

Unsupervised Domain Adaptation with Adversarial Residual Transform Networks

Guanyu Cai, Yuqin Wang, Lianghua He, and Mengchu Zhou, *Fellow, IEEE*

Abstract—Domain adaptation is widely used in learning problems lacking labels. Recent studies show that deep adversarial domain adaptation models can make remarkable improvements in performance, which include symmetric and asymmetric architectures. However, the former has poor generalization ability whereas the latter is very hard to train. In this paper, we propose a novel adversarial domain adaptation method named Adversarial Residual Transform Networks (ARTNs) to improve the generalization ability, which directly transforms the source features into the space of target features. In this model, residual connections are used to share features and adversarial loss is reconstructed, thus making the model more generalized and easier to train. Moreover, a special regularization term is added to the loss function to alleviate a vanishing gradient problem, which enables its training process stable. A series of experiments based on Amazon review dataset, digits datasets and Office-31 image datasets are conducted to show that the proposed ARTN can be comparable with the methods of the state-of-the-art.

Index Terms—Adversarial neural networks, unsupervised domain adaptation, residual connections, transfer learning.

I. INTRODUCTION

DEEP neural networks trained on large-scale labeled datasets could achieve excellent performance across varieties of tasks, such as sentiment analysis [1], [2], image classification [3]–[5] and semantic segmentation [6]. Yet they usually fail to generalize well on novel tasks because the transferability of features decreases as the distance between the base and target tasks increases [7]. A convincing explanation is that there exists a domain shift between training data and testing one [8], [9]. To alleviate the negative effect caused by a domain shift, domain adaptation (DA) is proposed to utilize labeled data from a source domain to generalize models generalize well on a target domain [10], [11].

Domain adaptation, which is a field belonging to transfer learning, has long been utilized to make it possible to exploit

the knowledge learned in one specific domain to effectively improve the performance in a related but different domain. Earlier methods of DA aim to learn domain-invariant feature representations from data by jointly minimizing a distance metric that actually measures the adaptability between a pair of source and target domains, such as Transfer Component Analysis [12], Geodesic Flow Kernel [13], and Transfer Kernel Learning [14]. In order to learn transferable features well, researchers apply deep neural networks to DA models [15]–[17]. A feature extractor neural network is trained by reducing “distance” between distributions of two different domains, on the assumption that the classifier trained by source data also works well in a target domain. In this kind of methods, Maximum Mean Discrepancy (MMD) loss is widely used for mapping different distributions [18]. For example, Deep Adaptation Networks (DAN) [19], Joint Adaptation Networks [20] and Residual Transfer Networks [21] apply MMD loss to several layers whereas Large Scale Detection through Adaptation [22] adds a domain adaptation layer that is updated based on MMD loss.

Recently, the idea of Generative Adversarial Networks (GANs) [23], [24] has been widely applied to DA. The methods of using GANs [25], [26] to transform source images to target ones are proposed and their classifiers are trained with the generated target images. However, when distributions of source and target domains are totally different, adversarial training has poor performance because of a gradient vanishing phenomenon. Alternative methods train GANs on features of source and target domains. Their generator is acted as a feature extractor, and discriminator as a domain classifier. There are symmetric and asymmetric adaptation architectures in adversarial domain adaptation, which can effectively adapt source and target distributions. The former’s features in the source and target domains are generated from the same network [27], [28], while the latter’s from different networks [29]. It is well-recognized that the former is poor at generalization whereas the latter is difficult to train.

To solve the above problems, in this work, we propose a novel feature-shared model for adversarial domain adaptation, which achieves the flexibility of asymmetric architecture and can be easily trained. In the proposed framework as shown in Fig. 1, a weight-shared feature extractor distills features from different domains, and a feature-shared transform network maps features from the source domain to the space of target features. Adversarial learning is completed with the losses from the label and domain classifiers. Note that we design residual connections between the extractor and the network to ease the learning of distribution mapping by sharing features.

Manuscript received July 16, 2018; revised December 11, 2018, and May 19, 2019; accepted August 2, 2019. This work was supported by the National Natural Science Foundation of China under Grant 61772369, 61773166, 61771144 and 61871004, Joint Funds of the National Science Foundation of China (U18092006), Shanghai Municipal Science and Technology Committee of Shanghai Outstanding Academic Leaders Plan (19XD1434000), Projects of International Cooperation of Shanghai Municipal Science and Technology Committee (19490712800). This paper is partially supported by the National Key R&D Program 2018YFB1004701, by the Fundamental Research Funds for the Central Universities. (Corresponding author: Lianghua He).

G. Cai, Y. Wang and L. He are with the Department of Computer Science and Technology, Tongji University, Shanghai 201804, China (email: caiguanyu@tongji.edu.cn; wangyuqin@tongji.edu.cn; he-lianghua@tongji.edu.cn).

M. Zhou is with the Department of Electrical and Computer Engineering, New Jersey Institute of Technology, Newark, NJ 07102 USA, and also with the Institute of Systems Engineering, Macau University of Science and Technology, Macau 999078, China (e-mail: zhou@njit.edu).

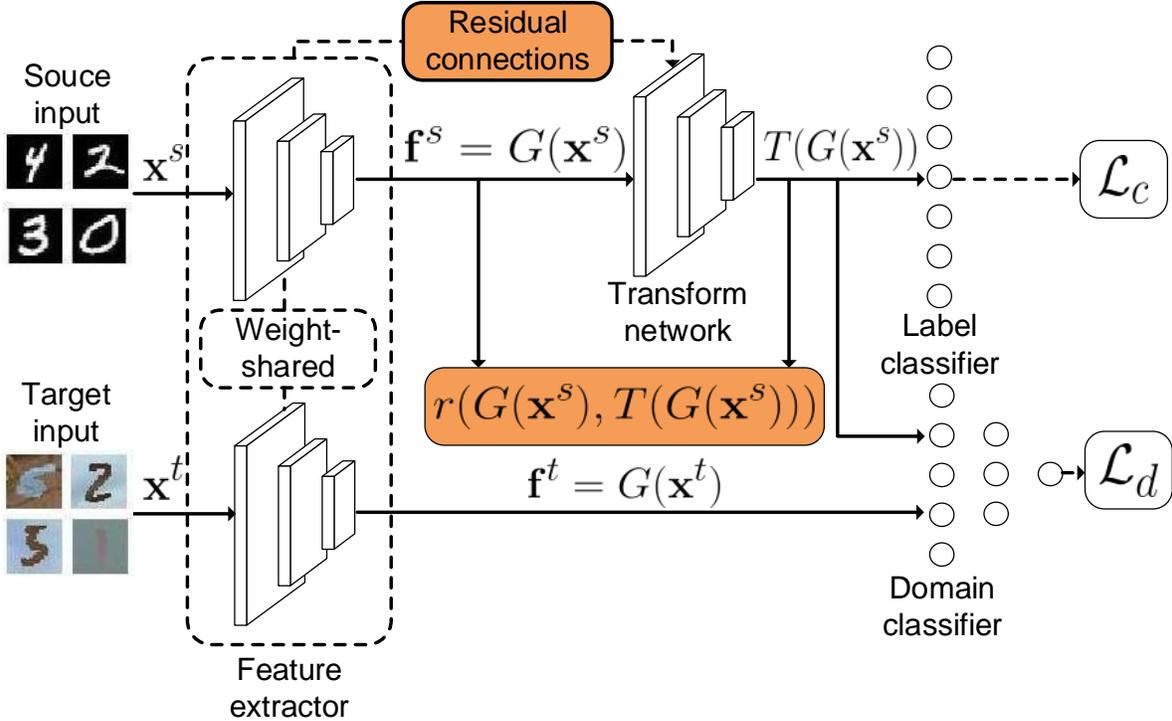


Fig. 1. The architecture of the proposed model. First feature extractor G distills the feature representations of source and target samples. Then, transform network T projects source features $f^s = G(x^s)$ to the space of target features $f^t = G(x^t)$. Finally, the label classifier C is trained with the fake target features $T(G(x^s))$ and predicts the labels of target features $G(x^t)$ during a test period. In addition, domain classifier D is trained to distinguish fake target features $T(G(x^s))$ and real target features $G(x^t)$, which can minimize the discrepancy between source and target domains through adversarial optimization. The regularization term $r(G(x^s), T(G(x^s)))$ measures the distance between $f^s = G(x^s)$ and $T(G(x^s))$.

In addition, in order to avoid getting stuck into local minima, we construct a regularization term to ensure that the model at least knows a vague, if not exact, direction to match different distributions and overcome a gradient vanishing problem. The main contributions of this work are as follows:

- 1) A novel adversarial model that learns a non-linear mapping from a source domain to a target one is proposed. By using the features generated from the source domain, this model owns high generalization ability in the target domain.
- 2) During training, concise regularization that ensures the model to select the shortest path from all the transfer paths is constructed, thereby helping stabilize an adversarial optimization process.
- 3) This work extensively evaluates the proposed method on several standard benchmarks. The results demonstrate that our model outperforms the state-of-the-art methods in accuracy. Notably, our model can maintain excellent generalization and anti-noise abilities.

Section II reviews some related work on unsupervised domain adaptation. In Section III, the proposed method is described. Several experiments are reported in Section IV. Section V concludes this paper.

II. RELATED WORK

DA has been extensively studied in recent years. Several studies give theoretical analyses of error bound when training and testing data are drawn from different distributions [8], [9]. This means that it is feasible to utilize the knowledge across domains. Moreover, the most important problem in DA is how to reduce the discrepancy between the source and target domains. Therefore, many methods modify a classifier to match different distributions [30]–[33]. However, the performances of these shallow methods are limited. Recently, because deep neural networks can learn feature representations including information identification, and these features are also transferable [7], they are widely used in DA methods. Most of them focus on how to measure the distance between different domains and how to design a network structure, and achieve remarkable performance.

A. Traditional Domain Adaptation

A general idea is to reweight instances in the source domain, where instances similar to the target distribution are considered with more importance. This kind of method has proven effective in adaptation for the differences between the source and target distributions. In detail, the calculated weight is taken as the loss coefficient of each source instance. Some

methods reweight instances with direct importance estimation algorithms, such as [34]–[37] and [32]. Other methods reweight instances with noisy labels to reduce the marginal and conditional shifts [38], [39].

Another DA idea is to explicitly make source and target distributions similar. Many statistic characteristics are chosen to be the metrics to align subspaces of different distributions. MMD, which measures expectation difference in reproducing kernel Hilbert space (RKHS) of source and target distributions, is widely used in many methods [12], [14], [30], [33]. In order to align more complex statistic characteristics, the studies [40]–[42] calculate statistic moments of different order to match different distributions, which are easy to implement with low computational complexity. Instead of exploiting statistic characteristics, some methods [13], [43]–[45] utilize manifold learning methods to transform a source distribution to a target one. In these methods, feature spaces are refined into low-dimensional spaces and feature distortion is thus avoided.

B. Deep Domain Adaptation

Recent, development of deep neural networks promotes deep domain adaptation. In [46], the experimental results demonstrate that features of deep neural networks instead of hand-crafted ones alleviate the negative influences of a domain shift even without any adaptations. However, a main limitation of using pre-trained deep features is that they severely restrict the range of application. Later, a number of methods combine statistic characteristics with deep neural networks to a unified framework, which greatly improve performance on different tasks. In [19]–[22], MMD is embedded in deep convolutional neural networks. In [42] and [47], high order moments are utilized to align feature spaces of source and target domains.

Instead of designing fancy regularizers, some methods design special architectures to minimize the discrepancy between source and target domains. In [48], a Siamese architecture is introduced to adapt pairs of source and target instances. In [49], [50], auto-encoders are suggested to learn the transformation from a source domain to a target one.

Some methods choose adversarial loss to learn manifest invariant factors underlying different populations with a domain discriminator subnetwork. In these models, deep features are learned to confuse a domain discriminator such that they could capture the most discriminative information about classification instead of characteristics of domains. Domain-Adversarial Neural Networks (DANN) [28] consist of a symmetric feature generator, label discriminator, and domain discriminator. The whole model can be directly optimized via a gradient reversal algorithm. Deep Reconstruction-Classification Networks [51] also take adversarial learning and add a reconstruction step for target images. Adversarial Discriminative Domain Adaptation (ADDA) [29] uses an asymmetric feature generator that is trained alternatively with a domain classifier.

The above-mentioned domain adversarial networks fall into two categories. Some methods, such as domain confusion networks [27] and DAN [28], share weights between source and target feature extractors. They use the same network to learn representations from different inputs, which learns

a symmetric transformation to utilize the transferability of features generated from deep neural networks and reduces parameters in the model. Other methods construct two networks for source and target domains, respectively [25], [26], [29]. They can learn an asymmetric transformation, allowing networks to learn parameters for each domain individually. In theory, asymmetric transformation can lead to more effective adaptation [52].

Adversarial domain adaptation has also been explored in generative adversarial networks (GANs). Coupled Generative Adversarial Networks (CoGANs) [25] apply GANs to DA. Two GANs are trained to generate source and target images, respectively. Pixel-Level Domain Adaptation [26] uses a conditional GAN model to synthesize target images to facilitate training a label classifier. Methods based on GANs can improve the performance of digits datasets, but their downside is a difficult training process as caused by gradient vanishing when facing more natural image datasets according to [53]. In this work, we focus on learning the mapping of different feature spaces instead of synthesizing target images, and propose a discriminative model aiming to adapt distinct domains.

III. ADVERSARIAL RESIDUAL TRANSFORM NETWORKS

In this section, we describe the details of our proposed model. We first define unsupervised domain adaptation and preliminary domain adversarial networks, and then demonstrate the key innovations of our model, which can well handle the problems encountered by the previous models. At last, we give a complete algorithm of matching the distributions of target and source domains using our model.

A. Definitions

When it comes to a machine learning task, a domain D corresponds to four parts: feature space \mathcal{X} , label space \mathcal{Y} , marginal probability distribution $P(\mathbf{X})$ and conditional probability distribution $P(\mathbf{X}|\mathbf{Y})$, where $\mathbf{X} \in \mathcal{X}$, $\mathbf{Y} \in \mathcal{Y}$. Subscript s and t are used to denote the source and target distributions. In a traditional machine learning task, training data are drawn from source domain D_s and testing data are drawn from target domain D_t , where their marginal and conditional probability distributions are the same ($P_s(\mathbf{X}^s) = P_t(\mathbf{X}^t)$, $P_s(\mathbf{X}^s|\mathbf{Y}^s) = P_t(\mathbf{X}^t|\mathbf{Y}^t)$). Thus, models trained in the source domain are feasible to the target one. However, in unsupervised DA, these assumptions are not valid, which leads us to a more difficult problem as follows.

Given a source domain as $D_s\{\mathbf{x}_i^s, \mathbf{y}_i^s\}_{i=1}^{n_s}$, where n_s is the number of source domain samples, \mathbf{x}_i^s is the i th instance in D_s and \mathbf{y}_i^s is the label of \mathbf{x}_i^s . Similarly, a target domain is denoted as $D_t\{\mathbf{x}_i^t\}_{i=1}^{n_t}$, where n_t is the number of target domain samples, \mathbf{x}_i^t is the i th instance in D_t . The source and target domains are drawn from distribution $P_s(\mathbf{X}^s)$ and $P_t(\mathbf{X}^t)$, respectively, which are different. In most cases, conditional probability distributions are also different ($P_s(\mathbf{X}^s|\mathbf{Y}^s) \neq P_t(\mathbf{X}^t|\mathbf{Y}^t)$). The goal is to learn a feature extractor G_t and a classifier C_t for D_t . G_t distills feature representations $\mathbf{f}^t = G_t(\mathbf{x}^t)$ from target samples, and C_t correctly predicts

the labels of target samples receiving $\mathbf{f}^t = G_t(\mathbf{x}^t)$. Because of lacking annotations in D_t , DA learns G_s and C_s with samples from D_s , and tries to adapt them to be useful in D_t .

B. Adversarial Domain Adaptation

To solve an unsupervised DA problem, a number of methods have been proposed. Among the most effective ones is adversarial domain adaptation. This work aims to modify this kind of framework to improve its generalization and anti-noise ability. In DA problems, it is difficult to train G_t and C_t for a target domain without labels. However, because there exists a correlation between the source and target domains, it is common to utilize G_s and C_s to predict labels of target samples. In order to make G_s and C_s valid in a target domain, adversarial DA models are usually used to train a feature extractor G , a label classifier C and a domain classifier D for both domains. In details, these models set $G = G_t = G_s$, and $C = C_t = C_s$, which means the feature extractor and label classifier are used for both source and target domains. Specifically, D also receives feature representations from G and is trained to minimize the discrepancy between source and target feature distributions: $G(\mathbf{x}^s)$ and $G(\mathbf{x}^t)$. An adversarial training procedure is a minimax two-player game [23]. One player D learns to distinguish whether features are from a source or target domain, whereas the other G tries to generate domain-invariant features. They have contradictory optimization objectives, and their objectives are optimized alternately in this minmax game. To train the whole network in an end-to-end way, DANN [28] adopts the following loss function:

$$\mathcal{L}(\theta_d, \theta_g, \theta_c) = \frac{1}{n_s} \sum_{\mathbf{x}_i \in D_s} \mathcal{L}_c(C(G(\mathbf{x}_i)), y_i) - \frac{\lambda}{n} \sum_{\mathbf{x}_i \in D_s \cup D_t} \mathcal{L}_d(D(G(\mathbf{x}_i)), d_i) \quad (1)$$

where $n = n_s + n_t$, and \mathcal{L}_c and \mathcal{L}_d denote losses of C and D , respectively. λ is a trade-off parameter between \mathcal{L}_c and \mathcal{L}_d . θ_d , θ_g and θ_c are the parameters of D , G and C , respectively. y_i and d_i denote the class and domain labels of images. After convergence, optimal parameters $\hat{\theta}_d$, $\hat{\theta}_c$ and $\hat{\theta}_g$ can deliver a saddle point given as:

$$\hat{\theta}_d = \arg \min_{\theta_d} \mathcal{L}_d(\theta_d, \theta_g) \quad (2)$$

$$\hat{\theta}_c = \arg \min_{\theta_c} \mathcal{L}_c(\theta_g, \theta_c) \quad (3)$$

$$\hat{\theta}_g = \arg \min_{\theta_g} \mathcal{L}(\theta_d, \theta_g, \theta_c) \quad (4)$$

In such framework, a DA model can be trained in an end-to-end way. The intuitive idea behind this model is that with the minmax two-player game going, D and G strengthen each other. When the training procedure converges, the features of different domains generated from G are very hard to distinguish by D . In this condition, features are domain-invariant and the feature distributions of different domains are adapted.

Theoretically, adversarial domain adaptation is based on \mathcal{H} -divergence in [8], [9]. However, it is almost impossible to apply \mathcal{H} -divergence in real-world algorithms. Because it is defined in a binary classification problem and requires a global search in all hypothesis. In [9], an approximate algorithm is given. Given a generalization error ϵ of discriminating between source and target instances, \mathcal{H} -divergence is computed as:

$$\hat{d}_{\mathcal{A}} = 2(1 - 2\epsilon) \quad (5)$$

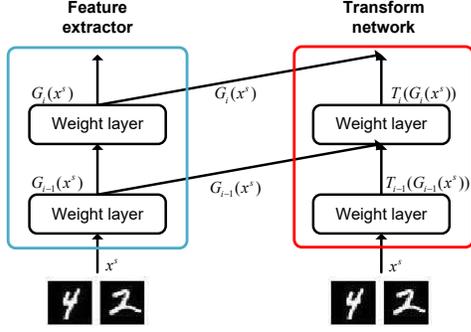
The value of $\hat{d}_{\mathcal{A}}$ is called the *Proxy A-distance* (PAD). In adversarial domain adaptation, the domain classifier composed of neural networks is trained to directly decrease PAD. Cooperating with the domain classifier, the feature extractor learns domain-invariant features from different domains, implying that discrepancy between source and target distributions is decreased. Several models based on this kind of architecture have achieved the top performances in different visual tasks [28], [51].

C. Residual Connections

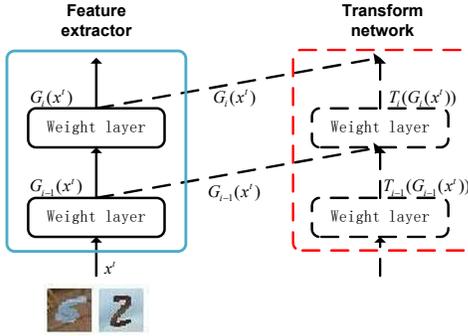
The proposed method does not rely on only feature extractor G to map different distributions. Instead, we construct an adversarial residual transform network (ARTN) T to project source features $\mathbf{f}^s = G(\mathbf{x}^s)$ to the space of target features. The network is trained to generate fake target features $T(G(\mathbf{x}^s))$, which are in the same distribution as real target features $\mathbf{f}^t = G(\mathbf{x}^t)$. Then, we use the fake target features $T(G(\mathbf{x}^s))$ and corresponding labels \mathbf{y}^s to train a classifier C for the target domain. After training, the labels of target samples are predicted by C .

In previous unsupervised DA methods, the weights of feature extractor G for source and target domains are shared [19]–[22]. However, regarding matching different distributions, the generalization ability of asymmetric transformation is better than that of symmetric one [29]. If the networks are trained to capture domain-invariant information from source features and utilize them to classify target samples, there would be a boost to their generalization ability. However, the asymmetric architecture proposed in [29] is hard to obtain such enhancement and the feature extractor for a target domain is easy to collapse, because there exists no relationship between the feature extractors of source and target domains. In order to make our model learn domain-invariant information and avoid diverging during its training, we propose a transform network that builds connections between source and target domain features.

The detailed architecture of residual connections between a feature extractor and a transform network is shown in Fig. 2. The weight-tied feature extractor G is trained to capture representations from source samples \mathbf{x}^s and target samples \mathbf{x}^t . The transform network stacks a few layers by using the same architecture with a feature extractor. Unlike the symmetric transformation, the proposed network shares features with the feature extractor instead of parameters. Our network is also different from asymmetric transformation where two networks have no relationship. We add residual connections between the feature extractor and the transform network to share features.



(a) Source samples



(b) Target samples

Fig. 2. Residual connections between the feature extractor and transform network. When the inputs are source samples, the feature extractor and transform network are both activated (the solid line represents that this network is activated, whereas the dashed line is opposite) and the features distilled from a feature extractor would be conveyed to a transform network. When the inputs are target samples, only the feature extractor is activated. The features are not shared between the feature extractor and transform network.

Therefore, with a carefully designed architecture, our model is able to alleviate the drawbacks of both symmetric and asymmetric models, which is never seen in the literature to our best knowledge.

Theoretically, by denoting a desired underlying mapping between source and target distributions as M and letting $G(x^t) = M(G(x^s))$, we intend to train transform network T to fit a mapping of $T(G(x^s)) = M(G(x^s)) - G(x^s)$. In details, for an N layer transform network, the i th layer of a transform network where $i < N$ is defined as:

$$T_i(G_i(x^s)) = T_i(T_{i-1}(G_{i-1}(x^s))) + G_i(x^s) \quad (6)$$

where $T_i(\cdot)$ and $G_i(\cdot)$ denote the i th layer of the transform network and feature extractor individually. The inputs of T are original data, and the output is $T_N(G_N(x^s))$.

Please note that residual connections in different methods tend to realize different purposes. Firstly, the effect of residual connections in our paper is different from others. For example, residual connections in ResNet [3] are used to shorten their training process by making gradients flow well. However, residual connections in U-Net [54] are used to enable pre-

cise localization for segmentation. Although these papers all utilize skip connections, they still make novel contributions. In our paper, residual connections are used to capture semantic information from a source domain. This modification has not been seen in other domain adaptation methods and the effect is also different from papers in other research areas. Secondly, the detailed modification is different from skip connections in other papers. Skip connections in ResNet and U-Net are constructed in a single network across different layers whereas ours are constructed in the same layer across different networks.

D. Vanishing Gradient Problem in Adversarial Training

The detailed theoretical derivation and training process of adversarial DA has been described in [28]. Yet there exists a vanishing gradient problem in adversarial training. In this section, its theoretical analysis is presented.

Once we adopt a transform network in adversarial DA and utilize cross entropy loss function for D , the adversarial nets-based DA framework of D and G needs the following minmax optimization:

$$\min_{G, T} \max_D \mathcal{L}(D, G, T) = \mathbb{E}_{\mathbf{x} \sim P_s} [\log D(T(G(\mathbf{x}))) + \mathbb{E}_{\mathbf{x} \sim P_t} [\log(1 - D(G(\mathbf{x}))) \quad (7)$$

where maximizing the loss with respect to D yields a tighter lower bound on the true domain distribution divergence, whereas minimizing the loss with respect to G and T minimizes the distribution divergence in the feature space.

For any given G and T , the optimal D^* is obtained at:

$$D^*(\mathbf{z}) = \frac{P_s(\mathbf{z})}{P_s(\mathbf{z}) + P_t(\mathbf{z})} \quad (8)$$

where \mathbf{z} is the sample in the feature space. For $\mathbf{z} \sim P_s$, $\mathbf{z} = T(G(\mathbf{x}))$, while for $\mathbf{z} \sim P_t$, $\mathbf{z} = G(\mathbf{x})$. Similar to [23], we give the proof as follows.

Proof. For any given G and T , the training criterion for D is to maximize $\mathcal{L}(D, G, T)$:

$$\begin{aligned} \max_D \mathcal{L}(D, G, T) &= \int_{\mathbf{x}} P_s(\mathbf{x}) \log D(T(G(\mathbf{x}))) + \\ &P_t(\mathbf{x}) \log(1 - D(G(\mathbf{x}))) d\mathbf{x} \\ &= \int_{\mathbf{z}} P_s(\mathbf{z}) \log D(\mathbf{z}) + \\ &P_t(\mathbf{z}) \log(1 - D(\mathbf{z})) d\mathbf{z} \end{aligned} \quad (9)$$

We take the partial differential of $\mathcal{L}(D, G, T)$ with respect to D , and achieve its maximum in $[0, 1]$ at $D^*(\mathbf{z}) = \frac{P_s(\mathbf{z})}{P_s(\mathbf{z}) + P_t(\mathbf{z})}$.

Given D^* , the minmax optimization can now be reformulated as:

$$\begin{aligned}
\min_{G,T} \mathcal{L}(D^*, G, T) &= \mathbb{E}_{\mathbf{z} \sim P_s} [\log D^*(\mathbf{z})] + \\
&\mathbb{E}_{\mathbf{z} \sim P_t} [\log(1 - D^*(\mathbf{z}))] \\
&= \mathbb{E}_{\mathbf{z} \sim P_s} \left[\log \frac{P_s(\mathbf{z})}{P_s(\mathbf{z}) + P_t(\mathbf{z})} \right] + \\
&\mathbb{E}_{\mathbf{z} \sim P_t} \left[\log \frac{P_t(\mathbf{z})}{P_s(\mathbf{z}) + P_t(\mathbf{z})} \right] \\
&= \mathbb{E}_{\mathbf{z} \sim P_s} \left[\log \frac{2P_s(\mathbf{z})}{P_s(\mathbf{z}) + P_t(\mathbf{z})} \right] + \\
&\mathbb{E}_{\mathbf{z} \sim P_t} \left[\log \frac{2P_t(\mathbf{z})}{P_s(\mathbf{z}) + P_t(\mathbf{z})} \right] - 2\log 2 \\
&= 2 \cdot JSD(P_s || P_t) - 2\log 2 \quad (10)
\end{aligned}$$

where $JSD(\cdot)$ is the Jensen-Shannon divergence. Since the Jensen-Shannon divergence between two distributions is always non-negative, and zero if they are equal, $\mathcal{L}^* = -2\log 2$ is the global minimum of $\mathcal{L}(D, G, T)$ where the only solution is $P_s = P_t$. In this case, the distributions of source and target domains are the same and the goal of DA is well achieved.

However, in practice, adversarial DA remains remarkably difficult to train. It is sensitive to the initialization of parameters and its training process tends to be unstable, i.e., $\mathcal{L}(D, G, T)$ does not converge. These problems are caused by a vanishing gradient phenomenon. In theory, Jensen-Shannon divergence measures the difference between source and target distributions. By minimizing it, source and target distributions in the feature space tend to be the same. However, if we utilize a gradient descent algorithm to optimize $\mathcal{L}(D, G, T)$ which is the most common algorithm for neural networks, Jensen-Shannon divergence is difficult to converge because its gradient is easily stuck into zero, to be proved next.

According to [53], P_s and P_t can be regarded as two distributions that have support contained in two closed manifolds \mathcal{M} and \mathcal{N} that do not have full dimension, respectively. P_s and P_t are continuous in their respective manifolds, which means that if a set A with measure 0 in \mathcal{M} , then $P_s(A) = 0$. In this case, $JSD(P_s || P_t) = \log 2$ for almost any P_s and P_t . We need to use Lemma 3.1 [53] to prove it:

Lemma 3.1: Let \mathcal{M} and \mathcal{P} be two regular submanifolds that do not perfectly align and do not have full dimension. Let $\mathcal{L} = \mathcal{M} \cap \mathcal{P}$. If \mathcal{M} and \mathcal{P} do not have boundary, then \mathcal{L} is also a manifold. and has strictly lower dimension than both of \mathcal{M} and \mathcal{P} . If they have boundary, \mathcal{L} is a union of at most 4 strictly lower dimensional manifolds. In both cases, \mathcal{L} has measure 0 in both \mathcal{M} and \mathcal{P}

Proof. By Lemma 3.1, we know that $\mathcal{L} = \mathcal{M} \cap \mathcal{P}$ has strictly lower dimension than both \mathcal{M} and \mathcal{P} do, such that

$$P_s(\mathcal{L}) = 0 \text{ and } P_t(\mathcal{L}) = 0.$$

$$\begin{aligned}
2 \cdot JSD(P_s || P_t) &= \int_{\mathbf{z}} P_s(\mathbf{z}) \log \frac{2P_s(\mathbf{z})}{P_s(\mathbf{z}) + P_t(\mathbf{z})} + \\
&P_t(\mathbf{z}) \log \frac{2P_t(\mathbf{z})}{P_s(\mathbf{z}) + P_t(\mathbf{z})} d\mathbf{z} \\
&= \int_{\mathbf{z} \in \mathcal{M} \setminus \mathcal{L}} P_s(\mathbf{z}) \log \frac{2P_s(\mathbf{z})}{P_s(\mathbf{z}) + P_t(\mathbf{z})} + \\
&P_t(\mathbf{z}) \log \frac{2P_t(\mathbf{z})}{P_s(\mathbf{z}) + P_t(\mathbf{z})} d\mathbf{z} \\
&+ \int_{\mathbf{z} \in \mathcal{N} \setminus \mathcal{L}} P_s(\mathbf{z}) \log \frac{2P_s(\mathbf{z})}{P_s(\mathbf{z}) + P_t(\mathbf{z})} + \\
&P_t(\mathbf{z}) \log \frac{2P_t(\mathbf{z})}{P_s(\mathbf{z}) + P_t(\mathbf{z})} d\mathbf{z} \\
&+ \int_{\mathbf{z} \in \mathcal{L}} P_s(\mathbf{z}) \log \frac{2P_s(\mathbf{z})}{P_s(\mathbf{z}) + P_t(\mathbf{z})} + \\
&P_t(\mathbf{z}) \log \frac{2P_t(\mathbf{z})}{P_s(\mathbf{z}) + P_t(\mathbf{z})} d\mathbf{z} \\
&+ \int_{\mathbf{z} \in (\mathcal{M} \cup \mathcal{N})^c} P_s(\mathbf{z}) \log \frac{2P_s(\mathbf{z})}{P_s(\mathbf{z}) + P_t(\mathbf{z})} + \\
&P_t(\mathbf{z}) \log \frac{2P_t(\mathbf{z})}{P_s(\mathbf{z}) + P_t(\mathbf{z})} d\mathbf{z} \quad (11)
\end{aligned}$$

where $(\mathcal{M} \cup \mathcal{N})^c$ is the complement of $(\mathcal{M} \cup \mathcal{N})$. For $\mathbf{z} \in \mathcal{M} \setminus \mathcal{L}$, $P_s(\mathbf{z}) = 1$ and $P_t(\mathbf{z}) = 0$. Similarly, for $\mathbf{z} \in \mathcal{N} \setminus \mathcal{L}$, $P_t(\mathbf{z}) = 1$ and $P_s(\mathbf{z}) = 0$. When $\mathbf{z} \in (\mathcal{M} \cup \mathcal{L})^c$ and $\mathbf{z} \in \mathcal{L}$, $P_s(\mathbf{z})$ and $P_t(\mathbf{z})$ are equal to zero. Therefore, $JSD(P_s || P_t) = \log 2$.

Note that when $JSD(P_s || P_t)$ is a constant, gradients for all the parameters in an adversarial DA network are zeros. Therefore, if a gradient descent algorithm is adopted, the vanishing gradient problem appears and divergence between source and target domains is difficult and sometimes impossible to be minimized.

E. Regularizer Based on Transport Theory

Once we have parametrized G and T , we employ adversarial loss to adapt different distributions. The architecture modification requires us to revise our loss function. Instead of measuring the distance between source features $\mathbf{f}^s = G(\mathbf{x}^s)$ and target features $\mathbf{f}^t = G(\mathbf{x}^t)$ generated from one feature extractor, the proposed model lets domain classifier D discriminate source features $\mathbf{f}^s = T(G(\mathbf{x}^s))$ from the transform network and target features $\mathbf{f}^t = G(\mathbf{x}^t)$ from the feature extractor. Thus, the loss function is modified from (1) into:

$$\begin{aligned}
\mathcal{L}(\theta_d, \theta_g, \theta_c, \theta_t) &= \frac{1}{n_s} \sum_{\mathbf{x}_i \in D_s} \mathcal{L}_c(C(T(G(\mathbf{x}_i))), y_i) - \\
&\frac{\lambda}{n} \left(\sum_{\mathbf{x}_i \in D_s} \mathcal{L}_s(D(T(G(\mathbf{x}_i))), d_i^s) + \right. \\
&\left. \sum_{\mathbf{x}_i \in D_t} \mathcal{L}_t(D(G(\mathbf{x}_i))), d_i^t) \right) \quad (12)
\end{aligned}$$

where d_i^s and d_i^t denote the domain labels of the i th source and target samples, respectively. \mathcal{L}_s and \mathcal{L}_t denote the domain loss of source and target samples, respectively. θ_t denotes the

parameters of T . This objective function replaces $G(\mathbf{x}_i)$ in (1) with $T(G(\mathbf{x}_i))$, indicating that our model uses features generated from transform network T to be the input of label classifier C and domain classifier D .

As our proof, if we optimize $\mathcal{L}(\theta_d, \theta_g, \theta_c, \theta_t)$ as general adversarial DA framework, a vanishing gradient problem would disturb us. To address this issue, we add a regularization term to the loss function based on the optimal transport problem as defined by Monge [44]. DA's goal is to find a mapping from a source domain to a target one, while the optimal transport problem gives a solution that transfers one distribution to another. Therefore, that problem can be represented in the form of Monge's formulation of the optimal transport problem [44], [45]. If we denote the probability measures over P_s and P_t as μ_s and μ_t , respectively, Monge's formulation of the optimal transport problem is:

$$T_0 = \arg \min_T \int_{\mathbf{x} \in P_s} r(\mathbf{x}, T(\mathbf{x})) d\mu(\mathbf{x}), s.t. T\#(\mu_s) = \mu_t \quad (13)$$

where $r(\cdot)$ denotes some kind of distance metric, T denotes a transport mapping from P_s to P_t , and T_0 is the optimal solution of T . $T\#(\mu_s)$ denotes the push forward of μ_s by a measurable function T . \mathbf{x} denotes the samples drawn from P_s . DA's goal is to find a transport mapping T_0 satisfying $T\#(\mu_s) = \mu_t$, which means that a transformation from source distribution P_s to target distribution P_t should be found. Specifically, in our model, we use transform network T to fit the transport mapping to meet $T\#(\mu_s) = \mu_t$ via adversarial training. By fitting $r(\cdot)$, we can construct a regularization term that measures the distance between $G(\mathbf{x}^s)$ and $T(G(\mathbf{x}^s))$. In our model, according to our empirical evaluation results, $r(\cdot)$ is the cosine distance between them:

$$r(G(\mathbf{x}^s), T(G(\mathbf{x}^s))) = -\frac{\langle G(\mathbf{x}^s) \cdot T(G(\mathbf{x}^s)) \rangle}{|G(\mathbf{x}^s)| \cdot |T(G(\mathbf{x}^s))|} \quad (14)$$

where $\langle \cdot \rangle$ denotes an inner product, and $|\cdot|$ denotes L_2 norm.

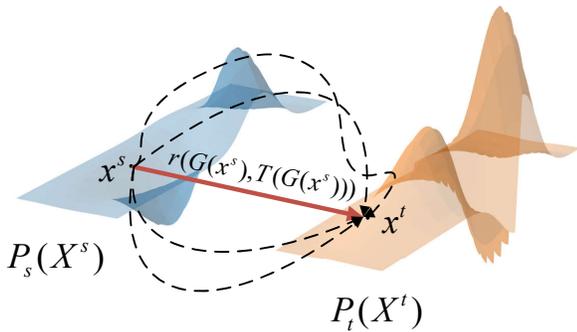


Fig. 3. When transferring features from the source to target domain, the regularization term proposed forces our model choosing the shortest path (red line).

For a transport problem, there are usually a few practical paths as shown in Fig. 3. The optimal transport theory seeks the most efficient way of transforming one distribution of mass to another, just like the red line in Fig. 3. In detail,

$\int r(\mathbf{x}, T(\mathbf{x})) d\mu(\mathbf{x})$ in (13), which indicates the expected cost of transportation, is to be minimized. If it reaches the minimum, the most efficient path would be found. Once we refine the optimal transport theory into unsupervised DA, the regularization term $r(G(\mathbf{x}^s), T(G(\mathbf{x}^s)))$ leads our model to select the most efficient way of transforming a source to target distribution. Specifically, the term attempts to constrain the distance between the features before and after the transformation. If we regard the distance as the cost of transformation, similar to the cost of transportation, the term attempts to select the shortest path from a number of transfer paths that map a source to target distribution.

On one hand, when the distributions of source and target domains are totally different, domain classifier D can so easily distinguish samples from different domains that \mathcal{L}_s and \mathcal{L}_t backpropagate very small gradients. In this situation, the regularization term $r(G(\mathbf{x}^s), T(G(\mathbf{x}^s)))$ can still provide gradients to the target mapping. On the other hand, when D directs updating parameters, the term would constrain the range of updating to prevent features from changing too rapidly, because it constrains differences between the features before and after the transformation. Thus, the stability of a training procedure of the proposed model is guaranteed via such added regularization term.

Consequently, our objective function becomes

$$\begin{aligned} \mathcal{L}(\theta_d, \theta_g, \theta_c, \theta_t) = & \frac{1}{n_s} \sum_{\mathbf{x}_i \in D_s} \mathcal{L}_c(C(T(G(\mathbf{x}_i))), y_i) - \\ & \frac{\lambda}{n} \left(\sum_{\mathbf{x}_i \in D_s} \mathcal{L}_s(D(T(G(\mathbf{x}_i))), d_i^s) + \right. \\ & \left. \sum_{\mathbf{x}_i \in D_t} \mathcal{L}_t(D(G(\mathbf{x}_i)), d_i^t) \right) + \\ & \beta \cdot r(G(\mathbf{x}^s), T(G(\mathbf{x}^s))) \end{aligned} \quad (15)$$

where β denotes the coefficient parameter of a regularization term. In addition, the optimization problem is to find parameters $\hat{\theta}_g, \hat{\theta}_d, \hat{\theta}_c$ and $\hat{\theta}_t$, where $\hat{\theta}_c, \hat{\theta}_g, \hat{\theta}_d$ and $\hat{\theta}_t$ satisfy:

$$\hat{\theta}_d = \arg \min_{\theta_d} (\mathcal{L}_s(\theta_d, \theta_g, \theta_t) + \mathcal{L}_t(\theta_d, \theta_g)) \quad (16)$$

$$\hat{\theta}_c = \arg \min_{\theta_c} \mathcal{L}_c(\theta_g, \theta_c, \theta_t) \quad (17)$$

$$\hat{\theta}_g = \arg \min_{\theta_g} \mathcal{L}(\theta_d, \theta_g, \theta_c, \theta_t) \quad (18)$$

$$\hat{\theta}_t = \arg \min_{\theta_t} (\mathcal{L}_c(\theta_g, \theta_c, \theta_t) + \mathcal{L}_s(\theta_d, \theta_g, \theta_t) + r(\theta_g, \theta_t)) \quad (19)$$

In this case, (10) is reformulated into:

$$\min_{G, T} \mathcal{L}(D^*, G, T) = 2 \cdot JSD(P_s || P_t) - 2 \log 2 + r(\theta_g, \theta_t) \quad (20)$$

where $\frac{\partial r(\theta_g, \theta_t)}{\partial \theta_g}$ and $\frac{\partial r(\theta_g, \theta_t)}{\partial \theta_t}$ would not be zeros. Since, $JSD(P_s || P_t) = 0$ appears easily, this regularization term provides gradients for parameters in G and T , thereby alleviating the adverse effects of the vanishing gradient problem.

F. Framework of Proposed Method

In a training period for our model, we have two stages. In the first one, feature extractor G and transform network T receive labeled source samples from $D_s\{\mathbf{x}_i^s, y_i^s, d_i^s\}_{i=1}^{n_s}$, and outputs \mathbf{f}^s and $T(\mathbf{f}^s)$. With class labels y^s and domain labels d^s , \mathcal{L}_c is computed by label classifier C , and \mathcal{L}_s is computed by domain classifier D . At the same time, the regularization term $r(G(\mathbf{x}^s), T(G(\mathbf{x}^s)))$ is also obtained according to \mathbf{f}^s and $T(\mathbf{f}^s)$. In the second stage, G receives unlabeled samples from $D_t\{\mathbf{x}_i^t, d_i^t\}_{i=1}^{n_t}$, and outputs \mathbf{f}^t . Similarly, \mathcal{L}_t is computed by domain classifier D . At last, all the above losses are multiplied by their corresponding coefficients, and then the model is optimized using these losses.

As for optimizing adversarial networks, previous studies have carried out a number of explorations [28], [29]. In [29], an iterative optimization strategy is proposed, where a feature extractor and domain classifier update their parameters iteratively. Specifically, it alternates between k steps of optimizing a domain classifier and one step of optimizing a feature extractor. This is the most common training strategy which is also widely used in GANs [23], [24]. One of the obstacles to it is that tuning the hyperparameter k . Unsuitable k may cause a failure of model training. As a result we have to tune this hyperparameter for each model carefully. Instead, in [28], the proposed gradient reversal layer (GRL) replaces iterative optimization. During forward propagation, GRL has no difference from normal layers, whereas during backpropagation, GRL reverses the gradient from the subsequent layer, multiplies it by a coefficient γ and passes it to the previous layer. Based on a large number of experiments, [28] adjusts γ using the following formula: $\gamma = \frac{2}{1+e^{-10p}} - 1$, where p is the training progress linearly changing from 0 to 1. In the implementation of ARTN, we choose GRL to optimize our model. With this strategy, there is no need to tune the hyperparameter k , and parameters of the feature extractor and domain classifier are updated in one backpropagation.

Algorithm 1 provides the pseudo-code of our proposed learning procedure. With stochastic gradient descent (SGD), parameters θ_d , θ_c and θ_g are updated. When the loss converges, the training stops.

If G , T and D have enough capacity, and at each loop of Algorithm 1, D is allowed to reach its optimal D^* given G and T , and P_t is updated so as to improve the following criterion:

$$\mathbb{E}_{\mathbf{x} \sim P_s}[\log D^*(T(G(\mathbf{x})))] + \mathbb{E}_{\mathbf{x} \sim P_t}[\log(1 - D^*(G(\mathbf{x})))] \quad (21)$$

Then P_t converges to P_s . Similar to [23], we give a brief proof as follows.

Proof. Consider that $U(P_t, D) = \int_{\mathbf{z}} P_s(\mathbf{z}) \log D(\mathbf{z}) + P_t(\mathbf{z}) \log(1 - D(\mathbf{z})) d\mathbf{z}$ as a function of P_t . Note that $U(P_t, D)$ is convex in P_t . The subderivatives of a supremum of convex functions include the derivative of the function at the point where the maximum is attained. In other words, if $f(x) = \sup_{\alpha \in \mathcal{A}} f_\alpha(x)$ and $f_\alpha(x)$ is convex in x for every α , then $\partial f_\beta(x) \in \partial f$ if $\beta = \operatorname{argsup}_{\alpha \in \mathcal{A}} f_\alpha(x)$. This is equivalent to computing a gradient descent update for P_t at D^* given the corresponding G and T . $\sup_D U(P_t, D)$ is convex in P_t

Algorithm 1 Learning Procedure of ARTN

Input:

Labeled source samples $D_s\{\mathbf{x}_i^s, y_i^s, d_i^s\}_{i=1}^{n_s}$
 Unlabeled target samples $D_t\{\mathbf{x}_i^t, d_i^t\}_{i=1}^{n_t}$
 Learning rate α , Coefficient parameters λ, β

Output:

Model parameters $\{\theta_d, \theta_g, \theta_c, \theta_t\}$

```

1: while not converge do
2:   for  $i$  from 1 to  $n_s$  do
3:      $\mathbf{f}_i^s = G(\mathbf{x}_i^s)$ 
4:      $\mathcal{L}_c = \text{crossentropy}(C(T(\mathbf{f}_i^s)), y_i^s)$ 
5:      $\mathcal{L}_s = \text{crossentropy}(D(T(\mathbf{f}_i^s)), d_i^s)$ 
6:      $\text{reg} = r(\mathbf{f}_i^s, T(\mathbf{f}_i^s))$ 
7:   end for
8:   for  $i$  from 1 to  $n_t$  do
9:      $\mathbf{f}_i^t = G(\mathbf{x}_i^t)$ 
10:     $\mathcal{L}_t = \text{crossentropy}(D(\mathbf{f}_i^t), d_i^t)$ 
11:  end for
12:   $\mathcal{L}_d \leftarrow \mathcal{L}_s + \mathcal{L}_t$ 
13:   $\theta_d \leftarrow \theta_d - \alpha \cdot \frac{\partial \mathcal{L}_d}{\partial \theta_d}$ 
14:   $\theta_c \leftarrow \theta_c - \alpha \cdot \frac{\partial \mathcal{L}_c}{\partial \theta_c}$ 
15:   $\theta_g \leftarrow \theta_g - \alpha \cdot \frac{\partial (\mathcal{L}_c - \lambda \mathcal{L}_d + \beta \text{reg})}{\partial \theta_g}$ 
16:   $\theta_t \leftarrow \theta_t - \alpha \cdot \frac{\partial (\mathcal{L}_c - \lambda \mathcal{L}_s + \beta \text{reg})}{\partial \theta_t}$ 
17: end while

```

with a unique global optimum as proven in (10). Hence, with sufficiently small updates of P_t , P_t converges to P_s .

IV. EXPERIMENTS

In order to evaluate the effectiveness of the proposed method, we test the proposed ARTN for unsupervised DA in several experiments that are recognized to be difficult. For the first experiment, we test our model in a sentiment analysis task. Second, to test its performance when source and target domains are relatively similar, the model is evaluated on several digits datasets. Third, to test it when facing a large discrepancy between source and target domains, the model is evaluated on a natural image dataset. Fourth, to test its anti-noise and generalization abilities, we test it when target images are added with varying noise. Fifth, to test the effectiveness of regularization in the proposed method, we compare the performance of ARTN with and without regularization on a natural image dataset. Finally, we investigate the effects of parameter λ on the performance of the proposed method. In all experiments, we implement models with Pytorch, and employ the learning strategy GRL mentioned in Section III, which reverses and propagates gradients to the feature extractors.

A. Sentiment Analysis

We use the **Amazon reviews** dataset with the same pre-processing used in mSDA [55] and DANN [28]. It contains reviews of four different categories of products: Books, DVDs, Kitchen Appliances and Electronics, which means that this dataset includes four domains and we can set up twelve domain adaptation tasks across them. Reviews are

encoded in 5 000 dimensional feature vectors of unigrams and bigrams, and labels are binary: 0 if a product is ranked up to 3 stars, and 1 if it is ranked 4 or 5 stars. In all twelve tasks, we use 2000 labeled source samples and 2000 unlabeled target samples to train our model. In a testing period, we test our model on separate target test sets (between 3000 and 6000 examples). To evaluate the effectiveness of our model, we compare it with DANN [28], DAN [19], Central Moment Discrepancy (CMD) for Domain-Invariant Representation Learning [47], Variational Fair Autoencoder (VFAE) [56] and the model with no adaptation. The results are directly cited from the original publication [28].

In this experiment, we use the same neural network as DANN [28]. Both domain and label classifiers consist of just one layer with 100 hidden units followed by the final output layer. Because there is only one hidden layer in the neural network, we build just one residual connection. ReLU activation function and batch normalization are employed. We choose SGD as the optimizer with its learning rate 0.001 and momentum 0.9. Parameters λ is set to 0.5, and β is set to 0.1. The batch size is set to 128. All the results are recorded after convergence.

Results are shown in Table I. The accuracy of ARTN is the highest in three out of twelve domain adaptation tasks. The accuracy of CMD-based model is the highest in six tasks and VFAE achieves the highest accuracy in three tasks. Therefore, in the experiment of sentiment analysis, ARTN is comparable with VFAE and slightly worse than CMD.

B. Digits



Fig. 4. Samples of digits dataset. The first to last rows correspond to MNIST, MNIST-M, SVHN and SYN NUMS.

In order to evaluate the performance when the discrepancy between source and target domains is relatively small, we experimentally test ARTN in several pairs of unsupervised domain adaptation tasks whose images are from the **MNIST**, **MNIST-M**, **SVHN** and **SYN NUMS** digits datasets. All these datasets consist of 10 classes, and we use the full training sets in all datasets. Example images from each dataset are shown in Fig. 4. In this experiment, we set three transfer tasks: **MNIST**→**MNIST-M**, **SVHN**→**MNIST**, and **SYN NUMS**→**SVHN**. As is shown in Fig. 4, images in **SYN NUMS** and **SVHN** are similar, whereas images in **MNIST** are much different from the other digits datasets.

We choose several unsupervised DA approaches to compare with the proposed one. CORAL [42], CMD [47] and DAN [19] rely on the distance metric between source and target distributions. DANN [28], CoGAN [25], Domain Transfer Network

(DTN) [57], CYCADA [58] and ADDA [29] are based on adversarial learning.

For **MNIST**→**MNIST-M**, we use a simple modified LeNet [59]. As for a domain classifier, we stack two fully connected layers: one layer with 100 hidden units followed by the final output layer. Each hidden unit uses a ReLU activation function. For **SVHN**→**MNIST** and **SYN NUMS**→**SVHN**, we use a three-layer convolutional network as a feature extractor, and a three-layer fully connected network as a domain classifier. In all tasks, batch normalization is employed. We employ SGD with 0.01 learning rate and the momentum 0.9. λ is set to 1, and β is set to 0.2. The batch size is set to 128. Prediction accuracy in the target domain is reported after convergence.

Results of our digits experiment are shown in Table II. Note that the basic networks for DTN and CYCADA are different from others, and the accuracy with no adaptation is included in the bracket. In **MNIST**→**MNIST-M**, the proposed model’s accuracy is 85.6% which outperforms the best of other methods by 0.6%. In **SYN NUMS**→**SVHN**, its accuracy achieves 89.1%, which is comparable to DANN’s. In **SYN NUMS**→**SVHN**, the accuracy of ARTN is 85.8%, which is just lower than CYCADA’s. Note that, CYCADA achieves the higher accuracy with better basic networks. For a fair comparison, the improvements compared with the adaptation-free models of ARTN, DTN, CyCADA pixel only and CyCADA pixel+feat are 30.9%, 8.3%, 3.2% and 23.3%, respectively. It is obvious that ARTN brings a bigger boost to the adaptation-free model. Totally, in two of three tasks, our approach outperforms other methods, and in the task whose source and target datasets are similar, it can achieve the same competitive results as the others.

C. Image Classification



Fig. 5. Samples of Office-31 dataset. The first to last rows correspond to AMAZON, DSLR and WEBCAM.

We further evaluate our model in a more complex setting. The proposed model is tested on a natural image dataset **Office-31**, which is a standard benchmark for visual domain adaptation, comprising 4,110 images and 31 categories collected from three domains: **AMAZON (A)**, images downloaded from amazon.com) with 2,817 images, **DSLR (D)**, high-resolution images captured by a digital SLR camera) with

TABLE I

CLASSIFICATION ACCURACY PERCENTAGE OF SENTIMENT ANALYSIS EXPERIMENT AMONG ALL TWELVE TASKS. THE FIRST COLUMN CORRESPONDS TO THE PERFORMANCE IF NO ADAPTION IS IMPLEMENTED. THE PROPOSED METHOD OUTPERFORMS THE OTHERS IN THREE OF TWELVE TASKS.

SOURCE→TARGET	SOURCE ONLY	DAN [19]	DANN [28]	ARTN	CMD [47]	VFAE [56]
BOOKS→DVD	78.7	79.6	78.4	81.4	80.5	79.9
BOOKS→ELECTRONICS	71.4	75.8	73.3	77.5	78.7	79.2
BOOKS→KITCHEN	74.5	78.7	77.9	78.8	81.3	81.6
DVD→BOOKS	74.6	78.0	72.3	78.8	79.5	75.5
DVD→ELECTRONICS	72.4	76.6	75.4	77.0	79.7	78.6
DVD→KITCHEN	76.5	79.6	78.3	79.3	83.0	82.2
ELECTRONICS→BOOKS	71.1	73.3	71.3	72.4	74.4	72.7
ELECTRONICS→DVD	71.9	74.8	73.8	73.9	76.3	76.5
ELECTRONICS→KITCHEN	84.4	85.7	85.4	86.4	86.0	85.0
KITCHEN→BOOKS	69.9	74.0	70.9	73.8	75.6	72.0
KITCHEN→DVD	73.4	76.3	74.0	75.7	77.5	73.3
KITCHEN→ELECTRONICS	83.3	84.4	84.3	86.1	85.4	83.8

TABLE II

CLASSIFICATION ACCURACY PERCENTAGE OF DIGITS CLASSIFICATIONS AMONG MNIST, MNIST-M, SVHN AND SYN NUMS. THE FIRST ROW CORRESPONDS TO THE PERFORMANCE IF NO ADAPTION IS IMPLEMENTED. THE PROPOSED METHOD OUTPERFORMS THE OTHERS IN TWO OF THREE TASKS WHEN IT COMES TO IMPROVEMENT COMPARED WITH THE BASIC NETWORK. IN ADDITION, THE RESULTS ARE CITED FROM LITERATURE.

METHOD	MNIST→MNIST-M	SYN NUMS→SVHN	SVHN→MNIST
SOURCE ONLY	51.4	86.7	54.9
CORAL [42]	57.7	85.2	63.1
DAN [19]	76.9	88.0	71.1
DANN [28]	76.7	91.1	73.9
CMD [47]	85.0	85.5	84.5
CoGAN [25]	-	-	DIVERGE
ADDA [29]	-	-	76.0
DTN [57]	-	-	84.4(76.1)
CYCADA PIXEL ONLY [58]	-	-	70.3(67.1)
CYCADA PIXEL+FEAT [58]	-	-	90.4(67.1)
ARTN	85.6	89.1	85.8

TABLE III

CLASSIFICATION ACCURACY PERCENTAGE OF EXPERIMENT ON THE OFFICE-31 DATASET. THE FIRST COLUMN CORRESPONDS TO THE PERFORMANCE IF NO ADAPTION IS IMPLEMENTED. THE SECOND TO LAST COLUMNS CORRESPOND TO THE PERFORMANCE OF DIFFERENT DA METHODS AND THE PROPOSED METHOD.

Method	DSLRL→AMAZON	WEBCAM→AMAZON	AMAZON→WEBCAM	AMAZON→DSLRL
AlexNet	51.1	49.8	61.6	63.8
DDC [27]	52.1	52.2	61.8	64.4
Deep CORAL [42]	52.8	51.5	66.4	66.8
DAN [19]	54.0	53.1	68.5	67.0
InceptionBN	60.1	57.9	70.3	70.5
LSSA [60]	57.8	57.8	67.7	71.3
CORAL [42]	59.0	60.2	70.9	71.9
AdaBN [61]	59.8	57.4	74.2	73.1
VGG16	58.2	57.8	67.6	73.9
CMD [47]	63.8	63.3	77.0	79.6
ResNet34	57.5	55.5	68.4	68.9
DANN [28]	58.1	56.3	73.7	75.3
ARTN	60.9	61.0	76.2	76.1

498 images and WEBCAM (**W**, low-resolution images captured by a Web camera) with 795 images. Samples of **Office-31** dataset are shown in Fig. 5. In order to test the generalization ability of different methods, we focus on the most difficult

four tasks [19]: **A→D**, **A→W**, **D→A** and **W→A**. In **A→W** and **A→D**, models are easier to train because images in source domain **A** are adequate. In **W→A** and **D→A**, there are only hundreds of images in the source domain but about 2,900

images in the target one. Thus models are very difficult to train. In addition, we test our model without regularization to analyze how the regularization term of our model affects its performance. In this experiment, we evaluate the effectiveness of our approach by comparing it with different models trained on the **Office-31** dataset. Note that some of the methods, such as DDC [27], Deep CORAL [42] and DAN [19], are based on AlexNet, some of the methods, such as LSSA [60], CORAL [42] and AdaBN [61], are based on InceptionBN, and CMD [47] is based on VGG16. Results of these methods are cited from original papers. Moreover, we implement DANN [28] and the model with no adaptation to be the baselines.

Because of lacking sufficient images, we implement our model based on ResNet34 [3] which is pre-trained on an ImageNet dataset, and fine-tune the model on Office-31. Different from the digits experiment, we build a residual connection for every three layers in ResNet34 instead of every layer. As for the domain classifier, we use a network with three fully connected layers. In addition, we replace the last layer of ResNet34 with a three-layer fully connected network, and use it to predict the labels of inputs. In all tasks, we employ the same SGD and parameter setting as before except that λ is set to 0.6 and batch size is 40. All the prediction accuracy results are recorded after training for 30 epochs.

Results of the experiment on Office-31 are shown in Table III. In $\mathbf{D} \rightarrow \mathbf{A}$, $\mathbf{W} \rightarrow \mathbf{A}$, $\mathbf{A} \rightarrow \mathbf{W}$ and $\mathbf{A} \rightarrow \mathbf{D}$, the proposed model achieves the accuracy of 60.9%, 61.0%, 76.2% and 76.1%, respectively. Thus, in all four tasks, the proposed model achieves the second highest accuracy. Note that these methods are based on different basic networks, besides the accuracy, improvement compared with the corresponding basic network is a fairer metric. The improvements of ARTN in four tasks are 3.4%, 5.5%, 7.8% and 7.2%, respectively. CMD, which outperforms all the related state-of-the-art methods on all four tasks, achieves 5.6%, 5.5%, 9.4% and 5.7% improvements. ARTN achieves a higher improvement in $\mathbf{A} \rightarrow \mathbf{D}$ and same improvement in $\mathbf{W} \rightarrow \mathbf{A}$ compared with the state-of-art method, CMD. The intuitive interpretation of CMD’s excellent performance is that it minimizes the sum of differences of higher order central moments of the corresponding activation distributions. Higher-order statistics describe the differences between distributions more comprehensively, but they also incur significantly more computational overhead than such methods as our proposed one. In practice, the number of moments is pre-set to be no more than five. In summary, ARTN outperforms all the other methods except CMD.

D. Generalization Analysis

A generalization test is taken by adding Gaussian noise to images in a target domain. In this way, the discrepancy between source and target domains is larger and discriminative information in a target domain is more difficult to capture. In this experiment, we test the anti-noise and generalization abilities of our model based on the digits experiment. For images in the source domain, we follow the settings in MNIST \rightarrow MNIST-M, SYN NUMS \rightarrow SVHN and SVHN \rightarrow MNIST respectively,

however, for images in the target domain, we add varying Gaussian noise. For MNIST \rightarrow MNIST-M and SYN NUMS \rightarrow SVHN, the standard deviation of Gaussian noise is selected from $\{0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0\}$. That in SVHN \rightarrow MNIST is from $\{1.0, 1.5, 2.0, 2.5, 3.0\}$. The means of Gaussian noise in all tasks are 0. Results are plotted in Fig. 6. The baseline method is a model without adaptation. We also compare the proposed method with DANN [28].

Comparing the proposed model with the adaptation-free model, we can see that although noises are added to the test images, the proposed model exhibits a great advantage over the adaptation-free model. In MNIST \rightarrow MNIST-M, when the standard deviation is 0.4, the accuracy of the adaptation-free model is 45.83%, whereas ours improves it by 36.6%. When standard deviation is 1.0, the accuracy of the adaptation-free model is 24.55%, whereas ours improves is by 76.4%. Similar results appear in SYN NUMS \rightarrow SVHN and SVHN \rightarrow MNIST, where the rate of improvement generally shows an upward trend in the case of a gradual increase of noise. Therefore, as the discrepancy between source and target domains increases, the performance advantage of the proposed model is becoming more and more obvious in comparison with a adaptation-free model. At the same time, the improvement percentage of our model is higher than DANN in almost all tasks, which means that the proposed method has better anti-noise abilities than DANN. This result demonstrates that even if there exists noise in a target domain, the proposed model can maintain excellent generalization and anti-noise abilities.

E. Regularization Analysis

We next analyze how the regularization term of our model affects the performance of our model. We test our model without regularization on Office-31 by setting $\beta = 0$. In this way, \mathcal{L} consists of \mathcal{L}_c , \mathcal{L}_s and \mathcal{L}_t only. Except for the regularization term, this experiment has same settings as the image classification experiment does.

Results of this experiment are shown in Fig. 7. In $\mathbf{D} \rightarrow \mathbf{A}$, $\mathbf{W} \rightarrow \mathbf{A}$, $\mathbf{A} \rightarrow \mathbf{W}$ and $\mathbf{A} \rightarrow \mathbf{D}$, the proposed model without regularization achieves the accuracy of 59.5%, 59.8%, 76.0% and 75.9%, which is lower by 1.4%, 1.2%, 0.2% and 0.2% of the proposed one with such term, respectively. The model with regularization outperforms DANN and the model without regularization in all tasks, which demonstrates the effectiveness of regularization. In another word, the regularization term strengthens the generalization ability of the proposed model. It should be noted that in $\mathbf{D} \rightarrow \mathbf{A}$, $\mathbf{W} \rightarrow \mathbf{A}$ and $\mathbf{A} \rightarrow \mathbf{W}$, the proposed model without regularization still outperforms DANN. This means that the improvement is not only from the regularization but also the modification of its architecture.

Besides performance improvement, we analyze how the regularization term affects the gradients during training. Because displaying every gradient of a parameter is impossible, we calculate $\|\nabla_{\theta} \mathcal{L}(\theta)\|$ to capture the overall statistics which is a metric adopted in [53]. We record $\|\nabla_{\theta} \mathcal{L}(\theta)\|$ of the model with and without a regularization term on Office-31. Moreover, we record the minimum, maximum and standard deviation of $\|\nabla_{\theta} \mathcal{L}(\theta)\|$ during the training period. $\|\nabla_{\theta} \mathcal{L}(\theta)\|$ in $\mathbf{D} \rightarrow \mathbf{A}$,

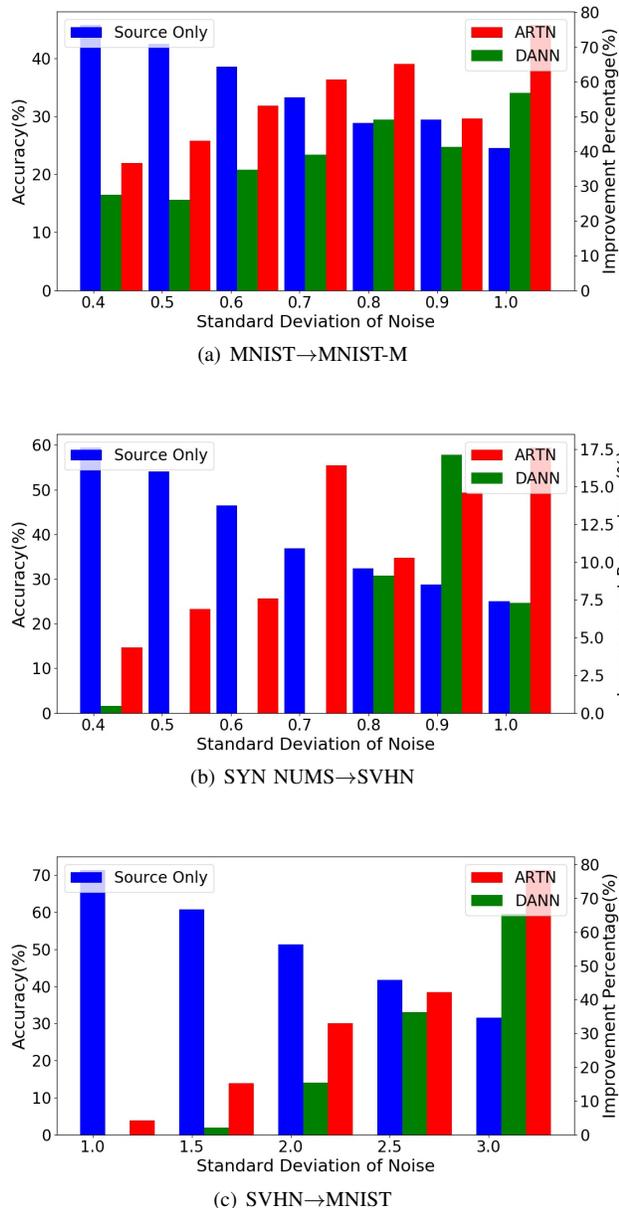


Fig. 6. Accuracy of the adaptation-free method and improvement of DANN and our model in MNIST→MNIST-M, SYN NUMS→SVHN and SVHN→MNIST, where we add gaussian noise to images in the target domain. X-axis represents the standard deviation of noise, and Y-axis represents the accuracy of adaptation-free method and improvement percentage of DANN and our model in the target domain.

$W \rightarrow A$, $A \rightarrow W$ and $A \rightarrow D$ are drawn in Fig. 8 and related statistical data are shown in Table IV. Please note that even if the gradient vanishing issue occurs during training, $\|\nabla_{\theta} \mathcal{L}(\theta)\|$ would not be very close to zero because it involves gradients of all parameters.

According to Fig. 8, especially in Fig. 8(a), (b) and (d), we can see that gradients of ARTN without a regularization term are easier to be unstable. There are more extreme large gradients in ARTN without a regularization term. The results are shown in Table IV. In all four tasks, ARTN with a regularization term gets smaller standard deviation than ARTN

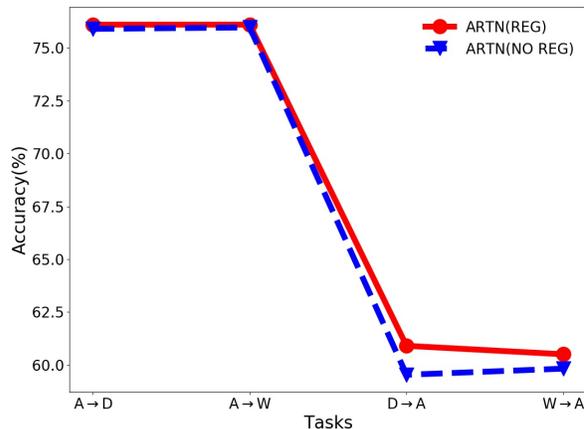


Fig. 7. Classification accuracy percentage of experiment on the Office-31 dataset. The red line corresponds to the proposed method with regularization and the blue line corresponds to the one without regularization. The regularization term shows a positive effect on the performance improvement.

TABLE IV
STATISTICAL DATA OF $\|\nabla_{\theta} \mathcal{L}(\theta)\|$ ARE RECORDED. BOTH THE MODELS WITH AND WITHOUT REGULARIZATION TERM ARE EVALUATED ON OFFICE-31. FOR EACH ROW, THE UPPER LINE CORRESPONDS TO ARTN WITH REGULARIZATION AND THE LOWER LINE CORRESPONDS TO ARTN WITHOUT REGULARIZATION.

Method	A→W	A→D	W→A	D→A
max	14.51	15.88	4.64	6.16
	18.50	20.35	4.81	7.02
min	3.24	2.96	2.07	2.32
	3.17	3.04	2.20	2.32
max-min	11.27	12.92	2.57	3.84
	15.33	17.31	2.61	4.70
std	1.25	1.36	0.38	0.59
	1.30	1.37	0.40	0.61

without it. This directly indicates that a regularization term promotes the stability of adversarial training in our model. In detail, the maximum and minimum gradients are also recorded. We find that maximum gradient of ARTN with a regularization term in four tasks are smaller than that of ARTN without it. Meanwhile, except in $W \rightarrow A$, the gap between maximum and minimum gradients suggests that ARTN with a regularization term is more stable than ARTN without it with a large margin. This fact can also be observed in Fig. 8. Thus, the effect of the proposed regularization term to stabilize adversarial training in our model is considered being verified in this experiment. The reason why the $W \rightarrow A$ case is an exception needs to be explored as future research.

F. Parameter Sensitivity

In this experiment, we investigate how parameter λ affects the performance of our model. In order to make the results convincing, we test our model on tasks $A \rightarrow W$, $A \rightarrow D$, $W \rightarrow A$, and $D \rightarrow A$ to acquire the variation of transfer classification performance as $\lambda \in \{0.4, 0.5, 0.6, 0.7, 0.8, 0.9\}$. Note that other

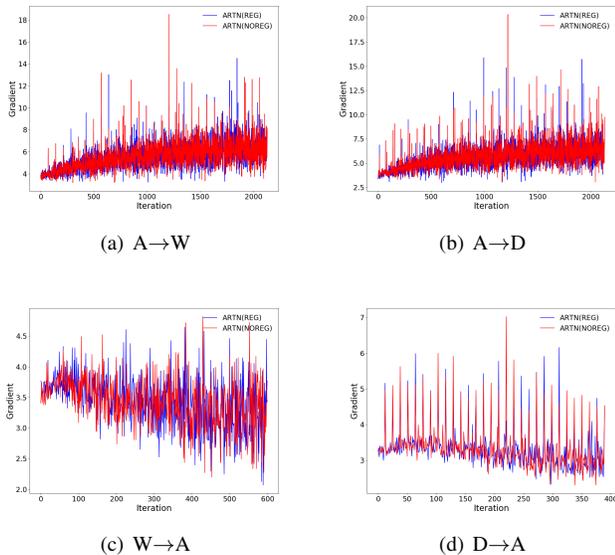


Fig. 8. $\|\nabla_{\theta} \mathcal{L}(\theta)\|$ in task $D \rightarrow A$, $W \rightarrow A$, $A \rightarrow W$ and $A \rightarrow D$. Red line corresponds to our model without regularization term while blue line corresponds to our model with regularization term.

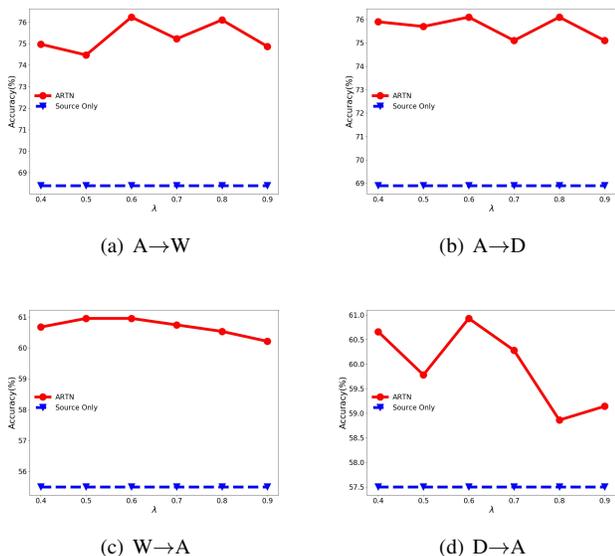


Fig. 9. Sensitivity of λ in task $A \rightarrow W$, $A \rightarrow D$, $W \rightarrow A$, and $D \rightarrow A$. Dashed lines show results of adaptation-free method.

settings are the same with those of the image classification experiment. In Fig. 9, a detailed illustration is given.

The results in three of four tasks, $A \rightarrow W$, $A \rightarrow D$ and $W \rightarrow A$ exhibit the same trend that the accuracy of ARTN is almost stable as λ varies. Only in $W \rightarrow A$, the accuracy fluctuates slightly with the variation of λ . Moreover, in the range of our settings, ARTN is always better than the model without adaptation and also better than most methods in Table III. This confirms the belief that ARTN is robust as λ changes, which means the proposed method is no need for tuning hyper parameters subtly.

V. CONCLUSION

We propose a novel unsupervised domain adaptation model based on adversarial learning. Different from previous adversarial adaptation models which rely on extracting domain-invariant representations, our model adds a feature-shared transform network to directly map features from a source domain to the space of target features. Furthermore, we add a regularization term to help strengthen its performance. Experimental results clearly demonstrate that the proposed model can match different domains effectively and is comparable with the state-of-the-art methods.

REFERENCES

- [1] Y. Kim, “Convolutional neural networks for sentence classification,” in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, 2014, pp. 1746–1751. [Online]. Available: <http://www.aclweb.org/anthology/D14-1181>
- [2] C. dos Santos and M. Gatti, “Deep convolutional neural networks for sentiment analysis of short texts,” in *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers*, 2014, pp. 69–78.
- [3] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [4] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in Neural Information Processing Systems*, 2012, pp. 1097–1105.
- [5] Z. Ren, K. Qian, Z. Zhang, V. Pandit, A. Baird, and B. Schuller, “Deep scalogram representations for acoustic scene classification,” *IEEE/CAA Journal of Automatica Sinica*, vol. 5, no. 3, pp. 662–669, 2018.
- [6] J. Long, E. Shelhamer, and T. Darrell, “Fully convolutional networks for semantic segmentation,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3431–3440.
- [7] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson, “How transferable are features in deep neural networks?” in *Advances in Neural Information Processing Systems*, 2014, pp. 3320–3328.
- [8] S. Ben-David, J. Blitzer, K. Crammer, A. Kulesza, F. Pereira, and J. W. Vaughan, “A theory of learning from different domains,” *Machine learning*, vol. 79, no. 1, pp. 151–175, 2010.
- [9] S. Ben-David, J. Blitzer, K. Crammer, and F. Pereira, “Analysis of representations for domain adaptation,” in *Advances in Neural Information Processing Systems*, 2007, pp. 137–144.
- [10] S. J. Pan and Q. Yang, “A survey on transfer learning,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 22, no. 10, pp. 1345–1359, 2010.
- [11] L. Shao, F. Zhu, and X. Li, “Transfer learning for visual categorization: A survey,” *IEEE Transactions on Neural Networks and Learning Systems*, vol. 26, no. 5, pp. 1019–1034, 2015.
- [12] S. J. Pan, I. W. Tsang, J. T. Kwok, and Q. Yang, “Domain adaptation via transfer component analysis,” *IEEE Transactions on Neural Networks*, vol. 22, no. 2, pp. 199–210, 2011.
- [13] B. Gong, Y. Shi, F. Sha, and K. Grauman, “Geodesic flow kernel for unsupervised domain adaptation,” in *Computer Vision and Pattern Recognition (CVPR)*, 2012, pp. 2066–2073.
- [14] M. Long, J. Wang, J. Sun, and S. Y. Philip, “Domain invariant transfer kernel learning,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 27, no. 6, pp. 1519–1532, 2015.
- [15] Y. Bengio, “Deep learning of representations for unsupervised and transfer learning,” in *Proceedings of ICML Workshop on Unsupervised and Transfer Learning*, 2012, pp. 17–36.
- [16] G. Mesnil, Y. Dauphin, X. Glorot, S. Rifai, Y. Bengio, I. Goodfellow, E. Lavoie, X. Muller, G. Desjardins, D. Warde-Farley *et al.*, “Unsupervised and transfer learning challenge: a deep learning approach,” in *Proceedings of the 2011 International Conference on Unsupervised and Transfer Learning workshop-Volume 27*, 2011, pp. 97–111.
- [17] G. Hinton, O. Vinyals, and J. Dean, “Distilling the knowledge in a neural network,” *arXiv preprint arXiv:1503.02531*, 2015.
- [18] A. Gretton, D. Sejdinovic, H. Strathmann, S. Balakrishnan, M. Pontil, K. Fukumizu, and B. K. Sriperumbudur, “Optimal kernel choice for large-scale two-sample tests,” in *Advances in Neural Information Processing Systems*, 2012, pp. 1205–1213.

- [19] M. Long, Y. Cao, J. Wang, and M. Jordan, "Learning transferable features with deep adaptation networks," in *International Conference on Machine Learning*, 2015, pp. 97–105.
- [20] M. Long, H. Zhu, J. Wang, and M. I. Jordan, "Deep transfer learning with joint adaptation networks," in *Proceedings of the 34th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, D. Precup and Y. W. Teh, Eds., vol. 70. International Convention Centre, Sydney, Australia: PMLR, 06–11 Aug 2017, pp. 2208–2217.
- [21] M. Long, H. Zhu, J. Wang, and M. I. Jordan, "Unsupervised domain adaptation with residual transfer networks," in *Proceedings of the 30th International Conference on Neural Information Processing Systems*, ser. NIPS'16. USA: Curran Associates Inc., 2016, pp. 136–144.
- [22] J. Hoffman, S. Guadarrama, E. S. Tzeng, R. Hu, J. Donahue, R. Girshick, T. Darrell, and K. Saenko, "LSDa: Large scale detection through adaptation," in *Advances in Neural Information Processing Systems*, 2014, pp. 3536–3544.
- [23] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in Neural Information Processing Systems*, 2014, pp. 2672–2680.
- [24] K. Wang, C. Gou, Y. Duan, Y. Lin, X. Zheng, and F.-Y. Wang, "Generative adversarial networks: introduction and outlook," *IEEE/CAA Journal of Automatica Sinica*, vol. 4, no. 4, pp. 588–598, 2017.
- [25] M.-Y. Liu and O. Tuzel, "Coupled generative adversarial networks," in *Advances in Neural Information Processing Systems*, 2016, pp. 469–477.
- [26] K. Bousmalis, N. Silberman, D. Dohan, D. Erhan, and D. Krishnan, "Unsupervised pixel-level domain adaptation with generative adversarial networks," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, vol. 1, no. 2, 2017, p. 7.
- [27] E. Tzeng, J. Hoffman, T. Darrell, and K. Saenko, "Simultaneous deep transfer across domains and tasks," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 4068–4076.
- [28] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky, "Domain-adversarial training of neural networks," *Journal of Machine Learning Research*, vol. 17, no. 59, pp. 1–35, 2016.
- [29] E. Tzeng, J. Hoffman, K. Saenko, and T. Darrell, "Adversarial discriminative domain adaptation," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2017, pp. 2962–2971.
- [30] Y. Chen, S. Song, S. Li, L. Yang, and C. Wu, "Domain space transfer extreme learning machine for domain adaptation," *IEEE Transactions on Cybernetics*, pp. 1–14, 2018.
- [31] S. Mehrkanoon and J. A. K. Suykens, "Regularized semipaired kernel cca for domain adaptation," *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–15, 2018.
- [32] S. Khalighi, B. Ribeiro, and U. J. Nunes, "Importance weighted import vector machine for unsupervised domain adaptation," *IEEE Transactions on Cybernetics*, vol. 47, no. 10, pp. 3280–3292, Oct 2017.
- [33] J. Li, K. Lu, Z. Huang, L. Zhu, and H. T. Shen, "Transfer independently together: A generalized framework for domain adaptation," *IEEE Transactions on Cybernetics*, pp. 1–12, 2018.
- [34] T. Kanamori, S. Hido, and M. Sugiyama, "A least-squares approach to direct importance estimation," *Journal of Machine Learning Research*, vol. 10, no. Jul, pp. 1391–1445, 2009.
- [35] J. Huang, A. Gretton, K. M. Borgwardt, B. Schölkopf, and A. J. Smola, "Correcting sample selection bias by unlabeled data," in *Advances in Neural Information Processing Systems*, 2007, pp. 601–608.
- [36] M. Sugiyama, S. Nakajima, H. Kashima, P. V. Buenau, and M. Kawanabe, "Direct importance estimation with model selection and its application to covariate shift adaptation," in *Advances in Neural Information Processing Systems*, 2008, pp. 1433–1440.
- [37] S. Li, S. Song, and G. Huang, "Prediction reweighting for domain adaptation," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 28, no. 7, pp. 1682–1695, 2017.
- [38] K. Zhang, B. Schölkopf, K. Muandet, and Z. Wang, "Domain adaptation under target and conditional shift," in *International Conference on Machine Learning*, 2013, pp. 819–827.
- [39] T. Liu and D. Tao, "Classification with noisy labels by importance reweighting," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 38, no. 3, pp. 447–461, 2016.
- [40] B. Fernando, A. Habrard, M. Sebban, and T. Tuytelaars, "Unsupervised visual domain adaptation using subspace alignment," in *Computer Vision (ICCV), 2013 IEEE International Conference on*. IEEE, 2013, pp. 2960–2967.
- [41] B. Sun and K. Saenko, "Subspace distribution alignment for unsupervised domain adaptation," in *BMVC*, 2015, pp. 24–1.
- [42] B. Sun, J. Feng, and K. Saenko, "Return of frustratingly easy domain adaptation," in *AAAI*, vol. 6, no. 7, 2016, p. 8.
- [43] L. Cheng and S. J. Pan, "Semi-supervised domain adaptation on manifolds," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 25, no. 12, pp. 2240–2249, 2014.
- [44] N. Courty, R. Flamary, D. Tuia, and A. Rakotomamonjy, "Optimal transport for domain adaptation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 9, pp. 1853–1865, 2017.
- [45] N. Courty, R. Flamary, A. Habrard, and A. Rakotomamonjy, "Joint distribution optimal transportation for domain adaptation," in *Advances in Neural Information Processing Systems*, 2017, pp. 3733–3742.
- [46] J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, and T. Darrell, "Decaf: A deep convolutional activation feature for generic visual recognition," in *International Conference on Machine Learning*, 2014, pp. 647–655.
- [47] W. Zellinger, T. Grubinger, E. Lughofer, T. Natschläger, and S. Saminger-Platz, "Central moment discrepancy (CMD) for domain-invariant representation learning," in *International Conference on Learning Representations*, 2017.
- [48] S. Chopra, R. Hadsell, and Y. LeCun, "Learning a similarity metric discriminatively, with application to face verification," in *Computer Vision and Pattern Recognition*, vol. 1, 2005, pp. 539–546.
- [49] X. Glorot, A. Bordes, and Y. Bengio, "Domain adaptation for large-scale sentiment classification: A deep learning approach," in *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, 2011, pp. 513–520.
- [50] M. Kan, S. Shan, and X. Chen, "Bi-shifting auto-encoder for unsupervised domain adaptation," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 3846–3854.
- [51] M. Ghifary, W. B. Kleijn, M. Zhang, D. Balduzzi, and W. Li, "Deep reconstruction-classification networks for unsupervised domain adaptation," in *European Conference on Computer Vision*. Springer, 2016, pp. 597–613.
- [52] A. Rozantsev, M. Salzmann, and P. Fua, "Beyond sharing weights for deep domain adaptation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018.
- [53] M. Arjovsky, S. Chintala, and L. Bottou, "Wasserstein generative adversarial networks," in *International Conference on Machine Learning*, 2017, pp. 214–223.
- [54] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.
- [55] M. Chen, Z. Xu, K. Q. Weinberger, and F. Sha, "Marginalized denoising autoencoders for domain adaptation," in *Proceedings of the 29th International Conference on Machine Learning*, ser. ICML'12. USA: Omnipress, 2012, pp. 1627–1634. [Online]. Available: <http://dl.acm.org/citation.cfm?id=3042573.3042781>
- [56] C. Louizos, K. Swersky, Y. Li, M. Welling, and R. Zemel, "The variational fair autoencoder," 2016.
- [57] Y. Taigman, A. Polyak, and L. Wolf, "Unsupervised cross-domain image generation," 2016.
- [58] J. Hoffman, E. Tzeng, T. Park, J.-Y. Zhu, P. Isola, K. Saenko, A. Efros, and T. Darrell, "CyCADA: Cycle-consistent adversarial domain adaptation," in *Proceedings of the 35th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, J. Dy and A. Krause, Eds., vol. 80. Stockholmsmssan, Stockholm Sweden: PMLR, 10–15 Jul 2018, pp. 1989–1998. [Online]. Available: <http://proceedings.mlr.press/v80/hoffman18a.html>
- [59] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [60] R. Aljundi, R. Emonet, D. Muselet, and M. Sebban, "Landmarks-based kernelized subspace alignment for unsupervised domain adaptation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 56–63.
- [61] Y. Li, N. Wang, J. Shi, J. Liu, and X. Hou, "Revisiting batch normalization for practical domain adaptation," 2017.