

Learning Multi-Tasks with Inconsistent Labels by using Auxiliary Big Task

Quan Feng^{a,b}, Songcan Chen^{a,b,*}

^a*College of Computer Science and Technology, Nanjing University of Aeronautics and Astronautics, Nanjing, Jiangsu, 211106, China*

^b*MIIT Key Laboratory of Pattern Analysis and Machine Intelligence, Nanjing University of Aeronautics and Astronautics, Nanjing, Jiangsu, 211106, China*

Abstract

Multi-task learning is to improve the performance of the model by transferring and exploiting common knowledge among tasks. Existing MTL works mainly focus on the scenario where label sets among multiple tasks (MTs) are usually the same, thus they can be utilized for learning across the tasks. While almost rare works explore the scenario where each task only has a small amount of training samples, and their label sets are just partially overlapped or even not. Learning such MTs is more challenging because of less correlation information available among these tasks. For this, we propose a framework to learn these tasks by jointly leveraging both abundant information from a learnt auxiliary big task with sufficiently many classes to cover those of all these tasks and the information shared among those partially-overlapped tasks. In our implementation of using the same neural network architecture of the learnt auxiliary task to learn individual tasks, the key idea is to utilize available label information to adaptively prune the hidden layer neurons of the auxiliary network to construct corresponding network for each task, while accompanying a joint learning across individual tasks. Our experimental results demonstrate its effectiveness in comparison with the state-of-the-art approaches.

Keywords: multi-task learning, inconsistent labels, auxiliary task

*Fully documented templates are available in the elsarticle package on CTAN.

*Corresponding author

Email address: `s.chen@nuaa.edu.cn` (Songcan Chen)

1. Introduction

Multi-task learning (MTL) is an approach of exploiting and transferring the relevant information among tasks to assist individual tasks obtain better generalization. Over the last few years, it has been proved to be effective in multiple different machine learning fields, such as object detection [1], image segmentation [2], image classification[3], natural language processing [4], speech recognition [5], drug discovery [6] and so on.

Currently, most existing MTL methods usually assume that learning tasks have the same label sets and use the same model [7],[8],[9]. Because these label sets contain abundant common knowledge and are transferred to each task to improve the learning performance of the MTL model [10]. However, there are more general situations in the real world, with only a small number of training samples in each task, and when their label sets overlap partially or even unoverlap, there would be less shared information between tasks, so learning such tasks will be more challenging. To meet the challenges, [11] uses a modulation and gating network to automatically adjust the shared characteristics among different tasks for the recommendation system. [12] learns various heterogeneous tasks by sharing similar convolutional kernels among multi-task networks. These methods aim to mine and use as much common knowledge hidden in the current tasks as possible, but for the above-mentioned general scenarios, these methods still leave an improved room in performance.

To achieve the above improvement, we re-focus on the two major issues affecting MTL. Firstly, *how to extract suitable knowledge from different tasks for current multiple tasks*. Since abundant knowledge exists in the nature of multiple tasks that directly affects the joint learning among tasks and useful knowledge for the current tasks can improve the performance of the whole MT model, a key is how to extract this knowledge to avoid notorious negative transfer [13]. For this reason, many methods have been proposed and can be divided into two types: non-deep methods and deep methods: 1) *the non-*

deep methods build on shallow models to learn the parameters involved, e.g., [14] extracts useful knowledge between tasks by regularizing a task-coupled kernel function (such as a support vector machine) for the user’s prediction of product selection. [15] obtains useful knowledge between tasks by learning the same covariance matrix to predict students’ test scores. 2) *The deep methods learn a shared representation from the individual task networks to improve the performance of the current tasks*. E.g., [16] designs a feature matching network (i.e., knowledge transfer) to capture shared features in different tasks. [17] uses a segmented attention head module to capture useful knowledge between tasks for depth estimation. [18] uses a two-level graph neural network to learn useful knowledge of different tasks to improve the performance of the MTL model. Secondly, *how to design an effective MTL sharing mechanism*. An effective sharing mechanism can increase the predictive performance of the MTL model by using useful knowledge between related tasks [19]. Inspired by this motivation, many classic MTs sharing mechanisms have been designed. According to whether the task’s feature/label spaces are consistent among tasks we can divide these mechanisms into two types: homogeneous task sharing and heterogeneous task sharing, as shown in Table 1. The homogeneous task sharing mechanisms can further be subdivided into 1) *hard sharing* based: the implementation of this type of methods assumes that all tasks share knowledge in the same hidden space. For example, [20] connects the aggregated features of specific layers between tasks for semantic segmentation and depth prediction of images. 2) *Soft sharing* based: the implementation of this type of methods assumes that all task models and parameters are independent, and the distance between model parameters is regularized to obtain similar parameters for joint learning. E.g., [21] uses the attention mechanism to share parameters in specific layers between different tasks to identify symptoms of depression. 3) *Mixed sharing* based: the implementation of this type of methods uses a special task strategy to select the layer of the multi-task network model can perform shared learning. Typically, [22] uses a specific task strategy to mix these common features with the current tasks for image semantic and normal segmentation.

The heterogeneous task sharing mechanism can likewise be further subdivided into 1) *Sparse sharing* based: the implementation of this type of methods is to form a sub-network appropriate for individual tasks from an overparameterized base network, and to extract the common knowledge from the overlapped parts of the sub-networks through the sparse strategy. For example, [23] extracts shared parameters as common knowledge to learn individual tasks by a mask in the overlapping part of the sub-networks. 2) *Gradient sharing* based: the implementation of this type of methods uses some similarities to measure the gradient difference between tasks and calculate the nonnegative weights in these tasks, thereby constructing a shared gradient. For example, [24] constructs a shared gradient to measure the gradient difference among individual tasks by cosine distance to predict hospital mortality. 3) *Hierarchical sharing* based: the implementation of this type of methods performs hierarchical sharing for different overlapping areas between multiple tasks. For example, [25] learns common knowledge from different levels of multiple task networks for natural language processing.

Unfortunately, most of the above works are designed for the scenario where the label sets are the same among tasks, rather than for the scenario where the label sets are partially overlapped or even unoverlapped. A few current works design various learning mechanisms for the latter scenario. However, such methods only capture useful knowledge among tasks, which is still difficult to effectively solve the scenario. For this, we propose a novel multi-task learning framework with the help of a big auxiliary deep network (DAMTL), whose intention is to use the auxiliary task(s) with abundant knowledge to assist learning given multiple tasks with partial, even unoverlapping label sets. Our framework is shown in Fig.1. Based on it, we first pre-train a big overparameterized auxiliary network which contains the class label information in all individual tasks; Next, we use a set of soft masks to selectively prune the neurons in the convolutional layers of this network to yield the corresponding network for each learning task. Finally, we jointly train all individual networks in an end-to-end manner.

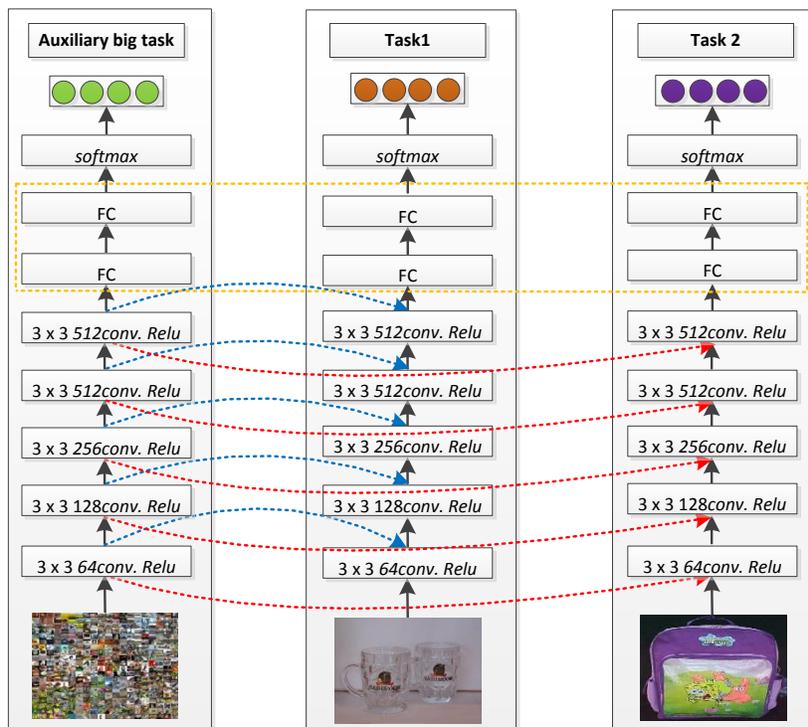


Figure 1: Our proposed DAMTL networks. The network consists of three identical independent task networks. The left side is the auxiliary big task, and the rest are individual tasks; \dashrightarrow and \dashrightarrow indicate the direction in which the knowledge of the auxiliary task is transferred to different tasks, and the yellow dotted box denotes the alignment layers.

In summary, our contribution can be summarized as below:

- 1) We present DAMTL framework and provide a new way to solve the problem of partial overlap or even unoverlap of label sets in MTL.
- 2) We design a novel knowledge extracting strategy that uses a set of soft masks to prune neurons in the convolutional layers of the auxiliary task network to extract knowledge for each task learning.
- 3) We propose a new alignment strategy that alleviates the possible class drift in the knowledge transfer from the auxiliary task to individual tasks in the DMATL network.
- 4) We conduct experiments on twelve public datasets and compare with the

Table 1: comparison of various MTL sharing mechanisms.

Sharing mechanism	Homogeneous	Heterogeneous	Supervised	Algorithms
Hard sharing	✓	×	✓	[26, 27]
Soft sharing	✓	×	✓	[27, 28]
Mixed sharing	✓	×	✓	[29, 30]
Sparse sharing	✓	✓	✓	[31, 23]
Gradient sharing	×	✓	✓	[24, 32]
Hierarchical sharing	×	✓	✓	[33, 34]

Homogeneous task, Heterogeneous task, Supervised learning, Representative algorithms

state-of-the-art methods to prove the effectiveness of our method.

The rest of this paper is arranged as follows. In Section 2, we briefly review related work in multi-task learning. In Section 3, we introduce the architecture of DAMTL, give the definition and some related theoretical application analysis. In the experimental stage of Section 4, we present image classification results on benchmark data sets. Finally, we conclude in Section 5. The code is available at <http://parnec.nuaa.edu.cn/3021/list.htm>.

2. Related Work

MTL has good performance in many applications, especially in the field of computer vision, so it has attracted a lot of attention in recent years. In this section, we briefly review the related works of MTL based on shared task features and MTL based on shared model parameters. Our work follows the latter research line.

2.1. MTL based on shared task features

The methods of this class usually assume that a common feature representation can be learned from individual tasks. According to the implementation manners, they likewise can roughly be divided into three sub-types:

1) *Selective sharing of task features*: for the tasks in the same subspace, they realize sharing by specifically regularizing the features among tasks. Typically, [35] uses the ℓ_2 norm to regularize the task weight matrices to extract shared features for the test score prediction of most school students. [36] uses the $\ell_{1,2}$ norm to regularize the weight matrices to extract shared features between tasks for learning multi-tasks with different feature dimensions. [37] uses $\ell_{2,1}$ norm to regularize the weight matrices of various modal tasks to jointly select common features for multi-modal classification of Alzheimer’s disease.

2) *Priori knowledge sharing of tasks*: for the tasks defined in the same subspace, they use the same prior knowledge among tasks to realize sharing. Typically, [38] embeds prior knowledge (i.e., pathological images with different magnification belong to the same subclass) into the feature extraction process among different tasks to verify the relationship between tasks and pathological image categories for fine-grained classification and pathophysiological image classification. [39] uses a kind of meta data (i.e., contextual attributes) as a priori knowledge to capture the relationship between different tasks for multiple tasks clustering. [40] uses the same subclass of the gland area as the prior information in the convolutional neural network to guide the network inference for pathological colon image analysis.

3) *Transformation sharing of task features*: for the tasks represented in the same subspace, they realize sharing by performing the nonlinear transformations of the original feature representation among tasks. Typically, [41] uses a set of non-linearly transformed feature sharing units for image semantic segmentation and normal estimation. [42] uses the feature adapter to learn the non-linear transformation of the tasks features to automatically evaluate the child’s speech ability.

2.2. MTL based on shared model parameters

The methods of this class usually associate different tasks with their partial model parameters or weights to realize sharing. According to their learning manners used, they can roughly be divided into three sub-types:

1) *Weighted sharing of weight matrices*: for the tasks represented in the same subspace, they realize sharing by weightedly combining a set of weight matrices among tasks. Typically, [43] weights the weight matrices among tasks for boundary classification of keywords. [44] partitions the weight matrices among tasks into common and private parts, then weights the common part for multi-label classification. [45] weights the weight matrices at the same spatial position in the pictures and transfers them to each task for image depth estimation, segmentation, and surface normal prediction.

2) *Common factor sharing via decomposing individual weight matrices*: for the weight matrix of each task model, they decompose these matrices into private and common parts, where the common part is used for sharing. Typically, [46] decomposes the weight matrices of multiple task models into common and private parts, and further uses the common part for visual target tracking. [47] sparsely decomposes the parameter tensor of the prediction model into multiple parameter matrices, and linearly combines the corresponding parameter matrices into a set of base matrices for sharing. [48] decomposes a collective matrix of drug-disease correlations to share the correlation matrix between them for drug discovery.

3) *Low-rank structure sharing of model weight matrices*: for the tasks represented in the same subspace, they capture the low-rank structure of the weight matrix among tasks by specifically regularizing to realize sharing. Typically, [49] uses feature tensor flattening of different tasks (i.e., a convex combination of matrix trace norms) to capture its low-rank structure for multi-task learning. [50] uses a set of low-rank matrices to capture the potential relationships between multiple tasks for Parkinson’s disease diagnosis. [51] uses a set of low-rank matrices constrained by the nuclear norm for target detection in hyper-spectral images.

3. Our Method

Most of the previously mentioned methods extract useful knowledge among tasks to make predictions. However, these methods are difficult to be further improved when faced with the partially overlapped, even unoverlapped labels among tasks. To overcome this difficulty, we try to leverage a big auxiliary task with abundant labels and class information to assist learning these tasks with limited data. As shown in Fig.2, our method mainly consists of three steps: 1) pre-training a big overparameterized auxiliary network, 2) selectively extracting the corresponding specific weight parameters for individual tasks from the auxiliary network, 3) transferring the weight parameters to these individual tasks to assist them learning. Specifically, Section 3.1 formally defines the problem. Section 3.2 details the proposed method, which extracts knowledge from the pre-trained auxiliary network through a soft making matrix, and then transfers them into individual tasks to form the corresponding networks. Finally, the whole DAMTL network is formulated.

3.1. Problem formulation

Given a big auxiliary task T_{aux} and a dataset $\mathcal{D}_{aux} = \{x_i, y_i\}_{i=1}^N$ containing N samples with $x_i \in \mathbb{R}^d$ and its associated label $y_i \in \{1, \dots, c\}$, where d and c are the numbers of dimensions and classes in the dataset \mathcal{D}_{aux} , respectively. Meanwhile we are given M individual tasks $\{T_j\}_{j=1}^M$, and corresponding training dataset $\mathcal{D}_j = \left\{x_k^j, y_k^j\right\}_{k=1}^{N_j}$ with N_j samples, $x_k^j \in \mathbb{R}^d$ and its associated label $y_k^j \in \{1, \dots, c_j\}$, where c_j is the number of classes in the dataset \mathcal{D}_j . We assume that the classset $\mathcal{C}_{T_{aux}}$ of the auxiliary task contains all the individual tasks classes \mathcal{C}_T , namely, $\mathcal{C}_{T_{aux}} = \mathcal{C}_{T_1} \cup \mathcal{C}_{T_2} \cup \dots, \mathcal{C}_{T_M}$, where \mathcal{C}_{T_i} and $\mathcal{C}_{T_j} (i \neq j)$ can be partially overlapped, or even unoverlapped. This makes DAMTL applicable under more general settings than most existing MTL methods.

3.2. Implementation

To learn the above tasks mentioned in Section 3.1, we first assume that this big auxiliary task network has L convolutional layers and H fully con-

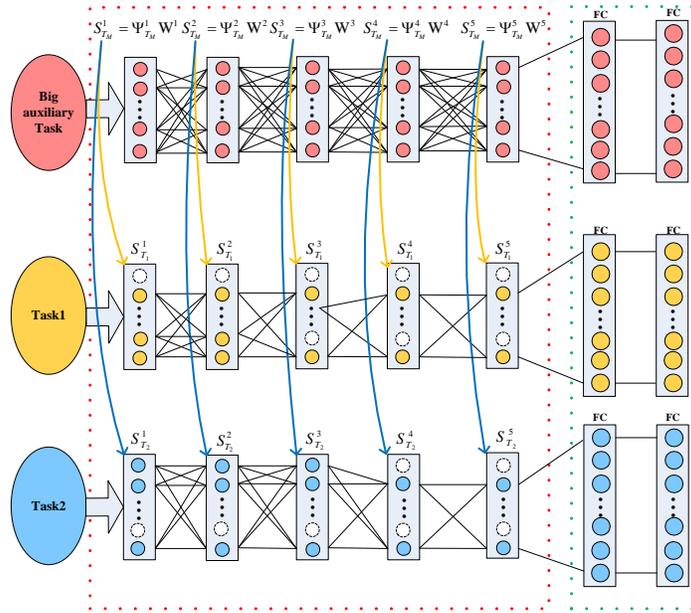


Figure 2: Detailed framework of the DMTAL network. The top of the framework is the auxiliary task network, and the rest are different task networks. All the filled circles of different colors denote neurons, while dashed circles are neurons that are pruned. The red and blue dotted boxes are the weight transfer layers and the alignment layers, respectively.

nected (FC) layers, and pre-train the network to obtain two sets of weights $W_{T_{aux}} = \{W_{aux}^1, W_{aux}^2, \dots, W_{aux}^L\}$ and $f_{T_{aux}} = \{f_{aux}^1, f_{aux}^2, \dots, f_{aux}^H\}$, where W_{aux}^l and f_{aux}^h represent the weights of the l -th convolutional layer and the weights of the h -th FC layer, respectively. Then, we manage to acquire necessary knowledge selectively from the trained big task to learn individual tasks by soft masking the weights of each corresponding convolutional layer, as shown in Fig 3. In what follows, we take the task T_j as a training example. To selectively extract specific knowledge layer-wisely from the big auxiliary task network, we introduce a soft masking matrix $\Psi_{T_j}^l$ as follows:

$$S_{T_j}^l = \Psi_{T_j}^l \odot W_{aux}^l, \quad (1)$$

here $S_{T_j}^l$ denotes the extracted knowledge from the auxiliary task, which will be transferred to the l -th convolutional layer in task T_j , $\Psi_{T_j}^l$ just takes non-negative value, and \odot denotes the Hardmard product. Then, we use Eq.(2) to activate the multiplication of $S_{T_j}^l$ and $F_{T_j}^{l-1}$ to realize the knowledge being transferred

$$F_{T_j}^l = \sigma \left(S_{T_j}^l F_{T_j}^{l-1} + b_{T_j}^l \right), \quad (2)$$

where $F_{T_j}^l$ is the feature representation of the l -th convolutional layer of the task T_j network, σ is the activation function, and $b_{T_j}^l$ is the bias vector.

Now optimizing the $S_{T_j}^l$ boils down to optimizing $\Psi_{T_j}^l$. Specifically, in order to ensure the prediction accuracy of the DAMTL network, we leverage *Conditional Maximum Mean Discrepancy* (CMMD) [52] to align the conditional probability distributions between T_{aux} and T_j in the FC layers to alleviate the possible class drift in the knowledge transfer from T_{aux} to the T_j , which is expressed as:

$$D(T_{aux}, T_j) = \sum_{c_j=1}^{C_{T_j}} \left\| \frac{1}{k_{c_j}^{T_j}} \sum_{k=1}^{k_{c_j}^{T_j}} f_{aux}^h(O_k^{h-1}) - \frac{1}{k_{c_j}^{T_j}} \sum_{k=1}^{k_{c_j}^{T_j}} f_{T_j}^h(O_k^{h-1}) \right\|_{\mathcal{H}}^2, \quad (3)$$

where (O_k^{h-1}) is the representation of the task T_j in the $(h-1)$ -th FC layer,

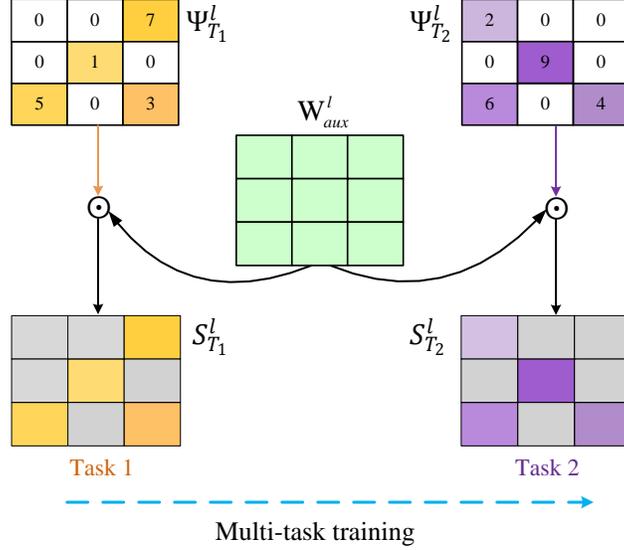


Figure 3: Knowledge extraction illustrations. The light green square in the middle represents the weight of the big auxiliary task network in the l -th convolutional layer. The squares in the upper left and upper right corners represent the soft masking matrix for different individual tasks, and the squares in the lower left and right corners represent the extracted knowledge.

f_{aux}^h and $f_{T_j}^h$ are the weights in the h -th FC layer, and $k_{c_j}^{T_j}$ is the number of samples of the c_j -th class in task T_j .

In practice, we use the inputs x_k ($k = 1, 2, \dots, n$) and the labels y_k ($k = 1, 2, \dots, n$) to minimize the following individual task loss

$$\mathcal{L}_{T_j}(\hat{y}_k^{T_j}, y_k^{T_j}) = -\sum_{k=1}^{N_j} y_k^{T_j} (\log \hat{y}_k^{T_j}) + \lambda_1 \sum_{l=1}^L \|\Psi_{T_j}^l\|_1 + \lambda_2 \sum_{h=1}^H D(T_{aux}, T_j), \quad (4)$$

where $\hat{y}_k^{T_j}$ is the predicted output, λ_1 and λ_2 are hyper-parameters, the first term is the cross-entropy loss of task T_j , and the second term uses the ℓ_1 norm to make the soft masking matrix $\Psi_{T_j}^l$ sparse. Finally, we conduct joint training of these tasks with the total loss

$$\mathcal{L}(\theta) = \sum_{j=1}^M \alpha_{T_j} \sum_{k=1}^{N_j} \mathcal{L}_{T_j}(\hat{y}_k^{T_j}, y_k^{T_j}), \quad (5)$$

where α_{T_j} is the hyper-parameter to be adjusted.

The whole process of the proposed method to solve DAMTL is summarized in Algorithm (1).

Algorithm 1: DAMTL

Notations: $W_{T_j} = \{W_{T_j}^l\}_{l=1}^L$ and $\Psi_{T_j} = \{\Psi_{T_j}^l\}_{l=1}^L$ denote the temporary variables and soft masking matrices in the task T_j network, respectively, $f_{T_j} = \{f_{T_j}^h\}_{h=1}^H$ are the weights of the FC layers in the task T_j network.

Input: $D_j, W_{T_{aux}}, \lambda_1, \lambda_2$

Output: $W_{T_j}, f_{T_j}, \Psi_{T_j}^l$

- 1 random initialization $\Psi_{T_j}^l$, initialize W_{T_j} with $W_{T_{aux}}$
- 2 **repeat**
- 3 **for** each convolution layer l **do**
- 4 $S_{T_j}^l \leftarrow (\Psi_{T_j}^l \odot W_{T_j}^l)$
- 5 $F_{T_j}^l = \sigma(S_{T_j}^l F_{T_j}^{l-1} + b_{T_j}^l)$
- 6 **end for**
- 7 **for** each fully connected layer h **do**
- 8 align the conditional probability distributions between T_{aux} and T_j by Eq.(3)
- 9 **end for**
- 10 Update the parameters $W_{T_j}^l, \Psi_{T_j}$ and f_{T_j} by the back-propagation algorithm
- 11 **until** convergence;

4. Experiments

In this section, we use the VGG network [53] as the base-model and the ImageNet dataset as a big auxiliary task to conduct two categories of experiments: 1) the label sets among tasks partially overlap; 2) the label sets among tasks do not overlap.

Table 2: Summarize statistics for datasets where part of the label sets overlap.

Datasets	classes	Size of image	Overlapping classes
Caltech-101	10	200 * 300	3
Caltech-256	10	371 * 326	3
Amazon	10	150 * 900/557 * 28	3
Webcam	10	200 * 150/900 * 557	3
Dlsr	10	200 * 150/900 * 557	3
Product	10	117 * 85/4384 * 2686	3

4.1. Datasets

We conduct experiments on the following data sets and divide 70% of the data as training set and the remaining 30% as testing set. The detailed information is shown in Table 2 and Table 3.

The **ImageNet Dataset**¹ is a computer vision dataset, which is used as an auxiliary task in the experiment. The dataset contains 21,841 categories and 14,197,122 images.

The **Office-Caltech Dataset**² is divided into two datasets, Office-Caltech10 and Office-Caltech31, and each dataset has 2,533 images, which are composed of three different subsets Dslr, Amazon and Webcam. We randomly select 10 categories in each subset to do experiment.

The **Office Home Dataset**³ is composed of Art, Clipart, Product, Real-World, and each subset covers 15500 images from 65 categories. We also randomly select 10 categories in each subset as the dataset.

The **Caltech-256 Dataset**⁴ is a dataset collected by the California Institute of Technology. It has 256 categories, and each category has more than 80 images and a total of 30,607 images. We also randomly select 10 categories in this data

¹<http://image-net.org/download>

²<https://people.eecs.berkeley.edu/~jhoffman/domainadapt/>

³<http://hemantdv.org/OfficeHome-Dataset/>

⁴http://www.vision.caltech.edu/Image_Datasets/Caltech256/

Table 3: Summarize the statistics of the datasets with unoverlapping label sets.

Datasets	classes	Size of image	Overlapping classes
Art	10	117 * 85/4384 * 2686	—
Real World	10	117 * 85/4384 * 2686	—
Caltech-101	10	200 * 300	—
Webcam	10	200 * 150/900 * 557	—
Amazon	10	200 * 150/900 * 557	—
Tiny ImageNet	10	64 * 64	—

set for experiment.

The **Tiny ImageNet Dataset**⁵ is a dataset collected by Stanford University, which is a subset of the ImageNet dataset. The dataset contains 200 categories and a total of 100,000 images, and the size of each image is $3 * 64 * 64$. We randomly select 10 categories from the Tiny ImageNet Dataset for experiments.

4.2. Comparison Methods

For evaluation, we use common single-task and multi-tasks network architectures to train each task separately/jointly, and its experimental results serve as our single-task and multi-tasks baseline. Simultaneously, we compare our proposed method with the following MTL methods:

Single task⁶ [53]: This method uses a single VGG network to learn the predictive model for each independent task.

Multi-task⁷: This method uses multiple identical VGG networks to jointly learn a multi-task prediction model.

Cross-Stitch⁸ [54]: This method uses a cross stitch unit to learn the com-

⁵<https://www.kaggle.com/c/tiny-imagenet>

⁶<https://github.com/machrisaa/tensorflow-vgg>

⁷https://github.com/luntai/VGG16_Keras_TensorFlow

⁸<https://github.com/helloyide/Cross-stitch-Networks-for-Multi-task-Learning>

mon features in two network feature layers for learning multiple tasks.

NDDR-CNN⁹ [55]: This method uses the NDDR module to automatically integrate the features of each sub-network layer for MTL.

MTAL¹⁰[12]: This method leverages the similarity between convolution kernels to capture common knowledge among multi-network for joint learning of various tasks.

DAMTL¹¹: This method uses the abound knowledge of a lager auxiliary task to help joint learning of multiple tasks.

4.3. Hyper-Parameter Tuning

In the contrasted deep neural network methods, we adjust the hidden units, learning rate, and the number of training steps in each layer according to the parameter settings of the corresponding reference. In DAMTL, we adjust the hyper-parameters in the same way. Specifically, for all experiments, we set the learning rate η to 0.01, α_{T_j} to 0.01, λ_1 and λ_2 are 0.9. In addition, in DAMTL network training, we select the rectified linear unit (ReLU) function as the activation function σ . All the deep learning models are implemented by Tensorflow.

4.4. Results of Model Performance:

We conduct experiments on the above-mentioned datasets and compare our methods with the state-of-the-arts (SOTAs), while analyzing the experimental results.

First, the performance of the DAMTL method shown in Table 4 and Table 5 is significantly better than other methods. In the PO-2 group of experiments, the DAMTL method is better than the MTAL methods.

Secondly, Table 4 shows that most of the MTL methods are better than the single-task learning method, which demonstrate that using the relationship

⁹<https://github.com/ethanygao/NDDR-CNN>

¹⁰<http://parnec.nuaa.edu.cn/3021/list.htm>

¹¹<http://parnec.nuaa.edu.cn/3021/list.htm>

Table 4: The performance comparison of various methods in the experiment of partially overlapping label sets between tasks. Among them, the bold numbers are the best classification results, and the underlined numbers are the sub-optimal classification results.

Methods	PO-1		PO-2		PO-3	
	Caltech-101	Caltech-256	Amazon	Webcam	Dlsr	Product
Single-Task	0.76 ± 0.037	0.45 ± 0.037	0.74 ± 0.036	0.63 ± 0.038	0.77 ± 0.031	0.63 ± 0.035
Multi-task	0.76 ± 0.050	0.53 ± 0.055	0.80 ± 0.046	0.69 ± 0.049	0.80 ± 0.046	<u>0.68 ± 0.049</u>
Cross-Stich	0.75 ± 0.050	0.53 ± 0.059	0.80 ± 0.046	0.67 ± 0.054	0.79 ± 0.042	0.65 ± 0.056
NDDR-CNN	0.75 ± 0.054	0.51 ± 0.055	0.79 ± 0.042	0.67 ± 0.044	0.75 ± 0.049	0.67 ± 0.050
MTAL	<u>0.78 ± 0.054</u>	<u>0.53 ± 0.051</u>	0.81 ± 0.042	<u>0.69 ± 0.048</u>	<u>0.80 ± 0.038</u>	0.66 ± 0.055
DAMTL	0.80 ± 0.050	0.54 ± 0.059	<u>0.80 ± 0.044</u>	0.70 ± 0.048	0.84 ± 0.042	0.74 ± 0.052

Table 5: Performance comparison of various methods in unoverlapping experiments of label sets between tasks. Among them, the bold numbers are the best classification results, and the underlined numbers are the sub-optimal classification results.

Methods	NO-1		NO-2		NO-3	
	Art	Real World	Caltech-101	Webcam	Amazon	T-ImagNet
Single-Task	0.60 ± 0.039	0.51 ± 0.043	0.76 ± 0.037	0.63 ± 0.038	0.77 ± 0.031	0.51 ± 0.048
Multi-task	0.62 ± 0.055	<u>0.53 ± 0.054</u>	0.78 ± 0.050	0.68 ± 0.050	0.78 ± 0.048	0.53 ± 0.058
Cross-Stich	<u>0.63 ± 0.057</u>	0.51 ± 0.056	<u>0.78 ± 0.048</u>	0.67 ± 0.050	0.79 ± 0.046	0.52 ± 0.062
NDDR-CNN	0.59 ± 0.056	0.52 ± 0.056	0.76 ± 0.051	0.66 ± 0.050	0.79 ± 0.044	0.46 ± 0.461
MTAL	0.61 ± 0.056	0.52 ± 0.058	0.73 ± 0.046	<u>0.69 ± 0.044</u>	<u>0.80 ± 0.042</u>	<u>0.54 ± 0.067</u>
DAMTL	0.67 ± 0.050	0.6 ± 0.047	0.80 ± 0.050	0.72 ± 0.044	0.81 ± 0.047	0.54 ± 0.065

between tasks to capture their useful information for interaction can promote the effectiveness of the joint learning of multiple tasks. In addition, we find that different multi-task learning methods have different performance results, which are caused by the differences among tasks. For example, the experimental results of the dataset Caltech-101 in PO-1 show that the Cross-Stich method and the NDDR-CNN method are lower than the accuracy of the single-task learning method.

Finally, the performance of DAMTL is better than that other methods in all datasets, which shows that extracting and transferring features from a big auxiliary task (which contains the class label information in all individual tasks) can help the joint learning of multiple tasks with partially overlap or even unoverlapping label sets and further improve DAMTL performance.

Also, Fig.4 shows the performance comparison of the mean and standard

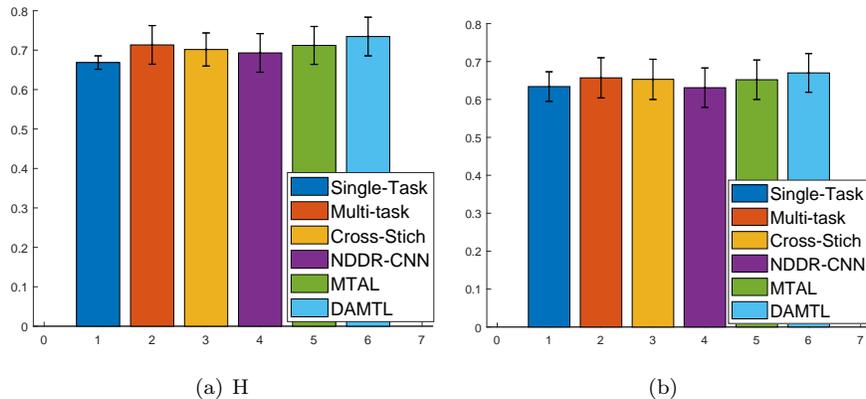


Figure 4: Performance comparison of various methods on mean and mean square error. Sub-Fig.(a) shows that the performance of various methods in the mean and mean square error of the scenario where the label sets between tasks are partially overlapped is various. Sub-Fig.(b) shows that the performance of various methods in the mean value and mean square error of the scenario where the label sets between tasks are completely unoverlapping is various.

deviation of the classification accuracy of various methods in the two categorical experiments. We observe that the overall performance of the DAMTL method is better than that other methods. The above experimental results are consistent with our theoretical analysis.

4.5. Time cost comparison

The result is shown in Fig.5 (a), we observe that when the label sets among tasks partially overlap, although the single-task method takes the shortest time, it does not use shared information, so the accuracy is the lowest. The NDDR-CNN method takes the longest time, but the accuracy is the second lowest (even worse than the multi-task benchmark). This indicates that the differences among tasks lead to adversarial interference in the learning process of this method. In addition, the multitasking, Cross-Stich, NDDR-CNN and MTAL methods are comparable in time cost, and our DAMTL method has the highest accuracy. As shown in Fig.5 (b), we observe that when the label sets among tasks do not overlap, the single-task method takes the shortest time, but its accuracy is the lowest. MTAL uses the convolution kernel sharing technology,

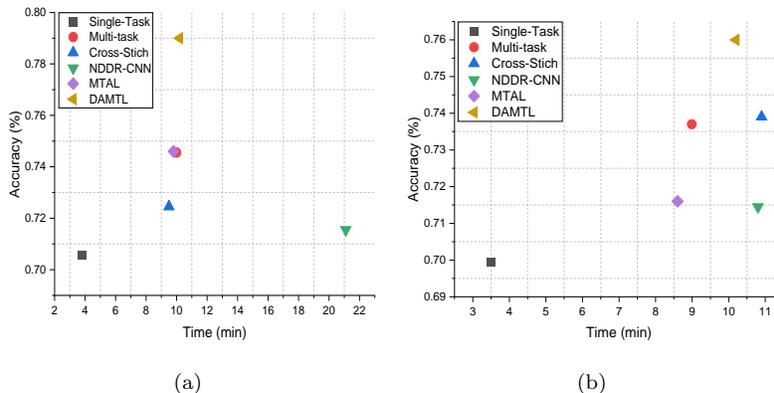


Figure 5: Comparison of the time cost (in minutes) for MTL on partially overlapping and unoverlapping label sets.

Table 6: Results of ablation study when task label sets partially overlap. $DAMTL_{Fa}$ represents the DAMTL network without feature alignment.

Methods	PO-1		PO-2		PO-3	
	Caltech-101	Caltech-256	Amazon	Webcam	Dlsr	Product
DAMTL	0.80 ± 0.050	0.54 ± 0.059	0.80 ± 0.044	0.70 ± 0.048	0.84 ± 0.042	0.74 ± 0.052
$DAMTL_{Fa}$	0.75 ± 0.053	0.50 ± 0.057	0.79 ± 0.046	0.65 ± 0.054	0.63 ± 0.049	0.76 ± 0.052

thus spending the second shortest time. In conclusion, our method has relatively low time cost and the highest accuracy in scenarios where the task label sets are partially overlapped or even unoverlapped.

4.6. Ablation Study

In this section, we demonstrate that the addition of Eq.(3) to DAMTL through ablation analysis is very important to achieve state-of-the-art performance. We first study the performance of multiple tasks with overlapping partial label sets without using Eq.(3). The performance of DAMTL network on Caltech-101, Caltech-256, Amazon, Webcam, Dlsr, and Product is shown on Table 6. Obviously, when the rest of the data set except Product, the performance of DAMTL is reduced. Then, we delete the Eq.(3) part for unoverlapping label sets between tasks. The performance of DAMTL network on Art, Real World, Caltech-101, Webcam, Amazon, and T-ImagNet is shown on Table 7.

Table 7: Results of ablation study when task label sets do not overlap. $DAMTL_{Fa}$ represents the DAMTL network without feature alignment.

Methods	NO-1		NO-2		NO-3	
	Art	Real World	Caltech-101	Webcam	Amazon	T-ImagNet
DAMTL	0.67 ± 0.050	0.6 ± 0.047	0.80 ± 0.050	0.72 ± 0.044	0.81 ± 0.047	0.54 ± 0.065
$DAMTL_{Fa}$	0.61 ± 0.048	0.48 ± 0.048	0.75 ± 0.044	0.64 ± 0.050	0.79 ± 0.041	0.50 ± 0.045

The performance of DAMT on all data sets is significantly degraded. Finally, we use Eq.(3) DAMTL network through ablation analysis which can effectively improve performance.

4.7. Model convergence analysis

In this section, we analyze the convergence of our proposed method on two benchmarks, i.e., Caltech-101 and Caltech-256 and the values of the objective function (5) with respect to iterations on the two datasets are shown in Fig.6 (a) and (b) respectively. From Fig.6 (a), we can see that the loss curves on the data sets Caltech-101 and Caltech-256 tend to converge after about 50 iterations. In Fig.6 (b), the loss curves on the data sets Art and Real-world stabilize after about 70 and 50 iterations, respectively.

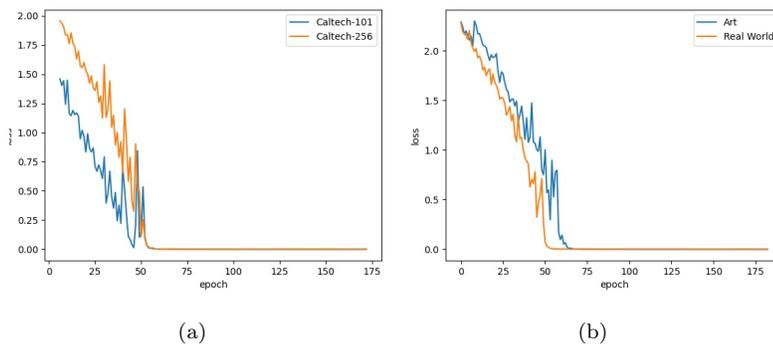


Figure 6: Sub-Fig.(a) shows the convergence curve of DAMTL when the label sets among tasks partially overlapping; sub-fig. Sub-Fig.(b) shows the convergence curve of DAMTL when the label sets among network tasks do not overlap.

5. Conclusion

In this work, we provide a deep multi-task learning framework DAMTL, which is used to deal with multi-tasks with partial or unoverlapping label sets among tasks. Compared with the previous MTL method, DAMTL leverages big auxiliary tasks to jointly learn multiple tasks with partially overlapping or unoverlapping label sets. In addition, the auxiliary strategies in DAMTL can be flexibly embedded in other deep multi-task learning frameworks or transfer learning frameworks. In order to evaluate the performance of DAMTL, we conduct experiments on twelve public datasets and compared state-of-the-art MTL methods. Experimental results show that the DAMTL framework has significant advantages. In summary, our work can enrich MTL research to a certain extent from two aspects: 1) a novel adaptive MT learning mechanism is used to deal with multiple tasks when the label sets are partially overlapped or even unoverlapped. 2) A new knowledge extraction strategy that uses a set of soft masking matrices to adaptively prune the hidden neurons in the auxiliary task network to extract specific knowledge that assist the current task learn to form a corresponding network for each task. However, in this work, we haven't solved the interpretable problems in MTL, and the following work will focus on such problems.

Acknowledgments

This work is supported by NSFC under Grant No. 61672281, and the Key Program of NSFC under Grant No. 61732006. The authors would like to thank Dr. Meng Cao and Yingyu Zhong for their generous help and beneficial discussions.

References

References

- [1] Y. Chen, D. Zhao, L. Lv, Q. Zhang, Multi-task learning for dangerous object detection in autonomous driving, *Information Sciences* 432 (2018)

559–571.

- [2] Q. Xu, Y. Zeng, W. Tang, W. Peng, T. Xia, Z. Li, F. Teng, W. Li, J. Guo, Multi-task joint learning model for segmenting and classifying tongue images using a deep neural network, *IEEE journal of biomedical and health informatics* 24 (9) (2020) 2481–2489.
- [3] Y. Jiang, Z. Deng, K.-S. Choi, F.-L. Chung, S. Wang, A novel multi-task task fuzzy classifier and its enhanced version for labeling-risk-aware multi-task classification, *Information Sciences* 357 (2016) 39–60.
- [4] M. J.Zhong, S. P.Swietojanski, J.Monteiro, J.Trmal, Y.Bengio, Multi-task self-supervised learning for robust speech recognition, in: *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 6989–6993. doi:10.1109/ICASSP40776.2020.9053569.
- [5] M. Ravanelli, J. Zhong, S. Pascual, P. Swietojanski, J. Monteiro, J. Trmal, Y. Bengio, Multi-task self-supervised learning for robust speech recognition, in: *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, IEEE, 2020, pp. 6989–6993.
- [6] Z. Zhao, J. Qin, Z. Gou, Y. Zhang, Y. Yang, Multi-task learning models for predicting active compounds, *Journal of Biomedical Informatics* 108 (2020) 103484.
- [7] S. Ruder, J. Bingel, I. Augenstein, A. Søgaard, Sluice networks: Learning what to share between loosely related tasks, *stat* 1050 (2017) 23.
- [8] J. Bingel, A. Søgaard, Identifying beneficial task relations for multi-task learning in deep neural networks, *arXiv e-prints* (2017) arXiv:1702.08303arXiv:1702.08303.
- [9] K. Hashimoto, C. Xiong, Y. Tsuruoka, R. Socher, A Joint Many-Task Model: Growing a Neural Network for Multiple NLP Tasks, *arXiv e-prints* (2016) arXiv:1611.01587arXiv:1611.01587.

- [10] M. Long, Z. Cao, J. Wang, P. S. Yu, Learning Multiple Tasks with Multilinear Relationship Networks, arXiv e-prints (2015) arXiv:1506.02117arXiv:1506.02117.
- [11] J. Ma, Z. Zhao, X. Yi, J. Chen, L. Hong, E. H. Chi, Modeling task relationships in multi-task learning with multi-gate mixture-of-experts, in: Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, 2018, pp. 1930–1939.
- [12] Q. Feng, S. Chen, Learning Twofold Heterogeneous Multi-Task by Sharing Similar Convolution Kernel Pairs, arXiv e-prints (2021) arXiv:2101.12431arXiv:2101.12431.
- [13] S. Wu, H. R. Zhang, C. Ré, Understanding and Improving Information Transfer in Multi-Task Learning, arXiv e-prints (2020) arXiv:2005.00944arXiv:2005.00944.
- [14] T. Evgeniou, M. Pontil, Regularized multi-task learning, in: Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining, 2004, pp. 109–117.
- [15] J. Honorio, D. Samaras, Multi-task learning of gaussian graphical models, in: ICML, 2010.
- [16] Q. Liu, X. Li, Z. He, N. Fan, D. Yuan, W. Liu, Y. Liang, Multi-task driven feature models for thermal infrared tracking, in: AAAI, 2020.
- [17] J. Wang, S. Zhang, Y. Wang, Z. Zhu, Learning efficient multi-task stereo matching network with richer feature information, Neurocomputing 421 (2021) 151–160.
- [18] P. Guo, C. Deng, L. Xu, X. Huang, Y. Zhang, Deep Multi-Task Augmented Feature Learning via Hierarchical Graph Neural Network, arXiv e-prints (2020) arXiv:2002.04813arXiv:2002.04813.

- [19] S. Vandenhende, S. Georgoulis, W. Van Gansbeke, M. Proesmans, D. Dai, L. Van Gool, Multi-Task Learning for Dense Prediction Tasks: A Survey, arXiv e-prints (2020) arXiv:2004.13379arXiv:2004.13379.
- [20] Z. Shen, C. Cui, J. Huang, J. Zong, M. Chen, Y. Yin, Deep adaptive feature aggregation in multi-task convolutional neural networks, in: Proceedings of the 29th ACM International Conference on Information & Knowledge Management, 2020, pp. 2213–2216.
- [21] S. Yadav, J. Chauhan, J. P. Sain, K. Thirunarayan, J. Schumm, Identifying depressive symptoms from tweets: Figurative language enabled multitask learning framework, in: Proceedings of the 28th International Conference on Computational Linguistics, 2020.
- [22] X. Sun, R. Panda, R. Feris, K. Saenko, AdaShare: Learning What To Share For Efficient Deep Multi-Task Learning, arXiv e-prints (2019) arXiv:1911.12423arXiv:1911.12423.
- [23] T. Sun, Y. Shao, X. Li, P. Liu, H. Yan, X. Qiu, X. Huang, Learning sparse sharing architectures for multiple tasks, in: Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 34, 2020, pp. 8936–8943.
- [24] S. Verboven, M. Hafeez Chaudhary, J. Berrevoets, W. Verbeke, HydaLearn: Highly Dynamic Task Weighting for Multi-task Learning with Auxiliary Tasks, arXiv e-prints (2020) arXiv:2008.11643arXiv:2008.11643.
- [25] V. Sanh, T. Wolf, S. Ruder, A hierarchical multi-task approach for learning embeddings from semantic tasks, in: Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 33, 2019, pp. 6949–6956.
- [26] J. Baxter, A bayesian/information theoretic model of learning to learn via multiple task sampling, Machine learning 28 (1) (1997) 7–39.
- [27] S. Ruder, J. Bingel, I. Augenstein, A. Søgaard, Latent multi-task architecture learning, in: Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 33, 2019, pp. 4822–4829.

- [28] G. Strezoski, N. v. Noord, M. Worring, Many task learning with task routing, in: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV), 2019.
- [29] C. Fernando, D. Banarse, C. Blundell, Y. Zwols, D. Ha, A. A. Rusu, A. Pritzel, D. Wierstra, Pathnet: Evolution channels gradient descent in super neural networks, CoRR, abs/1701.0873,2017.
- [30] G. Pironkov, S. U. Wood, S. Dupont, Hybrid-task learning for robust automatic speech recognition, Computer Speech & Language (2020) 101103.
- [31] P. Cao, X. Shan, D. Zhao, M. Huang, O. Zaiane, Sparse shared structure based multi-task learning for mri based cognitive performance prediction of alzheimer’s disease, Pattern Recognition 72 (2017) 219–235.
- [32] S. Lee, Y. Son, Multitask Learning with Single Gradient Step Update for Task Balancing, arXiv e-prints (2020) arXiv:2005.09910arXiv:2005.09910.
- [33] A. Søgaard, Y. Goldberg, Deep multi-task learning with low level tasks supervised at lower layers, in: ACL, 2016.
- [34] J. Fan, T. Zhao, Z. Kuang, Y. Zheng, J. Zhang, J. Yu, J. Peng, Hdmtl: Hierarchical deep multi-task learning for large-scale visual recognition, IEEE Transactions on Image Processing 26 (4) (2017) 1923–1938. doi:10.1109/TIP.2017.2667405.
- [35] W. Xue, Weighted feature-task-aware regularization learner for multitask learning, Pattern Analysis and Applications 23 (1) (2020) 253–263.
- [36] J. Zhang, J. Miao, K. Zhao, Y. Tian, Multi-task feature selection with sparse regularization to extract common and task-specific features, Neurocomputing 340 (2019) 76–89.
- [37] W. Shao, Y. Peng, C. Zu, M. Wang, D. Zhang, A. D. N. Initiative, et al., Hypergraph based multi-task feature selection for multimodal classification

of alzheimer’s disease, *Computerized Medical Imaging and Graphics* 80 (2020) 101663.

- [38] L. Li, X. Pan, H. Yang, Z. Liu, Y. He, Z. Li, Y. Fan, Z. Cao, L. Zhang, Multi-task deep learning for fine-grained classification and grading in breast cancer histopathological images, *Multimedia* (2020) 14509–15428.
- [39] Z. Zheng, Y. Wang, Q. Dai, H. Zheng, D. Wang, Metadata-driven task relation discovery for multi-task learning., in: *IJCAI*, 2019, pp. 4426–4432.
- [40] C. Yan, J. Xu, J. Xie, C. Cai, H. Lu, Prior-aware cnn with multi-task learning for colon images analysis, in: *2020 IEEE 17th International Symposium on Biomedical Imaging (ISBI)*, 2020, pp. 254–257. doi:10.1109/ISBI45749.2020.9098703.
- [41] I. Misra, A. Shrivastava, A. Gupta, M. Hebert, Cross-stitch networks for multi-task learning, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 3994–4003.
- [42] R. Duan, N. F. Chen, Unsupervised feature adaptation using adversarial multi-task training for automatic evaluation of children’s speech, *Proc. Interspeech* (2020) 3037–3041.
- [43] I. Augenstein, A. Søgaard, Multi-Task Learning of Keyphrase Boundary Classification, *arXiv e-prints* (2017) arXiv:1704.00514arXiv:1704.00514.
- [44] P. Rai, H. Daumé III, Infinite predictor subspace models for multitask learning, in: *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, 2010, pp. 613–620.
- [45] L. Zhou, Z. Cui, C. Xu, Z. Zhang, C. Wang, T. Zhang, J. Yang, Pattern-structure diffusion for multi-task learning, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.

- [46] Y. Wang, X. Luo, L. Ding, S. Fu, S. Hu, Multi-task non-negative matrix factorization for visual object tracking, *Pattern Analysis and Applications* 2020. doi:10.1007/s10044-019-00812-4.
- [47] J. Jeong, C. Jun, Sparse tensor decomposition for multi-task interaction selection, in: 2019 IEEE International Conference on Big Knowledge (ICBK), 2019, pp. 105–114. doi:10.1109/ICBK.2019.00022.
- [48] F. Huang, Y. Qiu, Q. Li, S. Liu, F. Ni, Predicting drug-disease associations via multi-task learning based on collective matrix factorization, *Frontiers in Bioengineering and Biotechnology* 8 (2020) 218.
- [49] Y. Zhang, Y. Zhang, W. Wang, Deep Multi-Task Learning via Generalized Tensor Trace Norm, *arXiv e-prints* (2020) arXiv:2002.04799arXiv:2002.04799.
- [50] Z. Chen, H. Lei, Y. Zhao, Z. Huang, X. Xiao, Y. Lei, E.-L. Tan, B. Lei, Template-oriented multi-task sparse low-rank learning for parkinson’s diseases diagnosis, in: *International Workshop on PRedictive Intelligence In MEdicine*, Springer, 2020, pp. 178–187.
- [51] X. Wu, X. Zhang, Y. Cen, Multi-task joint sparse and low-rank representation target detection for hyperspectral image, *IEEE Geoscience and Remote Sensing Letters* 16 (11) (2019) 1756–1760. doi:10.1109/LGRS.2019.2908196.
- [52] M. Long, J. Wang, G. Ding, J. Sun, P. S. Yu, Transfer joint matching for unsupervised domain adaptation, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 1410–1417.
- [53] K. Simonyan, A. Zisserman, Very Deep Convolutional Networks for Large-Scale Image Recognition, *arXiv e-prints* (2014) arXiv:1409.1556arXiv:1409.1556.
- [54] I. Misra, A. Shrivastava, A. Gupta, M. Hebert, Cross-stitch networks for multi-task learning (2016) 3994–4003.

- [55] Y. Gao, J. Ma, M. Zhao, W. Liu, A. L. Yuille, Nddr-cnn: Layerwise feature fusing in multi-task cnns by neural discriminative dimensionality reduction (2019) 3205–3214.