nonlinear mapping between source samples and target samples.

Since each row of  $\mathbf{A}$  indicates the importance of the corresponding source sample in reconstructing transformed target samples, we use the  $\ell_2$ -norm of each row of  $\mathbf{A}$  to measure the *relevance* between a source sample and the target samples. This leads to an  $\ell_{2,1}$ -norm regularization function  $\mathcal{R}(\mathbf{A}) = \|\mathbf{A}\|_{2,1}$  that enforces row sparsity on the transformation matrix  $\mathbf{A}$ .

Data transformations based on autoencoders only guarantee broad data structure preservation, which does not take pair-wise relationships between data points into consideration. Therefore, we need to include local geometric structures from the data into our objective function. Local geometric structures of the data often contain discriminative information of neighboring data point pairs [27–32], in which nearby data points are assumed to have similar representations. In order to characterize the local data structure, we construct a  $k$ -nearest neighbor ( $k$ NN) graph  $\mathbb{G}$  on the data space. The edge weight between two connected data points is determined by the similarity between those

two points. We define the adjacency matrix  $\mathbf{S}$  for the graph  $\mathbb{G}$  as follows: for a data point  $\mathbf{X}^{(i)}$ , its weight  $\mathbf{S}^{(i,j)} \neq 0$  if and only if  $\mathbf{X}^{(i)} \in \mathcal{N}_k(\mathbf{X}^{(j)})$  or  $\mathbf{X}^{(j)} \in \mathcal{N}_k(\mathbf{X}^{(i)})$ , where  $\mathcal{N}_k(\mathbf{X}^{(i)})$  denotes the  $k$ -nearest neighborhood set for  $\mathbf{X}^{(i)}$ ; otherwise,  $\mathbf{S}^{(i,j)} = 0$ . In this paper, we use cosine similarity to determine nonzero weights as  $\mathbf{S}^{(i,j)} = (\mathbf{X}^{(i)T} \mathbf{X}^{(j)}) / (\|\mathbf{X}^{(i)}\|_2 \|\mathbf{X}^{(j)}\|_2)$ . The Laplacian matrix  $\mathbf{L}$  of the graph  $\mathbb{G}$  is defined as  $\mathbf{L} = \mathbf{D} - \mathbf{S}$ , where  $\mathbf{D}$  is a diagonal matrix whose  $i^{\text{th}}$  element on the diagonal is defined as  $\mathbf{D}^{(i,i)} = \sum_{j=1}^n \mathbf{S}^{(i,j)}$ . With these definitions, we set up the following minimization objective for local data structure preservation:

$$\mathcal{G}(\Theta) = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \|\mathbf{Z}^{(i)} - \mathbf{Z}^{(j)}\|_2^2 \mathbf{S}^{(i,j)} = \text{Tr}(\mathbf{Z}\mathbf{L}\mathbf{Z}^T),$$

where  $\mathbf{Z}^{(i)} = f(\mathbf{W}_1 \mathbf{X}^{(i)} + \mathbf{b}_1)$  for  $i = 1, 2, \dots, n$ , and  $\mathbf{Z} = [\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}, \dots, \mathbf{Z}^{(n)}]$ .

The final objective function for source sample selection can be written in terms of the following minimization with respect to the parameters  $\Theta = [\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}_1, \mathbf{b}_2]$  and  $\mathbf{A}$ :

$$\{\hat{\Theta}, \hat{\mathbf{A}}\} = \arg \min_{\Theta, \mathbf{A}} \mathcal{L}(\Theta) + \mu \mathcal{C}(\Theta, \mathbf{A}) + \lambda \mathcal{R}(\mathbf{A}) + \gamma \mathcal{G}(\Theta), \quad (3)$$

where  $\mu$ ,  $\lambda$ , and  $\gamma$  are balance parameters.

#### 4.1.2. Optimization

The closed form solution of the optimization problem in (3) is hard to obtain due to the  $\ell_{2,1}$ -norm regularization term. We employ an alternating optimization scheme to solve this problem with  $\Theta$  and  $\mathbf{A}$  being iteratively updated, until the objective function value in (3) converges or a maximum number of iterations is reached.

When  $\mathbf{A}$  is fixed, (3) becomes

$$\hat{\Theta} = \arg \min_{\Theta} \mathcal{F}_1(\Theta) := \arg \min_{\Theta} \mathcal{L}(\Theta) + \mu \mathcal{C}(\Theta, \mathbf{A}) + \gamma \mathcal{G}(\Theta). \quad (4)$$

Following [33], we use a limited memory Broyden-Fletcher-Goldfarb-Shanno (L-BFGS) algorithm to solve (4). The L-BFGS algorithm is a type of quasi-Newton method that requires much fewer iterations to converge than first order methods

such as gradient descent. Furthermore, compared with other BFGS algorithms, the L-BFGS algorithm has low computational cost, making it possible to use the whole dataset for optimization and provide more stable performance than commonly used stochastic gradient descent algorithms. For example, the dimensionality of the parameter  $\Theta$  is the sum of the dimensionalities of  $\mathbf{W}_1 \in \mathbb{R}^{m \times d}$ ,  $\mathbf{W}_2 \in \mathbb{R}^{d \times m}$ ,  $\mathbf{b}_1 \in \mathbb{R}^m$ , and  $\mathbf{b}_2 \in \mathbb{R}^d$ , which is  $2md + d + m$ . Compared with the conventional BFGS algorithm, which requires computing and storing of  $(2md + d + m) \times (2md + d + m)$  Hessian matrices, the L-BFGS algorithm saves the past  $l$  updates of  $\Theta$  and corresponding gradients. Therefore, denoting the number of iterations in the optimization by  $t$ , the corresponding time complexity of L-BFGS is  $O(tlmd)$ . We refer readers to [34] for more details on L-BFGS algorithm, which we implement using the *minFunc* toolbox [35]. The solver requires the gradients of the objective function in (4) with respect to its parameters  $\Theta$ . The gradients for both  $\mathcal{L}(\Theta)$  and  $\mathcal{C}(\Theta, \mathbf{A})$  can be obtained through a back-propagation algorithm. We skip the details for the derivation of the gradients of both  $\mathcal{L}(\Theta)$  and  $\mathcal{C}(\Theta)$ , which are standard in the formulation of back-propagation for an autoencoder. The resulting gradients for  $\mathcal{L}(\Theta)$  are:

$$\frac{\partial \mathcal{L}(\Theta)}{\partial \mathbf{W}_1} = \frac{1}{n} \Delta_{\mathcal{L}_2} \mathbf{X}^T, \quad \frac{\partial \mathcal{L}(\Theta)}{\partial \mathbf{W}_2} = \frac{1}{n} \Delta_{\mathcal{L}_3} \mathbf{Y}^T, \quad \frac{\partial \mathcal{L}(\Theta)}{\partial \mathbf{b}_1} = \frac{1}{n} \Delta_{\mathcal{L}_2} \mathbf{1}, \quad \frac{\partial \mathcal{L}(\Theta)}{\partial \mathbf{b}_2} = \frac{1}{n} \Delta_{\mathcal{L}_3} \mathbf{1},$$

where each column of  $\Delta_{\mathcal{L}_2} \in \mathbb{R}^{m \times n}$  and  $\Delta_{\mathcal{L}_3} \in \mathbb{R}^{d \times n}$  contains the error term of the corresponding sample for the hidden layer and the output layer, respectively:

$$\Delta_{\mathcal{L}_3} = (\mathbf{h}(\mathbf{X}) - \mathbf{X}) \bullet \mathbf{h}(\mathbf{X}) \bullet (\mathbf{1} - \mathbf{h}(\mathbf{X})), \quad \Delta_{\mathcal{L}_2} = (\mathbf{W}_2^T \Delta_{\mathcal{L}_3}) \bullet \mathbf{Y} \bullet (\mathbf{1} - \mathbf{Y}),$$

with  $\bullet$  denoting the element-wise product operator. The gradients for  $\mathcal{C}(\Theta, \mathbf{A})$  are:

$$\begin{aligned} \frac{\partial \mathcal{C}(\Theta, \mathbf{A})}{\partial \mathbf{W}_1} &= \frac{1}{n_{\text{trg}}} \Delta_{\mathcal{C}_2} (\mathbf{X}_{\text{src}} \mathbf{A})^T, & \frac{\partial \mathcal{C}(\Theta, \mathbf{A})}{\partial \mathbf{W}_2} &= \frac{1}{n_{\text{trg}}} \Delta_{\mathcal{C}_3} \mathbf{Y}_{\text{trg}}^T, \\ \frac{\partial \mathcal{C}(\Theta, \mathbf{A})}{\partial \mathbf{b}_1} &= \frac{1}{n_{\text{trg}}} \Delta_{\mathcal{C}_2} \mathbf{1}, & \frac{\partial \mathcal{C}(\Theta, \mathbf{A})}{\partial \mathbf{b}_2} &= \frac{1}{n_{\text{trg}}} \Delta_{\mathcal{C}_3} \mathbf{1}. \end{aligned} \quad (5)$$

Both  $\Delta_{\mathcal{L}_2}^{(i)}$  and  $\Delta_{\mathcal{L}_3}^{(i)}$  in (4.1.2) play same roles as  $\Delta_{\mathcal{C}_2}^{(i)}$  and  $\Delta_{\mathcal{C}_3}^{(i)}$  in (5). Their definitions are:

$$\Delta_{\mathcal{C}_3} = (\mathbf{h}(\mathbf{X}_{\text{trg}}) - \mathbf{X}_{\text{src}} \mathbf{A}) \bullet \mathbf{h}(\mathbf{X}_{\text{trg}}) \bullet (\mathbf{1} - \mathbf{h}(\mathbf{X}_{\text{trg}})), \quad \Delta_{\mathcal{C}_2} = (\mathbf{W}_2^T \Delta_{\mathcal{C}_3}) \bullet \mathbf{Y}_{\text{trg}} \bullet (\mathbf{1} - \mathbf{Y}_{\text{trg}}),$$

where  $\mathbf{Y}_{\text{trg}} = \mathbf{f}(\mathbf{W}_1 \mathbf{X}_{\text{trg}} + \mathbf{b}_1)$ . The gradients of the graph term  $\mathcal{G}(\Theta) = \text{Tr}(\mathbf{Y}\mathbf{L}\mathbf{Y}^T)$  can be obtained in a straightforward fashion as follows:

$$\begin{aligned}\frac{\partial \mathcal{G}(\Theta)}{\partial \mathbf{W}_1} &= \frac{\partial \text{Tr}(\mathbf{Y}\mathbf{L}\mathbf{Y}^T)}{\partial \mathbf{Y}} \cdot \frac{\partial \mathbf{Y}}{\partial \mathbf{W}_1} = 2(\mathbf{Y}\mathbf{L} \bullet \mathbf{Y} \bullet (\mathbf{1} - \mathbf{Y})) \mathbf{X}^T, & \frac{\partial \mathcal{G}(\Theta)}{\partial \mathbf{W}_2} &= \mathbf{0}, \\ \frac{\partial \mathcal{G}(\Theta)}{\partial \mathbf{b}_1} &= \frac{\partial \text{Tr}(\mathbf{Y}\mathbf{L}\mathbf{Y}^T)}{\partial \mathbf{Y}} \cdot \frac{\partial \mathbf{Y}}{\partial \mathbf{b}_1} = 2(\mathbf{Y}\mathbf{L} \bullet \mathbf{Y} \bullet (\mathbf{1} - \mathbf{Y})) \mathbf{1}, & \frac{\partial \mathcal{G}(\Theta)}{\partial \mathbf{b}_2} &= \mathbf{0}.\end{aligned}$$

To conclude, the gradients of the objective function in (4) with respect to  $\Theta = [\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}_1, \mathbf{b}_2]$  can be written as

$$\begin{aligned}\frac{\partial \mathcal{F}_1(\Theta)}{\partial \mathbf{W}_1} &= \frac{1}{n} \Delta_{\mathcal{L}2} \mathbf{X}^T + \frac{\mu}{n_{\text{trg}}} \Delta_{\mathcal{C}2} (\mathbf{X}_{\text{src}} \mathbf{A})^T + 2\gamma (\mathbf{Y}\mathbf{L} \bullet \mathbf{Y} \bullet (\mathbf{1} - \mathbf{Y})) \mathbf{X}^T, \\ \frac{\partial \mathcal{F}_1(\Theta)}{\partial \mathbf{W}_2} &= \frac{1}{n} \Delta_{\mathcal{L}3} \mathbf{Y}^T + \frac{\mu}{n_{\text{trg}}} \Delta_{\mathcal{C}3} \mathbf{Y}_{\text{trg}}^T, \\ \frac{\partial \mathcal{F}_1(\Theta)}{\partial \mathbf{b}_1} &= \frac{1}{n} \Delta_{\mathcal{L}2} \mathbf{1} + \frac{\mu}{n_{\text{trg}}} \Delta_{\mathcal{C}2} \mathbf{1} + 2\gamma (\mathbf{Y}\mathbf{L} \bullet \mathbf{Y} \bullet (\mathbf{1} - \mathbf{Y})) \mathbf{1}, \\ \frac{\partial \mathcal{F}_1(\Theta)}{\partial \mathbf{b}_2} &= \frac{1}{n} \Delta_{\mathcal{L}3} \mathbf{1} + \frac{\mu}{n_{\text{trg}}} \Delta_{\mathcal{C}3} \mathbf{1}.\end{aligned}$$

When  $\Theta$  is fixed, (3) becomes

$$\begin{aligned}\hat{\mathbf{A}} &= \arg \min_{\mathbf{A}} \mathcal{F}_2(\mathbf{A}) := \arg \min_{\mathbf{A}} \mu \mathcal{C}(\Theta, \mathbf{A}) + \lambda \mathcal{R}(\mathbf{A}) \\ &= \arg \min_{\mathbf{A}} \frac{\mu}{2n_{\text{trg}}} \|\mathbf{X}_{\text{src}} \mathbf{A} - \mathbf{h}(\mathbf{X}_{\text{trg}}; \Theta)\|_F^2 + \lambda \|\mathbf{A}\|_{2,1}.\end{aligned}\tag{6}$$

Following [36], we calculate the derivative of  $\mathcal{F}_2(\mathbf{A})$  with respect to  $\mathbf{A}$  through

$$\frac{\partial \mathcal{F}_2(\mathbf{A})}{\partial \mathbf{A}} = \frac{\mu}{n_{\text{trg}}} [\mathbf{X}_{\text{src}}^T \mathbf{X}_{\text{src}} \mathbf{A} - \mathbf{X}_{\text{src}}^T \mathbf{h}(\mathbf{X}_{\text{trg}})] + \lambda \mathbf{A} \mathbf{U},$$

where  $\mathbf{U} \in \mathbb{R}^{n_{\text{src}} \times n_{\text{src}}}$  is a diagonal matrix whose  $i^{\text{th}}$  element on the diagonal is

$$\mathbf{U}^{(i,i)} = (2\|\mathbf{A}_{(i)}\|_2)^{-1}.\tag{7}$$

We add a small constant  $\epsilon$  to each element in  $\mathbf{A}$  to avoid overflow; thus  $\|\mathbf{A}^{(i)}\|_2$  is nonzero for each  $i$ . In this way,  $\mathcal{F}_2(\mathbf{A})$  becomes

$$\mathcal{S}(\mathbf{A}, \mathbf{U}) = \frac{\mu}{2n_{\text{trg}}} \|\mathbf{X}_{\text{src}} \mathbf{A} - \mathbf{h}(\mathbf{X}_{\text{trg}}; \Theta)\|_F^2 + \lambda \text{Tr}(\mathbf{A}^T \mathbf{U} \mathbf{A}).\tag{8}$$

Therefore, when  $\mathbf{U}$  is fixed, the optimal value of  $\mathbf{A}$  can be obtained through

$$\hat{\mathbf{A}} = (\mu \mathbf{X}_{\text{src}}^T \mathbf{X}_{\text{src}} + n_{\text{trg}} \lambda \mathbf{U})^{-1} \mu \mathbf{X}_{\text{src}}^T \mathbf{h}(\mathbf{X}_{\text{trg}}).\tag{9}$$

We can update  $\mathbf{U}$  through (7) when  $\mathbf{A}$  is fixed and update  $\mathbf{A}$  through (9) when  $\mathbf{U}$  is fixed with an iterative scheme until the value of  $\mathcal{F}_2(\mathbf{A})$  converges.

#### 4.1.3. Convergence Analysis

The optimization of GASTL is based on an alternative scheme. When  $\mathbf{A}$  is fixed, an L-BFGS algorithm is employed to optimize (4). The convergence analysis of L-BFGS algorithm can be found in [34]. Following [36], we show the convergence behavior when  $\Theta$  is fixed as follows.

**Proposition 1.** *For two variables  $a$  and  $b$ , if  $b$  is positive, then the following inequality holds:*

$$\frac{a^2}{2b} - a \geq \frac{b^2}{2b} - b \quad (10)$$

*Proof.* For two arbitrary variables  $a$  and  $b$ , we have

$$a^2 - 2ab + b^2 \geq 0 \Leftrightarrow a^2 - 2ab \geq b^2 - 2b^2$$

If  $b$  is positive, then we get (10).  $\square$

**Proposition 2.** *When  $\Theta$  is fixed, the optimization procedure of (6) is non-increasing over iteration by employing the optimization procedure in Section 4.1.2.*

*Proof.* When  $\mathbf{U}$  is fixed as  $\mathbf{U}^t$  in the  $t^{\text{th}}$  iteration, minimizing (8) with respect to  $\mathbf{A}$  is a convex optimization problem; thus leading to the following inequality:

$$\begin{aligned} & \frac{\mu}{2n_{\text{trg}}} \|\mathbf{X}_{\text{src}} \mathbf{A}^{t+1} - \mathbf{h}(\mathbf{X}_{\text{trg}}; \Theta)\|_F^2 + \lambda \text{Tr} [(\mathbf{A}^{t+1})^T \mathbf{U}^t \mathbf{A}^{t+1}] \\ & \leq \frac{\mu}{2n_{\text{trg}}} \|\mathbf{X}_{\text{src}} \mathbf{A}^t - \mathbf{h}(\mathbf{X}_{\text{trg}}; \Theta)\|_F^2 + \lambda \text{Tr} [(\mathbf{A}^t)^T \mathbf{U}^t \mathbf{A}^t] \end{aligned} \quad (11)$$

The trace item on the left side can be written as

$$\begin{aligned} & \text{Tr}((\mathbf{A}^{t+1})^T \mathbf{U}^t \mathbf{A}^{t+1}) \\ & = \|\mathbf{A}^{t+1}\|_{2,1} + \text{Tr}((\mathbf{A}^{t+1})^T \mathbf{U}^t \mathbf{A}^{t+1}) - \|\mathbf{A}^{t+1}\|_{2,1} \\ & = \|\mathbf{A}^{t+1}\|_{2,1} + \sum_{i=1}^{n_{\text{src}}} \left( \frac{\|\mathbf{A}_{(i)}^{t+1}\|_2^2}{2\|\mathbf{A}_{(i)}^t\|_2} - \|\mathbf{A}_{(i)}^{t+1}\|_2 \right). \end{aligned}$$

Likewise, the trace item on the right side can be written as

$$\text{Tr}((\mathbf{A}^t)^T \mathbf{U}^t \mathbf{A}^t) = \|\mathbf{A}^t\|_{2,1} + \sum_{i=1}^{n_{\text{src}}} \left( \frac{\|\mathbf{A}_{(i)}^t\|_2^2}{2\|\mathbf{A}_{(i)}^t\|_2} - \|\mathbf{A}_{(i)}^t\|_2 \right).$$

Therefore, (11) becomes

$$\begin{aligned}
& \frac{\mu}{2n_{\text{trg}}} \|\mathbf{X}_{\text{src}} \mathbf{A}^{t+1} - \mathbf{h}(\mathbf{X}_{\text{trg}}; \Theta)\|_F^2 + \lambda \|\mathbf{A}^{t+1}\|_{2,1} + \lambda \sum_{i=1}^{n_{\text{src}}} \left( \frac{\|\mathbf{A}_{(i)}^{t+1}\|_2^2}{2\|\mathbf{A}_{(i)}^t\|_2} - \|\mathbf{A}_{(i)}^{t+1}\|_2 \right) \\
& \leq \frac{\mu}{2n_{\text{trg}}} \|\mathbf{X}_{\text{src}} \mathbf{A}^t - \mathbf{h}(\mathbf{X}_{\text{trg}}; \Theta)\|_F^2 + \lambda \|\mathbf{A}^t\|_{2,1} + \lambda \sum_{i=1}^{n_{\text{src}}} \left( \frac{\|\mathbf{A}_{(i)}^t\|_2^2}{2\|\mathbf{A}_{(i)}^t\|_2} - \|\mathbf{A}_{(i)}^t\|_2 \right)
\end{aligned} \tag{12}$$

It is easy to obtain the following inequality based on the result of Proposition

1:

$$\sum_{i=1}^{n_{\text{src}}} \left( \frac{\|\mathbf{A}_{(i)}^{t+1}\|_2^2}{2\|\mathbf{A}_{(i)}^t\|_2} - \|\mathbf{A}_{(i)}^{t+1}\|_2 \right) \geq \sum_{i=1}^{n_{\text{src}}} \left( \frac{\|\mathbf{A}_{(i)}^t\|_2^2}{2\|\mathbf{A}_{(i)}^t\|_2} - \|\mathbf{A}_{(i)}^t\|_2 \right) \tag{13}$$

Plugging (13) into (12) leads to the following inequality:

$$\begin{aligned}
& \frac{\mu}{2n_{\text{trg}}} \|\mathbf{X}_{\text{src}} \mathbf{A}^{t+1} - \mathbf{h}(\mathbf{X}_{\text{trg}}; \Theta)\|_F^2 + \lambda \|\mathbf{A}^{t+1}\|_{2,1} \\
& \leq \frac{\mu}{2n_{\text{trg}}} \|\mathbf{X}_{\text{src}} \mathbf{A}^t - \mathbf{h}(\mathbf{X}_{\text{trg}}; \Theta)\|_F^2 + \lambda \|\mathbf{A}^t\|_{2,1}
\end{aligned} \tag{14}$$

□

**Proposition 3.** *The alternating optimization process of (3) with respect to  $\Theta$  and  $\mathbf{A}$  is convergent.*

*Proof.* When  $\mathbf{A}$  is fixed as  $\mathbf{A}^t$  in the  $t^{\text{th}}$  iteration, an L-BFGS algorithm is used to optimize (3) with regard to  $\Theta$ . Since L-BFGS algorithms are convergent,<sup>7</sup> we have

$$\mathcal{F}(\Theta^{t+1}, \mathbf{A}^t) \leq \mathcal{F}(\Theta^t, \mathbf{A}^t). \tag{15}$$

When  $\Theta$  is fixed as  $\Theta^{t+1}$  in the  $(t+1)^{\text{th}}$  iteration, we can get

$$\mathcal{F}(\Theta^{t+1}, \mathbf{A}^{t+1}) \leq \mathcal{F}(\Theta^{t+1}, \mathbf{A}^t) \tag{16}$$

based on Proposition 2. By combining (15) and (16), we have

$$\mathcal{F}(\Theta^{t+1}, \mathbf{A}^{t+1}) \leq \mathcal{F}(\Theta^t, \mathbf{A}^t)$$

Therefore, the optimization process of (3) is convergent. □

---

<sup>7</sup>We refer readers to [34] for more details on the convergence analysis of L-BFGS algorithms

## 4.2. Classifier Training

The next step is to use source samples combined with target samples to train a classifier, which can then be applied to unseen samples in the target domain for classification. Since source samples have different relevance levels with the target domain, we propose a scheme to assign weights to source samples that reflect their relevance. Source domain samples with large weights are kept while others are discarded. Subsequently, a classifier is trained using both target samples and selected source samples. For each source sample, pseudo-labels that indicate the transferability of the source sample to different target classes are used as true labels during classifier training, where transferability reflects the possibility of transferring a source sample to a target domain class [37]. The transferability values of the source samples are stored in a matrix  $\mathbf{Tr} \in \mathbb{R}^{n_{\text{src}} \times C_{\text{trg}}}$ , where  $C_{\text{trg}}$  is the cardinality of  $\mathcal{Y}_{\text{trg}}$ . Two pseudo-labeling schemes are proposed for comparison. Additionally, source sample weights are taken into consideration in classifier training. For each pseudo-labeling scheme, we evaluate both soft and hard classification with the softmax classifier.

### 4.2.1. Source Domain Sample Reweighting

As mentioned in Section 4.1.1, the  $\ell_2$ -norm value of each row of  $\mathbf{A}$  can be used to measure the relevance between the corresponding source sample and target samples. We propose a scheme to assign a weight to each source sample based on the corresponding row in  $\mathbf{A}$ . The weight for a source sample  $\mathbf{X}_{\text{src}}^{(i)}|_{i=1}^{n_{\text{src}}}$  is set proportional to the  $\ell_2$ -norm of the corresponding row in  $\mathbf{A}$ . That is, for a source sample  $\mathbf{X}_{\text{src}}^{(i)}$ , its weight vector  $\mathbf{Wt} \in \mathbb{R}^{n_{\text{src}}}$  has entries

$$\mathbf{Wt}(i) = \|\mathbf{A}_{(i)}\|_2 / \max_j (\|\mathbf{A}_{(j)}\|_2). \quad (17)$$

Note that the vector  $\mathbf{Wt}$  is normalized with the maximum entry value being 1. In addition, during classifier training, all target training samples are given weight 1.

#### 4.2.2. Pseudo-Labeling

We propose two transferability measure schemes and assign pseudo-labels to source samples based on transferability values.

**Scheme A:** For a given source sample  $\mathbf{X}_{\text{src}}^{(i)}$ , its transferability to a target class  $\mathbf{c}^{(j)}$  is measured by the square of the  $\ell_2$ -norm of a sub-vector consisting of the elements in the corresponding row of  $\mathbf{A}$  that belong to the target samples of  $\mathbf{c}^{(j)}$ . That is,

$$\mathbf{Tr}^{(i, \mathbf{c}^{(j)})} = \|\mathbf{A}^{(i, \mathbf{J}_{\mathbf{c}^{(j)}})}\|_2^2, \quad (18)$$

where  $\mathbf{J}_{\mathbf{c}^{(j)}}$  denotes the columns in  $\mathbf{A}$  that correspond to class  $\mathbf{c}^{(j)}$ .

**Scheme B:** By following [37], we adopt the isometric Gaussian probability [38] computed on the hidden layer representation of the trained single-layer autoencoder as the transferability of a given source sample  $\mathbf{X}_{\text{src}}^{(i)}$  to a target class  $\mathbf{c}^{(j)}$ . More concretely, the transferability is

$$\mathbf{Tr}^{(i, \mathbf{c}^{(j)})} = \mathcal{N}(\mathbf{Z}_{\text{src}}^{(i)} | \bar{\mathbf{Z}}_{\text{trg}}^{\mathbf{c}^{(j)}}, \sigma^2 \mathbf{I}), \quad (19)$$

where  $\mathbf{Z}_{\text{src}}^{(i)} \in \mathbb{R}^m$  is the hidden layer representation of the source sample  $\mathbf{X}_{\text{src}}^{(i)}$ , and  $\bar{\mathbf{Z}}_{\text{trg}}^{\mathbf{c}^{(j)}} \in \mathbb{R}^m$  is the mean of the hidden layer representation of target samples belonging to class  $\mathbf{c}^{(j)}$ . As such, the transferability is measured by the probability that the source sample belongs to a target class given the auxiliary information  $\bar{\mathbf{Z}}_{\text{trg}}^{\mathbf{c}^{(j)}}$ .

The pseudo-labels of source samples consist of a matrix  $\mathbf{L} \in \mathbb{R}^{n_{\text{src}} \times C_{\text{trg}}}$ . We assign pseudo-labels to the source samples based on their transferability values to different target domain classes. Given a source sample  $\mathbf{X}_{\text{src}}^{(i)}$ , for a *hard classifier*, we set

$$\mathbf{L}^{(i, j)} = \begin{cases} 1, & j = \operatorname{argmax}_k \mathbf{Tr}^{(i, \mathbf{c}^{(k)})}, \\ 0, & \text{otherwise.} \end{cases} \quad (20)$$

For a *soft classifier*, the normalized transferability values are used as pseudo-labels so that pseudo-labels reflect the likelihood of transferring source samples

to target domain classes:

$$\mathbf{L}^{(i,j)} = \mathbf{Tr}^{(i,\mathbf{c}^{(j)})} / \sum_{k=1}^{C_{\text{trg}}} \mathbf{Tr}^{(i,\mathbf{c}^{(k)})}. \quad (21)$$

Compared with hard classifiers, soft classifiers may help improve knowledge transfer performance since they are able to capture the relationship between each single source sample and multiple target categories instead of only one category. This is especially necessary for image classification tasks since there usually exist commonalities between image categories.

#### 4.2.3. Softmax Classifier Training with Weighted Samples

We employ a softmax classifier due to its simplicity and capability to do soft classification. The training data weights are included in classifier training, which leads to the following cost function

$$J(\Theta_{\mathbf{c}}) = -\frac{1}{n} \left[ \sum_{i=1}^n \mathbf{Wt}(i) \sum_{j=1}^{C_{\text{trg}}} \mathbf{L}^{(i,j)} \log \frac{\exp(\Theta_{\mathbf{c}}^{(j)T} \mathbf{X}^{(i)})}{\sum_{l=1}^{C_{\text{trg}}} \exp(\Theta_{\mathbf{c}}^{(l)T} \mathbf{X}^{(i)})} \right], \quad (22)$$

where  $\Theta_{\mathbf{c}} = [\Theta_{\mathbf{c}}^{(1)}, \Theta_{\mathbf{c}}^{(2)}, \dots, \Theta_{\mathbf{c}}^{(C_{\text{trg}})}]$  is the classifier parameter to be optimized. We use an L-BFGS algorithm to compute the optimal value of  $\Theta$ . The gradients needed for optimization are given by

$$\frac{\partial J(\Theta_{\mathbf{c}})}{\partial \Theta_{\mathbf{c}}^{(j)}} = -\frac{1}{n} \sum_{i=1}^n \left[ \mathbf{Wt}(i) \mathbf{X}^{(i)} \left( \mathbf{L}^{(i,j)} - \frac{\exp(\Theta_{\mathbf{c}}^{(j)T} \mathbf{X}^{(i)})}{\sum_{l=1}^{C_{\text{trg}}} \exp(\Theta_{\mathbf{c}}^{(l)T} \mathbf{X}^{(i)})} \right) \right].$$

The procedure for GASTL is summarized in Algorithm 1.

#### 4.3. Algorithm Analysis

In this section, we discuss the time complexity and the limitations of GASTL.

*Time Complexity:* The time complexity for the construction of a  $k$ NN graph is  $O(dn^2)$ , where  $d$  is the data dimensionality and  $n = n_{\text{src}} + n_{\text{trg}}$ . When  $\mathbf{A}$  is fixed, the time complexity of using L-BFGS algorithm to optimize (4) is  $O(tlmd)$ , where  $t$  is the number of iterations for parameter updating and  $l$  is the number of steps stored in memory; when  $\Theta$  is fixed, the time complexity to optimize (6) is  $O(n_{\text{src}}^3 + dn_{\text{src}}^2)$ , where  $O(n_{\text{src}}^3)$  results from matrix

---

**Algorithm 1** GASTL Algorithm

---

**Inputs:** Source dataset  $\mathbf{X}_{\text{src}}$ ; target dataset  $\mathbf{X}_{\text{trg}}$ ; graph neighborhood size  $k$  ;  
autoencoder hidden layer size  $m$ ; balance parameters  $\lambda$  and  $\gamma$ .

**Outputs:** Softmax classifier parameter  $\Theta_{\mathbf{c}}$ .

**Stage 1:** *Relevance Measure*

- 1: Construct a  $k$ -nearest neighbor graph  $\mathbb{G}$  on a combined dataset  $\mathbf{X} = [\mathbf{X}_{\text{src}} \ \mathbf{X}_{\text{trg}}]$ ;
- 2: Calculate the autoencoder parameter  $\Theta$  and transformation matrix  $\mathbf{A}$  by optimizing the objective function (3) with the alternating scheme described in Section 4.1.2;

**Stage 2:** *Classifier Training*

- 3: Calculate relevance weights for source samples according to the weighting scheme (17);
  - 4: Calculate the transferability of each source sample to each target class according to the transferability measure scheme (18) or (19);
  - 5: Assign target domain class labels to source samples according to the hard pseudo-labeling scheme (20) or the soft pseudo-labeling scheme (21);
  - 6: Construct a softmax classifier from the target data and labels, source data, relevance weights, and pseudo-labels by optimizing the cost function (22) to get classifier parameters  $\Theta_{\mathbf{c}}$ .
-

inversion and  $O(dn_{\text{src}}^2)$  results from matrix multiplication. The time complexity of source sample weighting is  $O(n_{\text{src}})$ , which results from finding maximum value among the  $\ell_2$ -norm values of rows in  $\mathbf{A}$ . The time complexity of pseudo-labeling is  $O(n_{\text{src}}C_{\text{trg}}^2)$ . The time complexity of the softmax classifier training is  $O(t_c l_c C_{\text{trg}} d)$ , where  $t_c$  is the number of iterations for parameter updating and  $l_c$  is the number of steps stored in memory in the classifier training step. Since the number of target classes  $C_{\text{trg}}$  is expected to be much smaller than the number of source and target samples  $n_{\text{src}}$  and  $n_{\text{trg}}$ , the time complexity of GASTL is  $O(dn^2 + t l m d + n_{\text{src}}^3 + t_c l_c C_{\text{trg}} d)$ .

*Limitation:* The analysis above indicates that the time complexity of GASTL is highly dependent on the number of samples, which will usually be dominated by the number of source samples. This is a potential bottleneck for GASTL to tackle with extremely large source datasets.

## 5. Experiments

In this section, we evaluate the knowledge transfer performance of GASTL by comparing it with other relevant state-of-the-art transfer learning techniques. To be more specific, we first select  $p$  source samples which are the most relevant to the target domain, and then use those selected source samples combined with labeled target samples to train a classifier. The classification rates and mean f1-scores on target testing samples are then used as metrics to evaluate knowledge transfer performance.

### 5.1. Dataset Preparation

Next, we provide information on the datasets used in our experiments.

Table 1: Details of datasets used in our experiment.

Dataset	Features	Samples	Classes	Type
Caltech101 (SIFTBOW)	1,000	3,000	100	Image
Caltech101 (VGG-19)	4,096	3,000	100	Image
IMDB	3,000	6,500	2	Text
Twitter	3,000	6,500	2	Text

- Dataset Information:** We employ one visual dataset (Caltech101<sup>8</sup>) and two natural language datasets (IMDB<sup>9</sup> and Twitter,<sup>10</sup> both for sentiment analysis). In order to eliminate the side effects caused by imbalanced classes, we set the number of samples from each class to be the same within each dataset through random selection.
- Feature Extraction:** In many cases, raw data cannot be used for knowledge transfer due to possible dimensionality inconsistencies. Therefore, it is necessary to do feature extraction on each dataset to make knowledge transfer feasible. For Caltech101, we employ both the SIFTBOW feature<sup>11</sup> proposed in [39] and the output of the last fully connected layer of the pre-trained VGG-19 model [40]. For both IMDB and Twitter, we use the method from [41]<sup>12</sup> to do feature extraction. We denote this feature as WORD2VEC in the sequel.

<sup>8</sup>Dataset downloaded from [http://www.vision.caltech.edu/Image\\_Datasets/Caltech101/](http://www.vision.caltech.edu/Image_Datasets/Caltech101/). The Caltech101 dataset contains both a “Faces“ and “Faces easy“ class, with each consisting of different versions of the same human face images. However, the images in “Faces“ contain more complex backgrounds. To avoid confusion between these two similar classes of images, we do not include the “Faces easy“ images in our experiments. Therefore, we keep 100 classes for Caltech101.

<sup>9</sup>Dataset downloaded from <https://drive.google.com/file/d/OB8yp1g0BCztyNOJaMDVoeXhHwM8/>.

<sup>10</sup>Dataset downloaded from <https://www.kaggle.com/c/twitter-sentiment-analysis2/data>

<sup>11</sup>Codes downloaded from <http://files.is.tue.mpg.de/pgehler/projects/iccv09/>.

<sup>12</sup>Codes downloaded from [https://github.com/yoonykim/CNN\\_sentence](https://github.com/yoonykim/CNN_sentence)

- **Source/Target Split:** For Caltech101, we randomly separate the 100 classes into 5 groups with 20 classes in each group. Five independent self-taught learning experiments were conducted on Caltech101: in each experiment samples in one group are used as target samples and those in the remaining four groups are used as source samples. For each class in the target domain, 15 samples were used for training and 15 samples were used for testing. For the two natural language datasets, we used one as source and the other one as target. That is, when IMDB was used as target, then Twitter was used as source and vice versa. For computational convenience, we did not use the entire datasets when either IMDB and Twitter is used as the source. We randomly selected 3,000 samples as source for both IMDB and Twitter. For each class in the target domain, 10, 100 and 1000 samples were used for training and 750 samples were used for testing.

## 5.2. Experimental Setup

We performed classification on the target testing samples in order to evaluate the effectiveness of the self-taught learning algorithms and two instance-based transfer methods. The three self-taught learning methods are STL [9], RDSTL [16], and S-Low [17], introduced in Section 3.1. The two instance-based transfer learning methods are KMM [22] and MLS [23]. Both KMM and MLS were tailored to the scenario of self-taught learning. The computational complexity of GASTL and the five competitors are listed in Table 2. We also compute the classification performance without knowledge transfer.

Both GASTL and the compared algorithms include parameters to adjust. In this experiment, we fix some parameters and tune others through a “grid search” strategy. For algorithms requiring source sample selection, we select the number of source samples  $p \in \{10, 20, 30, \dots, 100, 150, 200, 250, \dots, 500, 1000, 1500, n_{\text{src}}\}$ . In GASTL, the range of hidden layer sizes is set to  $m \in \{10, 50, 100, 200\}$ , while the balance parameters are given ranges of  $\lambda \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1\}$  and  $\gamma \in \{0, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}\}$ . The value of  $\mu$  is set to 1. The number of nearest neighbors in a  $k$ NN graph is set to 5. The value of  $\sigma^2$  in (19) is set to 1.

Table 2: Time complexity of GASTL and five competing methods. In this table,  $d$  denotes data dimensionality,  $n_{\text{src}}$  denotes number of source samples,  $n_{\text{trg}}$  denotes number of target samples, and  $n = n_{\text{src}} + n_{\text{trg}}$ ,  $t$  denotes number of iterations for optimization and  $t_c$  denotes number of iterations for softmax classifier optimization used in GASTL and KMM and MLS that are tailored to the scenario of self-taught learning.  $m$  denotes the autoencoder hidden layer size for GASTL and dictionary learning size for STL, RDSTL, and S-Low.

Method	Time Complexity
GASTL	$\mathbf{O}(dn^2 + tlm d + n_{\text{src}}^3 + +t_c l_c C_{\text{trg}} d)$
STL	$\mathbf{O}(d^3 + d^2 n + dm)$
RDSTL	$\mathbf{O}(tnmd C_{\text{trg}})$
S-Low	$\mathbf{O}(n_{\text{src}} n_{\text{trg}} + tnmd)$
KMM	$\mathbf{O}(n^2 + +t_c l_c C_{\text{trg}} d)$
MLS	$\mathbf{O}(n + +t_c l_c C_{\text{trg}} d)$

For the optimization of GASTL with L-BFGS, we set the number of iterations  $t_1$  and  $t_2$  to be 400 and the number of storing updates  $l_1$  and  $l_2$  to be 100.

All three self-taught learning methods (STL, RDSTL, and S-Low) are based on dictionary learning, where the resulting sparse code vectors were used as features for classification. For these three methods, we first performed PCA on training data due to three reasons: (i) the features listed in Table 1 have high dimensionalities and require large dictionary sizes, which would cause prohibitive training time; (ii) we noticed that the performances of these three methods were highly dependent on the choice of parameters, which means a smaller feature dimensionality would significantly reduce the training time needed due to parameter search; (iii) we did not see systematical differences on the performance when raw features and PCA features were employed for model learning, respectively. Following [9], we kept the number of principal components to preserve approximately 96% of the training sample variance.

For both KMM<sup>13</sup> and MLS<sup>14</sup>, we first obtained weights for the source samples, which are also assigned pseudo-labels. Between the two pseudo-labeling schemes, only Scheme B is applicable since Scheme A is dependent on the transformation matrix  $\mathbf{A}$ , and these two domain adaptation methods do not generate one. We use the features listed in Table 1 instead of autoencoder activations as we did in GASTL. Subsequent steps were exactly the same as those for GASTL.

### 5.3. Parameter Sensitivity

We study the performance variation of GASTL with respect to the hidden layer size  $m$  and the two balance parameters  $\lambda$  and  $\gamma$  as measured by the classification accuracy on target testing samples. We show the results on all four datasets.

We first study the parameter sensitivity of GASTL with respect to the hidden layer size  $m$ . Due to limited space, we only present a small portion of our experimental results in Fig. 2,<sup>15</sup> where “Soft” and “Hard” refer to whether a soft or hard classifier is used, while “A” and “B” refer to the pseudo-labeling schemes. Note that when  $m$  is fixed, the number of source samples used for knowledge transfer  $p$  and the values of  $\lambda$  and  $\gamma$  are adjustable. The classification results in Fig. 2 are the highest classification accuracy among all available combinations of  $p$ ,  $\lambda$ , and  $\gamma$  with  $m$  fixed to specific values. The results show that the performance of GASTL is not too sensitive to hidden layer size on the given datasets.

We also study the parameter sensitivity of GASTL with respect to the balance parameters  $\lambda$  and  $\gamma$ , under a fixed hidden layer size  $m$ . In order to do this, we record the classification accuracy corresponding to each combination of

---

<sup>13</sup>Codes downloaded from <http://www.gatsby.ucl.ac.uk/~gretton/covariateShiftFiles/covariateShiftSoftware.html>

<sup>14</sup>Codes downloaded from <https://github.com/jindongwang/transferlearning/tree/master/code>

<sup>15</sup>For Caltech101, we use the results of one set out of five with SIFTBOW features as autoencoder inputs. For both IMDB and Twitter, we use the results with each target dataset having 10 training samples per category.

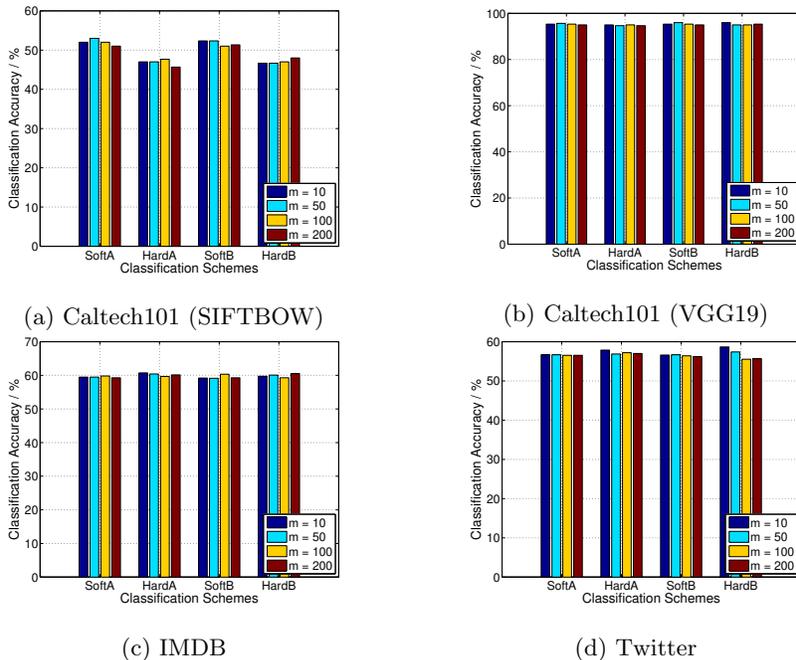


Figure 2: Performance of GASTL in classification as a function of the hidden layer size  $m$  of the autoencoder. Classification accuracy (%) is used as the evaluation metric.

$\lambda$  and  $\gamma$ , resulting in a  $5 \times 5$  matrix of performance measurement. Note that each entry in this matrix corresponds to the highest classification value over the various values of  $p$  tested. We then calculate the mean and standard deviation of these 25 elements, and the parameter sensitivity can be evaluated through the ratio between the standard deviation and the mean value. We choose the hidden layer size  $m = 10$  as Fig. 2 shows that the performance of GASTL is not sensitive to the value of  $m$ . The results are listed in Table 3, which shows that the performance of GASTL is quite stable with respect to the balance parameters  $\lambda$  and  $\gamma$  for Caltech101, IMDB, and Twitter.

#### 5.4. Performance Comparison

We present the classification accuracy results of GASTL and baselines on all datasets in Tables 4 to 6, corresponding to Caltech101 (SIFTBOW), Caltech101

Table 3: Performance stability of GASTL in classification with respect to balance parameters  $\lambda$  and  $\gamma$ . Classification accuracy mean (%) and standard deviation (%) are presented.

Dataset \ Scheme	Caltech101 (S)	Caltech101 (V)	IMDB	Twitter
SoftA	48.92 $\pm$ 1.53	94.48 $\pm$ 0.47	58.22 $\pm$ 0.60	55.79 $\pm$ 0.47
HardA	44.20 $\pm$ 1.17	92.84 $\pm$ 1.45	57.91 $\pm$ 0.98	55.41 $\pm$ 1.13
SoftB	48.75 $\pm$ 1.65	94.44 $\pm$ 0.51	58.08 $\pm$ 0.54	55.76 $\pm$ 0.47
HardB	44.27 $\pm$ 1.45	94.39 $\pm$ 1.20	57.79 $\pm$ 0.86	55.06 $\pm$ 1.29

(VGG-19), IMDB and Twitter, respectively.<sup>16</sup> respectively. Note that “Target Only” denotes the method that performs classifier training on target samples without knowledge transfer. In Tables 4 and 5, each column corresponds to one subset, while in Table 6, each column corresponds to one training sample number in each target class. We highlight the best two performances in each experiment given that we find in many cases the best two (or even more) performance are very close to each other.

We first provide an overall description on the comparison between GASTL and the competitors on each dataset. We can find that in Tables 4 and 5 the best performances are claimed by GASTL methods, and in Table 6, RDSTL is comparable to GASTL in a few cases. In other words, GASTL provides the best overall performance. We can also find that the classification performance of the two pseudo-labeling schemes are quite similar to each other in almost every case in the five tables.

Our next analysis focuses on the comparison between performance generated by soft classifiers and hard classifiers. In Table 4, it is obvious that GASTL methods with soft classifiers provide the best performance. For image datasets

---

<sup>16</sup>In the sequel, Caltech101 (S) and Caltech101 (V) are used to denote Caltech101 (SIFT-BOW) and Caltech101 (VGG-19). Results for both IMDB and Twitter are demonstrated in Table 6

Table 4: Performance of GASTL and competing feature selection algorithms in classification on Caltech101 with SIFTBOW features. Classification accuracy (%) is used as the evaluation metric.

Method \ Set ID	1	2	3	4	5
Target Only	42.33	61.67	42.33	48.33	46.00
STL	46.00	64.33	48.33	46.33	56.00
RDSTL	36.67	51.33	37.00	39.67	42.00
S-Low	35.00	51.00	36.67	34.33	36.67
KMM-Soft	48.67	66.00	47.33	52.67	56.33
KMM-Hard	43.67	63.67	43.33	47.67	48.33
MLS-Soft	47.33	66.33	47.67	51.67	53.33
MLS-Hard	45.33	62.00	42.00	46.00	47.67
GASTL-SoftA	<b>53.00</b>	<b>67.67</b>	<b>51.67</b>	<b>56.00</b>	<b>58.67</b>
GASTL-HardA	47.67	64.33	47.67	48.67	52.33
GASTL-SoftB	<b>52.33</b>	<b>68.00</b>	<b>51.33</b>	<b>53.33</b>	<b>58.67</b>
GASTL-HardB	48.00	64.00	46.33	49.33	51.00

such as Caltech101, it is unusual for the relevance between one source sample and a particular target class to be much larger than for other target classes. A soft classifier is able to characterize the relationship between a source sample and each target class during pseudo-labeling, while a hard classifier only selects the most similar class to each source sample and ignores other target classes, which may degrade knowledge transfer performance due to the possible useful information from other classes. This can also be validated by the results of KMM and MLS in Table 4, which shows the advantages of soft classifiers over hard classifiers. However, the classification rates listed in Table 5 are quite large and similar to each other. Therefore, these results cannot provide significant information on validating the advantages of the soft classifier over the hard classifier on image datasets. On the other hand, the performance increase brought

Table 5: Performance of GASTL and competing feature selection algorithms in classification on Caltech101 with VGG19 features. Classification accuracy (%) is used as the evaluation metric.

Method \ Set ID	1	2	3	4	5
Target Only	94.33	95.67	93.00	94.33	93.67
STL	95.00	95.33	92.00	94.67	94.33
RDSTL	91.67	93.00	85.67	90.00	89.33
S-Low	91.33	90.67	87.00	89.67	89.33
KMM-Soft	95.33	96.00	93.33	93.67	94.67
KMM-Hard	94.33	95.33	92.33	94.00	93.67
MLS-Soft	94.00	95.33	92.67	94.00	93.67
MLS-Hard	95.00	95.67	93.00	94.33	93.67
GASTL-SoftA	95.67	<b>97.00</b>	93.67	95.00	<b>96.00</b>
GASTL-HardA	95.00	<b>97.00</b>	<b>94.67</b>	<b>95.67</b>	<b>96.00</b>
GASTL-SoftB	<b>96.00</b>	<b>97.00</b>	93.67	94.67	<b>96.00</b>
GASTL-HardB	<b>96.00</b>	96.67	<b>94.67</b>	<b>95.67</b>	95.67

by knowledge transfer also depends on the difficulty of the classification problem. For example, the performance increase of Caltech101 using the SIFTBOW features is much larger than using the VGG-19 features.

In Table 6, hard classifiers consistently provide slightly better performance than soft classifiers for GASTL methods, while for KMM and MLS, the differences are smaller. According to our experimental setup, IMDB and Twitter play interchangeable roles as source and target. In each experiment both source and target domains share a label space with two labels (“positive sentiment” and “negative sentiment”). Therefore, in this case it is better to use hard classifier than soft classifier since the two labels indicate two mutually exclusive and largely distinguishable categories.

Table 6: Performance of GASTL and competing feature selection algorithms in classification on IMDB and Twitter. Classification accuracy (%) is used as the evaluation metric. “TS” denotes number of training samples in each target class. “Twitter  $\rightarrow$  IMDB” denotes knowledge transfer from Twitter to IMDB, while “IMDB  $\rightarrow$  Twitter” denotes the opposite direction for knowledge transfer.

Method \ TS	Twitter $\rightarrow$ IMDB			IMDB $\rightarrow$ Twitter		
	10	100	1000	10	100	1000
Target Only	57.60	68.20	73.27	55.47	58.93	77.80
STL	58.80	69.53	73.73	53.67	58.93	76.93
RDSTL	<b>60.53</b>	<b>70.47</b>	73.47	<b>58.60</b>	59.93	63.13
S-Low	59.93	64.87	73.07	55.87	59.47	64.47
KMM-Soft	57.13	69.13	73.53	56.27	61.20	78.27
KMM-Hard	58.00	68.60	73.33	57.47	60.33	78.47
MLS-Soft	56.40	68.47	72.87	56.27	59.40	77.93
MLS-Hard	56.80	68.80	72.53	57.27	60.40	78.13
GASTL-SoftA	59.80	69.33	74.60	56.67	61.93	78.20
GASTL-HardA	<b>60.73</b>	<b>70.47</b>	<b>77.67</b>	57.87	<b>62.73</b>	<b>78.53</b>
GASTL-SoftB	60.33	69.27	74.47	56.67	<b>62.00</b>	78.20
GASTL-HardB	<b>60.53</b>	<b>70.60</b>	<b>77.67</b>	<b>58.67</b>	61.93	<b>78.60</b>

## 5.5. Discussion

### 5.5.1. Effect of Local Data Structure Preservation

As mentioned in Section 4.1.1, local data structure preservation provides similar representation for nearby data points. Intuitively, local data structure preservation applied in the hidden layer of the autoencoder is likely to improve knowledge transfer performance because it is able to reduce hidden layer representation distortion as it is involved in data reconstruction and pseudo-labeling. In order to measure the effect of local data structure preservation on knowledge transfer, we compare the classification performance when  $\gamma = 0$  with the optimal one. In Fig. 3, the comparisons on one set of Caltech101 with

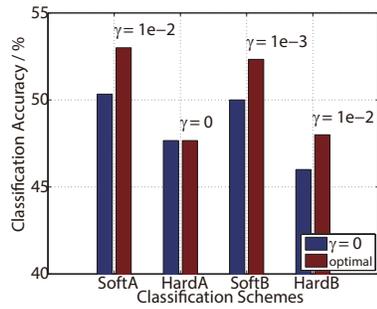
both SIFTBOW and VGG19 as the input to autoencoder and both IMDB and Twitter with each target dataset having 10 training samples per category are displayed. We find that in most cases setting  $\gamma = 0$  cannot achieve the optimal performance. Exceptions appear in the cases of “HardA” on Caltech101 with SIFTBOW features and both “HardA” and “HardB” on Twitter. Due to our observations on complete comparisons, classification rates obtained for  $\gamma = 0$  are predominantly lower than those obtained for  $\gamma \neq 0$ . Therefore, we regard these counted exceptions as outliers. Nonetheless, we found that the advantages in classification accuracy contributed by local data structure preservation are not obvious on Twitter. One possible explanation is that the local data structures in the WORD2VEC feature space cannot provide discriminative information for the Twitter dataset. Therefore, local data structure preservation negatively affected the knowledge transfer performance reflected by classification accuracy on unlabeled target samples.

### 5.5.2. Effect of Source Sample Selection

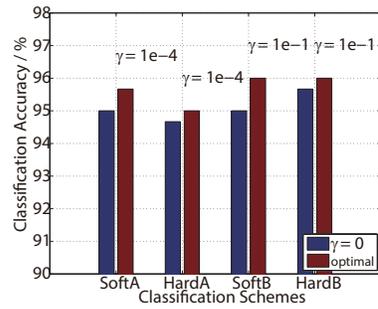
We claim that transferring knowledge from source samples indiscriminately may cause negative transfer since there is no guarantee that all source samples have sufficient relevance to the target domain. In order to demonstrate the advantage of source sample selection, we compare the classification accuracy for three different cases: the optimal value of  $p$  found with GASTL,  $p = 0$  (i.e. no transfer learning), and  $p = n_{src}$  (i.e. no sample selection). The results shown in Fig. 3 demonstrate not only that significant performance gains are obtained via GASTL, but also that in several cases the blind consideration of all source samples can in fact result in negative transfer, as seen by the reduced performance obtained with  $p = n_{src}$  versus  $p = 0$ .

## 6. Conclusions and Future Work

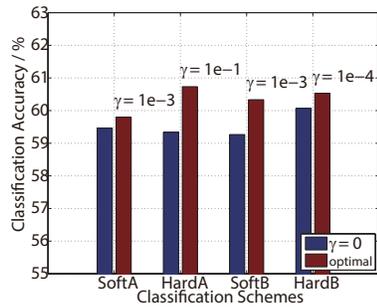
In this paper, we propose a graph and autoencoder based self-taught learning (GASTL) method. The main innovations in our self-taught learning methodology with respect to the literature can be summarized as (a) leveraging relevance



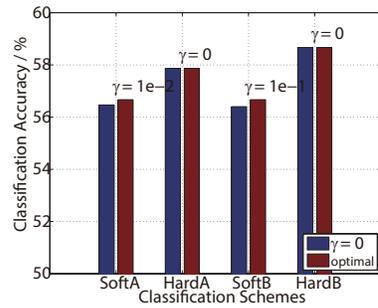
(a) Caltech101 (S)



(b) Caltech101 (V)

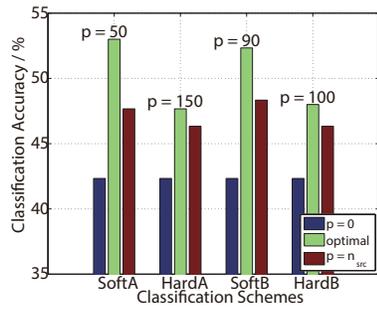


(c) IMDB

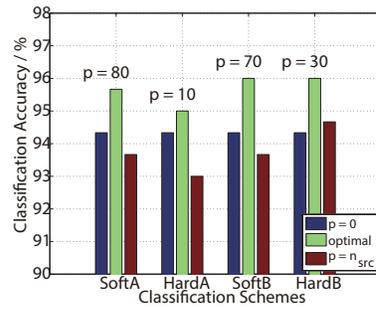


(d) Twitter

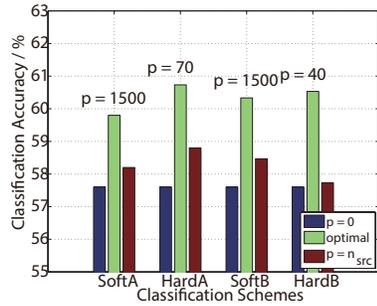
Figure 3: Comparison of GASTL performance when  $\gamma = 0$  and optimal GASTL performance. Classification accuracy (%) is used as the evaluation metric. Optimal values for  $\gamma$  are shown.



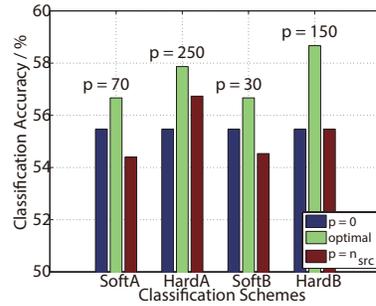
(a) Caltech101 (S)



(b) Caltech101 (V)



(c) IMDB



(d) Twitter

Figure 4: Comparison of GASTL performance when all source samples are used for classifier training and optimal GASTL performance. Classification accuracy (%) is used as the evaluation metric. Optimal values for  $p$  are shown.

metrics to select a subset of source samples in transfer learning; (b) considering cross-domain relevance for classifier training; and (c) developing our method for hard as well as soft classification problems. With our proposed framework, we decrease negative transfer and improve knowledge transfer performance in many scenarios. Experimental results demonstrate the advantages of GASTL versus methods in the literature.

For future work, we will focus on extending our work from a single-layer autoencoder-based design to one based on deep neural networks [42, 43] as well as reducing time complexity. We also plan to integrate discriminative information of target samples into our framework.

### Acknowledgment

We thank Dr. Sheng Li for providing the code of paper [17].

This research is supported in part by the Nanyang Assistant Professorship (NAP); AISG-GC-2019-003; NRF-NRFI05-2019-0002; NTU-SDU-CFAIR (NSC-2019-011); and NTU-WeBank JRI (NWJ-2019-007).

### References

- [1] T. G. Dietterich, Approximate Statistical Tests for Comparing Supervised Classification Learning Algorithms, *Neural Comput.* 10 (7) (1998) 1895–1923.
- [2] J. Cohen, P. Cohen, S. G. West, L. S. Aiken, *Applied Multiple Regression/Correlation Analysis for the Behavioral Sciences*, Routledge, 2013.
- [3] C. Hong, J. Yu, J. Zhang, X. Jin, K.-H. Lee, Multi-Modal Face Pose Estimation with Multi-Task Manifold Deep Learning, *IEEE Trans. Ind. Inf.* 15 (7) (2018) 3952–3961.
- [4] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al., Imagenet Large Scale Visual Recognition Challenge, *Int. J. Comput. Vis.* 115 (3) (2015) 211–252.

- [5] A. Rozantsev, M. Salzmann, P. Fua, Beyond Sharing Weights for Deep Domain Adaptation, *IEEE Trans. Pattern Anal. Mach. Intell.* 41 (4) (2018) 801–814.
- [6] X. Zhu, X. Wu, Class Noise Handling for Effective Cost-Sensitive Learning by Cost-Guided Iterative Classification Filtering, *IEEE Trans. Knowl. Data Eng.* 18 (10) (2006) 1435–1440.
- [7] K. Nigam, A. K. McCallum, S. Thrun, T. Mitchell, Text Classification from Labeled and Unlabeled Documents Using EM, *Mach. Learn.* 39 (2–3) (2000) 103–134.
- [8] S. J. Pan, Q. Yang, A Survey on Transfer Learning, *IEEE Trans. Knowl. Data Eng.* 22 (10) (2010) 1345–1359.
- [9] R. Raina, A. Battle, H. Lee, B. Packer, A. Y. Ng, Self-Taught Learning: Transfer Learning from Unlabeled Data, in: *Int. Conf. Mach. Learn.*, 2007, pp. 759–766.
- [10] W. Dai, Q. Yang, G.-R. Xue, Y. Yu, Self-Taught Clustering, in: *Int. Conf. Mach. Learn.*, 2008, pp. 200–207.
- [11] J. Kuen, K. M. Lim, C. P. Lee, Self-Taught Learning of a Deep Invariant Representation for Visual Tracking via Temporal Slowness Principle, *Pattern Recognit.* 48 (10) (2015) 2964–2982.
- [12] L. Bazzani, A. Bergamo, D. Anguelov, L. Torresani, Self-Taught Object Localization with Deep Networks, in: *IEEE Winter Conf. Appl. Comput. Vis.*, 2016, pp. 1–9.
- [13] Z. Jie, Y. Wei, X. Jin, J. Feng, W. Liu, Deep Self-Taught Learning for Weakly Supervised Object Localization, *arXiv preprint arXiv:1704.05188*.
- [14] R. Kemker, C. Kanan, Self-Taught Feature Learning for Hyperspectral Image Classification, *IEEE Trans. Geosci. Remote Sens.* 55 (5) (2017) 2693–2705.

- [15] P. He, P. Jia, S. Qiao, S. Duan, Self-Taught Learning Based on Sparse Autoencoder for E-Nose in Wound Infection Detection, *Sensors* 17 (10) (2017) 2279.
- [16] H. Wang, F. Nie, H. Huang, Robust and Discriminative Self-Taught Learning, in: *Int. Conf. Mach. Learn.*, 2013, pp. 298–306.
- [17] S. Li, K. Li, Y. Fu, Self-Taught Low-Rank Coding for Visual Learning, *IEEE Trans. Neural. Netw. Learn. Syst.* 29 (3) (2018) 645–656.
- [18] M. T. Rosenstein, Z. Marx, L. P. Kaelbling, T. G. Dietterich, To Transfer or Not to Transfer, in: *Adv. Neural Inf. Proc. Syst. Workshop on Transfer Learning*, Vol. 898, 2015, pp. 1–4.
- [19] L. Duan, D. Xu, S.-F. Chang, Exploiting Web Images for Event Recognition in Consumer Videos: A Multiple Source Domain Adaptation Approach, in: *IEEE Conf. Comp. Vis. and Pattern Recognit.*, 2012, pp. 1338–1345.
- [20] L. Yang, L. Jing, J. Yu, M. K. Ng, Learning Transferred Weights from Co-Occurrence Data for Heterogeneous Transfer Learning, *IEEE Trans. Neural Networks Learn. Syst.* 27 (11) (2015) 2187–2200.
- [21] Y. Zhu, X. Hu, Y. Zhang, P. Li, Transfer Learning with Stacked Reconstruction Independent Component Analysis, *Knowledge-Based Syst.* 152 (2018) 100–106.
- [22] J. Huang, A. Gretton, K. M. Borgwardt, B. Schölkopf, A. J. Smola, Correcting Sample Selection Bias by Unlabeled Data, in: *Adv. Neural Inf. Process. Syst.*, 2007, pp. 601–608.
- [23] R. Aljundi, R. Emonet, D. Muselet, M. Sebban, Landmarks-Based Kernelized Subspace Alignment for Unsupervised Domain Adaptation, in: *IEEE Conf. Comp. Vis. and Pattern Recognit.*, 2015, pp. 56–63.
- [24] R. K. Sanodiya, J. Mathew, A Framework for Semi-Supervised Metric Transfer Learning on Manifolds, *Knowledge-Based Syst.* 176 (2019) 1–14.

- [25] X. Hu, J. Pan, P. Li, H. Li, W. He, Y. Zhang, Multi-Bridge Transfer Learning, *Knowledge-Based Syst.* 97 (2016) 60–74.
- [26] M. Shao, D. Kit, Y. Fu, Generalized Transfer Subspace Learning through Low-Rank Constraint, *Int. J. Comput. Vis.* 109 (1–2) (2014) 74–93.
- [27] D. Cai, C. Zhang, X. He, Unsupervised Feature Selection for Multi-Cluster Data, in: *ACM Int. Conf. Knowl. Dis. Data Mining*, 2010, pp. 333–342.
- [28] R. Shang, W. Wang, R. Stolkin, L. Jiao, Subspace Learning-Based Graph Regularized Feature Selection, *Knowl.-Based Syst.* 112 (2016) 152–165.
- [29] R. Shang, W. Wang, R. Stolkin, L. Jiao, Non-Negative Spectral Learning and Sparse Regression-Based Dual-Graph Regularized Feature Selection, *IEEE Trans. Cybern.* 48 (2) (2017) 793–806.
- [30] R. Shang, Y. Meng, C. Liu, L. Jiao, A. M. G. Esfahani, R. Stolkin, Unsupervised Feature Selection Based on Kernel Fisher Discriminant Analysis and Regression Learning, *Mach. Learn.* 108 (4) (2019) 659–686.
- [31] J. Yu, D. Tao, M. Wang, Y. Rui, Learning to Rank Using User Clicks and Visual Features for Image Retrieval, *IEEE Trans. Cybern.* 45 (4) (2014) 767–779.
- [32] C. Hong, J. Yu, D. Tao, M. Wang, Image-Based Three-Dimensional Human Pose Recovery by Multiview Locality-Sensitive Sparse Retrieval, *IEEE Trans. Ind. Electron.* 62 (6) (2014) 3742–3751.
- [33] Q. V. Le, J. Ngiam, A. Coates, A. Lahiri, B. Prochnow, A. Y. Ng, On Optimization Methods for Deep Learning, in: *Int. Conf. Mach. Learn.*, 2011, pp. 265–272.
- [34] D. C. Liu, J. Nocedal, On the Limited Memory BFGS Method for Large Scale Optimization, *Math. Program.* 45 (1) (1989) 501–528.

- [35] M. Schmidt, minFunc: Unconstrained Differentiable Multivariate Optimization in Matlab, Available at: <http://www.cs.ubc.ca/~schmidtm/Software/minFunc.html> (2005).
- [36] C. Hou, F. Nie, X. Li, D. Yi, Y. Wu, Joint Embedding Learning and Sparse Regression: A Framework for Unsupervised Feature Selection, *IEEE Trans. Cybern.* 44 (6) (2014) 793–804.
- [37] Y. Guo, G. Ding, J. Han, Y. Gao, Zero-Shot Learning with Transferred Samples, *IEEE Trans. Image Proc.* 26 (7) (2017) 3277–3290.
- [38] R. Socher, M. Ganjoo, C. D. Manning, A. Ng, Zero-shot Learning through Cross-Modal Transfer, in: *Adv. Neural Inf. Proc. Syst.*, 2013, pp. 935–943.
- [39] P. Gehler, S. Nowozin, On Feature Combination for Multiclass Object Classification, in: *IEEE Int. Conf. Comput. Vis.*, 2009, pp. 221–228.
- [40] K. Simonyan, A. Zisserman, Very Deep Convolutional Networks for Large-Scale Image Recognition, in: *Int. Conf. Learn. Represent.*, 2015.
- [41] Y. Kim, Convolutional Neural Networks for Sentence Classification, in: *Conf. Empir. Methods Nat. Lang. Process.*, 2014, pp. 1746–1751.
- [42] J. Yu, X. Yang, F. Gao, D. Tao, Deep Multimodal Distance Metric Learning Using Click Constraints for Image Ranking, *IEEE Trans. Cybern.* 47 (12) (2016) 4014–4024.
- [43] C. Hong, J. Yu, J. Wan, D. Tao, M. Wang, Multimodal Deep Autoencoder for Human Pose Recovery, *IEEE Trans. Image Process.* 24 (12) (2015) 5659–5670.