

Discriminative Adversarial Domain Adaptation

Hui Tang, Kui Jia*

South China University of Technology
eehuitang@mail.scut.edu.cn, kuijia@scut.edu.cn

Abstract

Given labeled instances on a source domain and unlabeled ones on a target domain, unsupervised domain adaptation aims to learn a task classifier that can well classify target instances. Recent advances rely on domain-adversarial training of deep networks to learn domain-invariant features. However, due to an issue of mode collapse induced by the separate design of task and domain classifiers, these methods are limited in aligning the joint distributions of feature and category across domains. To overcome it, we propose a novel adversarial learning method termed Discriminative Adversarial Domain Adaptation (DADA). Based on an integrated category and domain classifier, DADA has a novel adversarial objective that encourages a mutually inhibitory relation between category and domain predictions for any input instance. We show that under practical conditions, it defines a minimax game that can promote the joint distribution alignment. Except for the traditional closed set domain adaptation, we also extend DADA for extremely challenging problem settings of partial and open set domain adaptation. Experiments show the efficacy of our proposed methods and we achieve the new state of the art for all the three settings on benchmark datasets.

Introduction

Many machine learning tasks are advanced by large-scale learning of deep models, with image classification (Rusakovsky et al. 2015) as one of the prominent examples. A key factor to achieve such advancements is the availability of massive labeled data on the domains of the tasks of interest. For many other tasks, however, training instances on the corresponding domains are either difficult to collect, or their labeling costs prohibitively. To address the scarcity of labeled data for these *target* tasks/domains, a general strategy is to leverage the massively available labeled data on related *source* ones via domain adaptation (Pan and Yang 2010). Even though the source and target tasks share the same label space (i.e. closed set domain adaptation), domain adaptation still suffers from the shift in data distributions. The main objective of domain adaptation is thus to learn domain-invariant features, so that task classifiers learned from the

source data can be readily applied to the target domain. In this work, we focus on the unsupervised setting where training instances on the target domain are completely unlabeled.

Recent domain adaptation methods are largely built on modern deep architectures. They rely on great model capacities of these networks to learn hierarchical features that are empirically shown to be more transferable across domains (Yosinski et al. 2014; Zhang, Tang, and Jia 2018). Among them, those based on domain-adversarial training (Ganin et al. 2016; Wang et al. 2019) achieve the current state of the art. Based on the seminal work of DANN (Ganin et al. 2016), they typically augment a classification network with an additional domain classifier. The domain classifier takes features from the feature extractor of the classification network as inputs, which is trained to differentiate between instances from the two domains. By playing a minimax game (Goodfellow et al. 2014), adversarial training aims to learn domain-invariant features.

Such domain-adversarial networks can largely reduce the domain discrepancy. However, the separate design of task and domain classifiers has the following shortcomings. Firstly, feature distributions can only be aligned to a certain level, since model capacity of the feature extractor could be large enough to compensate for the less aligned feature distributions. More importantly, given practical difficulties of aligning the source and target distributions with high granularity to the category level (especially for complex distributions with multi-mode structures), the task classifier obtained by minimizing the empirical source risk cannot well generalize to the target data due to an issue of mode collapse (Kurmi and Namboodiri 2019; Tran et al. 2019), i.e., the joint distributions of feature and category are not well aligned across the source and target domains.

Recent methods (Kurmi and Namboodiri 2019; Tran et al. 2019) take the first step to address the above shortcomings by jointly parameterizing the task and domain classifiers into an integrated one. To further push this line, based on such a classifier, we propose a novel adversarial learning method termed *Discriminative Adversarial Domain Adaptation (DADA)*, which encourages a *mutually inhibitory* relation between its domain prediction and category prediction for any input instance, as illustrated in Figure 1. This dis-

*Corresponding author.

criminative interaction between category and domain predictions underlies the ability of DADA to reduce domain discrepancy at both the feature and category levels. Intuitively, the adversarial training of DADA mainly conducts competition between the domain neuron (output) and the true category neuron (output). Different from the work (Tran et al. 2019) whose mechanism to align the joint distributions is rather implicit, DADA enables explicit alignment between the joint distributions, thus improving the classification of target data. Except for closed set domain adaptation, we also extend DADA for partial domain adaptation (Cao et al. 2018b), i.e. the target label space is subsumed by the source one, and open set domain adaptation (Saito et al. 2018c), i.e. the source label space is subsumed by the target one. Our main contributions can be summarized as follows.

- We propose in this work a novel adversarial learning method, termed DADA, for closed set domain adaptation. Based on an integrated category and domain classifier, DADA has a novel adversarial objective that encourages a *mutually inhibitory* relation between category and domain predictions for any input instance, which can promote the joint distribution alignment across domains.
- For more realistic partial domain adaptation, we extend DADA by a reliable category-level weighting mechanism, termed DADA-P, which can significantly reduce the negative influence of outlier source instances.
- For more challenging open set domain adaptation, we extend DADA by balancing the joint distribution alignment in the shared label space with the classification of outlier target instances, termed DADA-O.
- Experiments show the efficacy of our proposed methods and we achieve the new state of the art for all the three adaptation settings on benchmark datasets.

Related Works

Closed Set Domain Adaptation After the seminal work of DANN (Ganin et al. 2016), ADDA (Tzeng et al. 2017) proposes an untied weight sharing strategy to align the target feature distribution to a fixed source one. SimNet (Pinheiro 2018) replaces the standard FC-based cross-entropy classifier by a similarity-based one. MADA (Pei et al. 2018) and CDAN (Long et al. 2018b) integrate the discriminative category information into domain-adversarial training. VADA (Shu et al. 2018) reduces the cluster assumption violation to constrain domain-adversarial training. Some methods (Wang et al. 2019; Wen et al. 2019) focus on transferable regions to learn domain-invariant features and task classifier. TAT (Liu et al. 2019) enhances the discriminability of features to guarantee the adaptability. Some methods (Saito et al. 2018b; 2018a; Lee et al. 2019) utilize category predictions from two task classifiers to measure the domain discrepancy. The most related works (Kurmi and Namboodiri 2019; Tran et al. 2019) to us propose joint parameterization of the task and domain classifiers, which implicitly align the joint distributions. Differently, our proposed DADA makes the joint distribution alignment more explicit, thus promoting classification on the target domain.

Partial Domain Adaptation The work (Zhang et al. 2018a) weights each source instance by its importance to the target domain based on one domain classifier, and then trains another domain classifier on target and weighted source instances. The works (Cao et al. 2018a; 2018b) reduce the contribution of outlier source instances to the task or domain classifiers by utilizing category predictions. Differently, DADA-P weights the proposed source discriminative adversarial loss by a reliable category confidence.

Open Set Domain Adaptation Previous research (Jain, Scheirer, and Boulton 2014) proposes to reject an instance as the unknown category by threshold filtering. The work (Saito et al. 2018c) proposes to utilize adversarial training for both domain adaptation and unknown outlier detection. Differently, DADA-O balances the joint distribution alignment in the shared label space with the outlier rejection.

Method

Given $\{(\mathbf{x}_i^s, y_i^s)\}_{i=1}^{n_s}$ of labeled instances sampled from the source domain \mathcal{D}_s , and $\{\mathbf{x}_j^t\}_{j=1}^{n_t}$ of unlabeled instances sampled from the target domain \mathcal{D}_t , the objective of unsupervised domain adaptation is to learn a feature extractor $G(\cdot)$ and a task classifier $C(\cdot)$ such that the expected target risk $\mathbb{E}_{(\mathbf{x}^t, y^t) \sim \mathcal{D}_t}[\mathcal{L}_{cls}(C(G(\mathbf{x}^t)), y^t)]$ is low for a certain classification loss function $\mathcal{L}_{cls}(\cdot)$. The domains \mathcal{D}_s and \mathcal{D}_t are assumed to have different distributions. To achieve a low target risk, a typical strategy is to learn $G(\cdot)$ and $C(\cdot)$ by minimizing the sum of the source risk and some notion of *distance* between the source and target domain distributions, inspired by domain adaptation theories (Ben-David et al. 2007; 2010). This strategy is based on a simple rationale that the source risk would become a good indicator of the target risk when the distance between the two distributions is getting closer. While most of existing methods use distance measures based on the marginal distributions, it is arguably better to use those based on the joint distributions.

The above strategy is generally implemented by domain-adversarial learning (Ganin et al. 2016; Wang et al. 2019), where separate task classifier $C(\cdot)$ and domain classifier $D(\cdot)$ are typically stacked on top of the feature extractor $G(\cdot)$. As discussed before, this type of design has the following shortcomings: (1) model capacity of $G(\cdot)$ could be large enough to make $D(G(\mathbf{x}^s))$ and $D(G(\mathbf{x}^t))$ hardly differentiable for any instance, even though the marginal feature distributions are not well aligned; (2) more importantly, it is difficult to align the source and target distributions with high granularity to the category level (especially for complex distributions with multi-mode structures), and thus $C(\cdot)$ obtained by minimizing the empirical source risk cannot perfectly generalize to the target data due to an issue of mode collapse, i.e. the joint distributions are not well aligned.

To alleviate the above shortcomings, inspired by semi-supervised learning methods based on GANs (Salimans et al. 2016; Dai et al. 2017), the recent work (Tran et al. 2019) proposes joint parameterization of $C(\cdot)$ and $D(\cdot)$ into an integrated one $F(\cdot)$. Suppose the classification task of interest has K categories, $F(\cdot)$ is formed simply by augmenting the last FC layer of $C(\cdot)$ with one additional neuron.

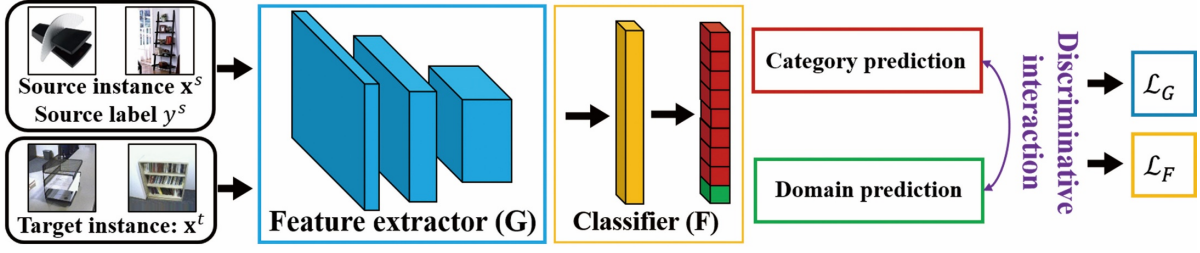


Figure 1: (Best viewed in color.) Discriminative Adversarial Domain Adaptation (DADA), which includes a feature extractor $G(\cdot)$ and an integrated category and domain classifier $F(\cdot)$. The blue and orange colors denote $G(\cdot)$ and $F(\cdot)$, and the losses applied to them, respectively. Note that DADA explicitly establishes a discriminative interaction between category and domain predictions. Please refer to the main text for how the adversarial training objective of DADA is defined.

Denote $\mathbf{p}(\mathbf{x}) \in [0, 1]^{K+1}$ as the output vector of class probabilities of $F(G(\mathbf{x}))$ for an instance \mathbf{x} , and $p_k(\mathbf{x})$, $k \in \{1, \dots, K+1\}$, as its k^{th} element. The k^{th} element of the conditional probability vector $\bar{\mathbf{p}}(\mathbf{x})$ is written as follows

$$\bar{p}_k(\mathbf{x}) = \begin{cases} \frac{p_k(\mathbf{x})}{1 - p_{K+1}(\mathbf{x})}, & k = 1, 2, \dots, K \\ 0, & k = K + 1 \end{cases} \quad (1)$$

For ease of subsequent notations, we also write $p_k^s = p_k(\mathbf{x}^s)$ and $p_k^t = p_k(\mathbf{x}^t)$. Then, such a network is trained by the classification-aware adversarial learning objective

$$\begin{aligned} \min_F & -\frac{1}{n_s} \sum_{i=1}^{n_s} \log p_{y_i^s}(\mathbf{x}_i^s) - \frac{1}{n_t} \sum_{j=1}^{n_t} \log p_{K+1}(\mathbf{x}_j^t) \\ \max_G & \frac{1}{n_s} \sum_{i=1}^{n_s} \log \bar{p}_{y_i^s}(\mathbf{x}_i^s) + \lambda \frac{1}{n_t} \sum_{j=1}^{n_t} \log(1 - p_{K+1}(\mathbf{x}_j^t)), \end{aligned} \quad (2)$$

where λ balances category classification and domain adversarial losses. The mechanism of this objective to align the joint distributions across domains is rather implicit.

To make it more explicit, based on the integrated classifier $F(\cdot)$, we propose a novel adversarial learning method termed *Discriminative Adversarial Domain Adaptation (DADA)*, which explicitly enables a discriminative interplay of predictions among the domain and K categories for any input instance, as illustrated in Figure 1. This discriminative interaction underlies the ability of DADA to promote the joint distribution alignment, as explained shortly.

Discriminative Adversarial Learning

To establish a direct interaction between category and domain predictions, we propose a novel source discriminative adversarial loss that is tailored to the design of the integrated classifier $F(\cdot)$. The proposed loss is inspired by the principle of binary cross-entropy loss. It is written as

$$\mathcal{L}^s(G, F) = -\frac{1}{n_s} \sum_{i=1}^{n_s} [(1 - p_{K+1}(\mathbf{x}_i^s)) \log p_{y_i^s}(\mathbf{x}_i^s) + p_{K+1}(\mathbf{x}_i^s) \log(1 - p_{y_i^s}(\mathbf{x}_i^s))]. \quad (3)$$

Intuitively, the proposed loss (3) establishes a mutually inhibitory relation between $p_{y^s}(\mathbf{x}^s)$ of the prediction on the

true category of \mathbf{x}^s , and $p_{K+1}(\mathbf{x}^s)$ of the prediction on the domain of \mathbf{x}^s . We first discuss how the proposed loss (3) works during adversarial training, and we show that under practical conditions, minimizing (3) over the classifier $F(\cdot)$ has the effects of discriminating among task categories while distinguishing the source domain from the target one, and maximizing (3) over the feature extractor $G(\cdot)$ can discriminatively align the source domain to the target one.

Discussion We first write the gradient formulas of \mathcal{L}^s on any source instance \mathbf{x}^s w.r.t. p_{y^s} and p_{K+1}^s as

$$\begin{aligned} \nabla_{p_{y^s}^s} \mathcal{L}^s &= \frac{\partial \mathcal{L}^s}{\partial p_{y^s}^s} = \frac{p_{y^s}^s p_{K+1}^s - (1 - p_{y^s}^s)(1 - p_{K+1}^s)}{p_{y^s}^s (1 - p_{y^s}^s)}, \\ \nabla_{p_{K+1}^s} \mathcal{L}^s &= \frac{\partial \mathcal{L}^s}{\partial p_{K+1}^s} = \log \frac{p_{y^s}^s}{1 - p_{y^s}^s}. \end{aligned}$$

Since both $p_{y^s}^s$ and p_{K+1}^s are among the $K+1$ output probabilities of the classifier $F(G(\mathbf{x}^s))$, we always have $p_{y^s}^s \leq 1 - p_{K+1}^s$ and $p_{K+1}^s \leq 1 - p_{y^s}^s$, suggesting $\nabla_{p_{y^s}^s} \mathcal{L}^s \leq 0$. When the loss (3) is minimized over $F(\cdot)$ via stochastic gradient descent (SGD), we have the update $p_{y^s}^s \leftarrow p_{y^s}^s - \eta \nabla_{p_{y^s}^s} \mathcal{L}^s$ where η is the learning rate, and since $\nabla_{p_{y^s}^s} \mathcal{L}^s \leq 0$, $p_{y^s}^s$ increases; when it is maximized over $G(\cdot)$ via stochastic gradient ascent (SGA), we have the update $p_{y^s}^s \leftarrow p_{y^s}^s + \eta \nabla_{p_{y^s}^s} \mathcal{L}^s$, and since $\nabla_{p_{y^s}^s} \mathcal{L}^s \leq 0$, $p_{y^s}^s$ decreases. Then, we discuss the change of p_{K+1}^s in two cases: (1) in case of $p_{y^s}^s > 0.5$ that guarantees $\nabla_{p_{K+1}^s} \mathcal{L}^s > 0$, when minimizing the loss (3) over $F(\cdot)$ by SGD update $p_{K+1}^s \leftarrow p_{K+1}^s - \eta \nabla_{p_{K+1}^s} \mathcal{L}^s$, we have decreased p_{K+1}^s , and when maximizing it over $G(\cdot)$ by SGA update $p_{K+1}^s \leftarrow p_{K+1}^s + \eta \nabla_{p_{K+1}^s} \mathcal{L}^s$, we have increased p_{K+1}^s ; (2) in case of $p_{y^s}^s < 0.5$ that guarantees $\nabla_{p_{K+1}^s} \mathcal{L}^s < 0$, when minimizing the loss (3) over $F(\cdot)$ by SGD update, we have increased p_{K+1}^s , and when maximizing it over $G(\cdot)$ by SGA update, we have decreased p_{K+1}^s , as shown in Figure 2.

For discriminative adversarial domain adaptation, we expect that (1) when minimizing the proposed loss (3) over $F(\cdot)$, task categories of the source domain is discriminative and the source domain is distinctive from the target one, which can be achieved when $p_{y^s}^s$ increases and p_{K+1}^s decreases; (2) when maximizing it over $G(\cdot)$, the source domain is aligned to the target one while retains discriminability, which can be achieved when $p_{y^s}^s$ decreases and p_{K+1}^s

| Cases | $\min_F \mathcal{L}^s$ | | $\max_G \mathcal{L}^s$ | |
|-------------------|------------------------|-------------|------------------------|-------------|
| | $p_{y^s}^s$ | p_{K+1}^s | $p_{y^s}^s$ | p_{K+1}^s |
| $p_{y^s}^s > 0.5$ | ↑ | ↓ | ↓ | ↑ |
| $p_{y^s}^s < 0.5$ | ↑ | ↑ | ↓ | ↓ |

Figure 2: Changes of $p_{y^s}^s$ and p_{K+1}^s when minimizing and maximizing the loss (3) in the two cases.

increases in the case of $p_{y^s}^s > 0.5$. To meet the expectations, the condition of $p_{y^s}^s > 0.5$ for all source instances should be always satisfied. This is practically achieved by pre-training DADA on the labeled source data using a K -way cross-entropy loss, and maintaining in the adversarial training of DADA the same supervision signal. We present in the supplemental material empirical evidence on benchmark datasets that shows the efficacy of our used scheme.

To achieve the joint distribution alignment, the explicit interplay between category and domain predictions for any target instance should also be created. Motivated by recent works (Pei et al. 2018; Long et al. 2018b) which alleviate the issue of mode collapse by aligning each instance to several most related categories, we propose a target discriminative adversarial loss based on the design of the integrated classifier $F(\cdot)$, by using the conditional category probabilities to weight the domain predictions. It is written as

$$\begin{aligned} \mathcal{L}_F^t(G, F) &= -\frac{1}{n_t} \sum_{j=1}^{n_t} \sum_{k=1}^K \bar{p}_k(\mathbf{x}_j^t) \log \hat{p}_{K+1}^k(\mathbf{x}_j^t) \\ \mathcal{L}_G^t(G, F) &= \frac{1}{n_t} \sum_{j=1}^{n_t} \sum_{k=1}^K \bar{p}_k(\mathbf{x}_j^t) \log(1 - \hat{p}_{K+1}^k(\mathbf{x}_j^t)), \end{aligned} \quad (4)$$

where the k^{th} element of the domain prediction vector $\hat{\mathbf{p}}^k$ for the k^{th} category is written as follows

$$\hat{p}_{k'}^k(\mathbf{x}) = \begin{cases} \frac{p_{k'}(\mathbf{x})}{p_k(\mathbf{x}) + p_{K+1}(\mathbf{x})}, & k' = k, K+1 \\ 0, & \text{otherwise} \end{cases}. \quad (5)$$

An intuitive explanation for our proposed (4) is provided in the supplemental material.

Established knowledge from cluster analysis (Nalewajski 2012) indicates that we can estimate clusters with a low probability of error only if the conditional entropy is small. To this end, we adopt the entropy minimization principle (Grandvalet and Bengio 2005), which is written as

$$\mathcal{L}_{em}^t(G, F) = \frac{1}{n_t} \sum_{j=1}^{n_t} \mathcal{H}(\bar{\mathbf{p}}(\mathbf{x}_j^t)), \quad (6)$$

where $\mathcal{H}(\cdot)$ computes the entropy of a probability vector. Combining (3), (4), and (6) gives the following minimax problem of our proposed DADA

$$\begin{aligned} \min_F \mathcal{L}_F &= \lambda(\mathcal{L}^s + \mathcal{L}_F^t) - \mathcal{L}_{em}^t \\ \max_G \mathcal{L}_G &= \lambda(\mathcal{L}^s + \mathcal{L}_G^t) - \mathcal{L}_{em}^t, \end{aligned} \quad (7)$$

where λ is a hyper-parameter that trade-offs the adversarial domain adaptation objective with the entropy minimization one in the unified optimization problem. Note that in the minimization problem of (7), \mathcal{L}_{em}^t serves as a regularizer for learning $F(\cdot)$ to avoid the trivial solution (i.e. all instances are assigned to the same category), and in the maximization problem of (7), it helps learn more target-discriminative features, which can alleviate the negative effect of adversarial feature adaptation on the adaptability (Liu et al. 2019).

By optimizing (7), the joint distribution alignment can be enhanced. This ability comes from the better use of discriminative information from both the source and target domains. Concretely, DADA constrains the domain classifier so that it clearly/explicitly knows the classification boundary, thus reducing false alignment between different categories. By deceiving such a strong domain classifier, DADA can learn a feature extractor that better aligns the two domains. *We also theoretically prove in the supplemental material that DADA can better bound the expected target error.*

Extension for Partial Domain Adaptation

Partial domain adaptation is a more realistic setting, where the target label space is subsumed by the source one. The false alignment between the outlier source categories and the target domain is unavoidable. To address it, existing methods (Cao et al. 2018a; Zhang et al. 2018a; Cao et al. 2018b) utilize the category or domain predictions, to decrease the contribution of source outliers to the training of task or domain classifiers. Inspired by these ideas, we extend DADA for partial domain adaptation by using a reliable category-level weighting mechanism, which is termed DADA-P.

Concretely, we average the conditional probability vectors $\bar{\mathbf{p}}(\mathbf{x}^t) \in [0, 1]^K$ over all target data and then normalize the averaged vector $\bar{\mathbf{c}} \in [0, 1]^K$ by dividing its largest element. The category weight vector $\mathbf{c} \in [0, 1]^K$ with c_k as its k^{th} element is derived by a convex combination of the normalized vector and an all-ones vector $\mathbf{1}$, as follows

$$\begin{aligned} \bar{\mathbf{c}} &= \frac{1}{n_t} \sum_{j=1}^{n_t} \bar{\mathbf{p}}(\mathbf{x}_j^t) \\ \mathbf{c} &= \lambda \frac{\bar{\mathbf{c}}}{\max(\bar{\mathbf{c}})} + (1 - \lambda)\mathbf{1}, \end{aligned} \quad (8)$$

where $\lambda \in [0, 1]$ is to suppress the detection noise of outlier source categories in the early stage of training. Then, we apply the category weight vector \mathbf{c} to the proposed discriminative adversarial loss for any source instance, leading to

$$\begin{aligned} \mathcal{L}^s(G, F) &= -\frac{1}{n_s} \sum_{i=1}^{n_s} c_{y_i^s} [(1 - p_{K+1}(\mathbf{x}_i^s)) \log p_{y_i^s}(\mathbf{x}_i^s) \\ &\quad + p_{K+1}(\mathbf{x}_i^s) \log(1 - p_{y_i^s}(\mathbf{x}_i^s))]. \end{aligned} \quad (9)$$

Since predicted probabilities on the outlier source categories are more likely to increase when minimizing $-\mathcal{L}_{em}^t$ over $F(\cdot)$, which incurs negative transfer. To avoid it, we minimize \mathcal{L}_{em}^t over $F(\cdot)$ and the objective of DADA-P is

$$\begin{aligned} \min_F \mathcal{L}_F &= \lambda(\mathcal{L}^s + \mathcal{L}_F^t) + \mathcal{L}_{em}^t \\ \max_G \mathcal{L}_G &= \lambda(\mathcal{L}^s + \mathcal{L}_G^t) - \mathcal{L}_{em}^t. \end{aligned} \quad (10)$$

By optimizing it, DADA-P can simultaneously alleviate negative transfer and promote the joint distribution alignment across domains in the shared label space.

Extension for Open Set Domain Adaptation

Open set domain adaptation is a very challenging setting, where the source label space is subsumed by the target one. We denominate the shared category and all unshared categories between the two domains as the “known category” and “unknown category” respectively. The goal of open set domain adaptation is to correctly classify any target instance as the known or unknown category. The false alignment between the known and unknown categories is inevitable. To this end, the work (Saito et al. 2018c) proposes to make a pseudo decision boundary for the unknown category, which enables the feature extractor to reject some target instances as outliers. Inspired by this work, we extend DADA for open set domain adaptation by training the classifier to classify all target instances as the unknown category with a small probability q , which is termed DADA-O. Assuming the predicted probability on the unknown category as the K^{th} element of $\mathbf{p}(\mathbf{x}^t)$, i.e., $p_K(\mathbf{x}^t)$, the modified target adversarial loss when minimized over the integrated classifier $F(\cdot)$ is

$$\begin{aligned} \mathcal{L}_F^t(G, F) = & \\ & - \frac{1}{n_t} \sum_{j=1}^{n_t} q \log p_K(\mathbf{x}_j^t) - (1 - q) \log p_{K+1}(\mathbf{x}_j^t), \end{aligned} \quad (11)$$

where $0 < q < 0.5$. When maximized over the feature extractor $G(\cdot)$, we still use the discriminative loss \mathcal{L}_G^t in (4). Replacing \mathcal{L}_F^t in (7) with (11) gives the overall adversarial objective of DADA-O, which can achieve a balance between domain adaptation and outlier rejection.

We utilize all target instances to obtain the concept of “unknown”, which is very helpful for the classification of unknown target instances as the unknown category but can cause the misclassification of known target instances as the unknown category. This issue can be alleviated by selecting an appropriate q . If q is too small, the unknown target instances cannot be correctly classified; if q is too large, the known target instances can be misclassified. By choosing an appropriate q , the feature extractor can separate the unknown target instances from the known ones while aligning the joint distributions in the shared label space.

Experiments

Datasets and Implementation Details

Office-31 (Saenko et al. 2010) is a popular benchmark domain adaptation dataset consisting of 4,110 images of 31 categories collected from three domains: Amazon (**A**), Webcam (**W**), and DSLR (**D**). We evaluate on six settings.

Syn2Real (Peng et al. 2018) is the largest benchmark. Syn2Real-C has over 280K images of 12 shared categories in the combined training, validation, and testing domains. The 152,397 images on the training domain are synthetic ones by rendering 3D models. The validation and test domains comprise real images, and the validation one has

55,388 images. We use the training domain as the source domain and validation one as the target domain. For partial domain adaptation, we choose images of the first 6 categories (in alphabetical order) in the validation domain as the target domain and form the setting: **Synthetic 12** \rightarrow **Real 6**. For open set domain adaptation, we evaluate on Syn2Real-O, which includes two domains. The training/synthetic domain uses synthetic images from the 12 categories of Syn2Real-C as “known”. The validation/real domain uses images of the 12 categories from the validation domain of Syn2Real-C as “known”, and 50k images from 69 other categories as “unknown”. We use the training and validation domains of Syn2Real-O as the source and target domains respectively.

Implementation Details We follow standard evaluation protocols for unsupervised domain adaptation (Ganin et al. 2016; Wang et al. 2019): we use all labeled source and all unlabeled target instances as the training data. For all tasks of Office-31 and **Synthetic 12** \rightarrow **Real 6**, based on ResNet-50 (He et al. 2016), we report the classification result on the target domain of mean(\pm standard deviation) over three random trials. For other tasks of Syn2Real, we evaluate the accuracy of each category based on ResNet-101 and ResNet-152 (for closed and open set domain adaptation respectively). For each base network, we use all its layers up to the second last one as the feature extractor $G(\cdot)$, and set the neuron number of its last FC layer as $K + 1$ to have the integrated classifier $F(\cdot)$. Exceptionally, we follow the work (Peng et al. 2018) and replace the last FC layer of ResNet-152 with three FC layers of 512 neurons. All base networks are pre-trained on ImageNet (Russakovsky et al. 2015). We firstly pre-train them on the labeled source data, and then fine-tune them on both the labeled source data and unlabeled target data via adversarial training, where we maintain the same supervision signal as the pre-training.

We follow DANN (Ganin et al. 2016) to use the SGD training schedule: the learning rate is adjusted by $\eta_p = \frac{\eta_0}{(1+\alpha p)^\beta}$, where p denotes the process of training iterations that is normalized to be in $[0, 1]$, and we set $\eta_0 = 0.0001$, $\alpha = 10$, and $\beta = 0.75$; the hyper-parameter λ is initialized at 0 and is gradually increased to 1 by $\lambda_p = \frac{2}{1+\exp(-\gamma p)} - 1$, where we set $\gamma = 10$. We empirically set $q = 0.1$. We implement all our methods by **PyTorch**. The code will be available at <https://github.com/huitangtang/DADA-AAAI2020>.

Analysis

Ablation Study We conduct ablation studies on Office-31 to investigate the effects of key components of our proposed DADA based on ResNet-50. Our ablation studies start with the very baseline termed “No Adaptation” that simply fine-tunes a ResNet-50 on the source data. To validate the mutually inhibitory relation enabled by DADA, we use DANN (Ganin et al. 2016) and DANN-CA (Tran et al. 2019) respectively as the second and third baselines. To investigate how the entropy minimization principle helps learn more target-discriminative features, we remove the entropy minimization loss (6) from our main minimax problem (7), denoted as “DADA (w/o em)”. To know effects of the proposed source and target discriminative adversarial losses (3) and (4), we

Table 1: Ablation studies using Office-31 based on ResNet-50. Please refer to the main text for how they are defined.

| Methods | A \rightarrow W | D \rightarrow W | W \rightarrow D | A \rightarrow D | D \rightarrow A | W \rightarrow A | Avg |
|------------------------|--------------------------------|--------------------------------|---------------------------------|--------------------------------|--------------------------------|--------------------------------|-------------|
| No Adaptation | 79.9 \pm 0.3 | 96.8 \pm 0.4 | 99.5 \pm 0.1 | 84.1 \pm 0.4 | 64.5 \pm 0.3 | 66.4 \pm 0.4 | 81.9 |
| DANN | 81.2 \pm 0.3 | 98.0 \pm 0.2 | 99.8 \pm 0.0 | 83.3 \pm 0.3 | 66.8 \pm 0.3 | 66.1 \pm 0.3 | 82.5 |
| DANN-CA | 85.4 \pm 0.4 | 98.2 \pm 0.2 | 99.8 \pm 0.0 | 87.1 \pm 0.4 | 68.5 \pm 0.2 | 67.6 \pm 0.3 | 84.4 |
| DADA (w/o em + w/o td) | 91.0 \pm 0.2 | 98.7 \pm 0.1 | 100.0\pm0.0 | 90.8 \pm 0.2 | 70.9 \pm 0.3 | 70.2 \pm 0.3 | 86.9 |
| DADA (w/o em) | 91.8 \pm 0.1 | 99.0 \pm 0.1 | 100.0\pm0.0 | 92.5 \pm 0.3 | 72.8 \pm 0.2 | 72.3 \pm 0.3 | 88.1 |
| DADA | 92.3\pm0.1 | 99.2\pm0.1 | 100.0\pm0.0 | 93.9\pm0.2 | 74.4\pm0.1 | 74.2\pm0.1 | 89.0 |

Table 2: Results for closed set domain adaptation on Office-31 based on ResNet-50. Note that SimNet is implemented by an **unknown** framework; MADA and DANN-CA are implemented by **Caffe**; all the other methods are implemented by **PyTorch**.

| Methods | A \rightarrow W | D \rightarrow W | W \rightarrow D | A \rightarrow D | D \rightarrow A | W \rightarrow A | Avg |
|------------------------------------|--------------------------------|--------------------------------|---------------------------------|--------------------------------|--------------------------------|--------------------------------|-------------|
| No Adaptation (He et al. 2016) | 79.9 \pm 0.3 | 96.8 \pm 0.4 | 99.5 \pm 0.1 | 84.1 \pm 0.4 | 64.5 \pm 0.3 | 66.4 \pm 0.4 | 81.9 |
| DAN (Long et al. 2018a) | 81.3 \pm 0.3 | 97.2 \pm 0.0 | 99.8 \pm 0.0 | 83.1 \pm 0.2 | 66.3 \pm 0.0 | 66.3 \pm 0.1 | 82.3 |
| DANN (Ganin et al. 2016) | 81.2 \pm 0.3 | 98.0 \pm 0.2 | 99.8 \pm 0.0 | 83.3 \pm 0.3 | 66.8 \pm 0.3 | 66.1 \pm 0.3 | 82.5 |
| ADDA (Tzeng et al. 2017) | 86.2 \pm 0.5 | 96.2 \pm 0.3 | 98.4 \pm 0.3 | 77.8 \pm 0.3 | 69.5 \pm 0.4 | 68.9 \pm 0.5 | 82.9 |
| MADA (Pei et al. 2018) | 90.0 \pm 0.1 | 97.4 \pm 0.1 | 99.6 \pm 0.1 | 87.8 \pm 0.2 | 70.3 \pm 0.3 | 66.4 \pm 0.3 | 85.2 |
| VADA (Shu et al. 2018) | 86.5 \pm 0.5 | 98.2 \pm 0.4 | 99.7 \pm 0.2 | 86.7 \pm 0.4 | 70.1 \pm 0.4 | 70.5 \pm 0.4 | 85.4 |
| DANN-CA (Tran et al. 2019) | 91.35 | 98.24 | 99.48 | 89.94 | 69.63 | 68.76 | 86.2 |
| GTA (Sankaranarayanan et al. 2018) | 89.5 \pm 0.5 | 97.9 \pm 0.3 | 99.8 \pm 0.4 | 87.7 \pm 0.5 | 72.8 \pm 0.3 | 71.4 \pm 0.4 | 86.5 |
| MCD (Saito et al. 2018b) | 88.6 \pm 0.2 | 98.5 \pm 0.1 | 100.0\pm0.0 | 92.2 \pm 0.2 | 69.5 \pm 0.1 | 69.7 \pm 0.3 | 86.5 |
| CDAN+E (Long et al. 2018b) | 94.1\pm0.1 | 98.6 \pm 0.1 | 100.0\pm0.0 | 92.9 \pm 0.2 | 71.0 \pm 0.3 | 69.3 \pm 0.3 | 87.7 |
| TADA (Wang et al. 2019) | 94.3 \pm 0.3 | 98.7 \pm 0.1 | 99.8 \pm 0.2 | 91.6 \pm 0.3 | 72.9 \pm 0.2 | 73.0 \pm 0.3 | 88.4 |
| SymNets (Zhang et al. 2019) | 90.8 \pm 0.1 | 98.8 \pm 0.3 | 100.0\pm0.0 | 93.9\pm0.5 | 74.6\pm0.6 | 72.5 \pm 0.5 | 88.4 |
| TAT (Liu et al. 2019) | 92.5 \pm 0.3 | 99.3\pm0.1 | 100.0\pm0.0 | 93.2 \pm 0.2 | 73.1 \pm 0.3 | 72.1 \pm 0.3 | 88.4 |
| DADA | 92.3 \pm 0.1 | 99.2 \pm 0.1 | 100.0\pm0.0 | 93.9\pm0.2 | 74.4 \pm 0.1 | 74.2\pm0.1 | 89.0 |

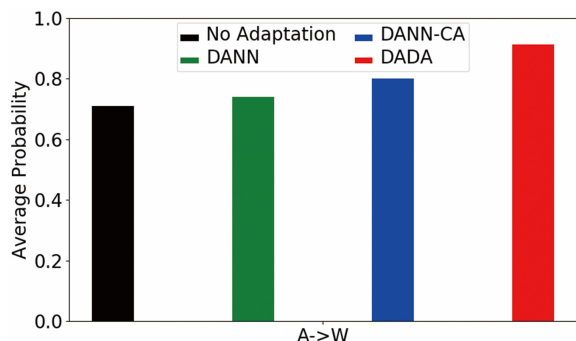


Figure 3: Average probability on the true category over all target instances by task classifiers of different methods.

remove both (6) and (4) from (7), denoted as “DADA (w/o em + w/o td)”.

Results in Table 1 show that although DANN improves over “No Adaptation”, its result is much worse than DANN-CA, verifying the efficacy of the design of the integrated classifier $F(\cdot)$. “DADA (w/o em + w/o td)” improves over DANN-CA and “DADA (w/o em)” improves over “DADA (w/o em + w/o td)”, showing the efficacy of our proposed discriminative adversarial learning. DADA significantly outperforms DANN and DANN-CA, confirming the efficacy of the proposed mutually inhibitory relation between the category and domain predictions in aligning the joint distribu-

tions of feature and category across domains. Table 1 also confirms that entropy minimization is helpful to learn more target-discriminative features.

Quantitative Comparison To compare the efficacy of different methods in reducing domain discrepancy at the category level, we visualize the average probability on the true category over all target instances by task classifiers of No Adaptation, DANN, DANN-CA, and DADA on $A \rightarrow W$ in Figure 3. Note that here we use labels of the target data for the quantization of category-level domain discrepancy. Figure 3 shows that our proposed DADA gives the predicted probability on the true category of any target instance a better chance to approach 1, meaning that target instances are more likely to be correctly classified by DADA, i.e., a better category-level domain alignment.

Results

Closed Set Domain Adaptation We compare in Tables 2 and 3 our proposed method with existing ones on Office-31 and Syn2Real-C based on ResNet-50 and ResNet-101 respectively. Whenever available, results of existing methods are quoted from their respective papers or the recent works (Pei et al. 2018; Long et al. 2018b; Liu et al. 2019; Saito et al. 2018b). Our proposed DADA outperforms existing methods, testifying the efficacy of DADA in aligning the joint distributions of feature and category across domains.

Partial Domain Adaptation We compare in Table 5 our proposed method to existing ones on Syn2Real-C based on

Table 3: Results for closed set domain adaptation on Syn2Real-C based on ResNet-101. Note that all compared methods are based on **PyTorch** implementation.

| Methods | plane | bcycl | bus | car | horse | knife | mcycl | person | plant | sktbrd | train | truck | mean |
|--------------------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| No Adaptation (He et al. 2016) | 55.1 | 53.3 | 61.9 | 59.1 | 80.6 | 17.9 | 79.7 | 31.2 | 81.0 | 26.5 | 73.5 | 8.5 | 52.4 |
| DANN (Ganin et al. 2016) | 81.9 | 77.7 | 82.8 | 44.3 | 81.2 | 29.5 | 65.1 | 28.6 | 51.9 | 54.6 | 82.8 | 7.8 | 57.4 |
| DAN (Long et al. 2018a) | 87.1 | 63.0 | 76.5 | 42.0 | 90.3 | 42.9 | 85.9 | 53.1 | 49.7 | 36.3 | 85.8 | 20.7 | 61.1 |
| MCD (Saito et al. 2018b) | 87.0 | 60.9 | 83.7 | 64.0 | 88.9 | 79.6 | 84.7 | 76.9 | 88.6 | 40.3 | 83.0 | 25.8 | 71.9 |
| GPDA (Kim et al. 2019) | 83.0 | 74.3 | 80.4 | 66.0 | 87.6 | 75.3 | 83.8 | 73.1 | 90.1 | 57.3 | 80.2 | 37.9 | 73.3 |
| ADR (Saito et al. 2018a) | 87.8 | 79.5 | 83.7 | 65.3 | 92.3 | 61.8 | 88.9 | 73.2 | 87.8 | 60.0 | 85.5 | 32.3 | 74.8 |
| DADA | 92.9 | 74.2 | 82.5 | 65.0 | 90.9 | 93.8 | 87.2 | 74.2 | 89.9 | 71.5 | 86.5 | 48.7 | 79.8 |

Table 4: Results for open set domain adaptation on Syn2Real-O based on ResNet-152. *Known* indicates the mean classification result over the known categories whereas *Mean* also includes the unknown category. The table below shows the results when the Known-to-Unknown Ratio in the target domain is set to 1 : 10. All compared methods are based on **PyTorch** implementation.

| Known-to-Unknown Ratio = 1 : 1 | | | | | | | | | | | | | | | |
|---------------------------------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|
| Methods | plane | bcycl | bus | car | horse | knife | mcycl | person | plant | sktbrd | train | truck | unk | Known | Mean |
| No Adaptation (He et al. 2016) | 49 | 20 | 29 | 47 | 62 | 27 | 79 | 3 | 37 | 19 | 70 | 1 | 62 | 36 | 38 |
| DAN (Long et al. 2018a) | 51 | 40 | 42 | 56 | 68 | 24 | 75 | 2 | 39 | 30 | 71 | 2 | 75 | 41 | 44 |
| DANN (Ganin et al. 2016) | 59 | 41 | 16 | 54 | 77 | 18 | 88 | 4 | 44 | 32 | 68 | 4 | 61 | 42 | 43 |
| AODA (Saito et al. 2018c) | 85 | 71 | 65 | 53 | 83 | 10 | 79 | 36 | 73 | 56 | 79 | 32 | 87 | 60 | 62 |
| DADA-O | 88 | 76 | 76 | 64 | 79 | 46 | 91 | 62 | 52 | 63 | 86 | 8 | 55 | 66 | 65 |
| Known-to-Unknown Ratio = 1 : 10 | | | | | | | | | | | | | | | |
| AODA (Saito et al. 2018c) | 80 | 63 | 59 | 63 | 83 | 12 | 89 | 5 | 61 | 14 | 79 | 0 | 69 | 51 | 52 |
| DADA-O | 77 | 63 | 75 | 71 | 38 | 33 | 92 | 58 | 47 | 50 | 89 | 1 | 50 | 58 | 57 |

Table 5: Results for partial domain adaptation on Syn2Real-C based on ResNet-50. Note that all compared methods are based on **PyTorch** implementation.

| Methods | Synthetic 12→Real 6 |
|--------------------------------|---------------------|
| No Adaptation (He et al. 2016) | 45.26 |
| DAN (Long et al. 2018a) | 47.60 |
| DANN (Ganin et al. 2016) | 51.01 |
| RTN (Long et al. 2016) | 50.04 |
| PADA (Cao et al. 2018b) | 53.53 |
| DADA-P | 69.06 |

ResNet-50. Results of existing methods are quoted from the work (Cao et al. 2018b). Our proposed DADA-P substantially outperforms all comparative methods by +15.53%, showing the effectiveness of DADA-P on reducing the negative influence of source outliers while promoting the joint distribution alignment in the shared label space.

Open Set Domain Adaptation We compare in Table 4 our proposed method with existing ones on Syn2Real-O based on ResNet-152. Results of existing methods are quoted from the recent work (Peng et al. 2018). Our proposed DADA-O outperforms all comparative methods in both evaluation metrics of Known and Mean, showing the efficacy of DADA-O in both aligning joint distributions of the known instances and identifying the unknown target instances. It is noteworthy that DADA-O improves over the state-of-the-art method AODA by a large margin when the known-to-unknown ratio in the target domain is much smaller than 1, i.e. the false alignment between the known source and unknown target instances will be much more serious. This observation confirms the efficacy of DADA-O.

We provide more results and analysis for the three problem settings in the supplemental material.

Conclusion

We propose a novel adversarial learning method termed Discriminative Adversarial Domain Adaptation (DADA) to overcome the limitation in aligning the joint distributions of feature and category across domains, which is due to an issue of mode collapse induced by the separate design of task and domain classifiers. Based on an integrated task and domain classifier, DADA has a novel adversarial objective that encourages a mutually inhibitory relation between the category and domain predictions, which can promote the joint distribution alignment. Unlike previous methods, DADA explicitly enables a discriminative interaction between category and domain predictions. Except for closed set domain adaptation, we also extend DADA for more challenging problem settings of partial and open set domain adaptation. Experiments on benchmark datasets testify the efficacy of our proposed methods for all the three settings.

Acknowledgments

This work is supported in part by the National Natural Science Foundation of China (Grant No.: 61771201), the Program for Guangdong Introducing Innovative and Entrepreneurial Teams (Grant No.: 2017ZT07X183), and the Guangdong R&D key project of China (Grant No.: 2019B010155001).

References

- Ben-David, S.; Blitzer, J.; Crammer, K.; and Pereira, F. 2007. Analysis of representations for domain adaptation. In Schölkopf, B.; Platt, J. C.; and Hoffman, T., eds., *Advances in Neural Information Processing Systems*. MIT Press. 137–144.
- Ben-David, S.; Blitzer, J.; Crammer, K.; Kulesza, A.; Pereira, F.; and Vaughan, J. W. 2010. A theory of learning from different domains. *Machine Learning* 79(1):151–175.

- Bousmalis, K.; Trigeorgis, G.; Silberman, N.; Krishnan, D.; and Erhan, D. 2016. Domain separation networks. In *Advances in Neural Information Processing Systems*.
- Busto, P. P.; Iqbal, A.; and Gall, J. 2018. Open set domain adaptation for image and action recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 1–1.
- Cao, Z.; Long, M.; Wang, J.; and Jordan, M. I. 2018a. Partial transfer learning with selective adversarial networks. In *Computer Vision and Pattern Recognition*.
- Cao, Z.; Ma, L.; Long, M.; and Wang, J. 2018b. Partial adversarial domain adaptation. In *European Conference on Computer Vision*.
- Chen, Q.; Liu, Y.; Wang, Z.; Wassell, I.; and Chetty, K. 2018. Re-weighted adversarial adaptation network for unsupervised domain adaptation. In *Computer Vision and Pattern Recognition*.
- Chen, C.; Chen, Z.; Jiang, B.; and Jin, X. 2019a. Joint domain alignment and discriminative feature learning for unsupervised deep domain adaptation. In *Association for the Advancement of Artificial Intelligence*.
- Chen, C.; Xie, W.; Xu, T.; Huang, W.; Rong, Y.; Ding, X.; Huang, Y.; and Huang, J. 2019b. Progressive feature alignment for unsupervised domain adaptation. *Computer Vision and Pattern Recognition*.
- Courty, N.; Flamary, R.; Habrard, A.; and Rakotomamonjy, A. 2017. Joint distribution optimal transportation for domain adaptation. In Guyon, I.; Luxburg, U. V.; Bengio, S.; Wallach, H.; Fergus, R.; Vishwanathan, S.; and Garnett, R., eds., *Advances in Neural Information Processing Systems 30*. Curran Associates, Inc. 3730–3739.
- Cybenko, G. 1989. Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals and Systems* 2(4):303–314.
- Dai, Z.; Yang, Z.; Yang, F.; Cohen, W. W.; and Salakhutdinov, R. R. 2017. Good semi-supervised learning that requires a bad gan. In Guyon, I.; Luxburg, U. V.; Bengio, S.; Wallach, H.; Fergus, R.; Vishwanathan, S.; and Garnett, R., eds., *Advances in Neural Information Processing Systems 30*. Curran Associates, Inc. 6510–6520.
- Ganin, Y.; Ustinova, E.; Ajakan, H.; Germain, P.; Larochelle, H.; Laviolette, F.; Marchand, M.; and Lempitsky, V. 2016. Domain-adversarial training of neural networks. *J. Mach. Learn. Res.* 17(1):2096–2030.
- Ghifary, M.; Kleijn, W. B.; Zhang, M.; Balduzzi, D.; and Li, W. 2016. Deep reconstruction-classification networks for unsupervised domain adaptation. In *The European Conference on Computer Vision*.
- Goodfellow, I.; Bengio, Y.; and Courville, A. 2016. *Deep Learning*. MIT Press. <http://www.deeplearningbook.org>.
- Goodfellow, I.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A.; and Bengio, Y. 2014. Generative adversarial nets. In Ghahramani, Z.; Welling, M.; Cortes, C.; Lawrence, N. D.; and Weinberger, K. Q., eds., *Advances in Neural Information Processing Systems 27*. Curran Associates, Inc. 2672–2680.
- Grandvalet, Y., and Bengio, Y. 2005. Semi-supervised learning by entropy minimization. In Saul, L. K.; Weiss, Y.; and Bottou, L., eds., *Advances in Neural Information Processing Systems 17*. MIT Press. 529–536.
- Griffin, G.; Holub, A.; and Perona, P. 2007. Caltech-256 object category dataset. *CalTech Report*.
- Haeusser, P.; Frerix, T.; Mordvintsev, A.; and Cremers, D. 2017. Associative domain adaptation. In *International Conference on Computer Vision*.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *Computer Vision and Pattern Recognition*.
- Hoffman, J.; Tzeng, E.; Park, T.; Zhu, J.-Y.; Isola, P.; Saenko, K.; Efros, A.; and Darrell, T. 2018. CyCADA: Cycle-consistent adversarial domain adaptation. In *International Conference on Machine Learning*.
- Hornik, K.; Stinchcombe, M.; and White, H. 1989. Multilayer feedforward networks are universal approximators. *Neural Networks* 2(5):359–366.
- Hull, J. J. 1994. A database for handwritten text recognition research. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 16(5):550–554.
- Jain, L. P.; Scheirer, W. J.; and Boulton, T. E. 2014. Multi-class open set recognition using probability of inclusion. In *European Conference on Computer Vision*.
- Kim, M.; Sahu, P.; Gholami, B.; and Pavlovic, V. 2019. Unsupervised visual domain adaptation: A deep max-margin gaussian process approach. In *Computer Vision and Pattern Recognition*.
- Krizhevsky, A.; Sutskever, I.; and Hinton, G. E. 2012. ImageNet classification with deep convolutional neural networks. In Pereira, F.; Burges, C. J. C.; Bottou, L.; and Weinberger, K. Q., eds., *Advances in Neural Information Processing Systems 25*. Curran Associates, Inc. 1097–1105.
- Kurmi, V. K., and Namboodiri, V. P. 2019. Looking back at labels: A class based domain adaptation technique. *ArXiv abs/1904.01341*.
- Lecun, Y.; Bottou, L.; Bengio, Y.; and Haffner, P. 1998. Gradient-based learning applied to document recognition. *Proceedings of the IEEE* 86(11):2278–2324.
- Lee, C.-Y.; Batra, T.; Baig, M. H.; and Ulbricht, D. 2019. Sliced wasserstein discrepancy for unsupervised domain adaptation. In *Computer Vision and Pattern Recognition*.
- Li, S.; Song, S.; and Wu, C. 2018. Layer-wise domain correction for unsupervised domain adaptation. *Frontiers of Information Technology and Electronic Engineering* 19:91–103.
- Liu, M.-Y., and Tuzel, O. 2016. Coupled generative adversarial networks. In *Advances in Neural Information Processing Systems*.
- Liu, H.; Long, M.; Wang, J.; and Jordan, M. 2019. Transferable adversarial training: A general approach to adapting deep classifiers. In *International Conference on Machine Learning*.
- Long, M.; Zhu, H.; Wang, J.; and Jordan, M. I. 2016. Unsupervised domain adaptation with residual transfer networks. In *Advances in Neural Information Processing Systems*. Curran Associates Inc.
- Long, M.; Zhu, H.; Wang, J.; and Jordan, M. I. 2017. Deep transfer learning with joint adaptation networks. In *International Conference on Machine Learning*.
- Long, M.; Cao, Y.; Cao, Z.; Wang, J.; and Jordan, M. I. 2018a. Transferable representation learning with deep adaptation networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 1–1.
- Long, M.; CAO, Z.; Wang, J.; and Jordan, M. I. 2018b. Conditional adversarial domain adaptation. In Bengio, S.; Wallach, H.; Larochelle, H.; Grauman, K.; Cesa-Bianchi, N.; and Garnett, R., eds., *Advances in Neural Information Processing Systems 31*. Curran Associates, Inc. 1640–1650.
- Luo, Z.; Zou, Y.; Hoffman, J.; and Fei-Fei, L. F. 2017. Label efficient learning of transferable representations across domains and tasks. In Guyon, I.; Luxburg, U. V.; Bengio, S.; Wallach, H.; Fergus, R.; Vishwanathan, S.; and Garnett, R., eds., *Advances in Neural Information Processing Systems 30*. Curran Associates, Inc. 165–177.

- Mansour, Y.; Mohri, M.; and Rostamizadeh, A. 2009. Domain adaptation: Learning bounds and algorithms. *COLT 2009 - The 22nd Conference on Learning Theory*.
- Morerio, P.; Cavazza, J.; and Murino, V. 2018. Minimal-entropy correlation alignment for unsupervised deep domain adaptation. In *International Conference on Learning Representations*.
- Nalewajski, R. F. 2012. *Elements of Information Theory*. Berlin, Heidelberg: Springer Berlin Heidelberg. 371–395.
- Netzer, Y.; Wang, T.; Coates, A.; Bissacco, A.; Wu, B.; and Ng, A. Y. 2011. Reading digits in natural images with unsupervised feature learning. In *NIPS Workshop on Deep Learning and Unsupervised Feature Learning 2011*.
- Pan, S. J., and Yang, Q. 2010. A survey on transfer learning. *IEEE Transactions on Knowledge and Data Engineering* 22:1345–1359.
- Pan, Y.; Yao, T.; Li, Y.; Wang, Y.; Ngo, C.-W.; and Mei, T. 2019. Transferrable prototypical networks for unsupervised domain adaptation. *Computer Vision and Pattern Recognition*.
- Pei, Z.; Cao, Z.; Long, M.; and Wang, J. 2018. Multi-adversarial domain adaptation. In *Association for the Advancement of Artificial Intelligence*.
- Peng, X.; Usman, B.; Saito, K.; Kaushik, N.; Hoffman, J.; and Saenko, K. 2018. Syn2real: A new benchmark for synthetic-to-real visual domain adaptation. *ArXiv abs/1806.09755*.
- Pinheiro, P. O. 2018. Unsupervised domain adaptation with similarity learning. In *Computer Vision and Pattern Recognition*.
- Russakovsky, O.; Deng, J.; Su, H.; Krause, J.; Satheesh, S.; Ma, S.; Huang, Z.; Karpathy, A.; Khosla, A.; Bernstein, M.; et al. 2015. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision* 115(3):211–252.
- Russo, P.; Carlucci, F. M.; Tommasi, T.; and Caputo, B. 2018. From source to target and back: Symmetric bi-directional adaptive gan. In *Computer Vision and Pattern Recognition*.
- Saenko, K.; Kulis, B.; Fritz, M.; and Darrell, T. 2010. Adapting visual category models to new domains. In *European Conference on Computer Vision*.
- Saito, K.; Ushiku, Y.; Harada, T.; and Saenko, K. 2018a. Adversarial dropout regularization. In *International Conference on Learning Representations*.
- Saito, K.; Watanabe, K.; Ushiku, Y.; and Harada, T. 2018b. Maximum classifier discrepancy for unsupervised domain adaptation. In *Computer Vision and Pattern Recognition*.
- Saito, K.; Yamamoto, S.; Ushiku, Y.; and Harada, T. 2018c. Open set domain adaptation by backpropagation. In *European Conference on Computer Vision*.
- Salimans, T.; Goodfellow, I.; Zaremba, W.; Cheung, V.; Radford, A.; Chen, X.; and Chen, X. 2016. Improved techniques for training gans. In Lee, D. D.; Sugiyama, M.; Luxburg, U. V.; Guyon, I.; and Garnett, R., eds., *Advances in Neural Information Processing Systems* 29. Curran Associates, Inc. 2234–2242.
- Sankaranarayanan, S.; Balaji, Y.; Castillo, C. D.; and Chellappa, R. 2018. Generate to adapt: Aligning domains using generative adversarial networks. In *Computer Vision and Pattern Recognition*.
- Shu, R.; Bui, H.; Narui, H.; and Ermon, S. 2018. A DIRT-t approach to unsupervised domain adaptation. In *International Conference on Learning Representations*.
- Tran, L.; Sohn, K.; Yu, X.; Liu, X.; and Chandraker, M. K. 2019. Gotta adapt 'em all: Joint pixel and feature-level domain adaptation for recognition in the wild. In *Computer Vision and Pattern Recognition*.
- Tzeng, E.; Hoffman, J.; Zhang, N.; Saenko, K.; and Darrell, T. 2014. Deep domain confusion: Maximizing for domain invariance. *CoRR abs/1412.3474*.
- Tzeng, E.; Hoffman, J.; Darrell, T.; and Saenko, K. 2015. Simultaneous deep transfer across domains and tasks. In *International Conference on Computer Vision*.
- Tzeng, E.; Hoffman, J.; Saenko, K.; and Darrell, T. 2017. Adversarial discriminative domain adaptation. In *Computer Vision and Pattern Recognition*.
- van der Maaten, L., and Hinton, G. 2008. Visualizing data using t-sne. *Journal of Machine Learning Research* 9 (Nov):2579–2605.
- Venkateswara, H.; Eusebio, J.; Chakraborty, S.; and Panchanathan, S. 2017. Deep hashing network for unsupervised domain adaptation. *Computer Vision and Pattern Recognition* 5385–5394.
- Wang, X.; Li, L.; Ye, W.; Long, M.; and Wang, J. 2019. Transferable attention for domain adaptation. In *Association for the Advancement of Artificial Intelligence*.
- Wen, J.; Liu, R.; Zheng, N.; Zheng, Q.; Gong, Z.; and Yuan, J. 2019. Exploiting local feature patterns for unsupervised domain adaptation. In *Association for the Advancement of Artificial Intelligence*.
- Xie, S.; Zheng, Z.; Chen, L.; and Chen, C. 2018. Learning semantic representations for unsupervised domain adaptation. In Dy, J., and Krause, A., eds., *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, 5423–5432. Stockholm, Sweden: PMLR.
- Yosinski, J.; Clune, J.; Bengio, Y.; and Lipson, H. 2014. How transferable are features in deep neural networks? In Ghahramani, Z.; Welling, M.; Cortes, C.; Lawrence, N. D.; and Weinberger, K. Q., eds., *Advances in Neural Information Processing Systems*. Curran Associates, Inc. 3320–3328.
- Zhang, J.; Ding, Z.; Li, W.; and Ogunbona, P. 2018a. Importance weighted adversarial nets for partial domain adaptation. In *Computer Vision and Pattern Recognition*.
- Zhang, W.; Ouyang, W.; Li, W.; and Xu, D. 2018b. Collaborative and adversarial network for unsupervised domain adaptation. In *Computer Vision and Pattern Recognition*.
- Zhang, Y.; Tang, H.; Jia, K.; and Tan, M. 2019. Domain-symmetric networks for adversarial domain adaptation. In *Computer Vision and Pattern Recognition*.
- Zhang, Y.; Tang, H.; and Jia, K. 2018. Fine-grained visual categorization using meta-learning optimization with sample selection of auxiliary data. In *The European Conference on Computer Vision*.
- Zou, H.; Zhou, Y.; Yang, J.; Liu, H.; Das, H. P.; and Spanos, C. J. 2019. Consensus adversarial domain adaptation. In *Association for the Advancement of Artificial Intelligence*.

We provide an intuitive explanation for our proposed loss (4) in Section A. We theoretically prove that our proposed method can better bound the expected target error than existing ones in Section B. We provide more results and analysis on benchmark datasets of Digits, Office-31, Office-Home, and ImageNet-Caltech for closed set, partial, and open set domain adaptation in Section C. We present empirical evidence on benchmark datasets of digits that shows the efficacy of our used training scheme in Section D. We will release the code soon.

A Intuitive Explanation for Our Proposed Loss (4)

We denote the output vector of class scores of $F(G(\mathbf{x}))$ before the final softmax operation for an instance \mathbf{x} as $\mathbf{o}(\mathbf{x}) \in \mathbb{R}^{K+1}$, and its k^{th} element as $o_k(\mathbf{x})$, $k \in \{1, \dots, K+1\}$. We denote the output vector of class probabilities of $F(G(\mathbf{x}))$ after the final softmax operation for an instance \mathbf{x} as $\mathbf{p}(\mathbf{x}) \in [0, 1]^{K+1}$, and its k^{th} element as $p_k(\mathbf{x})$, $k \in \{1, \dots, K+1\}$. We write $p_k(\mathbf{x})$, $k \in \{1, \dots, K+1\}$ as

$$p_k(\mathbf{x}) = \frac{\exp(o_k(\mathbf{x}))}{\sum_{k'=1}^{K+1} \exp(o_{k'}(\mathbf{x}))}. \quad (\text{A.1})$$

We always have $\sum_{k=1}^{K+1} p_k(\mathbf{x}) = 1$ for any instance \mathbf{x} . When maximized over the feature extractor $G(\cdot)$, the adversarial loss on an unlabeled target instance \mathbf{x}^t (cf. objective (2) in Section **Discriminative Adversarial Domain Adaptation** in the paper) is written as

$$\begin{aligned} l^t(G, F) &= \log(1 - p_{K+1}(\mathbf{x}^t)) \\ &= \log\left(\sum_{k=1}^K p_k(\mathbf{x}^t)\right) \\ &= \log\left(\frac{\sum_{k=1}^K \exp(o_k(\mathbf{x}^t))}{\sum_{k'=1}^{K+1} \exp(o_{k'}(\mathbf{x}^t))}\right). \end{aligned} \quad (\text{A.2})$$

We write the gradient formulas of l^t w.r.t. $o_k(\mathbf{x})$, $k \in \{1, \dots, K\}$ as

$$\begin{aligned} \nabla_{o_k(\mathbf{x}^t)} &= \frac{\partial l^t}{\partial o_k(\mathbf{x}^t)} \\ &= \frac{\sum_{k'=1}^{K+1} \exp(o_{k'}(\mathbf{x}^t)) \cdot \exp(o_k(\mathbf{x}^t)) \cdot \exp(o_{K+1}(\mathbf{x}^t))}{\sum_{k=1}^K \exp(o_k(\mathbf{x}^t)) \left(\sum_{k'=1}^{K+1} \exp(o_{k'}(\mathbf{x}^t))\right)^2}, \end{aligned} \quad (\text{A.3})$$

where $\nabla_{o_k(\mathbf{x}^t)}$, $k \in \{1, \dots, K\}$ differ in the term of $\exp(o_k(\mathbf{x}^t))$, meaning that they are proportional to the class scores of $o_k(\mathbf{x}^t)$. In other words, the higher the class score is (i.e., the higher the class probability is), the stronger gradient the corresponding category neuron back-propagates, suggesting that the target instance is aligned to several most confident/related categories on the source domain. Such a mechanism to align the joint distributions of feature and category across domains is rather implicit. To make it more explicit, our proposed target discriminative adversarial loss (cf. loss (4) in Section **Discriminative Adversarial Learning** in the paper) uses the conditional probabilities to weight the category-wise domain predictions. By such a design, the discriminative adversarial training on the target data explicitly conducts the competition between the domain neuron (output) and the most confident category neuron (output) as the discriminative adversarial training on the

source data does, thus promoting the category-level domain alignment. This is what we mean by the *mutually inhibitory relation between the category and domain predictions* for any input instance.

This intuitive explanation manifests that the adversarial training of DADA clearly and explicitly utilizes the discriminative information of the target domain, thus improving the alignment of joint distributions of feature and category across domains.

B Generalization Error Analysis for Our Proposed DADA

We prove that our proposed DADA can better bound the expected target error than existing domain adaptation methods (Ganin et al. 2016; Tzeng et al. 2017; Pei et al. 2018; Pinheiro 2018; Zhang et al. 2018b; Shu et al. 2018; Long et al. 2018b; Wang et al. 2019; Wen et al. 2019; Tran et al. 2019), taking the similar formalism of theoretical results of domain adaptation (Ben-David et al. 2007; 2010).

For all hypothesis spaces introduced below, we assume them of finite effective size, i.e., finite VC dimension, so that the following distance measures defined over these spaces can be estimated from finite instances (Ben-David et al. 2010). We consider a fixed representation function $G(\cdot)$ from the instance set X to the feature space Z , i.e., $\mathbf{z} = G(\mathbf{x})$, and a hypothesis space \mathcal{H} for the K -category task classifier $C(\cdot)$ from the feature space Z to the label space Y , i.e., $C \in \mathcal{H}$ (Ganin et al. 2016). Note that $\mathbf{y} \in Y$ is the K -dimensional one-hot vector for any label y . Denote the marginal feature distribution and the joint distribution of feature and category by P_Z^s and $P_{Z,Y}^s$ for the source domain \mathcal{D}_s , and similarly P_Z^t and $P_{Z,Y}^t$ for the target domain \mathcal{D}_t , respectively. Let $\epsilon_s(C) = \mathbb{E}_{(\mathbf{z}, \mathbf{y}) \sim P_{Z,Y}^s} \mathbb{I}[C(\mathbf{z}) \neq \mathbf{y}]$ be the expected source error of a hypothesis $C \in \mathcal{H}$ w.r.t. the joint distribution $P_{Z,Y}^s$, where $\mathbb{I}[a]$ is the indicator function which is 1 if predicate a is true, and 0 otherwise. Similarly, $\epsilon_t(C) = \mathbb{E}_{(\mathbf{z}, \mathbf{y}) \sim P_{Z,Y}^t} \mathbb{I}[C(\mathbf{z}) \neq \mathbf{y}]$ denotes the expected target error of C w.r.t. the joint distribution $P_{Z,Y}^t$. Let $C^* = \operatorname{argmin}_{C \in \mathcal{H}} [\epsilon_s(C) + \epsilon_t(C)]$ be the ideal joint hypothesis that explicitly embodies the notion of adaptability (Ben-David et al. 2010). Let $\epsilon_s(C, C^*) = \mathbb{E}_{(\mathbf{z}, \mathbf{y}) \sim P_{Z,Y}^s} \mathbb{I}[C(\mathbf{z}) \neq C^*(\mathbf{z})]$ and $\epsilon_t(C, C^*) = \mathbb{E}_{(\mathbf{z}, \mathbf{y}) \sim P_{Z,Y}^t} \mathbb{I}[C(\mathbf{z}) \neq C^*(\mathbf{z})]$ be the disagreement between hypotheses C and C^* w.r.t. the joint distributions $P_{Z,Y}^s$ and $P_{Z,Y}^t$ respectively. Specified by the two works (Ben-David et al. 2007; 2010), the probabilistic bound of the expected target error $\epsilon_t(C)$ of the hypothesis C is given by the sum of the expected source error $\epsilon_s(C)$, the combined error $[\epsilon_s(C^*) + \epsilon_t(C^*)]$ of the ideal joint hypothesis C^* , and the distribution discrepancy across data domains, as the follow

$$\begin{aligned} \epsilon_t(C) &\leq \\ \epsilon_s(C) + [\epsilon_s(C^*) + \epsilon_t(C^*)] + |\epsilon_s(C, C^*) - \epsilon_t(C, C^*)|. \end{aligned} \quad (\text{B.1})$$

For domain adaptation to be possible, a natural assumption is that there exists the ideal joint hypothesis $C^* \in \mathcal{H}$ so that the combined error $[\epsilon_s(C^*) + \epsilon_t(C^*)]$ is small. The ideal joint hypothesis C^* may not be unique, since in practice we always have the same error obtained by two different machine learning models. Denote a set of ideal joint hypotheses by \mathcal{H}^* , which is a subset of \mathcal{H} , i.e., $\mathcal{H}^* \subset \mathcal{H}$. Based on this assumption, domain adaptation aims to reduce the domain discrepancy $|\epsilon_s(C, C^*) - \epsilon_t(C, C^*)|$. Let $\mathbf{c} = C(\mathbf{z})$ be the proxy of the label vector \mathbf{y} of \mathbf{z} , for every pair of $(\mathbf{z}, \mathbf{y}) \sim P_{Z,Y}^s \cup P_{Z,Y}^t$. Denote the thus obtained proxies of the joint distributions $P_{Z,Y}^s$ and $P_{Z,Y}^t$ by $P_{Z,C}^s = (\mathbf{z}, C(\mathbf{z}))_{\mathbf{z} \sim P_Z^s}$ and $P_{Z,C}^t = (\mathbf{z}, C(\mathbf{z}))_{\mathbf{z} \sim P_Z^t}$, respectively (Courty et al. 2017). Then, by definition, $\epsilon_s(C, C^*) =$

$\mathbb{E}_{(\mathbf{z}, \mathbf{y}) \sim P_{Z, Y}^s} \mathbb{I}[C(\mathbf{z}) \neq C^*(\mathbf{z})] = \mathbb{E}_{(\mathbf{z}, \mathbf{c}) \sim P_{Z, C}^s} \mathbb{I}[\mathbf{c} \neq C^*(\mathbf{z})]$, and similarly $\epsilon_t(C, C^*) = \mathbb{E}_{(\mathbf{z}, \mathbf{y}) \sim P_{Z, Y}^t} \mathbb{I}[C(\mathbf{z}) \neq C^*(\mathbf{z})] = \mathbb{E}_{(\mathbf{z}, \mathbf{c}) \sim P_{Z, C}^t} \mathbb{I}[\mathbf{c} \neq C^*(\mathbf{z})]$. Based on the two joint distribution proxies, we have the domain discrepancy

$$\begin{aligned} & |\epsilon_s(C, C^*) - \epsilon_t(C, C^*)| \\ &= |\mathbb{E}_{(\mathbf{z}, \mathbf{y}) \sim P_{Z, Y}^s} \mathbb{I}[C(\mathbf{z}) \neq C^*(\mathbf{z})] \\ &\quad - \mathbb{E}_{(\mathbf{z}, \mathbf{y}) \sim P_{Z, Y}^t} \mathbb{I}[C(\mathbf{z}) \neq C^*(\mathbf{z})]| \quad (\text{B.2}) \\ &= |\mathbb{E}_{(\mathbf{z}, \mathbf{c}) \sim P_{Z, C}^s} \mathbb{I}[\mathbf{c} \neq C^*(\mathbf{z})] \\ &\quad - \mathbb{E}_{(\mathbf{z}, \mathbf{c}) \sim P_{Z, C}^t} \mathbb{I}[\mathbf{c} \neq C^*(\mathbf{z})]|. \end{aligned}$$

Inspired by the two works (Long et al. 2018b; Mansour, Mohri, and Rostamizadeh 2009), we next introduce four definitions of the distance measure that can upper bound the domain discrepancy.

Definition 1. Let $\mathcal{F}_{\mathcal{H}^*} = \{F(C^*(\mathbf{z}), \mathbf{c}) = \mathbb{I}[\mathbf{c} \neq C^*(\mathbf{z})] | C^* \in \mathcal{H}^*\}$ be a (loss) difference hypothesis space over the joint variable of $(C^*(\mathbf{z}), \mathbf{c})$, where $F : (C^*(\mathbf{z}), \mathbf{c}) \mapsto \{0, 1\}$ computes the empirical 0-1 classification loss of the task classifier $C^* \in \mathcal{H}^*$ for any input pair of $(\mathbf{z}, \mathbf{c}) \sim P_{Z, C}^s \cup P_{Z, C}^t$. Then, the $\mathcal{F}_{\mathcal{H}^*}$ -distance between two distributions $P_{Z, C}^s$ and $P_{Z, C}^t$, is defined as

$$\begin{aligned} & d_{\mathcal{F}_{\mathcal{H}^*}}(P_{Z, C}^s, P_{Z, C}^t) \\ &\triangleq \sup_{F \in \mathcal{F}_{\mathcal{H}^*}, C^* \in \mathcal{H}^*} |\mathbb{E}_{(\mathbf{z}, \mathbf{c}) \sim P_{Z, C}^s} F(C^*(\mathbf{z}), \mathbf{c}) \\ &\quad - \mathbb{E}_{(\mathbf{z}, \mathbf{c}) \sim P_{Z, C}^t} F(C^*(\mathbf{z}), \mathbf{c})| \quad (\text{B.3}) \\ &= \sup_{C^* \in \mathcal{H}^*} |\mathbb{E}_{(\mathbf{z}, \mathbf{c}) \sim P_{Z, C}^s} \mathbb{I}[\mathbf{c} \neq C^*(\mathbf{z})] \\ &\quad - \mathbb{E}_{(\mathbf{z}, \mathbf{c}) \sim P_{Z, C}^t} \mathbb{I}[\mathbf{c} \neq C^*(\mathbf{z})]|. \end{aligned}$$

Definition 2. Let \mathcal{F} be a (loss) difference hypothesis space, which contains a class of functions $F : (\mathbf{z}, \mathbf{c}) \mapsto \{0, 1\}$ over the joint variable of $(\mathbf{z}, \mathbf{c}) \sim P_{Z, C}^s \cup P_{Z, C}^t$. Then, the \mathcal{F} -distance between two distributions $P_{Z, C}^s$ and $P_{Z, C}^t$, is defined as

$$\begin{aligned} & d_{\mathcal{F}}(P_{Z, C}^s, P_{Z, C}^t) \\ &\triangleq \sup_{F \in \mathcal{F}} |\mathbb{E}_{(\mathbf{z}, \mathbf{c}) \sim P_{Z, C}^s} F(\mathbf{z}, \mathbf{c}) - \mathbb{E}_{(\mathbf{z}, \mathbf{c}) \sim P_{Z, C}^t} F(\mathbf{z}, \mathbf{c})|. \quad (\text{B.4}) \end{aligned}$$

Definition 3. Let $\mathcal{F}_{\mathcal{H}} = \{F : (C'(\mathbf{z}), \mathbf{c}) \mapsto \{0, 1\} | C' \in \mathcal{H}\}$ be a (loss) difference hypothesis space over the joint variable of $(C'(\mathbf{z}), \mathbf{c})$, where $F(C'(\mathbf{z}), \mathbf{c})$ computes the empirical 0-1 classification loss of the task classifier $C' \in \mathcal{H}$ for any input pair of $(\mathbf{z}, \mathbf{c}) \sim P_{Z, C}^s \cup P_{Z, C}^t$. Then, the $\mathcal{F}_{\mathcal{H}}$ -distance between two distributions $P_{Z, C}^s$ and $P_{Z, C}^t$, is defined as

$$\begin{aligned} & d_{\mathcal{F}_{\mathcal{H}}}(P_{Z, C}^s, P_{Z, C}^t) \\ &\triangleq \sup_{F \in \mathcal{F}_{\mathcal{H}}, C' \in \mathcal{H}} |\mathbb{E}_{(\mathbf{z}, \mathbf{c}) \sim P_{Z, C}^s} F(C'(\mathbf{z}), \mathbf{c}) \\ &\quad - \mathbb{E}_{(\mathbf{z}, \mathbf{c}) \sim P_{Z, C}^t} F(C'(\mathbf{z}), \mathbf{c})|. \quad (\text{B.5}) \end{aligned}$$

Definition 4. Let \mathcal{D} be a (loss) difference hypothesis space, which contains a class of functions $D : \mathbf{z} \mapsto \{0, 1\}$ over $\mathbf{z} \sim P_Z^s \cup P_Z^t$. Then, the \mathcal{D} -distance between two distributions $P_{Z, C}^s$ and $P_{Z, C}^t$, is defined as

$$d_{\mathcal{D}}(P_Z^s, P_Z^t) \triangleq \sup_{D \in \mathcal{D}} |\mathbb{E}_{\mathbf{z} \sim P_Z^s} D(\mathbf{z}) - \mathbb{E}_{\mathbf{z} \sim P_Z^t} D(\mathbf{z})|. \quad (\text{B.6})$$

We are now ready to give an upper bound on the domain discrepancy in terms of the distance measures we have defined.

Theorem 1. The distribution discrepancy between the source and target domains $|\epsilon_s(C, C^*) - \epsilon_t(C, C^*)|$ can be upper bounded by the $\mathcal{F}_{\mathcal{H}^*}$ -distance, the $\mathcal{F}_{\mathcal{H}}$ -distance, the \mathcal{F} -distance, and the \mathcal{D} -distance as follows

$$\begin{aligned} & |\epsilon_s(C, C^*) - \epsilon_t(C, C^*)| \\ &\leq d_{\mathcal{F}_{\mathcal{H}^*}}(P_{Z, C}^s, P_{Z, C}^t) \\ &\leq d_{\mathcal{F}_{\mathcal{H}}}(P_{Z, C}^s, P_{Z, C}^t) \quad (\text{B.7}) \\ &\leq d_{\mathcal{F}}(P_{Z, C}^s, P_{Z, C}^t) \\ &\leq d_{\mathcal{D}}(P_Z^s, P_Z^t). \end{aligned}$$

Proof. Comparing (B.2) and (B.3), since $|\mathbb{E}_{(\mathbf{z}, \mathbf{c}) \sim P_{Z, C}^s} \mathbb{I}[\mathbf{c} \neq C^*(\mathbf{z})] - \mathbb{E}_{(\mathbf{z}, \mathbf{c}) \sim P_{Z, C}^t} \mathbb{I}[\mathbf{c} \neq C^*(\mathbf{z})]| \leq \sup_{C^* \in \mathcal{H}^*} |\mathbb{E}_{(\mathbf{z}, \mathbf{c}) \sim P_{Z, C}^s} \mathbb{I}[\mathbf{c} \neq C^*(\mathbf{z})] - \mathbb{E}_{(\mathbf{z}, \mathbf{c}) \sim P_{Z, C}^t} \mathbb{I}[\mathbf{c} \neq C^*(\mathbf{z})]|$, we have $|\epsilon_s(C, C^*) - \epsilon_t(C, C^*)| \leq d_{\mathcal{F}_{\mathcal{H}^*}}(P_{Z, C}^s, P_{Z, C}^t)$.

Since by definition the hypothesis space \mathcal{F} contains all functions that map (\mathbf{z}, \mathbf{c}) to $\{0, 1\}$, $F(C^*(\mathbf{z}), \mathbf{c})$ is also a function in \mathcal{F} that can be written as the form of functions in $\mathcal{F}_{\mathcal{H}^*}$. The hypothesis space $\mathcal{F}_{\mathcal{H}^*}$ is subsumed by \mathcal{F} , i.e., $\mathcal{F}_{\mathcal{H}^*} \subset \mathcal{F}$. Thus, we have $|\epsilon_s(C, C^*) - \epsilon_t(C, C^*)| \leq d_{\mathcal{F}_{\mathcal{H}^*}}(P_{Z, C}^s, P_{Z, C}^t) \leq d_{\mathcal{F}}(P_{Z, C}^s, P_{Z, C}^t)$.

Similarly, since $\mathcal{F}_{\mathcal{H}} \subset \mathcal{F}$, we have $d_{\mathcal{F}_{\mathcal{H}}}(P_{Z, C}^s, P_{Z, C}^t) \leq d_{\mathcal{F}}(P_{Z, C}^s, P_{Z, C}^t)$. Since by definition the ideal joint hypothesis set $\mathcal{H}^* \subset \mathcal{H}$, the hypothesis space $\mathcal{F}_{\mathcal{H}^*}$ is subsumed by $\mathcal{F}_{\mathcal{H}}$, i.e., $\mathcal{F}_{\mathcal{H}^*} \subset \mathcal{F}_{\mathcal{H}}$. Thus, we have $|\epsilon_s(C, C^*) - \epsilon_t(C, C^*)| \leq d_{\mathcal{F}_{\mathcal{H}^*}}(P_{Z, C}^s, P_{Z, C}^t) \leq d_{\mathcal{F}_{\mathcal{H}}}(P_{Z, C}^s, P_{Z, C}^t)$.

Since by definition the hypothesis space \mathcal{D} contains all functions that map \mathbf{z} to $\{0, 1\}$, $F(\mathbf{z}, \mathbf{c}) = F(\mathbf{z}, C(\mathbf{z}))$ is also a function in \mathcal{D} that can be written as the form of functions in \mathcal{F} . The hypothesis space \mathcal{F} is subsumed by \mathcal{D} , i.e., $\mathcal{F} \subset \mathcal{D}$. Thus, we have $d_{\mathcal{F}}(P_{Z, C}^s, P_{Z, C}^t) \leq d_{\mathcal{D}}(P_Z^s, P_Z^t)$.

These prove the inequality $|\epsilon_s(C, C^*) - \epsilon_t(C, C^*)| \leq d_{\mathcal{F}_{\mathcal{H}^*}}(P_{Z, C}^s, P_{Z, C}^t) \leq d_{\mathcal{F}_{\mathcal{H}}}(P_{Z, C}^s, P_{Z, C}^t) \leq d_{\mathcal{F}}(P_{Z, C}^s, P_{Z, C}^t) \leq d_{\mathcal{D}}(P_Z^s, P_Z^t)$. \square

Theorem 1 shows that the $\mathcal{F}_{\mathcal{H}^*}$ -distance can best upper bound the domain discrepancy $|\epsilon_s(C, C^*) - \epsilon_t(C, C^*)|$, but cannot be computable, since instances on the target domain for unsupervised domain adaptation are unlabeled; the $\mathcal{F}_{\mathcal{H}}$ -distance can better bound the domain discrepancy $|\epsilon_s(C, C^*) - \epsilon_t(C, C^*)|$ than the \mathcal{F} -distance and the \mathcal{D} -distance, and the hypothesis space $\mathcal{F}_{\mathcal{H}}$ can be implemented by conditioning the function $F(\mathbf{z}, \mathbf{c}) \in \mathcal{F}$ on the other one $C(\mathbf{z}) \in \mathcal{H}$; the \mathcal{F} -distance can tighter bound the domain discrepancy $|\epsilon_s(C, C^*) - \epsilon_t(C, C^*)|$ than the \mathcal{D} -distance, and the hypothesis space \mathcal{F} can be realized by taking as input both the feature representation \mathbf{z} and the category prediction \mathbf{c} ; the \mathcal{D} -distance can loosely bound the domain discrepancy $|\epsilon_s(C, C^*) - \epsilon_t(C, C^*)|$, and the hypothesis space \mathcal{D} can be instantiated by taking as input only the feature representation \mathbf{z} . Since existing deep domain adaptation methods are based on deep neuron networks, the inference of the hypothesis space $\mathcal{F}_{\mathcal{H}^*} \subset \mathcal{F}_{\mathcal{H}} \subset \mathcal{F} \subset \mathcal{D}$ is reasonable and realistic in that, for any given function, there must exist a feedforward neural network or multilayer perceptron, which can approximate it with arbitrarily small error (Hornik, Stinchcombe, and White 1989; Cybenko 1989), however, the effective model capacity is limited by the capabilities of the optimization algorithm (Goodfellow, Bengio, and Courville 2016).

Since these methods (Ganin et al. 2016; Tzeng et al. 2017; Pinheiro 2018; Zhang et al. 2018b; Shu et al. 2018) are based on a separate domain classifier that takes as input only the feature representation, they aim to measure and minimize the \mathcal{D} -distance. Since these methods (Pei et al. 2018; Long et al. 2018b; Wang et al. 2019;

Wen et al. 2019) are based on one or several conditional domain classifiers that take as input both the feature representation and the category prediction, they aim to measure and minimize the \mathcal{F} -distance. Since the recent work (Tran et al. 2019) and the proposed DADA unify the task and domain classifiers into an integrated one, i.e., conditioning the domain classifier on the task classifier, they aim to measure and minimize the $\mathcal{F}_{\mathcal{H}}$ -distance. The $\mathcal{F}_{\mathcal{H}}$ -distance can be upper bounded by the optimal solution of the integrated domain and task classifier $F(\cdot)$. In the meanwhile, the upper bound of $\mathcal{F}_{\mathcal{H}}$ -distance is minimized by learning a domain-invariant feature extractor $G(\cdot)$.

Furthermore, our proposed DADA can be intuitively formalized as category-regularized domain-adversarial training, since our proposed discriminative adversarial training can learn an integrated classifier $F(\cdot)$ that has explicit intra-domain discrimination and inter-domain indistinguishability, which may enable a better performed ideal joint hypothesis C^* . Consequently, the expected target error $\epsilon_t(C)$ can be better approximated by the expected source error $\epsilon_s(C)$. As verified above, our proposed DADA can formally better bound the expected target error than existing domain adaptation methods.

C Additional Results and Analysis

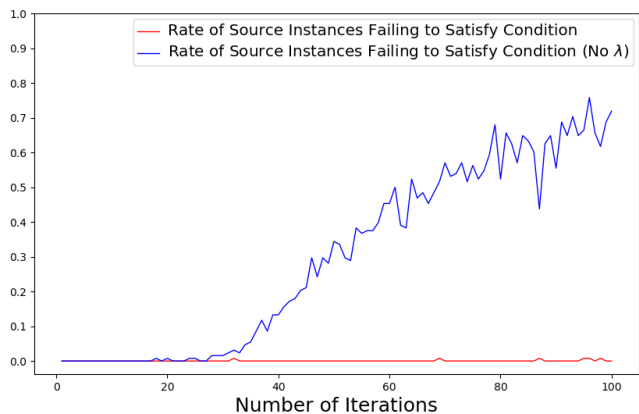
C.1 Datasets

Digits datasets of MNIST (Lecun et al. 1998), Street View House Numbers (SVHN) (Netzer et al. 2011), and USPS (Hull 1994) are popular. we follow ADR (Saito et al. 2018a) and evaluate on three adaptation settings of **SVHN**→**MNIST**, **MNIST**→**USPS**, and **USPS**→**MNIST**. For all adaptation settings, we adopt the same network architecture and experimental setting as ADR.

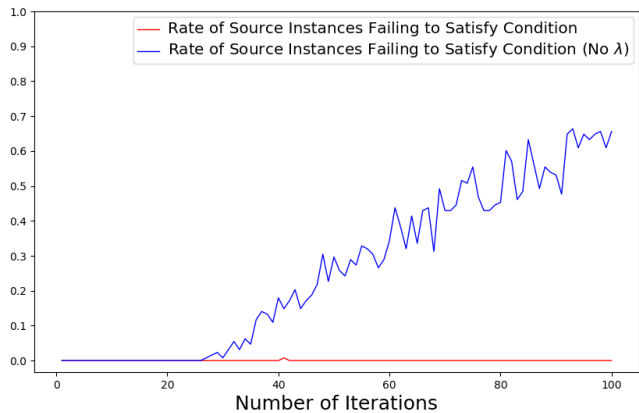
Office-31 (Saenko et al. 2010) is a benchmark domain adaptation dataset as introduced in Section **Datasets and Implementation Details** in the paper. For partial domain adaptation, we select images of 10 categories shared by Office-31 and Caltech-256 (Griffin, Holub, and Perona 2007) in each domain of Office-31 as the target domain. Note that the source domain here contains 31 categories and the target domain here contains 10 categories. For open set domain adaptation, we use the selected 10 categories as the known categories. In alphabetical order, 11 – 20 categories and 21 – 31 categories are used as the unknown categories in the source and target domains respectively. In this setting, an 11-category classification is performed.

Office-Home (Venkateswara et al. 2017) is a much more challenging benchmark dataset for domain adaptation, which includes 15,500 images of 65 object categories in office and home scenes, shared by four extremely distinct domains: Artistic images (**Ar**), Clip Art (**Cl**), Product images (**Pr**), and Real-World images (**Rw**). We build 12 adaptation settings: **Ar** → **Cl**, **Ar** → **Pr**, **Ar** → **Rw**, **Cl** → **Ar**, **Cl** → **Pr**, **Cl** → **Rw**, **Pr** → **Ar**, **Pr** → **Cl**, **Pr** → **Rw**, **Rw** → **Ar**, **Rw** → **Cl**, **Rw** → **Pr**. For partial domain adaptation, we choose images of the first 25 categories (in alphabetical order) in each domain of this dataset as target domains. Note that each source domain here contains 65 categories and each target domain here contains 25 categories.

ImageNet-Caltech is built from ImageNet (Russakovsky et al. 2015) that contains 1000 categories, and Caltech-256 (Griffin, Holub, and Perona 2007) that contains 256 categories. They share 84 common categories, thus we construct two adaptation settings: **I (1000)** → **C (84)**, and **C (256)** → **I (84)**. When ImageNet is used as the source domain, we use its training set; when it is used as the target domain, we use its validation set to prevent the model from the effect of pre-training on its training set.



(a) **A**→**D**



(b) **D**→**A**

Figure 4: An illustration for the effect of the λ on the rate of source instances failing to satisfy the condition in the early stage (e.g., the first 100 iterations) of adversarial training on the two adaptation settings of (a) **A**→**D** and (b) **D**→**A**.

C.2 Closed Set Domain Adaptation.

Effect of the λ We provide the empirical evidence on Office-31 (Saenko et al. 2010) based on ResNet-50 (He et al. 2016) for the effect of the hyper-parameter λ on keeping the source instances satisfying the condition of $p_{y^s}^s > 0.5$ (cf. Section **Discriminative Adversarial Learning** in the paper for its derivation) in the early stage of adversarial training of DADA in Figure 4, which shows that the rate of source instances failing to satisfy the condition rises rapidly in the early stage of adversarial training when the λ is not used.

Alternative Choice of Adversarial Loss for Target Instances For a target adversarial loss, when maximized over the feature extractor $G(\cdot)$, we have an alternative choice. In this section, we give further discussion and experiments to compare our used \mathcal{L}_G^t in loss (4) in the paper with this alternative.

Inspired by the works (Tzeng et al. 2015; Zhang et al. 2019), one may opt for a symmetric adversarial loss

$$\mathcal{L}_G^t(G, F) = \frac{1}{n_t} \sum_{j=1}^{n_t} \sum_{k=1}^K \bar{p}_k(\mathbf{x}_j^t) \left[\frac{1}{2} \log p_{K+1}(\mathbf{x}_j^t) + \frac{1}{2} \log(1 - p_{K+1}(\mathbf{x}_j^t)) \right], \quad (\text{C.1})$$

Table 6: Comparison on Office-31 based on ResNet-50 with an alternative choice of adversarial loss for target instances. Please refer to the main text for how this alternative is defined.

| Methods | A \rightarrow W | D \rightarrow W | W \rightarrow D | A \rightarrow D | D \rightarrow A | W \rightarrow A | Avg |
|---------|-----------------------|-----------------------|------------------------|-----------------------|-----------------------|-----------------------|-------------|
| DADA-DC | 90.4 \pm 0.1 | 98.7 \pm 0.1 | 100.0 \pm 0.0 | 92.5 \pm 0.3 | 72.5 \pm 0.2 | 73.0 \pm 0.3 | 87.9 |
| DADA | 92.3 \pm 0.1 | 99.2 \pm 0.1 | 100.0 \pm 0.0 | 93.9 \pm 0.2 | 74.4 \pm 0.1 | 74.2 \pm 0.1 | 89.0 |

Table 7: Analysis of robustness for different methods on benchmark datasets of MNIST (Lecun et al. 1998), SVHN (Netzer et al. 2011), and USPS (Hull 1994) based on modified LeNet.

| Methods | SVHN \rightarrow MNIST | MNIST \rightarrow USPS | USPS \rightarrow MNIST | Avg |
|--|----------------------------------|--------------------------|--------------------------|-------------|
| No Adaptation | 67.1 | 77.0 | 68.1 | 70.7 |
| DDC (Tzeng et al. 2014) | 68.1 \pm 0.3 | 79.1 \pm 0.5 | 66.5 \pm 3.3 | 71.2 |
| DANN (Ganin et al. 2016) | 73.9 | 77.1 \pm 1.8 | 73.0 \pm 0.2 | 74.7 |
| DRCN (Ghifary et al. 2016) | 82.0 \pm 0.1 | 91.8 \pm 0.09 | 73.7 \pm 0.04 | 82.5 |
| ADDA (Tzeng et al. 2017) | 76.0 \pm 1.8 | 89.4 \pm 0.2 | 90.1 \pm 0.8 | 85.2 |
| SBADA-GAN (Russo et al. 2018) | 76.1 | 97.6 | 95.0 | 89.6 |
| RAAN (Chen et al. 2018) | 89.2 | 89.0 | 92.1 | 90.1 |
| ADR (Saito et al. 2018a) | 94.1 \pm 1.37 | 91.3 \pm 0.65 | 91.5 \pm 3.61 | 92.3 |
| TPN (Pan et al. 2019) | 93.0 | 92.1 | 94.1 | 93.1 |
| CyCADA (Hoffman et al. 2018) | 90.4 \pm 0.4 | 95.6 \pm 0.2 | 96.5 \pm 0.1 | 94.2 |
| MCD (Saito et al. 2018b) | 96.2 \pm 0.4 | 94.2 \pm 0.7 | 94.1 \pm 0.3 | 94.8 |
| CADA (Zou et al. 2019) | 90.9 \pm 0.2 | 96.4 \pm 0.1 | 97.0 \pm 0.1 | 94.8 |
| DAN (Long et al. 2018a) | 71.1 | - | - | - |
| CoGAN (Liu and Tuzel 2016) | - | 91.2 \pm 0.8 | 89.1 \pm 0.8 | - |
| DSN (Bousmalis et al. 2016) | 82.7 | 91.3 | - | - |
| LDC (Li, Song, and Wu 2018) | 89.5 \pm 2.1 | - | - | - |
| MSTN (Xie et al. 2018) | 91.7 \pm 1.5 | 92.9 \pm 1.1 | - | - |
| PFAN (Chen et al. 2019b) | 93.9 \pm 0.8 | 95.0 \pm 1.3 | - | - |
| JDDA-C (Chen et al. 2019a) | 94.2 \pm 0.1 | - | 96.7 \pm 0.1 | - |
| MECA (Morerio, Cavazza, and Murino 2018) | 95.2 | - | - | - |
| ASSC (Haeusser et al. 2017) | 95.7 \pm 1.5 | - | - | - |
| DADA | 95.6 \pm 0.5 | 96.1 \pm 0.4 | 96.5 \pm 0.2 | 96.1 |

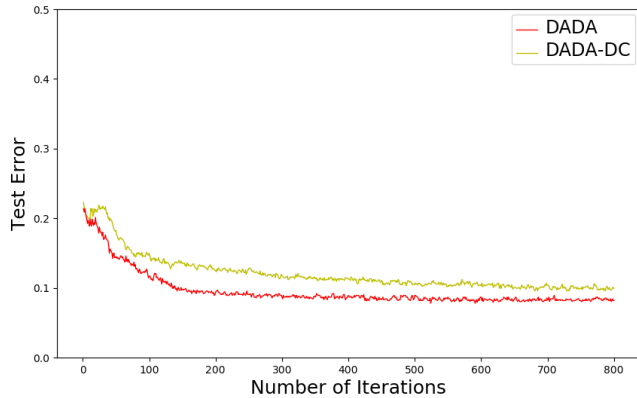


Figure 5: Convergence performance in terms of test error on the adaptation setting $\mathbf{A} \rightarrow \mathbf{W}$. Note that here we only show the convergence performance during adversarial training.

which when maximized over $G(\cdot)$, gives a confused prediction of $p_{K+1}(\mathbf{x}^t) = 0.5$. This result does not give category prediction $p_{y^t}(\mathbf{x}^t)$ on the *unknown* true category y^t of a target instance \mathbf{x}^t a

chance to approach 1. Thus, this alternative choice is sub-optimal.

In contrast, our used \mathcal{L}_G^t in loss (4) in the paper gives a prediction of $p_{K+1}(\mathbf{x}^t) = 0$ when maximized over $G(\cdot)$. This result gives $p_{y^t}(\mathbf{x}^t)$ a better chance to approach 1, i.e. $\bar{p}_{y^t}(\mathbf{x}^t)$ is more likely to approach 1. In other words, the target data are more likely to be correctly classified, which is enabled by our proposed *mutually inhibitory relation between the category and domain predictions*.

To compare the effectiveness of our used \mathcal{L}_G^t in loss (4) in the paper and this alternative choice, we conduct experiments on Office-31 (Saenko et al. 2010) based on ResNet-50 (He et al. 2016), by replacing \mathcal{L}_G^t in loss (4) in the paper with the domain confusion loss (C.1) in our main minimax problem (7) in the paper. We denote these this alternative as “DADA-DC”. Results in Table 6 and convergence performances in Figure 5 show advantages of our used \mathcal{L}_G^t in loss (4) in the paper.

Feature Visualization To visualize how different methods are effective at aligning learned features on the source and target domains, we use t-SNE embeddings (van der Maaten and Hinton 2008) to plot the output activations from the feature extractors of “No Adaptation”, DANN, DANN-CA, and DADA. Figure 6 gives the plotting, where samples are from the adaptation setting of $\mathbf{A} \rightarrow \mathbf{W}$ of Office-31 (Saenko et al. 2010) based on ResNet-50 (He et al. 2016). Figure 6 shows qualitative improvements of these meth-

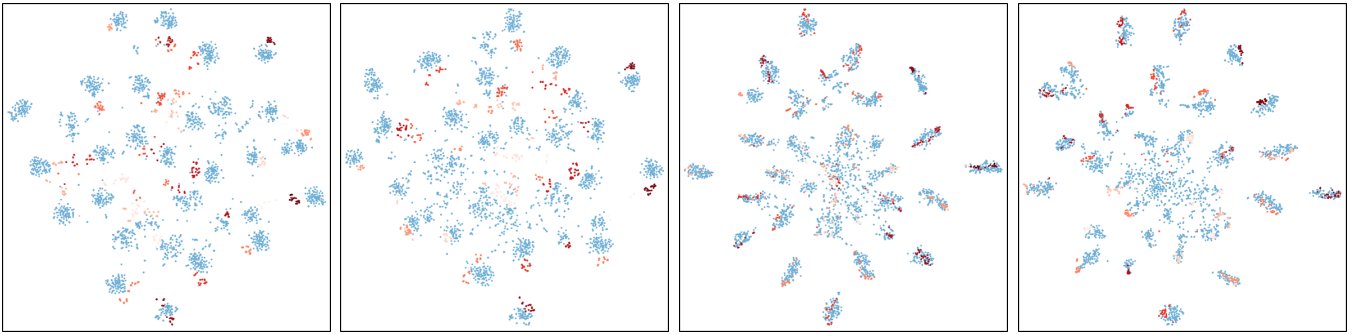


Figure 6: The t-SNE visualization of feature alignment between the source (blue) and target (red) domains by No Adaptation, DANN, DANN-CA, and DADA (from left to right). Samples of plotting are from the adaptation setting of $\mathbf{A} \rightarrow \mathbf{W}$ in Table 1 in the paper. Note that different degrees of the red color indicate different target categories.

ods at aligning features across data domains, i.e., the distribution of target samples (red) changes from the scattered state of DANN to multiple category-wise clusters of DADA, which are aligned with source samples (blue) of corresponding categories. Note that for the source domain, since we aim to achieve a balance between the transferability and discriminability (Liu et al. 2019), the categories are not perfectly separated. Since the transferability is enhanced, the target categories are well separated.

Digits To validate the robustness of our proposed DADA, we evaluate different methods on Digits datasets of MNIST (Lecun et al. 1998), SVHN (Netzer et al. 2011), and USPS (Hull 1994) based on modified LeNet in Table 7. Note that results of existing methods are quoted from their respective papers or the recent works (Saito et al. 2018a; 2018b). We follow these methods and report accuracies on the target test data in the format of $\text{mean} \pm \text{std}$ over five random trials. Our proposed DADA consistently achieves a good result on different adaptation settings, showing its excellent robustness.

C.3 Partial Domain Adaptation.

For each partial adaptation setting of Office-31, Office-Home, and ImageNet-Caltech, we follow the work (Cao et al. 2018b) to report the mean classification result on the target domain over three random trials.

Office-31 We compare in Table 8 our proposed method with existing ones on Office-31 based on ResNet-50 (He et al. 2016) pre-trained on ImageNet (Russakovsky et al. 2015). Results of existing methods are quoted from PADA (Cao et al. 2018b). Our proposed DADA-P outperforms all comparative methods by a large margin, showing the effectiveness of the adopted category-level weighting mechanism on reducing the negative influence of source outliers on adaptation settings with small source domain and small target domain, e.g., $\mathbf{A} \rightarrow \mathbf{W}$. Although PADA uses the same weighting mechanism, it performs much worse than our proposed DADA-P, suggesting the effectiveness of DADA-P on enhancing the positive influence of shared categories.

From the experimental results, several interesting observations can be derived. (1) Previous deep domain adaptation methods including those based on domain-adversarial training (e.g., DANN) and those based on MMD (e.g., DAN) perform much worse than the very baseline “No Adaptation”, showing the huge impact of negative transfer. Domain-adversarial training based methods aim to learn domain-invariant intermediate features to deceive the domain classifier, and MMD based methods aim to minimize the dis-

crepancy between data distributions of the source and target domains. Both of them align the whole source domain to the whole target one. However, in partial domain adaptation, since the source domain contains categories that do not exist in the target domain, i.e., outlier source categories, they will suffer false alignment between the outlier source categories and the target domain. This explains their poor performance in partial domain adaptation. (2) Among previous deep domain adaptation methods, RTN is the only one that performs better than “No Adaptation”. RTN exploits the entropy minimization principle (Grandvalet and Bengio 2005) to encourage the low-density separation of target categories. Its target classifier directly has access to the unlabeled target data and can amend itself to pass through the target low-density regions where the outlier source categories may exist, which alleviate the negative influence of source outliers. Nevertheless, PADA, which does not use the entropy minimization principle but a category-level weighting mechanism, performs much better than RTN, demonstrating that RTN still suffers negative transfer and may be not able to bridge such a large domain discrepancy caused by different label spaces. (3) Although our proposed DADA-P applies the same weighting mechanism as PADA, it performs much better than PADA. PADA has a separate design of task and domain classifiers and only aims to align marginal feature distributions, whereas our proposed DADA-P based on an integrated domain and task classifier, aims to promote the joint distribution alignment across domains. This explains the good performance of our proposed method in partial domain adaptation.

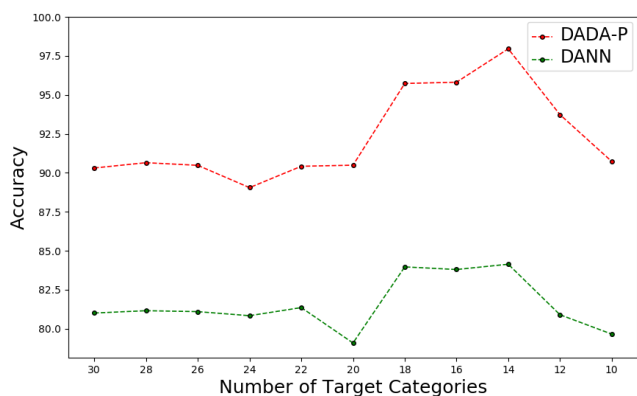
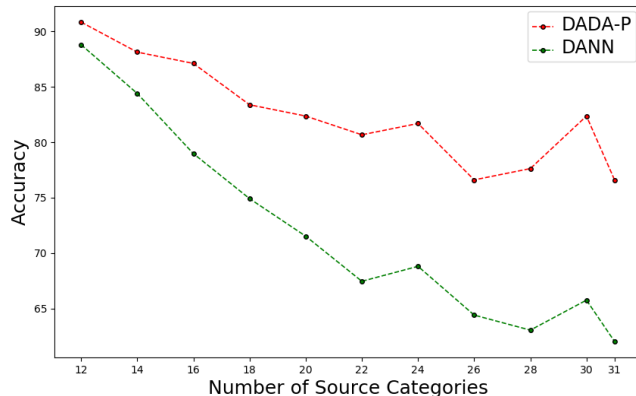
To investigate a wider spectrum of partial domain adaptation, we conduct experiments by varying the number of target categories. Figure 7 shows results for the baseline DANN (Ganin et al. 2016) and our proposed DADA-P on the partial adaptation setting $\mathbf{A} \rightarrow \mathbf{W}$ of Office-31 with a base network of ResNet-50. The source domain has always 31 categories, but the number of target categories varies from 30 to 10, i.e., $\{30, 28, 26, 24, 22, 20, 18, 16, 14, 12, 10\}$. As the number of target categories decreases, performances of the two methods have no evident decline in spite of the aggravation of negative transfer effect, since the difficulty of domain adaptation problem itself becomes smaller. We observe a sharp rise and a dramatic drop when the number of target categories decreases from 20 to 18 and from 14 to 12 respectively. One explanation is that the positive influence incurred by reducing the difficulty of domain adaptation problem itself is more (for the former observation) or less (for the latter one) than the negative influence caused by increasing the domain discrepancy. The results show that our proposed DADA-P performs much better than DANN in all settings. It is noteworthy

Table 8: Results for partial domain adaptation on Office-31 based on ResNet-50.

| Methods | A \rightarrow W | D \rightarrow W | W \rightarrow D | A \rightarrow D | D \rightarrow A | W \rightarrow A | Avg |
|-------------------------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|--------------|
| No Adaptation (He et al. 2016) | 54.52 | 94.57 | 94.27 | 65.61 | 73.17 | 71.71 | 75.64 |
| DAN (Long et al. 2018a) | 46.44 | 53.56 | 58.60 | 42.68 | 65.66 | 65.34 | 55.38 |
| DANN (Ganin et al. 2016) | 41.35 | 46.78 | 38.85 | 41.36 | 41.34 | 44.68 | 42.39 |
| ADDA (Tzeng et al. 2017) | 43.65 | 46.48 | 40.12 | 43.66 | 42.76 | 45.95 | 43.77 |
| RTN (Long et al. 2016) | 75.25 | 97.12 | 98.32 | 66.88 | 85.59 | 85.70 | 84.81 |
| JAN (Long et al. 2017) | 43.39 | 53.56 | 41.40 | 35.67 | 51.04 | 51.57 | 46.11 |
| Luo <i>et al.</i> (Luo et al. 2017) | 73.22 | 93.90 | 96.82 | 76.43 | 83.62 | 84.76 | 84.79 |
| PADA (Cao et al. 2018b) | 86.54 | 99.32 | 100.00 | 82.17 | 92.69 | 95.41 | 92.69 |
| DADA-P | 90.73 | 100.00 | 100.00 | 87.90 | 94.71 | 94.89 | 94.71 |

Table 9: Results for partial domain adaptation on Office-31 based on AlexNet.

| Methods | A \rightarrow W | D \rightarrow W | W \rightarrow D | A \rightarrow D | D \rightarrow A | W \rightarrow A | Avg |
|--|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|--------------|
| No Adaptation (Krizhevsky, Sutskever, and Hinton 2012) | 58.51 | 95.05 | 98.08 | 71.23 | 70.60 | 67.74 | 76.87 |
| DAN (Long et al. 2018a) | 56.52 | 71.86 | 86.78 | 51.86 | 50.42 | 52.29 | 61.62 |
| DANN (Ganin et al. 2016) | 49.49 | 93.55 | 90.44 | 49.68 | 46.72 | 48.81 | 63.11 |
| ADDA (Tzeng et al. 2017) | 70.68 | 96.44 | 98.65 | 72.90 | 74.26 | 75.56 | 81.42 |
| RTN (Long et al. 2016) | 66.78 | 86.77 | 99.36 | 70.06 | 73.52 | 76.41 | 78.82 |
| SAN (Cao et al. 2018a) | 80.02 | 98.64 | 100.00 | 81.28 | 80.58 | 83.09 | 87.27 |
| Zhang <i>et al.</i> (Zhang et al. 2018a) | 76.27 | 98.98 | 100.00 | 78.98 | 89.46 | 81.73 | 87.57 |
| DADA-P | 76.61 | 98.98 | 100.00 | 85.56 | 93.81 | 93.28 | 91.37 |

Figure 7: The accuracy curve of varying the number of target categories for the baseline DANN (Ganin et al. 2016) and our proposed DADA-P on the partial adaptation setting A \rightarrow W of Office-31 with a base network of ResNet-50.Figure 8: The accuracy curve of varying the number of source categories for the baseline DANN (Ganin et al. 2016) and our proposed DADA-P on the partial adaptation setting A \rightarrow W of Office-31 with a base network of AlexNet.

that the relative performance improvement becomes larger when the number of target categories decreases, testifying the superiority of our methods in reducing the influence of negative transfer. Thus, given a source domain, our methods can perform much better when applied to the target domain with unknown number of categories.

We compare in Table 9 our proposed method with existing ones on Office-31 based on AlexNet (Krizhevsky, Sutskever, and Hinton 2012) pre-trained on ImageNet. Results of existing methods are quoted from their respective papers or SAN (Cao et al. 2018a). Our proposed DADA-P achieves a much better result than all comparative methods, showing the efficacy of our methods with a shallower neuron network as the base network.

To investigate the influence of the number of outlier source categories on the performance, we conduct experiments by varying the number of source categories. Figure 8 shows results for the baseline DANN (Ganin et al. 2016) and our proposed DADA-P on the partial adaptation setting A \rightarrow W of Office-31 with a base network of AlexNet. The target domain has always 10 categories, but the number of source categories varies from 12 to 31, i.e., $\{12, 14, 16, 18, 20, 22, 24, 26, 28, 30, 31\}$. As the number of source categories increases, performances of the two methods have evident decline but also some rises, e.g., when the number of source categories increases from 22 to 24 and from 28 to 30. One explanation is that the positive influence incurred by increasing dis-

Table 10: Results for partial domain adaptation on Office-Home based on ResNet-50.

| Methods | Ar→Cl | Ar→Pr | Ar→Rw | Cl→Ar | Cl→Pr | Cl→Rw | Pr→Ar | Pr→Cl | Pr→Rw | Rw→Ar | Rw→Cl | Rw→Pr | Avg |
|--------------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| No Adaptation (He et al. 2016) | 38.57 | 60.78 | 75.21 | 39.94 | 48.12 | 52.90 | 49.68 | 30.91 | 70.79 | 65.38 | 41.79 | 70.42 | 53.71 |
| DAN (Long et al. 2018a) | 44.36 | 61.79 | 74.49 | 41.78 | 45.21 | 54.11 | 46.92 | 38.14 | 68.42 | 64.37 | 45.37 | 68.85 | 54.48 |
| DANN (Ganin et al. 2016) | 44.89 | 54.06 | 68.97 | 36.27 | 34.34 | 45.22 | 44.08 | 38.03 | 68.69 | 52.98 | 34.68 | 46.50 | 47.39 |
| RTN (Long et al. 2016) | 49.37 | 64.33 | 76.19 | 47.56 | 51.74 | 57.67 | 50.38 | 41.45 | 75.53 | 70.17 | 51.82 | 74.78 | 59.25 |
| PADA (Cao et al. 2018b) | 51.95 | 67.00 | 78.74 | 52.16 | 53.78 | 59.03 | 52.61 | 43.22 | 78.79 | 73.73 | 56.60 | 77.09 | 62.06 |
| DADA-P | 52.92 | 82.54 | 86.78 | 71.23 | 69.75 | 76.72 | 73.06 | 52.84 | 85.90 | 77.69 | 56.50 | 85.98 | 72.66 |

Table 11: Results for partial domain adaptation on ImageNet-Caltech based on ResNet-50.

| Methods | I→C | C→I | Avg |
|--------------------------------|--------------|--------------|--------------|
| No Adaptation (He et al. 2016) | 71.65 | 66.14 | 68.90 |
| DAN (Long et al. 2018a) | 71.57 | 66.48 | 69.03 |
| DANN (Ganin et al. 2016) | 68.67 | 52.97 | 60.82 |
| RTN (Long et al. 2016) | 72.24 | 68.33 | 70.29 |
| PADA (Cao et al. 2018b) | 75.03 | 70.48 | 72.76 |
| DADA-P | 80.94 | 76.91 | 78.93 |

criminative information of categories, especially those related to the target domain, is more than the negative influence caused by increasing the domain discrepancy. The results show that our proposed DADA-P significantly outperforms DANN in all settings. Particularly, the relative performance improvement is larger when the number of source categories is larger, demonstrating that our methods are more robust to the number of outlier source categories. Thus, for a given target task, our methods can have a much better performance when utilizing different source tasks.

Office-Home We compare in Table 10 our proposed method with existing ones on Office-Home based on ResNet-50. Results of existing methods are quoted from PADA (Cao et al. 2018b). Our proposed DADA-P significantly outperforms all comparative methods, showing the efficacy of DADA-P on adaptation settings with more categories in both the source and target domains and larger domain discrepancy between the two domains, e.g., $Cl \rightarrow Rw$.

ImageNet-Caltech We compare in Table 11 our proposed method with existing ones on ImageNet-Caltech based on ResNet-50. Results of existing methods are quoted from PADA (Cao et al. 2018b). Our proposed DADA-P outperforms all comparative methods by a large margin, showing the effectiveness of DADA-P on adaptation settings with large-scale source and target domains and a large number of categories in the two domains.

C.4 Open Set Domain Adaptation.

We compare in Table 12 our proposed method with existing ones on Office-31 based on AlexNet (Krizhevsky, Sutskever, and Hinton 2012) pre-trained on ImageNet (Russakovsky et al. 2015). Results of existing methods are quoted from AODA (Saito et al. 2018c). Our proposed DADA-O outperforms all comparative methods in both evaluation metrics of OS* and OS, showing the efficacy of DADA-O in both aligning distributions of the known instances across domains and identifying the unknown target instances as the unknown category for open set domain adaptation.

From the experimental results, we have some interesting observations. **(1)** DAN and DANN perform much worse than “No Adaptation”. DAN and DANN aim to align the whole marginal feature distributions across the source and target domains. If the

target domain contains unknown instances, false alignment between the known source instances and unknown target ones will occur, resulting in a sharp drop of the classification performance. **(2)** DANN performs worse than DAN, since DANN is better at aligning marginal feature distributions across data domains, leading to more serious false alignment. **(3)** ATI- λ and AODA can effectively reduce false alignment, since they have a good outlier rejection mechanism to recognize the unknown instances. **(4)** The results of all comparative methods on almost all adaptation settings are better in the evaluation metric OS than OS*, showing that many known target instances are classified as the unknown category. Since Open-set SVM is trained to detect outliers and the task classifier of AODA is trained to recognize all the target instances as the unknown category, they are inclined to classify the target instances as the unknown category. **(5)** For our proposed DADA-O, the results of all adaptation settings are better in the evaluation metric OS* than OS, since their classifiers are trained to classify all target instance as the unknown category with a small probability q , which can minimize the misclassification of the known target instances as the unknown category.

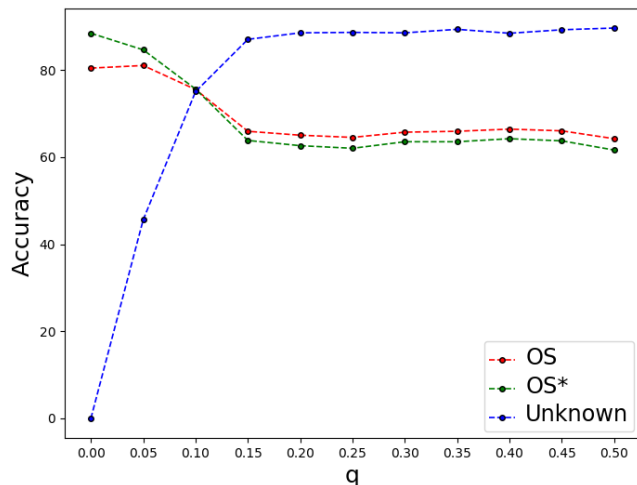


Figure 9: The accuracy curve of varying q for our proposed DADA-O on the open set adaptation setting $A \rightarrow W$ of Office-31 with a base network of AlexNet. The accuracy for unknown target instances is denoted by the blue line.

To investigate the influence of q on the performance, we conduct experiments by varying q . Figure 9 shows results for our proposed DADA on the open set adaptation setting $A \rightarrow W$ of Office-31 with a base network of AlexNet. As q increases, accuracies of OS and OS* decrease and the accuracy of Unknown increases, which means that the target instances are more likely classified as the unknown category. This confirms the statements we present in Section **Extension for Open Set Domain Adaptation** in the paper.

Table 12: Results for open set domain adaptation on Office-31 based on AlexNet. Note that all methods do not use unknown source instances. OS^* indicates the mean classification result over known categories whereas OS also includes the unknown category.

| Methods | A \rightarrow W | | D \rightarrow W | | W \rightarrow D | | A \rightarrow D | | D \rightarrow A | | W \rightarrow A | | Avg | |
|--|-------------------|-------------|-------------------|-------------|-------------------|-------------|-------------------|-------------|-------------------|-------------|-------------------|-------------|-------------|-------------|
| | OS | OS* | OS | OS* | OS | OS* | OS | OS* | OS | OS* | OS | OS* | OS | OS* |
| No Adaptation (Krizhevsky, Sutskever, and Hinton 2012) | 57.1 | 55.0 | 44.1 | 39.3 | 62.5 | 59.2 | 59.6 | 59.1 | 14.3 | 5.9 | 13.0 | 4.5 | 40.6 | 37.1 |
| DAN (Long et al. 2018a) | 41.5 | 36.2 | 34.4 | 28.4 | 62.0 | 58.5 | 47.8 | 44.3 | 9.9 | 0.9 | 11.5 | 2.7 | 34.5 | 28.5 |
| DANN (Ganin et al. 2016) | 31.0 | 24.3 | 33.6 | 27.3 | 49.7 | 44.8 | 40.8 | 35.6 | 10.4 | 1.5 | 11.5 | 2.7 | 29.5 | 22.7 |
| ATI- λ (Busto, Iqbal, and Gall 2018) | 65.3 | - | 82.2 | - | 92.7 | - | 72.0 | - | 66.4 | - | 71.6 | - | 75.0 | - |
| AODA (Saito et al. 2018c) | 70.1 | 69.1 | 94.4 | 94.6 | 96.8 | 96.9 | 76.6 | 76.4 | 62.5 | 62.3 | 82.3 | 82.2 | 80.4 | 80.2 |
| DADA-O | 75.5 | 75.6 | 91.2 | 93.0 | 93.3 | 94.4 | 82.7 | 83.9 | 73.5 | 74.8 | 71.1 | 71.6 | 81.2 | 82.2 |

When $q = 0$, the objective of the feature extractor is to align the whole source domain and the whole target domain, resulting in the misclassification of all unknown target instances as the known categories, as illustrated in Figure 9. This demonstrates that the model does not learn feature representations that can separate the unknown target instances from the known instances. To make a trade-off, we empirically set $q = 0.1$ for all open set adaptation settings.

D Investigation for Our Used Training Scheme

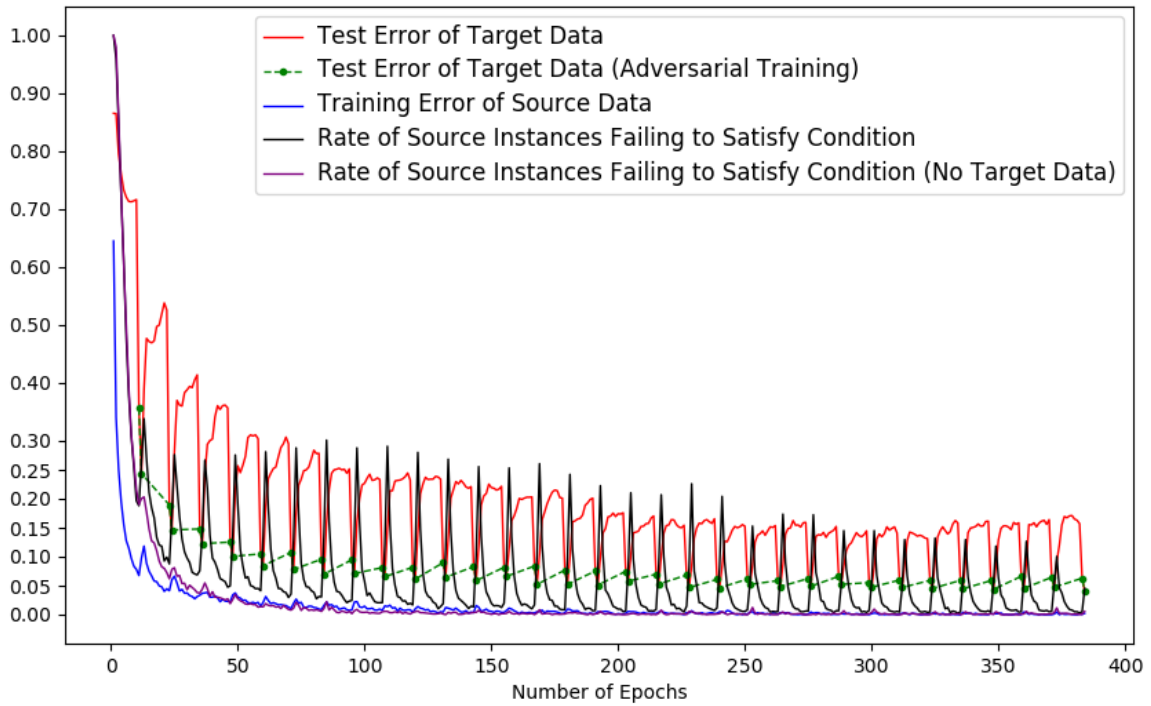
In this section, we investigate our used training scheme of pre-training DADA on the labeled source data and maintaining the same supervision signal in the adversarial training of DADA, on benchmark datasets of MNIST (Lecun et al. 1998) and USPS (Hull 1994), where two adaptation settings of **MNIST** \rightarrow **USPS** and **USPS** \rightarrow **MNIST** are built.

To always satisfy the condition of $p_{y^s}^s > 0.5$ discussed in Section **Discriminative Adversarial Learning** in the paper, we train DADA of $F(G(\cdot))$ by a well-designed scheme, which can be formulated as alternating the classification training on the labeled source data and the adversarial training of DADA on the labeled source data and unlabeled target data. We denote the number of training epochs or training iterations for classification training in each alternation respectively as T_{cls} and \hat{T}_{cls} , the number of training epochs or training iterations for adversarial training in each alternation respectively as T_{adv} and \hat{T}_{adv} , and the number of alternating the classification training and adversarial training as N_{alter} . For the two adaptation settings of **MNIST** \rightarrow **USPS** and **USPS** \rightarrow **MNIST**, T_{cls} , T_{adv} , and N_{alter} are respectively set to 10, 2, and 16, according to the rate of source instances failing to satisfy the condition; the hyper-parameter λ (cf. Section **Discriminative Adversarial Learning** in the paper for its definition) is not used, since T_{adv} is a quite small number. We investigate the efficacy of our used training scheme on keeping the condition satisfied by visualizing training processes on the two adaptation settings in Figure 10.

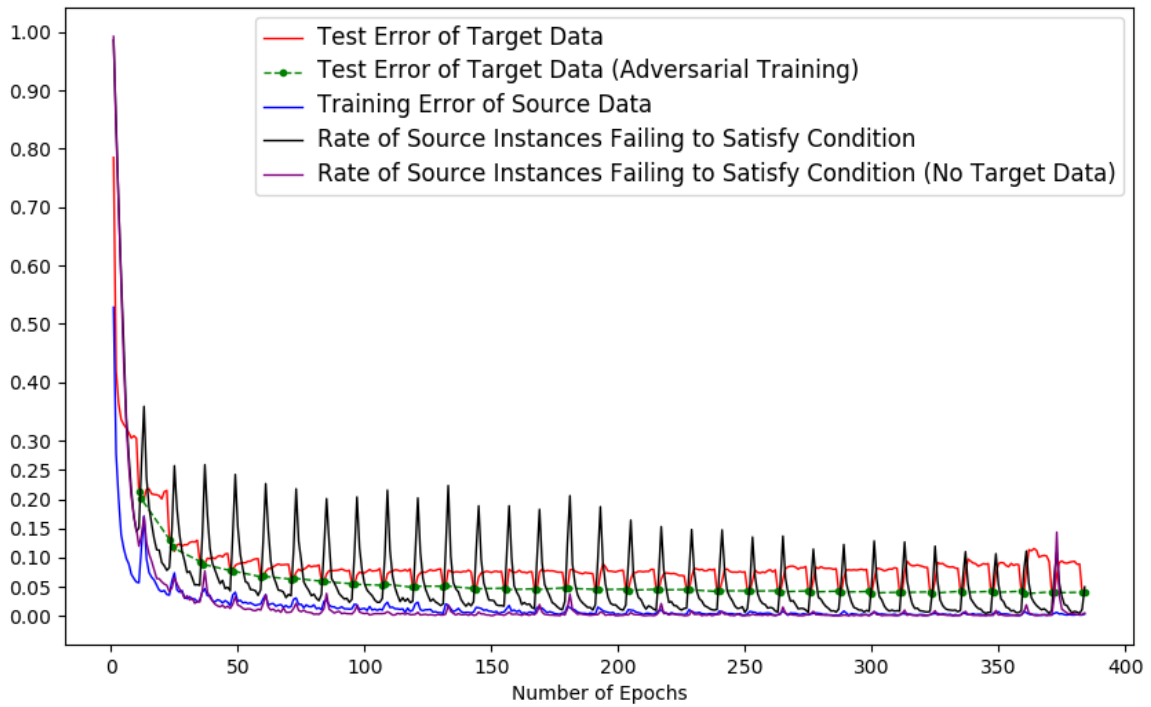
From Figure 10, we can obtain several interesting observations. (1) The classification training makes “Rate of Source Instances Failing to Satisfy Condition” fall into a valley whereas the adversarial training of DADA makes it rise to a peak, showing that a part of source instances change from satisfying the condition to not satisfying it during adversarial training. (2) “Rate of Source Instances Failing to Satisfy Condition (No Target Data)” is much lower than “Rate of Source Instances Failing to Satisfy Condition” at epochs of adversarial training, showing that the training of target data affects the source data and results in that a part of them do not satisfy the condition. (3) “Rate of Source Instances Failing to Satisfy Condition” declines to a very low value in an oscillatory manner,

showing the efficacy of this training scheme on keeping the condition satisfied. (4) “Training Error of Source Data” is low at epochs of adversarial training, showing that our proposed DADA has the same effect as classification training. (5) All valleys of “Test Error of Target Data” are derived from the adversarial training of DADA, showing the excellent efficacy of our proposed DADA in aligning the source and target domains. (6) At epochs of adversarial training, the lower “Rate of Source Instances Failing to Satisfy Condition” is, the more improvement of performance is obtained, showing the necessity of satisfying the condition. (6) The good performances of DADA on the two adaptation settings of **MNIST** \rightarrow **USPS** and **USPS** \rightarrow **MNIST**, which are very close to the perfect performance of 100%, confirm the efficacy of our proposed DADA in aligning the joint distributions of feature and category across the two domains.

For each closed set adaptation setting of Office-31, T_{cls} , \hat{T}_{adv} , and N_{alter} are respectively set to 200, 800, and 1. For the closed set adaptation setting of Syn2Real, \hat{T}_{cls} , \hat{T}_{adv} , and N_{alter} are respectively set to 2000, 1000, 1. For all these adaptation settings, the hyper-parameter λ is used.



(a) **MNIST**→**USPS** ($N_{alter} = 32$)



(b) **USPS**→**MNIST** ($N_{alter} = 32$)

Figure 10: Training processes in terms of the test error of the target data for each epoch, the test error of the target data for each epoch of adversarial training, the training error of the source data for each epoch, the rate of source instances failing to satisfy the condition for each epoch, and the rate of source instances failing to satisfy the condition for each epoch when no target data is used in the adversarial training, on the two adaptation settings of (a) **MNIST**→**USPS** and (b) **USPS**→**MNIST**.