

Unsupervised Multi-Target Domain Adaptation: An Information Theoretic Approach

Behnam Gholami
Pritish Sahu
Ognjen (Oggi) Rudovic
Konstantinos Bousmalis
Vladimir Pavlovic

bb510@cs.rutgers.edu
ps851@cs.rutgers.edu
orudovic@mit.edu
Konstantinos@google.com
vladimir@cs.rutgers.edu

Abstract

Unsupervised domain adaptation (**uDA**) models focus on pairwise adaptation settings where there is a single, labeled, source and a single target domain. However, in many real-world settings one seeks to adapt to multiple, but somewhat similar, target domains. Applying pairwise adaptation approaches to this setting may be suboptimal, as they fail to leverage shared information among multiple domains. In this work we propose an information theoretic approach for domain adaptation in the novel context of multiple target domains with unlabeled instances and one source domain with labeled instances. Our model aims to find a shared latent space common to all domains, while simultaneously accounting for the remaining private, domain-specific factors. Disentanglement of shared and private information is accomplished using a unified information-theoretic approach, which also serves to establish a stronger link between the latent representations and the observed data. The resulting model, accompanied by an efficient optimization algorithm, allows simultaneous adaptation from a single source to multiple target domains. We test our approach on three challenging publicly-available datasets, showing that it outperforms several popular domain adaptation methods.

1 Introduction

In real-world data, the training and test data often do not come from the same underlying distribution [40]. For instance, in the task of object recognition/classification from image data, this is may be due to the image noise, changes in the object view, etc., which induce different biases in the observed data sampled during the training and test stage. Consequently, assumptions made by traditional learning algorithms are often violated, resulting in degradation of the algorithms' performance during inference of test data. Domain Adaptation (**DA**) approaches (e.g., [16, 19, 23, 44]) aim to tackle this by transferring knowledge from a source domain (training data) to an unlabeled target domain (test data) to reduce the discrepancy between the source and target data distributions, typically by exploring domain-invariant data structures.

Existing **DA** methods can be divided into: (semi)supervised **DA**, and unsupervised **DA** [15]. The former assume that in addition to the labeled data of the source domain, some labeled data from the target domain are also available for training/adapting the classifiers. By contrast, the latter does not require any labels from the target domain but rather explores the similarity in the data distributions of the two domains. In this work, we focus on the unsupervised **DA** (**uDA**) scenario, which is more challenging due to the lack of correspondences in source and target labels.

Most works on **uDA** today focus on a single-source-single-target-domain scenario. However, in many real-world applications, unlabeled data may come from different domains, thus, with different statistical properties but with common task-related content. For instance, we may have access to images of the same class of objects (e.g., cars) recorded by various types of cameras, and/or under different camera views and at different times, rendering multiple different domains

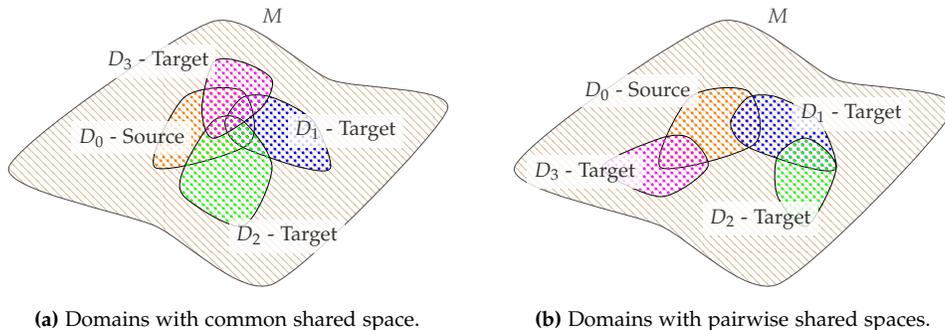


Figure 1: Illustration of domains with common (a) and pairwise-shared spaces (b). We tackle the domain adaptation task when all domains share a common task/space, which is then leveraged to transfer knowledge across multiple target domains.

(e.g., datasets). Likewise, facial expressions of emotions, such as joy and surprise, shown by different people and recorded under different views, result in multiple domains with varying data distributions. In most cases, these domains have similar *underlying* data distributions, which can be leveraged to build more effective and robust classifiers for tasks such as the object or emotion recognition across multiple datasets/domains. To this end, traditional **uDA** methods focus on the single-source-single-target **DA** scenario. However, in the presence of multiple domains, as typically encountered in real-world settings, this pair-wise adaptation approach may be suboptimal as it fails to leverage simultaneously the knowledge shared across multiple task-related domains.

Recently, Zhao et al. [48] showed that by having access to multiple source domains can facilitate better adaptation to a single target domain, when compared to the pair-wise **DA** approach. While this is intuitive due to the access to multiple *labelled* source domains, offering more adaptation flexibility for the target domain (i.e., by efficiently exploring the data labels across multiple source domains that are most related to the target domain), it comes at the expense of the data labelling in multiple source domains, which can be costly and time-consuming.

In either case, a single source domain or readily available multiple source domains, to the best of our knowledge, a simultaneous adaptation to *multiple* and *unlabelled* target domains remains an unexplored **DA** scenario. However, this **DA** scenario is important as we usually have access to multiple unlabeled domains; yet, the adaptation process is also more challenging due to the lack of supervision in the target domains. Still, multi-target **DA** can have advantages over a single-target **DA** when: (i) there is direct knowledge sharing between the source and multiple target domains (Fig. 1a), and (ii) the source and a target domain are related through another target domain (Fig. 1b). While this seems intuitive, it is critical how the data from multiple *unlabelled* target domains are leveraged within the multi-target **DA** approach, in order to improve its performance over the single target **DA** approaches and naive fusion of multiple target domains.

To this end, we propose a Multi-Target **DA**-Information-Theoretic-Approach (**MTDA-ITA**) for single-source-multi-target **DA**. We exploit a single source domain and focus on multiple target domains to investigate the effects of multi-target **DA**; however, the proposed model can easily be extended to multiple source domains. This approach leverages the data from multiple target domains to improve performance compared to individually learning from pairwise source-target domains. Specifically, we simultaneously factorize the information from each available target domain and learn separate subspaces for modeling the shared (i.e., correlated across the domains) and private (i.e., independent between the domains) subspaces of the data [38]. To this end, we employ deep learning to derive an information theoretic approach where we jointly maximize the

mutual information between the domain labels and private (domain-specific) features, while minimizing the mutual information between the the domain labels and the shared (domain-invariant) features. Consequently, the more robust feature representations are learned for each target domain by exploiting dependencies between multiple target domains. We show on benchmark datasets for **DA** that this approach leads to overall improved performance on each target domain, compared to independent **DA** for each pair of source-target domains, or the naive combination of multiple target domains, and state-of-the-art models applicable to the target task.

2 Preliminaries

2.1 Information Theory: Background

Let $\mathbf{x} = (x_1, x_2, \dots, x_n)$ denote a n -dimensional random variable with probability density function (pdf) given by $p(\mathbf{x})$. Shannon differential entropy is defined in the usual way as $H(\mathbf{x}) = -\mathbb{E}_{\mathbf{x}} [\ln p(\mathbf{x})]$ where \mathbb{E} denotes the expectation operator. Let $\mathbf{z} = (z_1, z_2, \dots, z_m)$ denote a m -dimensional random variable with pdf $p(\mathbf{z})$. Then mutual information between two random variables, \mathbf{x} and \mathbf{z} , is defined as $I(\mathbf{x}; \mathbf{z}) = H(\mathbf{x}) + H(\mathbf{z}) - H(\mathbf{x}, \mathbf{z})$. Mutual information can also be viewed as the reduction in uncertainty about one variable given another variable—i.e., $I(\mathbf{x}; \mathbf{z}) = H(\mathbf{x}) - H(\mathbf{x}|\mathbf{z}) = H(\mathbf{z}) - H(\mathbf{z}|\mathbf{x})$.

3 Method

In this section, we describe our proposed information theoretic approach that supports domain adaptation for multiple target domains simultaneously, finding factorized latent spaces that are non-redundant, and that can capture explicitly the shared (domain invariant) and the private (domain dependent) features of the data well suited for better generalization for domain adaptation.

3.1 Problem Formulation

Without loss of generalizability, we consider a multi-class (K -class) classification problem as the running example. Furthermore, let $(\mathbf{X}, \mathbf{Y}, \mathbf{D}) = \{(\mathbf{x}_i, \mathbf{y}_i, \mathbf{d}_i)\}_{i=0}^N$ be a collection of M domains (a labeled source domain, and $M - 1$ unlabeled target domains), where \mathbf{x}_i denotes the i -th sample, and $\mathbf{y}_i = [y_i^1, y_i^2, \dots, y_i^K]$ and $\mathbf{d}_i = [d_i^1, d_i^2, \dots, d_i^M]$ are the K -D and M -D encoding of the class and domain labels for \mathbf{x}_i , respectively. Note that the class labels are only available for the source samples.

The latent space representation of the data point \mathbf{x} is denoted as $\mathbf{z} = [\mathbf{z}_s, \mathbf{z}_p]$, where \mathbf{z}_s and \mathbf{z}_p are the (latent) shared and private features of the data point \mathbf{x} , respectively. Let \mathbf{z}_s and \mathbf{z}_p be some stochastic function of (\mathbf{x}, \mathbf{d}) parameterized by (θ_s, θ_p) , respectively, and \mathbf{y} be some stochastic function of \mathbf{z}_s parameterized by θ_c . We propose to maximize the following objective function:

$$\mathcal{L}(\theta_s, \theta_p, \theta_c; \mathbf{x}, \mathbf{y}, \mathbf{d}) = \lambda_r I(\mathbf{x}; \mathbf{z}) + \lambda_c I(\mathbf{y}; \mathbf{z}_s) + \lambda_d (I(\mathbf{d}; \mathbf{z}_p) - I(\mathbf{d}; \mathbf{z}_s)), \quad (3.1)$$

where $I(x; y)$ denotes the Mutual Information between the random variables \mathbf{x} and \mathbf{y} . λ_r, λ_c and λ_d denote the hyper-parameters controlling the weights of the objective terms. The proposed objective function (3.1) maximizes the three terms described below:

- $I(\mathbf{x}; \mathbf{z})$: encourages the latent features (both shared and private) to preserve information about the data samples (that can be used to reconstruct \mathbf{x} from \mathbf{z}).
- $I(\mathbf{y}; \mathbf{z}_s)$: enables to correctly predict the true class label of the samples out of their common shared features.

- $I(\mathbf{d}; \mathbf{z}_p) - I(\mathbf{d}; \mathbf{z}_s)$: encourages the latent private features to preserve the information about the domain label and penalizes the latent shared features to be domain informative. This not only reduces the redundancy in the shared and private features, but also, penalizes the redundancy of different private spaces, while preserving the shared information.

An additional term could be used to minimize the mutual information between the shared (\mathbf{z}_s) and private (\mathbf{z}_p) features. However, computing the mutual information (even approximating it) is intractable due to the highly complex joint distribution $p(\mathbf{z}_s, \mathbf{z}_p)$. Since we want \mathbf{z}_s and \mathbf{z}_p features to encode different aspects of \mathbf{x} , we enforce such constraint by jointly maximizing the term: $I(\mathbf{d}; \mathbf{z}_p) - I(\mathbf{d}; \mathbf{z}_s)$.

3.2 Optimization

The following lower bound for mutual information is derived using the non-negativity of KL-divergence [7]; i.e., $\sum_{\mathbf{x}} p(\mathbf{x}|\mathbf{z}) \ln \frac{p(\mathbf{x}|\mathbf{z})}{q(\mathbf{x}|\mathbf{z})} \geq 0$ gives:

$$I(\mathbf{x}; \mathbf{z}) \geq H(\mathbf{x}) + \mathbb{E}_{p(\mathbf{x}, \mathbf{z})} [\ln q(\mathbf{x}|\mathbf{z}; \phi)], \quad (3.2)$$

where $H(\mathbf{x})$ denotes the Shannon Entropy of the random variable \mathbf{x} . $q(\mathbf{x}|\mathbf{z}; \phi)$ is any arbitrary distribution parametrized by ϕ . We need a variational distribution $q(\mathbf{x}|\mathbf{z}; \phi)$ because the posterior distribution $p(\mathbf{x}|\mathbf{z}) = p(\mathbf{z}|\mathbf{x})p(\mathbf{x})/p(\mathbf{z})$ is intractable since the true data distribution $p(\mathbf{x})$ is assumed to be unknown. Similarly, we can derive lower bounds for $I(\mathbf{d}; \mathbf{z}_p) \geq H(\mathbf{d}) + \mathbb{E}_{p(\mathbf{d}, \mathbf{z}_p)} [\ln q(\mathbf{d}|\mathbf{z}_p; \psi)]$ and $I(\mathbf{d}; \mathbf{z}_s) \geq H(\mathbf{d}) + \mathbb{E}_{p(\mathbf{d}, \mathbf{z}_s)} [\ln q(\mathbf{d}|\mathbf{z}_s; \psi)]$, where $q(\mathbf{d}|\mathbf{z}_p; \psi)$ is any arbitrary distribution parametrized by ψ .¹ We further compute $I(\mathbf{y}; \mathbf{z}_s)$ as $I(\mathbf{y}; \mathbf{z}_s) = H(\mathbf{y}) + \mathbb{E}_{p(\mathbf{y}, \mathbf{z}_s)} [\ln p(\mathbf{y}|\mathbf{z}_s)]$.

Let next $E_s(\mathbf{x}; \theta_s)$ be a function parameterized by θ_s that maps a sample \mathbf{x} to the *shared* features \mathbf{z}_s , and $E_p(\mathbf{x}; \theta_p)$ be an analogous function which maps \mathbf{x} to \mathbf{z}_p , the features that are *private* to each domain (Fig. 2). We also define $F(\mathbf{z}_s, \mathbf{z}_p; \phi)$ as a decoding function mapping the concatenation of the latent features \mathbf{z}_s and \mathbf{z}_p to a sample reconstruction $\hat{\mathbf{x}}$, and $D(\mathbf{z}; \psi)$ as a decoding function mapping \mathbf{z}_s and \mathbf{z}_p to a M -dimensional probability vector: the predictions of the domain label $\hat{\mathbf{d}}$. Finally, $C(\mathbf{z}_s; \theta_c)$ is a task-specific function mapping \mathbf{z}_s to a K -dimensional probability vector of the class label $\hat{\mathbf{y}}$.

By representing $p(\mathbf{d}), p(\mathbf{x}), p(\mathbf{y})$ as the empirical distribution of a finite training set (e.g. $p(\mathbf{d}) = \frac{1}{N} \sum_{i=1}^N \delta(d - d_i)$) as in the case of variational autoencoders (VAE) [1, 35], parametrizing $p(\mathbf{z}_s|\mathbf{x})$, and $p(\mathbf{z}_p|\mathbf{x})$ as deterministic networks $E_s(\mathbf{x})$ and $E_p(\mathbf{x})$ respectively, and modeling the variational distributions as $\ln q(\mathbf{x}|\mathbf{z}; \psi) = \|\mathbf{x} - F(\mathbf{z}; \phi)\|_1$, $\ln q(\mathbf{d}|\mathbf{z}) = \mathbf{d}^\top \ln D(\mathbf{z}; \psi)$, and $\ln p(\mathbf{y}|\mathbf{z}_s) = \mathbf{y}^\top \ln C(\mathbf{z}_s; \theta_c)$, where $\|\cdot\|_1$ denotes the L_1 norm, the optimization task can be posed as a minimax saddle point problem, where we use adversarial training to maximize (3.1) w.r.t. the stochastic parameters $(\theta_s, \theta_p, \theta_c)$, and to minimize (3.1) w.r.t. the variational parameters (ϕ, ψ) , using Stochastic Gradient Descent (SGD).

Optimizing the parameters ϕ of the decoder F

$$\hat{\phi} = \arg \min_{\phi} \mathcal{L}_F = \frac{\lambda_r}{N} \sum_{i=1}^N \|\mathbf{x}_i - F(E_s(\mathbf{x}_i), E_p(\mathbf{x}_i))\|_1. \quad (3.3)$$

The decoder $F(\mathbf{z}_s, \mathbf{z}_p; \phi)$ is trained in such a way so as to minimize the difference between original input \mathbf{x} and its decoding from corresponding shared and private features via the decoder $F(\cdot)$.

¹Note that, for simplicity, we shared the parameters ψ between the approximate posterior distributions $q(\mathbf{d}|\mathbf{z}_s, \psi)$ and $q(\mathbf{d}|\mathbf{z}_p, \psi)$.

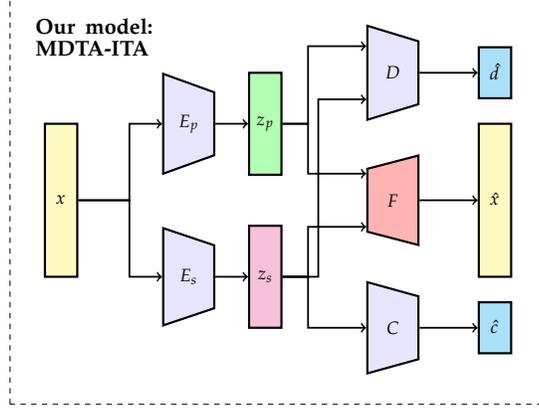


Figure 2: MDTA-ITA: The encoder $E_s(\mathbf{x})$ captures the feature representations (z_s) for a given input sample \mathbf{x} that are shared among domains. $E_p(\mathbf{x})$ captures domain-specific private features (z_p) using the *shared* private encoder. The shared decoder $F(z_p, z_s)$ learns to reconstruct the input sample by using both the private and shared features. The domain classifier D learns to correctly predict the domain labels of the actual samples from both their shared and private features while the classifier C learns to correctly predict the class labels from the shared features.

Optimizing the parameters ψ of the domain classifier D

$$\hat{\psi} = \arg \min_{\psi} \mathcal{L}_D = -\frac{\lambda_d}{N} \sum_{i=1}^N \mathbf{d}_i^\top \ln D(E_s(\mathbf{x}_i)) - \frac{\lambda_d}{N} \sum_{i=1}^N \mathbf{d}_i^\top \ln D(E_p(\mathbf{x}_i)). \quad (3.4)$$

$D(\mathbf{z}; \psi)$ can be considered as a classifier whose task is to distinguish between the shared/private features of the different domains. More precisely, the two terms in Eq. 3.4 encourage D to correctly predict the domain labels from the shared and private features, respectively.

Optimizing the parameters θ_c of the label classifier C

$$\hat{\theta}_c = \arg \min_{\theta_c} \{ -H(\mathbf{y}) - \mathbb{E}_{p(\mathbf{y}, \mathbf{z}_s)} [\ln p(\mathbf{y}|\mathbf{z}_s)] \}. \quad (3.5)$$

Since we have access to the source labels, $H(\mathbf{y})$ is a constant for source samples. we can approximate $H[\mathbf{y}]$ for the target samples using the output of the classifier C , leading to the following optimization problem:

$$\begin{aligned} \hat{\theta}_c = \arg \min_{\theta_c} \mathcal{L}_C = & -\frac{1}{N} \sum_{i=1}^{N_s} \mathbf{y}_i^\top \ln C(E_s(\mathbf{x}_i)) - \frac{\lambda_c}{N - N_s} \sum_{i=N_s+1}^N C(E_s(\mathbf{x}_i))^\top \ln C(E_s(\mathbf{x}_i)) \\ & + \frac{\lambda_c}{N - N_s} \sum_{i=N_s+1}^N C(E_s(\mathbf{x}_i))^\top \ln \left(\frac{1}{N - N_s} \sum_{i=N_s+1}^N C(E_s(\mathbf{x}_i)) \right), \end{aligned} \quad (3.6)$$

where N_s denotes the number of source samples. Intuitively, we enforce the classifier $C(\mathbf{z}_s; \theta_c)$ to correctly predict the class labels of the source samples by the first term in Eq. 3.6. We use the second term to minimize the entropy of $p(\mathbf{y}|\mathbf{z}_s)$ for the target samples; effectively, reducing the effects of “confusing” labels of target samples, as given by $p(\mathbf{y}|\mathbf{z}_s)$ that leads to decision boundaries occur far away from target data-dense regions in the feature space. The intuition behind the last

term is that if we minimize the entropy only(second term), we may arrive at a degenerate solution where every point x_t is assigned to the same class. Hence, the last term encourages the classifier $C(\cdot)$ to have balanced labeling for the target samples where it reaches its minimum, $\ln K$, when each class is selected with uniform probability.

Optimizing the parameter θ_p of the private encoder E_p

$$\hat{\theta}_p = \arg \min_{\theta_p} \mathcal{L}_p = \frac{\lambda_r}{N} \sum_{i=1}^N \|\mathbf{x}_i - F(E_s(\mathbf{x}_i), E_p(\mathbf{x}_i))\|_1 - \frac{\lambda_d}{N} \sum_{i=1}^N \mathbf{d}_i^\top D(E_p(\mathbf{x}_i)). \quad (3.7)$$

The first term in Eq. 3.7 encourages the private encoder $E_p(\mathbf{x}; \theta_p)$ to preserve the recovery ability of the private features. The second term, $E_p(\cdot)$, enforces distinct private features be produced for each domain by penalizing the representation redundancy in different private spaces. This, in turn, encourages moving this common information from multiple domains to their shared space.

Optimizing the parameter θ_s of the shared encoder E_s

$$\begin{aligned} \hat{\theta}_s = \arg \min_{\theta_s} \mathcal{L}_s = & \frac{\lambda_r}{N} \sum_{i=1}^N \|\mathbf{x} - F(E_s(\mathbf{x}_i), E_p(\mathbf{x}_i))\|_1 - \frac{\lambda_c}{N} \sum_{i=1}^{N_s} \mathbf{y}_i^\top \ln C(E_s(\mathbf{x}_i)) \\ & - \frac{\lambda_d}{N} \sum_{i=1}^N \mathbf{d}_i^\top \ln D(E_s(\mathbf{x}_i)) - \frac{\lambda_c}{N - N_s} \sum_{i=N_s+1}^N C(E_s(\mathbf{x}_i))^\top \ln C(E_s(\mathbf{x}_i)) \\ & + \frac{\lambda_c}{N - N_s} \sum_{i=N_s+1}^N C(E_s(\mathbf{x}_i))^\top \ln \left(\frac{1}{N - N_s} \sum_{i=N_s+1}^N C(E_s(\mathbf{x}_i)) \right). \end{aligned} \quad (3.8)$$

The first term in Eq. 3.8 encourages the shared encoder $E_s(\mathbf{x}; \theta_s)$ to preserve the recovery ability of the shared features. The second term is the source domain classification loss penalty that encourages E_s to produce discriminative features for the labeled source samples. The third term simulates the adversarial training by trying to fool the domain classifier $D(\cdot)$ when predicting the domain labels \mathbf{d} , given the shared features \mathbf{z}_s . The effect of this is two-fold: (i) the rendered shared features are more distinct from the corresponding private features, (ii) the shared features of different domains are encouraged to be similar to each other. The last two terms encourage $E_s(\cdot)$ to produce the shared features for target samples so that the classifier is confident on the unlabeled target data, driving the shared features away from the decision boundaries. To train our model, we alternate between updating the shared encoder $E_s(\cdot)$, the private encoder $E_p(\cdot)$, the decoder $F(\cdot)$, the classifier $C(\cdot)$, and the domain classifier $D(\cdot)$ using the SGD algorithm (see Algorithm 1 for more details).

4 Related Work

There has been extensive prior work on domain adaptation [15]. Recent papers have focused on transferring deep neural network representations from a labeled source dataset to an unlabeled target domain, where the main strategy is to find a feature space such that the confusion between source and target distributions in that space is maximized ([8, 10, 14, 33, 36, 37, 43, 46]). For this, it is critical to first define a measure of divergence between source and target distributions. For instance, several methods have used the Maximum Mean Discrepancy (MMD) loss for this purpose (e.g., [10, 28, 45]). MMD computes the norm of the difference between two domain means in the

Algorithm 1 MDTA-ITA Algorithm

Require: $\{\mathbf{X}, \mathbf{Y}, \mathbf{D}\}$: M domain datasets. $\lambda_r, \lambda_c, \lambda_d$: Model hyper-parameters. η : Learning rate.**Ensure:** $\theta_s, \theta_p, \theta_c, \phi, \psi$: Model parameters.

- 1: Initialize $\theta_s, \theta_p, \theta_c, \phi, \psi$;
 - 2: **repeat**
 - 3: Sample a mini-batch from each of source/target domain datasets.
 - 4: Update $\{\theta_s\}$ by minimizing \mathcal{L}_s in Eq.(3.8) through the gradient descent: $\theta_s = \theta_s - \eta \frac{\partial \mathcal{L}_s}{\partial \theta_s}$.
 - 5: Update $\{\theta_p\}$ by minimizing \mathcal{L}_p in Eq.(3.7) through the gradient descent: $\theta_p = \theta_p - \eta \frac{\partial \mathcal{L}_p}{\partial \theta_p}$.
 - 6: Update $\{\theta_c\}$ by minimizing \mathcal{L}_c in Eq.(3.6) through the gradient descent: $\theta_c = \theta_c - \eta \frac{\partial \mathcal{L}_c}{\partial \theta_c}$.
 - 7: Update $\{\phi\}$ by minimizing \mathcal{L}_ϕ in Eq.(3.4) through the gradient descent: $\phi = \phi - \eta \frac{\partial \mathcal{L}_\phi}{\partial \phi}$.
 - 8: Update $\{\psi\}$ by minimizing \mathcal{L}_ψ in Eq.(3.3) through the gradient descent: $\psi = \psi - \eta \frac{\partial \mathcal{L}_\psi}{\partial \psi}$.
 - 9: **until** Convergence;
 - 10: return $\{\theta_s, \theta_p, \theta_c, \phi, \psi\}$.
-

reproducing Kernel Hilbert Space (RKHS) induced by a pre-specified kernel. The Deep Adaptation Network (DAN) [29] applied MMD to layers embedded in a RKHS, effectively matching higher order statistics of the two distributions. The deep Correlation Alignment (CORAL) method [39] attempts to match the mean and covariance of the two distributions. Deep Transfer Network (DTN) [47] achieved source/target distribution alignment via two types of network layers based on MMD distance: the shared feature extraction layer, which learns a subspace that matches the marginal distributions of the source and the target samples, and the discrimination layer, which matches the conditional distributions by classifier transduction.

Recently proposed unsupervised DA methods ([8, 14, 33, 36, 37, 46]) operate by training deep neural networks using adversarial training, which allows the learning of feature representations that are simultaneously discriminative of source labels, and indistinguishable between the source and target domain. For instance, Ganin et al. [17] proposed a DA mechanism called Domain-Adversarial Training of Neural Networks (DANN), which enables the network to learn domain invariant representations in an adversarial way by adding a domain classifier and back-propagating inverse gradients. Adversarial Discriminative Domain Adaptation (ADDA) [42] learns a discriminative feature subspace using the source labels, followed by a separate encoding of the target data to this subspace using an asymmetric mapping learned through a domain-adversarial loss. Liu et al. [27] makes a shared-latent space assumption and proposes an unsupervised image-to-image translation (UNIT) framework based on Coupled GANs [26]. Another example is the pixel-level domain adaptation models that perform the distribution alignment not in the feature space but directly in raw pixel space. PixelDA [10] uses adversarial approaches to adapt source-domain images as if drawn from the target domain while maintaining the original content.

While these approaches have shown success in DA tasks with single source-target domains, they are not designed to leverage information from multiple domains simultaneously. More recently, Zhao et al. [48] introduced an adversarial framework called MDAN for multiple source single target domain adaptation where a domain classifier, induced by minimizing the H-divergence between multiple source and a target domain, is used to align their feature distributions in a shared space. Instead, in our approach we focus on multi-target DA where we perform adaptation of multiple *unlabelled* target domains. Although both our model and MDAN use the similar notion of the domain classifier to minimize the domain mismatch in shared space, the domain classifier

induced by our information-theoretic (IT) loss also acts to separate domains in the private space (see Eqs. 3.4, and 3.7 for more details), improving the essential reconstruction ability, similar to [9].

4.1 Connection to Information Theoretic Representation Learning

The idea of using information theoretic (IT) objectives for representation learning was originally introduced in Tishby & Zaslavsky [41]. Since their approach for optimizing the IT objective functions relied on the iterative Blahut Arimoto algorithm [41], it is not feasible to apply to DNN frameworks. Similar to our approach, there have been some recent works [3-5, 12, 32] to approximate the MI by applying variational bounds on MI, though not in the context of domain adaptation.

Mohamed & Rezende [32] utilized the variational bounds on MI, and apply it to deep neural networks in the context of reinforcement learning. Chalk et al. [12] and Amini et al. [4], developed the same variational lower bound in the context of Information Bottleneck (IB) principle [41], where the former applied it to sparse coding problems, and used the kernel trick to achieve nonlinear mappings, whereas the latter applied it to deep neural networks to handle large datasets thanks to the SGD algorithm. Achille & Soatto [2] proposed a variational bound on the IM in the context of IB, from the perspective of variational dropout and demonstrated its utility in learning disentangled representations for variational autoencoders.

The main difference between our method and the above methods is that these methods throw away the information in the data not related to the task by minimizing the mutual information between the data points and the latent representations that may lead to ignoring the individual characteristics (private features) of the datasets in a multiple dataset regime, whereas our method explicitly models what is unique to each domain (dataset) that improves the model’s ability to extract domain-invariant features.

In the unsupervised representation learning literature, our work is also related to the VAE-based models [11]. However, we propose to tackle the task using our IT approach using deterministic mappings instead of the traditional evidence lower bound (ELBO) optimization with stochastic mappings. One of the main drawbacks of ELBO-based approaches is that they can result in poor latent representation, when a powerful decoder effectively ignores the latent space [11]. In the spirit of recent works on improved representational learning [3], we seek to recover a good latent representation by replacing the ELBO objective with an IT-driven loss. In contrast to the unsupervised representation learning approaches, our setting also allows us to further improve the latent representation using the labeled data in the source domain while leveraging the sharing of dependencies across different target domains.

4.2 Connection to Multiple Domain Transfer Networks

Recent studies have shown remarkable success in multiple domain transfer (MDT) [6, 13, 20, 21] though not in the context of the image classification, rather in the context of image generation. Choi et al. [13] proposed **StarGAN**, a generative adversarial network capable of learning mappings among multiple domains in the context of image to image translation framework. The goal of **StarGAN** is to train a single generator G though this requires passing in a vector along with each input to the generator specifying the output domain desired, that learns mappings among multiple domains. To achieve this, G is trained to translate an input image x into an output image x' conditioned on the target domain label d , $G(x, d) \rightarrow x'$. Similar to our domain classifier module D , they introduce an auxiliary classifier that allows a single discriminator to control multiple domains.

Anoosheh et al. [6] introduced **ComboGAN**, which decouples the domains and networks from each other. Similar to our encoder/decoder modules, **ComboGAN**’s generator networks contain encoder/decoders assigning each encoder and decoder to a domain. They combine the encoders

and decoders of the trained model like building blocks, taking as input any domain and outputting any other. For example during inference, to transform an image x from an arbitrary domain X to x' from domain X' , they simply perform $x' = G_{X'X}(x) = Decoder'_{X'}(Encoder_X(x))$. The result of $Encoder_X(x)$ can even be cached when translating to other domains as not to repeat computation.

The main differences between the MDT methods and ours is that, unlike our method which does domain alignment in feature space, MDT methods adapt representations not in feature space but rather in raw pixel space; translating samples from one domain to the “style” of a other domains. This works well for limited domain shifts where the domains are similar in pixel-space, but can be too limiting for settings with larger domain shifts that results in poor performance in significant structural change of the samples in different domains.

4.3 Connection to Domain Separation Networks

The method closest to our work is Domain Separation Networks (DSN) [9], which use the notion of auto-encoders to explicitly separate the feature representations private to each source/target domain from those that are shared between the domains. Although extending DSN to multiple domains might seem trivial, DSN requires an autoencoder per domain, making the model impractical in the case of more than a couple of domains.

The overall loss of DSN consists of a reconstruction loss for each domain modeled by a shared decoder, a similarity loss such as MMD, which encourages domain invariance modeled by a shared encoder, and a dissimilarity loss modeled by two private encoders: one for the source domain and one for the target domain. While one could attempt to generalize DSN to multiple target domains by having individual per-target domain private encoders, doing so would prove problematic when the number of target domains is large — each private encoder would require a large “private” dataset to learn the private parameters. Precisely, for multiple (M) target domains, we could train a DSN model with one shared encoder, $M + 1$ private encoder (one for each domain), and one shared decoder. This leads to $M + 3$ models to train that implies the number of models increases linearly with the number of domains, as does the required training time. Second, DSN uses an orthogonality constraint among the shared and the private representations which may not be strong enough to remove redundancy and enforce disentangling among different private spaces. Precisely, DSN defines the loss via a soft subspace orthogonality constraint between the private and shared representation of each domain. However, it does not enforce the private representation of different domains to be different that may result in redundancy of different private spaces.

In addition, DSN enforces separation of spaces using the notion of Euclidean orthogonality, e.g., $\|z_s - z_p\|^2$. In case of multiple target domains, this would result in learning of all pairs of private spaces independently. To address those deficiencies, we first explicitly couple different private encoders into a single private encoder model, E_{θ_p} of Fig. 2, which allows us to generalize to an arbitrary number of target domains. To assure that the information among the private and shared spaces is not shared (i.e., “orthogonal”), we define an information-theoretic criteria enforced by a domain classifier, D_ψ of Fig. 2, which aims to segment the private space into clusters that correspond to individual target domains. By using D_ψ within the adversarial framework, MTDA-ITA learns simultaneously the shared and private features from different domains (see Fig. 5). We also show in Sec. 5 that our model performs better than the trivial extension of DSNs to the multi-domain case.

5 Experimental Results

We compare the proposed method with state-of-the-art methods on standard benchmark datasets: a digit classification task that includes 4 datasets: MNIST [24], MNIST-M [18], SVHN [34],

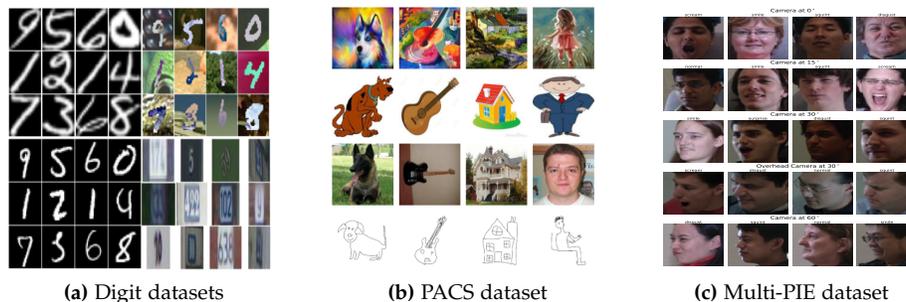


Figure 3: Exemplary images from different datasets. a) Digits datasets, b) PACS dataset (first row: Art-painting, second row: Cartoon, Third row: Photo, last row: Sketch), c) Multi-PIE dataset (each row corresponds to a different camera angle and each subject depicts an expression (“normal”, “smile”, “surprise”, “squint”, “disgust”, “scream”) at every camera position).

USPS [42], Multi-PIE expression recognition dataset², and PACS multi-domain image recognition benchmark [25], a new dataset designed for the cross-domain recognition problems. Figure 3 illustrates image samples from different datasets and domains. We evaluate the performance of all methods with classification accuracy metric.

We used ADAM [22] for training; the learning rate was set to 0.0002 and momentums to 0.5 and 0.999. We used batches of size 16 from each domain, and the input images were mean-centered/rescaled to $[-1, 1]$. The hyper-parameters are empirically set as $\lambda_r = 1.0$, $\lambda_c = 0.01$, $\lambda_d = 0.20$.

For the network architecture, our private/shared encoders consisted of three convolutional layers as the front-end and four basic residual blocks as the back-end. The decoder consisted of four basic residual blocks as the front-end and three transposed convolutional layers as the back-end. The discriminator and the classifier consisted of stacks of convolutional layers. We used ReLU for nonlinearity. TanH function is used as the activation function of the last layer in the decoder F for scaling the output pixels to $[-1, 1]$. The details of the networks are given in Appendix.

The quantitative evaluation involves a comparison of the performance of our model to previous work and to “Source Only” and “1-NN” baselines that do not use any domain adaptation. For “Source Only” baseline, we train models only on the unaltered source training data and evaluate on the target test data. We compare the proposed method MTDA-ITA with several related methods designed for pair-wise source-target adaptation: CORAL [39], DANN [17], ADDA [42], DTN [47], UNIT [27], PixelDA [10], and DSN [9]. We reported the results of two following baselines: (i) one is to combine all the target domains into a single one and train it using MTDA-ITA, which we denote as (c-MTDA-ITA). (ii) the other one is to train multiple MTDA-ITA separately, where each one corresponds to a source-target pair which we denote as (s-MTDA-ITA). For completeness, we reported the results of the competing methods by combining all the target domains into a single one (denoted by c-DTN, c-ADDA, and c-DSN) as well. We also extend DSN to multiple domains by adding multiple private encoder to it (denoted by mp-DSN) and contrast it with our model.

5.1 Digits Datasets

We combine four popular digits datasets (MNIST, MNIST-M, SVHN, and USPS) to build the multi-target domain dataset. All images were uniformly rescaled to 32×32 . We take each of

²<http://www.cs.cmu.edu/afs/cs/project/PIE/MultiPie/Multi-Pie/Home.html>

method	S → M	S → MM	S → U	M → S	M → MM	M → U	MM → S	MM → M	MM → U	U → S	U → M	U → MM
Source Only	62.10	40.43	39.90	30.29	55.98	78.30	40.00	84.46	80.43	23.41	50.64	41.45
r-NN	35.86	18.21	29.31	28.01	12.58	41.22	21.45	82.13	36.90	15.34	38.45	18.54
CORAL [39]	63.10	54.37	50.15	33.40	57.70	81.05	40.20	84.90	87.54	38.90	85.01	60.45
DANN [18]	73.80	61.05	62.54	35.50	77.40	81.60	51.80	61.05	85.34	35.50	77.40	61.60
ADDA [42]	77.68	64.23	64.10	30.04	91.47	90.51	40.64	92.82	80.70	41.23	90.10	56.21
c-ADDA	80.10	56.80	64.80	27.50	83.30	84.10	35.43	88.47	74.19	39.36	84.67	52.54
DTN [47]	81.40	63.70	60.12	40.40	85.70	85.80	48.80	88.80	90.68	42.43	89.04	55.78
c-DTN	82.10	59.30	56.87	38.32	80.90	79.31	44.21	83.60	84.98	39.75	85.04	48.86
PixelDA [10]	–	–	–	–	98.10*	94.10*	–	–	–	–	–	–
UNIT [27]	90.6*	–	–	–	–	92.90	–	–	–	–	90.60	–
DSN [9]	82.70	64.80	65.30	49.30	83.20	91.65	51.50	90.20	89.95	48.20	91.40	60.45
c-DSN	83.10	60.56	60.35	46.80	80.49	88.21	47.10	84.60	84.80	40.50	86.05	56.25
mp-DSN	83.40	61.00	58.10	47.35	79.30	78.45	47.15	85.51	83.24	38.30	87.40	55.47
s-MTDA-ITA	82.90	63.10	63.54	49.60	82.42	89.21	50.55	94.82	89.05	40.13	87.10	61.01
c-MTDA-ITA	79.20	59.90	63.70	45.30	77.12	87.47	47.32	90.20	90.01	41.10	85.35	60.31
MTDA-ITA	84.60	65.30	70.03	52.01	85.50	94.20	53.50	98.20	94.10	46.00	91.50	67.30

Table 1: Mean classification accuracy on digit classification. M: MNIST; MM: MNIST-M, S: SVHN, U: USPS. The best is shown in red. c-X: combining all target domains into a single one and train it using X. s-MTDA-ITA: training multiple MTDA-ITA where each one correspond to a source-target pair. mp-DSN: extended DSN with multiple private encoder. *UNIT trains with the extended SVHN (> 500K images vs ours 72K). *PixelDA uses ($\approx 1,000$) of labeled target domain data as a validation set for tuning the hyperparameters.

MNIST-M, SVHN, USPS, and MNIST as source domain in turn, and the rest as targets. We use all labeled source images and all unlabeled target images, following the standard evaluation protocol for unsupervised domain adaptation [18, 30].

We show the accuracy of different methods in Tab. 1. The results show that first of all **cMTDA-ITA** has worse performance than **sMTDA-ITA** and **MTDA-ITA**. We have similar observations for **ADDA**, **DTN**, and **DSN** that demonstrates a naive combination of different target datasets can sometimes even decrease the performance of the competing methods. Moreover, **MTDA-ITA** outperforms the state-of-the-art methods in most of domain transformations. The higher performance of **MTDA-ITA** compared to other methods is mainly attributed to the joint adaptation of related domains where each domain could benefit of other related domains. Furthermore, from the results obtained, we see that it is beneficial to use information coming from unlabeled target data (see Eq. 3.6 for updating the classifier C_{θ_c}) during the learning process, compared to when no data from target domain is used (See the ablation study section for more information). Indeed, using our scheme, we find a representation space in which embeds the knowledge from the target domain into the learned classifier. By contrast, the competing methods do not provide a principled way of sharing information across all domains, leading to overall lower performance. The results also verify the superiority of **MTDA-ITA** over **mp-DSN**. This can be due to (i) having multiple private encoders increase the number of parameters that may lead to model overfitting, (ii) superiority of the **MTDA-ITA**’s domain adversarial loss over the **DSN**’s MMD loss to separate the shared and private features, (iii) utilization of the unlabeled target data to regularize the classifier in **MTDA-ITA**.

5.2 Multi-PIE dataset

The **Multi-PIE** dataset includes face images of 337 individuals captured from different expressions, views, and illumination conditions (Fig. 3(c)). For this experiment, we use 5 different camera views (positions) C05, C08, C09, C13, and C14 as different domains (Fig. 3(c)) and the face expressions (normal, smile, surprise, squint, disgust, scream) as labels. Each domain contains 27120 images of size $64 \times 64 \times 3$. We used each view as the source domain, in turn, and the rest as targets. We expect the face inclination angle to reflect the complexity of transfer learning. Tables 2 and 3 shows the classification accuracy for C05, C08, C09, C13, and C14 as source domain. As can be seen,

method	Co5 → Co8	Co5 → Co9	Co5 → C13	Co5 → C14	C13 → Co5	C13 → Co8	C13 → Co9	C13 → C14	C14 → Co5	C14 → Co8	C14 → Co9	C14 → C13
Source Only	31.56	40.67	39.89	54.70	50.79	45.90	40.04	59.68	60.03	36.80	40.11	60.57
1-NN	27.28	31.22	33.66	47.04	33.21	37.01	34.45	48.79	47.44	28.24	30.86	44.86
CORAL [39]	36.55	38.60	40.60	55.29	54.89	48.90	40.30	68.90	59.98	40.63	40.80	65.11
DANN[18]	40.30	41.20	40.12	58.90	57.86	50.30	45.30	70.68	57.20	40.22	40.77	70.50
ADDA[42]	33.21	30.86	52.44	70.18	64.83	63.20	55.48	74.25	73.62	43.56	38.68	72.84
c-ADDA	46.88	36.38	39.14	65.41	59.20	30.70	53.20	68.33	65.88	30.60	45.34	64.30
DTN[47]	38.50	30.56	55.78	68.90	63.78	60.45	60.55	72.60	70.67	41.55	41.45	70.67
c-DTN	41.70	31.10	50.19	60.34	57.53	55.24	57.14	65.16	63.80	38.97	39.80	62.10
PixelDA[10]	44.93	44.75	45.18	46.88	45.68	44.95	44.45	90.50	46.28	45.89	44.45	69.15
UNIT[27]	44.47	44.47	44.47	44.51	44.14	44.47	44.21	44.47	43.03	44.44	44.47	44.47
DSN[9]	45.12	44.35	48.12	75.00	64.15	57.70	49.15	80.75	82.20	38.75	45.00	80.50
c-DSN	42.52	38.54	34.15	69.45	57.34	31.63	51.17	74.52	82.01	34.25	42.63	79.42
mp-DSN[42]	41.30	35.14	34.40	65.70	55.20	30.40	47.80	75.30	80.75	30.20	43.00	79.02
s-MTDA-ITA	44.40	44.60	47.65	80.20	70.10	58.90	58.10	80.12	82.05	45.90	52.67	81.60
c-MTDA-ITA	40.49	40.70	42.80	71.60	60.34	55.67	57.10	73.50	76.80	43.10	48.10	80.90
MTDA-ITA	49.01	48.23	53.13	84.29	78.40	66.70	70.30	85.49	87.20	61.40	60.05	86.70

Table 2: Mean classification accuracy on Multi-PIE classification. The best is shown in red.

method	Co8 → Co5	Co8 → Co9	Co8 → C13	Co8 → C14	Co9 → Co5	Co9 → Co8	Co9 → C13	Co9 → C14
Source Only	33.70	50.10	50.80	40.13	33.32	48.24	49.24	36.19
1-NN	28.75	35.39	39.79	32.13	26.82	35.30	34.26	28.41
CORAL [39]	35.89	55.79	60.00	40.67	35.89	51.56	50.45	40.67
DANN[18]	40.20	56.89	55.83	43.25	50.63	58.40	55.81	48.90
ADDA[42]	37.40	58.40	60.40	42.10	29.40	53.30	45.30	38.30
c-ADDA	41.60	39.65	50.00	46.25	45.01	52.14	37.43	43.26
DTN[47]	44.13	57.42	55.89	45.76	44.53	57.34	52.43	51.55
c-DTN	45.10	49.78	47.43	45.79	49.80	55.69	50.10	52.31
PixelDA[10]	46.45	44.33	44.87	46.83	45.63	16.37	45.43	47.00
UNIT[27]	43.88	43.99	44.47	44.47	44.47	43.95	44.64	44.47
DSN[9]	46.25	47.50	62.15	39.72	45.85	56.65	56.5	42.87
c-DSN	45.82	44.64	45.60	46.32	45.18	45.52	44.79	47.37
mp-DSN	42.19	44.70	42.47	40.50	45.00	43.80	45.79	42.39
s-MTDA-ITA	44.77	45.61	60.00	46.70	49.06	55.33	59.90	50.64
c-MTDA-ITA	44.35	42.67	58.90	44.32	46.74	54.11	56.89	49.64
MTDA-ITA	46.30	60.60	60.50	50.40	55.59	57.80	64.20	56.34

Table 3: Mean classification accuracy on Multi-PIE classification. The best (red).

MTDA-ITA achieves the best performances as well as the best scores in most settings that verifies the effectiveness of MTDA-ITA for multi-target domain adaptation. Clearly, with the increasing camera angle, the image structure changes up to a certain extent (the views become heterogeneous). However, our method produces better results even under such very challenging conditions.

5.3 PACS dataset

This dataset contains 9991 images ($227 \times 227 \times 3$ dimension) across 7 categories (‘dog’, ‘elephant’, ‘giraffe’, ‘guitar’, ‘house’, ‘horse’ and ‘person’) and 4 domains of different stylistic depictions (‘Photo’, ‘Art painting’, ‘Cartoon’ and ‘Sketch’). The very diverse depiction styles provide a significant gap between domains, coupled with the small number of data samples, making it extremely challenging for domain adaptation. Consequently, the dataset was originally used for multi-source to single target domain adaptation [25]. Instead, we tackle a significantly more challenging problem of single-source to multiple target adaptation. Tab. 4 shows the classification accuracy of various methods. MTDA-ITA consistently achieves the best performance for all transfer tasks. Evaluations were obtained by training all models (ADDA, DSN, and ours) from scratch on the PACS dataset. Note that the overall performance figures are low due to the extreme difficulty of the transfer task, induced by large differences among domains.

method	P \rightarrow A	P \rightarrow C	P \rightarrow S	A \rightarrow P	A \rightarrow C	A \rightarrow S
1-NN	15.28	18.16	25.60	22.70	19.75	22.70
ADDA [42]	24.35	20.12	22.45	32.57	17.68	18.90
DSN [9]	28.42	21.14	25.64	29.54	25.89	24.69
s-MTDA-ITA	28.02	21.64	26.24	31.06	25.09	25.89
c-MTDA-ITA	25.35	20.24	23.64	26.54	20.30	22.38
MTDA-ITA	31.40	23.05	28.24	35.74	27.00	28.90

Table 4: Mean classification accuracy on PACS dataset classification. A:Art-painting, C:Cartoon, S:Sketch, P:Photo. The best (red).

5.4 Ablation Studies

We performed an ablation study on the proposed model measuring impact of various terms on the model’s performance. To this end, we conducted additional experiments for the digit datasets with different components ablation, i.e., training without the reconstruction loss (denoted as **MTDA-woR**) by setting $\lambda_r = 0$, training without the classifier entropy loss (denoted as **MTDA-woE**) by setting $\lambda_c = 0$, training without the multi-domain separation loss (denoted as **MTDA-woD**) by setting $\lambda_d = 0$.

As can be seen from Fig. 4, disabling each of the above components leads to degraded performance. More precisely, the average drop by disabling the classifier entropy loss is $\approx 3.5\%$. Similarly, by disabling the reconstruction loss and the multi-domain separation loss, we have $\approx 4.5\%$ and $\approx 22\%$ average drop in performance, respectively. Clearly, by disabling the multi-domain separation loss, the accuracy drops significantly due to the severe data distribution mismatch between different domains. The figure also demonstrates that leveraging the unlabeled data from multiple target domains during training enhances the generalization ability of the model that leads to higher performance. In addition, the performance drop caused by removing the reconstruction loss, i.e., without the private encoder/decoder, indicates (i) the benefit of modeling the latent features as the combination of shared and private features, (ii) the ability of the model’s domain adversarial loss to effectively learn those features.

In order to examine the effect of the private features on the model’s classification performance, we took the **MTDA-ITA** and trained it without the private encoder (denoted as **MTDA-woP**). As Fig. 4 shows, without the private features, the model performed consistently worse ($\approx 2\%$ average drop in performance) in all scenarios. This demonstrates explicitly modeling what is unique to each domain can improve the model’s ability to extract domain-invariant features.

In summary, this ablation study showed that the individual components bring complimentary information to achieve the best classification results.

5.5 Feature Visualization

We use t-SNE [31] on Digit dataset to visualize shared and private feature representations from different domains. Fig. 5 shows shared and private features from source (SVHN) and target domains before (a),(b) and after adaptation (c),(d). **MTDA-ITA** significantly reduces the domain mismatch for the shared features (circle markers in Fig. 5d, strong mixing of domain labels in this cluster, Fig. 5c) and increases it for the private features (triangle markers, pure and well-separated domain clusters in Fig. 5c). This is partially due to the proposed multi-domain separation loss through the use of the domain classifier D , which penalizes the domain mismatch for the shared features and rewards the mismatch for the private features. Moreover, as supported by the quantitative results in Tab. 1, joint adaptation of related domains and the classifier, accomplished through the model, leads to superior class separability, compared to original features. This is

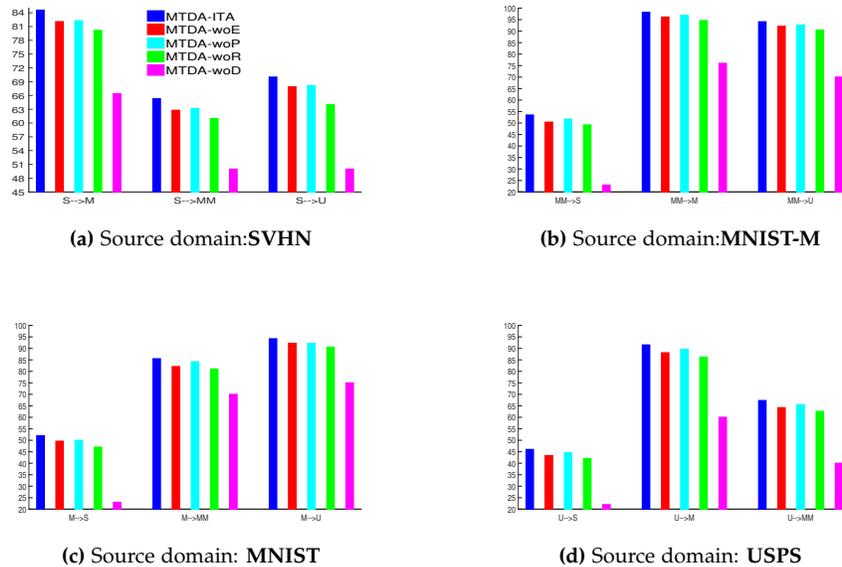


Figure 4: Ablation of MTDA-ITA on Digit dataset. We show that each component of our method, Reconstruction loss, Classifier entropy loss with separating shared/private features, contributes to the overall performance.

depicted in Fig. 5d, where the points in the shared space (large cluster) are grouped into class-specific subgroups (color indicates class label), while they are mixed in private spaces (smaller clusters). This is in contrast to Fig. 5b, where original features show no class-specificity.

We also show the learned shared and private features for the models MTDA-woE, MTDA-woP, MTDA-woR, and MTDA-woD, in Figs. 5e to 5l. Note that since the private encoder E_p is disabled for MTDA-woR, and MTDA-woP, no private features are depicted in Figs. 5g to 5j. The class label separation in the shared space for MTDA-woE, MTDA-woP, and MTDA-woR, Figs. 5f, 5h and 5j, is still evident but not as strong as in the full model, Fig. 5d, corroborating the small loss in classification accuracy observed in Fig. 4a. On the other hand, MTDA-woD has significant mixing of class labels in the shared space, Fig. 5l, more so than MTDA-woE, MTDA-woR, and MTDA-woP, implying worse classification prediction in Fig. 4a due to the severe mismatch between different domains.

6 Conclusion

This paper presented an information theoretic end-to-end approach to uDA in the context of a single source and multiple target domains that share a common task or properties. The proposed method learns feature representations invariant under multiple domain shifts and simultaneously discriminative for the learning task. This is accomplished by explicitly separating representations private to each domain and shared between source and target domains using a novel discrimination strategy. Our use of a single private domain encoder results in a highly scalable model, easily optimized using established back-propagation approaches. Results on three benchmark datasets for image classification show superiority of the proposed method compared to the state-of-the-art methods for unsupervised domain adaptation of visual domain categories.

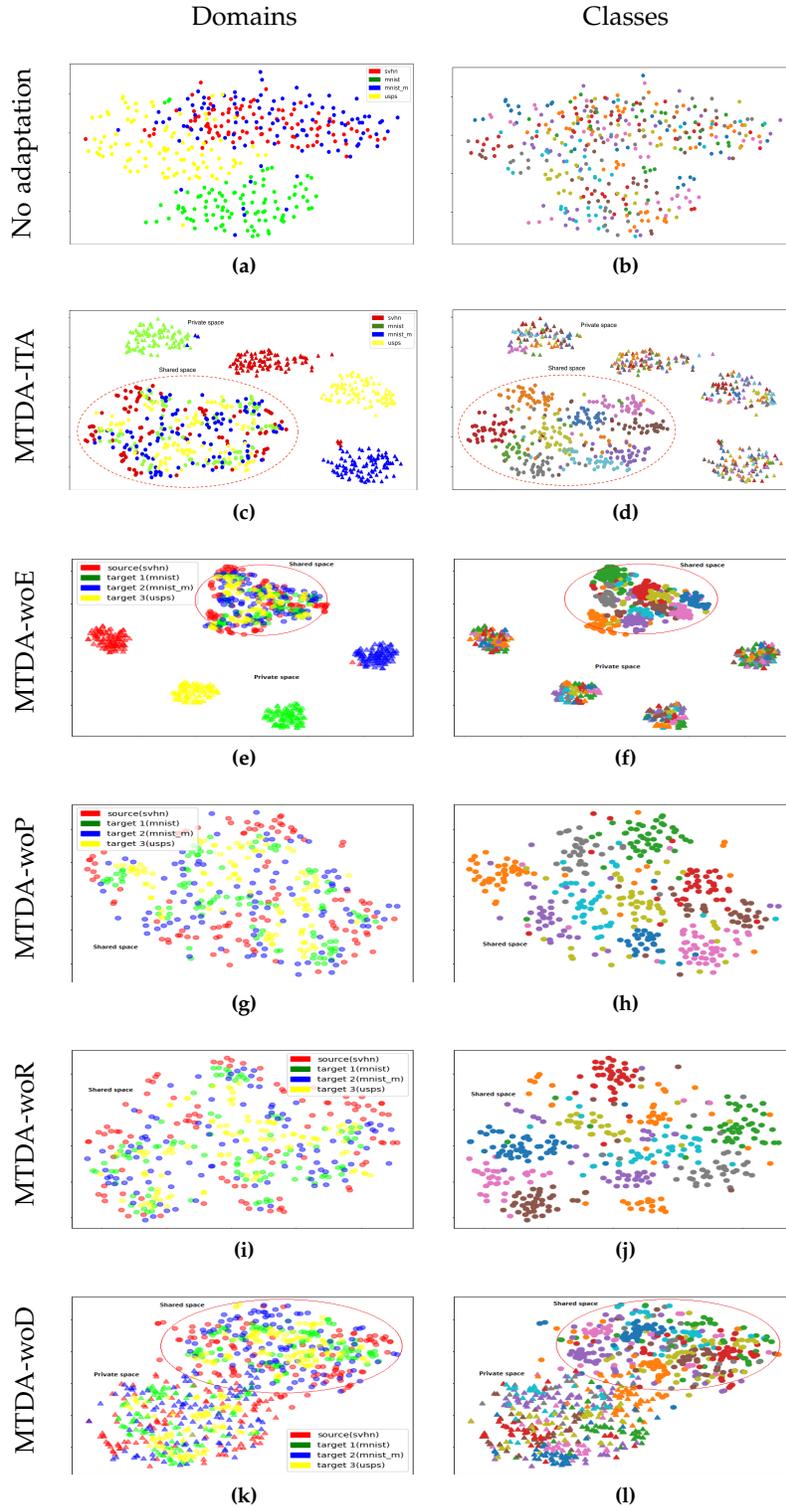


Figure 5: Feature visualization for embedding of digit datasets using t-SNE algorithm. The first and the second columns show the domains and classes, respectively, with color indicating domain and class membership. (a),(b) Original features. (c),(d) learned features for **MTDA-ITA** (triangle marker: private features, circle marker: shared features). Large clusters in the right column represent points from the shared space, while the smaller ones are from the private spaces. The remaining figures depict the learned features without: (e),(f) the classifier entropy loss, **MTDA-woE**; (g),(h) the private encoder, **MTDA-woP**; (i),(j) the reconstruction loss/decoder, **MTDA-woR**; and (k),(l) the multi-domain separation loss, **MTDA-woD**.

References

- [1] M. E. Abbasnejad, A. Dick, and A. van den Hengel. Infinite variational autoencoder for semi-supervised learning. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 781–790. IEEE, 2017.
- [2] A. Achille and S. Soatto. Information dropout: Learning optimal representations through noisy computation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018.
- [3] A. Alemi, B. Poole, I. Fischer, J. Dillon, R. A. Saurous, and K. Murphy. Fixing a broken elbo. In *International Conference on Machine Learning*, pages 159–168, 2018.
- [4] A. A. Alemi, I. Fischer, J. V. Dillon, and K. Murphy. Deep variational information bottleneck. *International Conference on Learning Representation (ICLR)*, 2016.
- [5] A. A. Alemi, I. Fischer, and J. V. Dillon. Uncertainty in the variational information bottleneck. *arXiv preprint arXiv:1807.00906*, 2018.
- [6] A. Anoosheh, E. Agustsson, R. Timofte, and L. Van Gool. Combogan: Unrestrained scalability for image domain translation. *ICLR Workshops*, 2017.
- [7] D. Barber and F. Agakov. The im algorithm: a variational approach to information maximization. In *Advances in Neural Information Processing Systems (NIPS)*, pages 201–208. MIT Press, 2003.
- [8] S. Benaim and L. Wolf. One-sided unsupervised domain mapping. In *Advances in Neural Information Processing Systems (NIPS)*, pages 752–762, 2017.
- [9] K. Bousmalis, G. Trigeorgis, N. Silberman, D. Krishnan, and D. Erhan. Domain separation networks. In *Advances in Neural Information Processing Systems (NIPS)*, pages 343–351, 2016.
- [10] K. Bousmalis, N. Silberman, D. Dohan, D. Erhan, and D. Krishnan. Unsupervised pixel-level domain adaptation with generative adversarial networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, page 7, 2017.
- [11] S. R. Bowman, L. Vilnis, O. Vinyals, A. M. Dai, R. Jozefowicz, and S. Bengio. Generating sentences from a continuous space. *arXiv preprint arXiv:1511.06349*, 2015.
- [12] M. Chalk, O. Marre, and G. Tkacik. Relevant sparse codes with variational information bottleneck. In *Advances in Neural Information Processing Systems (NIPS)*, pages 1957–1965, 2016.
- [13] Y. Choi, M. Choi, M. Kim, J.-W. Ha, S. Kim, and J. Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 1711, 2017.
- [14] N. Courty, R. Flamary, A. Habrard, and A. Rakotomamonjy. Joint distribution optimal transportation for domain adaptation. In *Advances in Neural Information Processing Systems (NIPS)*, pages 3733–3742, 2017.
- [15] G. Csurka. A comprehensive survey on domain adaptation for visual applications. pages 1–35. Springer, 2017.
- [16] B. Fernando, A. Habrard, M. Sebban, and T. Tuytelaars. Unsupervised visual domain adaptation using subspace alignment. In *IEEE International Conference on Computer Vision (ICCV)*, pages 2960–2967, 2013.

- [17] Y. Ganin and V. Lempitsky. Unsupervised domain adaptation by backpropagation. *International Conference on Machine Learning (ICML)*, 2015.
- [18] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky. Domain-adversarial training of neural networks. *The Journal of Machine Learning Research*, 17(1):2096–2030, 2016.
- [19] B. Gong, Y. Shi, F. Sha, and K. Grauman. Geodesic flow kernel for unsupervised domain adaptation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2066–2073, 2012.
- [20] G.-Y. Hao, H.-X. Yu, and W.-S. Zheng. Mixgan: Learning concepts from different domains for mixture generation. *arXiv preprint arXiv:1807.01659*, 2018.
- [21] H. Kameoka, T. Kaneko, K. Tanaka, and N. Hojo. Stargan-vc: Non-parallel many-to-many voice conversion with star generative adversarial networks. *arXiv preprint arXiv:1806.02169*, 2018.
- [22] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *International Conference on Learning Representation (ICLR)*, 2015.
- [23] E. Kodirov, T. Xiang, Z. Fu, and S. Gong. Unsupervised domain adaptation for zero-shot learning. In *IEEE International Conference on Computer Vision (ICCV)*, pages 2452–2460, 2015.
- [24] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [25] D. Li, Y. Yang, Y.-Z. Song, and T. M. Hospedales. Deeper, broader and artier domain generalization. In *IEEE International Conference on Computer Vision (ICCV)*, pages 5543–5551. IEEE, 2017.
- [26] M.-Y. Liu and O. Tuzel. Coupled generative adversarial networks. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems (NIPS)*, pages 469–477. Curran Associates, Inc., 2016.
- [27] M.-Y. Liu, T. Breuel, and J. Kautz. Unsupervised image-to-image translation networks. In *Advances in Neural Information Processing Systems (NIPS)*, pages 700–708, 2017.
- [28] M. Long, J. Wang, G. Ding, J. Sun, and P. S. Yu. Transfer joint matching for unsupervised domain adaptation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1410–1417, 2014.
- [29] M. Long, Y. Cao, J. Wang, and M. I. Jordan. Learning transferable features with deep adaptation networks. *International Conference on Machine Learning (ICML)*, 2015.
- [30] M. Long, J. Wang, Y. Cao, J. Sun, and S. Y. Philip. Deep learning of transferable representation for scalable domain adaptation. *IEEE Transactions on Knowledge and Data Engineering*, 28(8): 2027–2040, 2016.
- [31] L. v. d. Maaten and G. Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(Nov):2579–2605, 2008.
- [32] S. Mohamed and D. J. Rezende. Variational information maximisation for intrinsically motivated reinforcement learning. In *Advances in neural information processing systems (NIPS)*, pages 2125–2133, 2015.

- [33] S. Motiian, Q. Jones, S. Iranmanesh, and G. Doretto. Few-shot adversarial domain adaptation. In *Advances in Neural Information Processing Systems (NIPS)*, pages 6673–6683, 2017.
- [34] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng. Reading digits in natural images with unsupervised feature learning. In *NIPS workshop on deep learning and unsupervised feature learning*, volume 2011, page 5, 2011.
- [35] Y. Pu, W. Wang, R. Henao, L. Chen, Z. Gan, C. Li, and L. Carin. Adversarial symmetric variational autoencoder. In *Advances in Neural Information Processing Systems (NIPS)*, pages 4330–4339, 2017.
- [36] S.-A. Rebuffi, H. Bilen, and A. Vedaldi. Learning multiple visual domains with residual adapters. In *Advances in Neural Information Processing Systems (NIPS)*, pages 506–516, 2017.
- [37] K. Saito, Y. Ushiku, and T. Harada. Asymmetric tri-training for unsupervised domain adaptation. *International Conference on Machine Learning (ICML)*, 2017.
- [38] M. Salzman, C. H. Ek, R. Urtasun, and T. Darrell. Factorized orthogonal latent spaces. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 701–708, 2010.
- [39] B. Sun and K. Saenko. Deep coral: Correlation alignment for deep domain adaptation. In *European Conference on Computer Vision (ECCV)*, pages 443–450. Springer, 2016.
- [40] B. Sun, J. Feng, and K. Saenko. Return of frustratingly easy domain adaptation. *Thirtieth AAAI Conference on Artificial Intelligence*, 2016.
- [41] N. Tishby and N. Zaslavsky. Deep learning and the information bottleneck principle. In *IEEE Information Theory Workshop (ITW)*, pages 1–5. IEEE, 2015.
- [42] E. Tzeng, J. Hoffman, K. Saenko, and T. Darrell. Adversarial discriminative domain adaptation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, page 4, 2017.
- [43] H. Yan, Y. Ding, P. Li, Q. Wang, Y. Xu, and W. Zuo. Mind the class weight bias: Weighted maximum mean discrepancy for unsupervised domain adaptation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [44] D. Yoo, N. Kim, S. Park, A. S. Paek, and I. S. Kweon. Pixel-level domain transfer. In *European Conference on Computer Vision (ECCV)*, pages 517–532, 2016.
- [45] W. Zellinger, T. Grubinger, E. Lughofer, T. Natschläger, and S. Saminger-Platz. Central moment discrepancy (cmd) for domain-invariant representation learning. *International Conference on Learning Representation (ICLR)*, 2017.
- [46] J. Zhang, W. Li, and P. Ogunbona. Joint geometrical and statistical alignment for visual domain adaptation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [47] X. Zhang, F. X. Yu, S.-F. Chang, and S. Wang. Deep transfer network: Unsupervised domain adaptation. *arXiv preprint arXiv:1503.00591*, 2015.
- [48] H. Zhao, S. Zhang, G. Wu, J. P. Costeira, J. Moura, and G. J. Gordon. Multiple source domain adaptation with adversarial training of neural networks. *International Conference on Learning Representation (ICLR) workshops*, 05 2017.

Layer	Encoders (shared, private)
1	CONV-(N ₁₆ ,K ₇ ,S ₁), ReLU
2	CONV-(N ₃₂ ,K ₃ ,S ₂), ReLU
3	CONV-(N ₆₄ ,K ₃ ,S ₂), ReLU
4	RESBLK-(N ₆₄ ,K ₃ ,S ₁)
5	RESBLK-(N ₆₄ ,K ₃ ,S ₁)
6	RESBLK-(N ₆₄ ,K ₃ ,S ₁)
7	RESBLK-(N ₆₄ ,K ₃ ,S ₁)
Layer	Decoder
1	RESBLK-(N ₆₄ ,K ₃ ,S ₁)
2	RESBLK-(N ₆₄ ,K ₃ ,S ₁)
3	RESBLK-(N ₆₄ ,K ₃ ,S ₁)
4	RESBLK-(N ₆₄ ,K ₃ ,S ₁)
5	DCONV-(N ₃₂ ,K ₃ ,S ₂), ReLU
6	DCONV-(N ₁₆ ,K ₃ ,S ₂), ReLU
7	DCONV-(N ₁ ,K ₁ ,S ₁), TanH
Layer	Discriminator
1	CONV-(N ₄ ,K ₃ ,S ₁), ReLU
2	CONV-(N ₈ ,K ₃ ,S ₁), ReLU
3	CONV-(N ₁₆ ,K ₃ ,S ₁), ReLU
4	CONV-(N ₃₂ ,K ₃ ,S ₁), ReLU
5	CONV-(N ₁ ,K ₃ ,S ₁), ReLU
6	DENSE-(ND), Softmax
Layer	Classifier
1	CONV-(N ₄ ,K ₃ ,S ₁), ReLU
2	CONV-(N ₈ ,K ₃ ,S ₁), ReLU
3	CONV-(N ₁₆ ,K ₃ ,S ₁), ReLU
4	CONV-(N ₃₂ ,K ₃ ,S ₁), ReLU
5	CONV-(N ₁ ,K ₃ ,S ₁), ReLU
6	DENSE-(NC), Softmax

Table 5: Network architecture for the experiments.

A Network Architecture

The network architecture used for the experiments is given in Table 5. We use the following abbreviation for ease of presentation: N=Neurons, K=Kernel size, S=Stride size, D=Number of Domains, C=number of Classes. The transposed convolutional layer is denoted by DCONV. The residual basic block is denoted as RESBLK.