

Cross-domain error minimization for unsupervised domain adaptation

Yuntao Du, Yinghao Chen*, Fengli Cui *,
Xiaowen Zhang, and Chongjun Wang **

State Key Laboratory for Novel Software Technology at Nanjing University,
Nanjing University, Nanjing 210023, China,
duyuntao@smail.nju.edu.cn, yinghaochen48@gmail.com,
cuifengli1997@gmail.com, zhangxw@smail.nju.edu.cn,
chjwang@nju.edu.cn

Abstract. Unsupervised domain adaptation aims to transfer knowledge from a labeled source domain to an unlabeled target domain. Previous methods focus on learning domain-invariant features to decrease the discrepancy between the feature distributions as well as minimizing the source error and have made remarkable progress. However, a recently proposed theory reveals that such a strategy is not sufficient for a successful domain adaptation. It shows that besides a small source error, both the discrepancy between the feature distributions and the discrepancy between the labeling functions should be small across domains. The discrepancy between the labeling functions is essentially the **cross-domain errors** which are ignored by existing methods. To overcome this issue, in this paper, a novel method is proposed to integrate all the objectives into a unified optimization framework. Moreover, the incorrect pseudo labels widely used in previous methods can lead to error accumulation during learning. To alleviate this problem, the pseudo labels are obtained by utilizing structural information of the target domain besides source classifier and we propose a curriculum learning based strategy to select the target samples with more accurate pseudo-labels during training. Comprehensive experiments are conducted, and the results validate that our approach outperforms state-of-the-art methods.

Keywords: Transfer learning · Domain adaptation · Cross-domain errors.

1 Introduction

Traditional machine learning methods have achieved significant progress in various application scenarios [14, 33]. Training a model usually requires a large amount of labeled data. However, it is difficult to collect annotated data in some scenarios, such as medical image recognition [30] and automatic driving [42]. Such a case may lead to performance degradation for traditional machine learning methods. *Unsupervised domain adaptation* aims to overcome such challenge by transferring knowledge from a different but related domain (source domain) with labeled samples to a target domain with unlabeled samples [28]. And unsupervised domain adaptation based methods have achieved remarkable progress in many fields, such as image classification [45], automatic driving [42] and medical image precessing [30].

* equal contribution

** corresponding author

According to a classical theory of domain adaptation [1], the error of a hypothesis h in the target domain $\varepsilon_t(h)$ is bounded by three terms: the empirical error in the source domain $\hat{\varepsilon}_s(h)$, the distribution discrepancy across domains $d(\mathcal{D}_s, \mathcal{D}_t)$ and the ideal joint error λ^* :

$$\varepsilon_t(h) \leq \hat{\varepsilon}_s(h) + d(\mathcal{D}_s, \mathcal{D}_t) + \lambda^* \quad (1)$$

Note that $\mathcal{D}_s, \mathcal{D}_t$ denotes the source domain and the target domain, respectively. $\lambda^* = \varepsilon_s(h^*) + \varepsilon_t(h^*)$ is the ideal joint error and $h^* := \arg \min_{h \in \mathcal{H}} \varepsilon_s(h) + \varepsilon_t(h)$ is the ideal joint hypothesis. It is usually assumed that there is an ideal joint hypothesis h^* which can achieve good performance in both domains, making λ^* becoming a small and constant term. Therefore, besides minimizing the source empirical error, many methods focus on learning domain-invariant representations, i.e., intermediate features whose distributions are similar in the source and the target domain to achieve a small target error [6, 20, 29, 34, 36, 39, 44]. In shallow domain adaptation, *distribution alignment* is a widely used strategy for domain adaptation [21, 22, 27, 35, 38]. These methods assume that there exists a common space where the distributions of two domains are similar and they concentrate on finding a feature transformation matrix that projects the features of two domains into a common subspace with less distribution discrepancy.

Although having achieved remarkable progress, recent researches show that transforming the feature representations to be domain-invariant may inevitably distort the original feature distributions and enlarge the error of the ideal joint hypothesis [5, 18]. It reminds us that the error of the ideal joint error λ^* can not be ignored. However, it is usually intractable to compute the ideal joint error λ^* , because there are no labeled data in the target domain. Recently, a general and interpretable generalization upper bound without the pessimistic term λ^* for domain adaptation has been proposed in [47]:

$$\varepsilon_t(h) \leq \hat{\varepsilon}_s(h) + d(\mathcal{D}_s, \mathcal{D}_t) + \min\{E_{\mathcal{D}_s}[|f_s - f_t|], E_{\mathcal{D}_t}[|f_s - f_t|]\} \quad (2)$$

where f_s and f_t are the labeling functions (i.e., the classifiers to be learned) in both domains. The first two terms in Equ (2) are similar compared with Equ (1), while the third term is different. The third term measures the discrepancy between the labeling functions from the source and the target domain. Obviously, $E_{\mathcal{D}_s}[|f_s - f_t|] = \varepsilon_s(f_t)$ and $E_{\mathcal{D}_t}[|f_s - f_t|] = \varepsilon_t(f_s)$. As a result, the discrepancy between the labeling functions is essentially the **cross-domain errors**. Specifically, the cross-domain errors are the classification error of the source classifier in the target domain and the classification error of the target classifier in the source domain. Altogether, the newly proposed theory provides a sufficient condition for the success of domain adaptation: besides a small source error, not only the discrepancy between the feature distributions but also the cross-domain errors need to be small across domains, while the cross-domain errors are ignored by existing methods.

Besides, estimating the classifier errors is important for domain adaptation. Various classifiers such as k -NN, linear classifier and SVMs have been used in shallow domain adaptation [3, 21, 22, 27]. Recently, some methods adopt the *prototype classifier* [12] for classification in domain adaptation. The prototype classifier is a non-parametric classifier, where one class can be represented by one or more prototypes. And a sample can be classified according to the distances between the sample and the class prototypes.

In this paper, we propose a general framework named *Cross-Domain Error Minimization* (CDEM) based on the prototype classifier. CDEM aims to simultaneously learn domain-invariant features and minimize the cross-domain errors, besides minimizing the source

classification error. To minimize the cross-domain errors, we maintain a classifier for each domain separately, instead of assuming that there is an ideal joint classifier that can perform well in both domains. Moreover, we conduct discriminative feature learning for better classification. To sum up, as shown in Fig 1, there are four objectives in the proposed method. (i) Minimizing the classification errors in both domains to optimize the empirical errors. (ii) Performing distribution alignment to decrease the discrepancy between feature distributions. (iii) Minimizing the cross-domain errors to decrease the discrepancy between the labeling functions across domains. (iv) Performing discriminative learning to learn discriminative features. Note that the objectives (i), (ii) and (iv) have been explored in previous methods [27, 37, 38], while the objective (iii) is ignored by existing methods. We integrate the four objectives into a unified optimization problem to learn a feature transformation matrix via a closed-form solution. After transformation, the discrepancy between the feature distributions and the cross-domain errors will be small, and the source classifier can generalize well in the target domain.

Since the labels are unavailable in the target domain, we use *pseudo labels* instead in the learning process. Inevitably, there are some incorrect pseudo labels, which will cause error accumulation during learning [39]. To alleviate this problem, the pseudo labels of the target samples are obtained based on the structural information in the target domain and the source classifier, in this way, the pseudo labels are likely to be more accurate. Moreover, we propose to use *curriculum learning* [2] based strategy to select target samples with high prediction confidence during training. We regard the samples with high prediction confidence as "easy" samples and the samples with low prediction confidence as "hard" samples. The strategy is to learn the transformation matrix with "easy" samples at the early stage and with "hard" samples at the later stage. With the iterations going on, we gradually add more and more target samples to the training process.

Note that CDEM is composed of two processes: learning transformation matrix and selecting target samples. We perform these two processes in an alternative manner for better adaptation. Comprehensive experiments are conducted on three real-world object datasets. The results show that CDEM outperforms the state-of-the-art adaptation methods on most of the tasks (16 out of 24), which validates the substantial effects of simultaneously learning domain-invariant features and minimizing cross-domain errors for domain adaptation.

2 Related Work

Domain adaptation theory. The theory in [1] is one of the pioneering theoretical works in this field. A new statistics named $\mathcal{H}\Delta\mathcal{H}$ -divergence is proposed as a substitution of traditional distribution discrepancies (e.g. L_1 distance, KL-divergence) and a generalization error bound is presented. The theory shows that the target error is bounded by the source error and the distribution discrepancy across domains, so most domain adaptation methods aim to minimize the source error and reduce the distribution discrepancy across domains. A general class of loss functions satisfying symmetry and subadditivity are considered in [25] and a new generalization theory with respect to the newly proposed discrepancy distance is developed. A margin-aware generalization bound based on asymmetric margin loss is proposed in [25] and reveals the trade-off between generalization error and the choice of margin. Recently, a theory considering labeling functions is proposed in [46], which shows that the error of the target domain is bounded by three terms: the source error, the discrepancy in feature distributions and the discrepancy between the labeling functions across domains. The discrepancy between the labeling functions are essentially the cross-domain

errors which are ingored by existing methods. CDEM is able to optimize all the objectives simultaneously.

Domain adaptation algorithm. The mostly used shallow domain adaptation approaches include instance reweighting [3, 7, 16] and distribution alignment [22, 27, 34, 37, 38].

The instance reweighting methods assume that a certain portion of the data in the source domain can be reused for learning in the target domain and the samples in the source domain can be reweighted according to the relevance with the target domain. Tradaboost [7] is the most representative method which is inspired by Adaboost [41]. The source samples classified correctly by the target classifier have larger weight while the samples classified wrongly have less weight. LDML [16] also evaluates each sample and makes full use of the pivotal samples to filter out outliers. DMM [3] learns a transfer support vector machine via extracting domain-invariant feature representations and estimating unbiased instance weights to jointly minimize the distribution discrepancy. In fact, the strategy for selecting target samples based on *curriculum learning* can be regarded as a special case of instance reweighting, where the weight of selected samples is 1, while the weight of unselected samples is 0.

The distribution alignment methods assume that there exists a common space where the distributions of two domains are similar and focus on finding a feature transformation that projects features of two domains into another latent shared subspace with less distribution discrepancy. TCA [27] tries to align marginal distribution across domains, which learns a domain-invariant representation during feature mapping. Based on TCA, JDA [22] tries to align marginal distribution and conditional distribution simultaneously. Considering the balance between the marginal distribution and conditional distribution discrepancy, both BDA [37] and MEDA [38] adopt a balance factor to leverage the importance of different distributions. However, these methods all focus on learning domain-invariant features across domains and ignore the cross-domain errors. While our proposed method takes the cross-domain errors into consideration.

3 Motivation

3.1 Problem Definition

In this paper, we focus on unsupervised domain adaptation. There are a source domain $\mathcal{D}_s = \{(x_s^i, y_s^i)\}_{i=1}^{n_s}$ with n_s labeled source examples and a target domain $\mathcal{D}_t = \{x_t^j\}_{j=1}^{n_t}$ with n_t unlabeled target examples. It is assumed that the feature space and the label space are the same across domains, i.e., $\mathcal{X}_s = \mathcal{X}_t \in \mathbb{R}^d$, $\mathcal{Y}_s = \mathcal{Y}_t = \{1, 2, \dots, C\}$, while the source examples and target examples are drawn from different joint distributions $P(\mathcal{X}_s, \mathcal{Y}_s)$ and $Q(\mathcal{X}_t, \mathcal{Y}_t)$, respectively. The goal of CDEM is to learn a feature transformation matrix $P \in R^{d \times k}$, which projects the features of both domains into a common space to reduce the shift in the joint distribution across domains, such that the target error $\varepsilon_t(h) = E_{(x,y) \sim Q}[h(x) \neq y]$ can be minimized, where h is the classifier to be learned.

3.2 Main Idea

As shown in Fig 1 (a), there is a large discrepancy across domains before adaptation. Previous methods only focus on minimizing the source error and performing distribution alignment to reduce the domain discrepancy (Fig 1 (b-c)). As the new theory revealed [47], in addition to minimizing the source error and learning domain-invariant features, it is also important to minimize the cross-domain errors. As shown in Fig 1(d), although performing distribution alignment can reduce the domain discrepancy, the samples near the

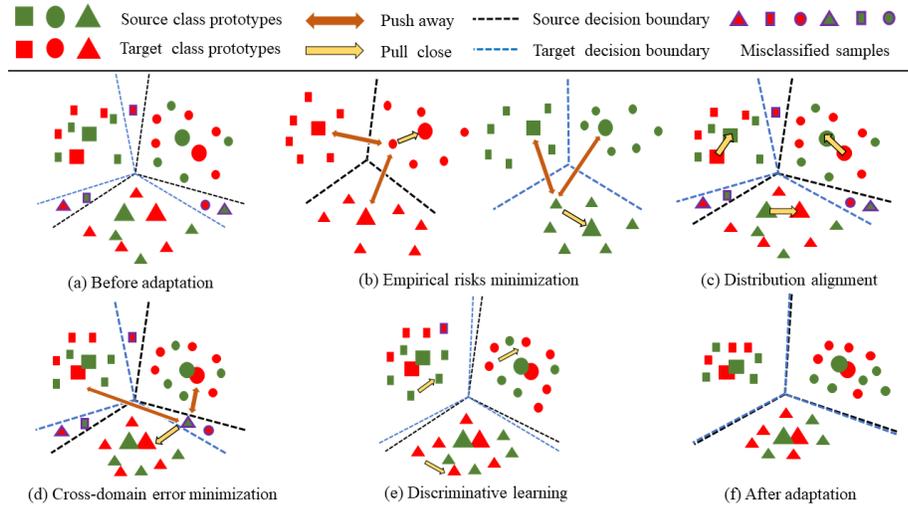


Fig. 1: An overview of the proposed method. In this paper, we use the prototype classifier as the basic classifier. There is one prototype in each class, and we choose the class center as the prototype. We use the distances between samples and prototypes to calculate the classification error. (a) Before adaptation, the classifier trained in the source domain can not generalize well in the target domain. (b-e) We aim to minimize empirical errors in both domains, perform distribution alignment to learn domain-invariant features, minimize the cross-domain errors to pull the decision boundaries across domains close, and perform discriminative learning to learn discriminative features. (f) After adaptation, the discrepancy across domains is reduced, so that the target samples can be classified correctly by the source classifier. Best viewed in color.

decision boundary are easy to be misclassified. Because performing distribution alignment only considers the discrepancy between the feature distributions, while the cross-domain errors are ignored. In the proposed method, minimizing the cross-domain errors can pull the decision boundaries across domains close, so that we can obtain a further reduced domain discrepancy. Moreover, we also perform discriminative learning to learn discriminative features (Fig 1(e)). Eventually, the domain discrepancy can be reduced and the classifier in the source domain can generalize well in the target domain (Fig 1(f)).

To sum up, we propose a general framework named *cross-domain error minimization* (CDEM), which is composed of four objectives:

$$h = \arg \min_{h \in \mathcal{H}} \sum_{i=1}^{n_s+n_t} l(h(x_i), y_i) + l_d(\mathcal{D}_s, \mathcal{D}_t) + l_f(\mathcal{D}_s, \mathcal{D}_t) + l_m(\mathcal{D}_s, \mathcal{D}_t) \quad (3)$$

where $l(h(x_i), y_i)$ is the classification errors in both domains. $l_d(\mathcal{D}_s, \mathcal{D}_t)$ and $l_f(\mathcal{D}_s, \mathcal{D}_t)$ represent the discrepancy between the feature distributions and the discrepancy between the labeling functions across domains, respectively. $l_m(\mathcal{D}_s, \mathcal{D}_t)$ is the discriminative objective to learn discriminative features. Note that CDEM is a shallow domain adaptation method and use the prototype classifier as the classifier, where no extra parameters are learned except the transformation matrix P . The framework is general and can generalize to other methods such as deep models.

As the labels in the target domain are unavailable, the *pseudo labels* for the target data are used for training instead. However, they are always some incorrect pseudo labels and

may lead to catastrophic error accumulation during learning. To alleviate this problem, we use the *curriculum learning* based strategy to select the target samples with more accurate pseudo labels which are obtained by taking advantage of source classifier and structural information of the target domain. With the iterations going on, we gradually add more and more target samples to the training process.

3.3 Classification Error

In this paper, we choose the prototype classifier [12] as the classifiers in both domains since the prototype classifier is a non-parametric classifier and is widely used in many tasks. As shown in Fig 1, we maintain one prototype for each class and adopt prototype matching for classification. The class centers $\{\mu_c\}_{c=1}^C$ are used as the prototype of each class in this paper. And we denote the classifier in the source domain as f_s and the classifier in the target domain as f_t . Given a training set $\mathcal{D} = \{x_i, y_i\}_{i=1}^{|\mathcal{D}|}$ with $|\mathcal{D}|$ samples, a sample $x \in \mathcal{D}$, the class center (prototype) for each class is defined as $\mu_c = \frac{1}{n_c} \sum_{x_i \in \mathcal{D}_c} x_i$, where $\mathcal{D}_c = \{x_i : x_i \in \mathcal{D}, y(x_i) = c\}$ and $n_c = |\mathcal{D}_c|$. We can derive the conditional probability of a given sample x belonging to class y as:

$$p(y|x) = \frac{\exp(-\|x - \mu_y\|)}{\sum_{c=1}^C \exp(-\|x - \mu_c\|)} \quad (4)$$

Assume the sample x belongs to class c , it is expected that the conditional probability $p(y|x)$ is close to $[0, 0, \dots, 1, \dots, 0]$, which is a C -dimensional one hot vector with the c -th dimension to be 1. Our goal is to pull the sample close to the center of c -th class while push the sample away from other $C - 1$ class centers. Note that instead of pushing samples directly away from $C - 1$ centers, we view the data of other $C - 1$ classes as a whole, and use the center of the $C - 1$ classes $\widehat{\mu}_c$ to calculate the distance. As a result, the algorithm complexity can be reduced and the proposed algorithm can be accelerated. The objective of minimizing classification error can be represented as,

$$\min \sum_{(x,c) \sim \mathcal{D}} \|x - \mu_c\|_2^2 - \beta \|x - \widehat{\mu}_c\|_2^2 \quad (5)$$

where $\widehat{\mu}_c = \frac{1}{n_c^*} \sum_{x_i \in \mathcal{D}/\mathcal{D}_c} x_i$ and $\widehat{\mu}_c$ is the center of all classes except class c in the training set, $n_c^* = |\mathcal{D}/\mathcal{D}_c|$, β is the regularization parameter.

4 Method

In this section, we will describe all the objectives and the method to select target samples separately.

4.1 Empirical Error Minimization

For classifying the samples correctly, the first objective of CDEM is to minimize the empirical errors in both domains. Since there are no labeled data in the target domain, we use the pseudo labels [22] instead. The empirical errors in both domains are represented as,

$$\begin{aligned} \sum_{i=1}^{n_s+n_t} l(h(x_i), y_i) &= \varepsilon_s(f_s) + \varepsilon_t(f_t) \\ &= \sum_{c=1}^C \sum_{x_i \in \mathcal{D}_{s,c}} \left(\|P^T(x_i - \mu_{s,c})\|_2^2 - \beta \|P^T(x_i - \widehat{\mu}_{s,c})\|_2^2 \right) \\ &\quad + \sum_{c=1}^C \sum_{x_j \in \mathcal{D}_{t,c}} \left(\|P^T(x_j - \mu_{t,c})\|_2^2 - \beta \|P^T(x_j - \widehat{\mu}_{t,c})\|_2^2 \right) \end{aligned} \quad (6)$$

where $\mathcal{D}_{s,c} = \{x_i : x_i \in \mathcal{D}_s, y(x_i) = c\}$ is the set of examples belonging to class c in the source domain and $y(x_i)$ is the true label of x_i . Correspondingly, $\mathcal{D}_{t,c} = \{x_j : x_j \in \mathcal{D}_t, \hat{y}(x_j) = c\}$ is the set of examples belonging to class c in the target domain, where $\hat{y}(x_j)$ is the pseudo label of x_j . $\mu_{s,c} = \frac{1}{n_{s,c}} \sum_{x_i \in \mathcal{D}_{s,c}} x_i$ and $\mu_{t,c} = \frac{1}{n_{t,c}} \sum_{x_j \in \mathcal{D}_{t,c}} x_j$ are the centers of c -th class in the source domain and the target domain respectively, where $n_{s,c} = |\mathcal{D}_{s,c}|$ and $n_{t,c} = |\mathcal{D}_{t,c}|$. Similarly, $\widehat{\mu}_{s,c} = \frac{1}{n_{s,c}^*} \sum_{x_i \in \mathcal{D}_s / \mathcal{D}_{s,c}} x_i$ and $\widehat{\mu}_{t,c} = \frac{1}{n_{t,c}^*} \sum_{x_j \in \mathcal{D}_t / \mathcal{D}_{t,c}} x_j$ are the centers of all classes except class c in the source domain and the target domain respectively, where $n_{s,c}^* = |\mathcal{D}_s / \mathcal{D}_{s,c}|$, $n_{t,c}^* = |\mathcal{D}_t / \mathcal{D}_{t,c}|$.

We further rewrite the first term of the objective function in Equ (6) as follows,

$$\begin{aligned} & \sum_{c=1}^C \sum_{x_i \in \mathcal{D}_{s,c}} \left(\|P^T(x_i - \mu_{s,c})\|_2^2 - \beta \|P^T(x_i - \widehat{\mu}_{s,c})\|_2^2 \right) \\ &= \sum_{c=1}^C \left((1 - \beta) \sum_{x_i \in \mathcal{D}_{s,c}} \|P^T(x_i - \mu_{s,c})\|_2^2 - \beta n_{s,c} \|P^T(\mu_{s,c} - \widehat{\mu}_{s,c})\|_2^2 \right) \quad (7) \\ &= (1 - \beta) \sum_{c=1}^C \sum_{x_i \in \mathcal{D}_{s,c}} \|P^T(x_i - \mu_{s,c})\|_2^2 - \beta \sum_{c=1}^C n_{s,c} \|P^T(\mu_{s,c} - \widehat{\mu}_{s,c})\|_2^2 \end{aligned}$$

Inspired by Linear Discriminant Analysis (LDA) [26] and follow previous method [17], we further transform the two terms, which can be considered as intra-class variance in Equ (7), into similar expressions as Equ (8).

$$\begin{aligned} & (1 - \beta) \sum_{c=1}^C \sum_{x_i \in \mathcal{D}_{s,c}} \|P^T(x_i - \mu_{s,c})\|_2^2 - \beta \sum_{c=1}^C n_{s,c} \|P^T(\mu_{s,c} - \widehat{\mu}_{s,c})\|_2^2 \\ &= \text{tr}(P^T X_s (I - Y_s (Y_s^T Y_s)^{-1} Y_s^T) X_s^T P) - \beta \sum_{c=1}^C n_{s,c} \text{tr}(P^T X_s \widehat{Q}_{s,c} X_s^T P) \quad (8) \end{aligned}$$

Where $X_s \in \mathbb{R}^{d \times n_s}$ and $Y_s \in \mathbb{R}^{n_s \times C}$ are the samples and labels in the source domain. $\text{tr}(\cdot)$ is the trace of a matrix. By using target samples $X_t \in \mathbb{R}^{d \times n_t}$ and pseudo labels $\hat{Y}_t \in \mathbb{R}^{n_t \times C}$, the same strategy is also used to transform the second term in Equ (6). Denote $X = X_s \cup X_t \in \mathbb{R}^{d \times (n_s + n_t)}$, the objective of minimizing empirical errors can be written as,

$$\varepsilon_s(f_s) + \varepsilon_t(f_t) = (1 - \beta) (\text{tr}(P^T X Q^Y X^T P) - \beta \sum_{c=1}^C \text{tr}(P^T X \widehat{Q}^c X^T P)) \quad (9)$$

where

$$Q^Y = \begin{bmatrix} I - Y_s (Y_s^T Y_s)^{-1} Y_s^T & \mathbf{0} \\ \mathbf{0} & I - \hat{Y}_t (\hat{Y}_t^T \hat{Y}_t)^{-1} \hat{Y}_t^T \end{bmatrix}, \widehat{Q}^c = \begin{bmatrix} n_{s,c} \widehat{Q}_{s,c} & \mathbf{0} \\ \mathbf{0} & n_{t,c} \widehat{Q}_{t,c} \end{bmatrix} \quad (10)$$

$$(\widehat{Q}_{s,c})_{ij} = \begin{cases} \frac{1}{n_{s,c} n_{s,c}}, & x_i, x_j \in \mathcal{D}_{s,c} \\ \frac{1}{n_{s,c}^* n_{s,c}^*}, & x_i, x_j \in \mathcal{D}_s / \mathcal{D}_{s,c} \\ -\frac{1}{n_{s,c} n_{s,c}^*}, & \text{otherwise} \end{cases}, (\widehat{Q}_{t,c})_{ij} = \begin{cases} \frac{1}{n_{t,c} n_{t,c}}, & x_i, x_j \in \mathcal{D}_{t,c} \\ \frac{1}{n_{t,c}^* n_{t,c}^*}, & x_i, x_j \in \mathcal{D}_t / \mathcal{D}_{t,c} \\ -\frac{1}{n_{t,c} n_{t,c}^*}, & \text{otherwise} \end{cases} \quad (11)$$

4.2 Distribution Alignment

As there are feature distribution discrepancy across domains, the second objective of CDEM is to learn domain-invariant features for decreasing the discrepancy between feature distributions across domains. Distribution alignment is a popular method in domain adaptation [22, 27, 38]. To reduce the shift between feature distributions across domains, we follow [19] and adopt *Maximum Mean Discrepancy* (MMD) as the distance measure to compute marginal distribution discrepancy $d_m(\mathcal{D}_s, \mathcal{D}_t)$ across domains based on the distance between the sample means of two domains in the feature embeddings:

$$d_m(\mathcal{D}_s, \mathcal{D}_t) = \left\| \frac{1}{n_s} \sum_{x_i \in \mathcal{D}_s} P^T x_i - \frac{1}{n_t} \sum_{x_j \in \mathcal{D}_t} P^T x_j \right\|^2 = \text{tr}(P^T X M_0 X^T P) \quad (12)$$

Based on the pseudo labels of the target data, we minimize the conditional distribution discrepancy $d_c(\mathcal{D}_s, \mathcal{D}_t)$ between domains:

$$d_c(\mathcal{D}_s, \mathcal{D}_t) = \sum_{c=1}^C \left\| \frac{1}{n_{s,c}} \sum_{x_i \in \mathcal{D}_{s,c}} P^T x_i - \frac{1}{n_{t,c}} \sum_{x_j \in \mathcal{D}_{t,c}} P^T x_j \right\|^2 = \sum_{c=1}^C \text{tr}(P^T X M_c X^T P) \quad (13)$$

where,

$$(M_0)_{ij} = \begin{cases} \frac{1}{n_s^2}, & x_i, x_j \in \mathcal{D}_s \\ \frac{1}{n_t^2}, & x_i, x_j \in \mathcal{D}_t \\ -\frac{1}{n_s n_t}, & \text{otherwise} \end{cases}, (M_c)_{ij} = \begin{cases} \frac{1}{n_{s,c}^2}, & x_i, x_j \in \mathcal{D}_{s,c} \\ \frac{1}{n_{t,c}^2}, & x_i, x_j \in \mathcal{D}_{t,c} \\ -\frac{1}{n_{s,c} n_{t,c}}, & \begin{cases} x_i \in \mathcal{D}_{s,c}, x_j \in \mathcal{D}_{t,c} \\ x_j \in \mathcal{D}_{s,c}, x_i \in \mathcal{D}_{t,c} \end{cases} \\ 0, & \text{otherwise} \end{cases} \quad (14)$$

Denote $M = M_0 + \sum_{c=1}^C M_c$, then the objective of distribution alignment is equal to:

$$l_d(\mathcal{D}_s, \mathcal{D}_t) = d_m(\mathcal{D}_s, \mathcal{D}_t) + d_c(\mathcal{D}_s, \mathcal{D}_t) = \text{tr}(P^T X M X^T P) \quad (15)$$

4.3 Cross-Domain Error Minimization

Although performing distribution alignment can pull the two domains close, it is not enough for a good adaptation across domains. The discrepancy between the labeling functions, which is essentially the cross-domain errors, is another factor leading to the domain discrepancy [47] while is ignored by existing methods. Thus, the third objective of CDEM is to minimize cross-domain errors, by which the decision boundaries across domains can be close and the samples near the decision boundaries can be classified correctly, achieving a further reduced domain discrepancy and better adaptation.

It is noticed that the cross-domain errors are the performance of the source classifier in the target domain and the performance of the target classifier in the source domain. As we use the prototype classifier, the cross-domain error in each domain is represented by the distances between the source samples (target samples) and the corresponding class centers in the target domain (source domain). For example, the cross-domain error in the source domain $\varepsilon_s(f_t)$ is the empirical error of applying the target classifier f_t to the source domain \mathcal{D}_s . Technically, the cross-domain errors in both domains are represented as,

$$\begin{aligned} l_f(\mathcal{D}_s, \mathcal{D}_t) &= \varepsilon_s(f_t) + \varepsilon_t(f_s) \\ &= \sum_{c=1}^C \sum_{x_i \in \mathcal{D}_{s,c}} \left(\|P^T(x_i - \mu_{t,c})\|_2^2 - \beta \|P^T(x_i - \widehat{\mu}_{t,c})\|_2^2 \right) \\ &\quad + \sum_{c=1}^C \sum_{x_j \in \mathcal{D}_{t,c}} \left(\|P^T(x_j - \mu_{s,c})\|_2^2 - \beta \|P^T(x_j - \widehat{\mu}_{s,c})\|_2^2 \right) \end{aligned} \quad (16)$$

Similar to the first objective, we transform the formula in Equ (16) as the following,

$$\begin{aligned} \varepsilon_s(f_t) + \varepsilon_t(f_s) = & (1 - \beta) \text{tr}(P^T X Q^Y X^T P) + \sum_{c=1}^C n^c \text{tr}(P^T X M_c X^T P) \\ & - \beta \sum_{c=1}^C \text{tr}(P^T X (n_{s,c} \widehat{Q}_{s,t}^c + n_{t,c} \widehat{Q}_{t,s}^c) X^T P) \end{aligned} \quad (17)$$

where

$$(\widehat{Q}_{s,t}^c)_{ij} = \begin{cases} \frac{1}{n_{s,c} n_{s,c}}, & x_i, x_j \in \mathcal{D}_{s,c} \\ \frac{1}{n_{t,c}^* n_{t,c}^*}, & x_i, x_j \in \mathcal{D}_{t,c} \\ -\frac{1}{n_{s,c} n_{t,c}^*}, & \begin{cases} x_i \in \mathcal{D}_{s,c}, x_j \in \mathcal{D}_{t,c} \\ x_j \in \mathcal{D}_{s,c}, x_i \in \mathcal{D}_{t,c} \end{cases} \\ 0 & \text{otherwise} \end{cases} \quad (\widehat{Q}_{t,s}^c)_{ij} = \begin{cases} \frac{1}{n_{t,c} n_{t,c}}, & x_i, x_j \in \mathcal{D}_{t,c} \\ \frac{1}{n_{s,c}^* n_{s,c}^*}, & x_i, x_j \in \mathcal{D}_{s,c} \\ -\frac{1}{n_{t,c} n_{s,c}^*}, & \begin{cases} x_i \in \mathcal{D}_{s,c}, x_j \in \mathcal{D}_{t,c} \\ x_j \in \mathcal{D}_{s,c}, x_i \in \mathcal{D}_{t,c} \end{cases} \\ 0 & \text{otherwise} \end{cases} \quad (18)$$

4.4 Discriminative Feature Learning

Learning domain-invariant features to reduce the domain discrepancy may harm the discriminability of the features [43]. So the fourth objective of CDEM is to perform discriminative learning to enhance the discriminability of the features [4]. To be specific, we resort to explore the structural information of all the samples to make the samples belonging to the same class close, which is useful for classification. Thus, the discriminative objective is,

$$l_m(\mathcal{D}_s, \mathcal{D}_t) = \sum_{x_i, x_j \in X} \|P^T x_i - P^T x_j\|_2^2 W_{ij} \quad (19)$$

where $W \in \mathbb{R}^{(n_s+n_t) \times (n_s+n_t)}$ is the similarity matrix, which is defined as follows,

$$W_{ij} = \begin{cases} 1, & y_i(\hat{y}_i) = y_j(\hat{y}_j) \\ 0, & y_i(\hat{y}_i) \neq y_j(\hat{y}_j) \end{cases} \quad (20)$$

This objective can be transformed as follows,

$$\begin{aligned} \sum_{x_i, x_j \in X} \|P^T x_i - P^T x_j\|_2^2 W_{ij} &= \text{tr} \left(\sum_{x_i \in X} P^T x_i B_{ii} x_i^T P - \sum_{x_i, x_j \in X} P^T x_i W_{ij} x_j^T P \right) \\ &= \text{tr}(P^T X B X^T P - P^T X W X^T P) = \text{tr}(P^T X L X^T P) \end{aligned} \quad (21)$$

Where, $L = B - W \in \mathbb{R}^{(n_s+n_t) \times (n_s+n_t)}$ is the laplacian matrix, and $B \in \mathbb{R}^{(n_s+n_t) \times (n_s+n_t)}$ is a diagonal matrix with $(B)_{ii} = \sum_j (W)_{ij}$.

4.5 Optimization

Combining the four objectives together, we get the following optimization problem,

$$\begin{aligned} L(p) &= \text{tr}(P^T X Q^Y X^T P) - \beta \sum_{c=1}^C \text{tr}(P^T X \widehat{Q}^c X^T P) + \lambda \text{tr}(P^T X M X^T P) \\ &\quad - \gamma \sum_{c=1}^C \text{tr}(P^T X (\widehat{Q}_{s,t}^c + \widehat{Q}_{t,s}^c) X^T P) + \eta \text{tr}(P^T X L X^T P) + \delta \|P\|_F^2 \\ &= \text{tr}(P^T X \Omega X^T P) + \delta \|P\|_F^2 \\ \text{s.t.} \quad & P^T X H X^T P = I \end{aligned} \quad (22)$$

where $\Omega = Q^Y + \lambda M + \eta L - \sum_{c=1}^C (\beta \widehat{Q}^c + \gamma \widehat{Q}_{s,t}^c + \gamma \widehat{Q}_{t,s}^c)$ and $H = \mathbf{I} - \frac{1}{n_s+n_t} \mathbf{1}$ is the centering matrix. According to the constrained theory, we denote $\Theta = \text{diag}(\theta_1, \dots, \theta_k) \in \mathbb{R}^{k \times k}$ as the Langrange multiplier, and derive the Langrange function for problem (22) as,

$$L = \text{tr}(P^T X \Omega X^T P) + \delta \|P\|_F^2 + \text{tr}((I - P^T X H X^T P) \Theta) \quad (23)$$

Setting $\frac{\partial L}{\partial P} = \mathbf{0}$, we get generalized eigendecomposition,

$$(X \Omega X^T + \delta I) P = X H X^T P \Theta \quad (24)$$

Finally, finding the optimal feature transformation matrix P is reduced to solving Equ (24) for the k smallest eigenvectors.

4.6 Selective target samples

To avoid the catastrophic error accumulation caused by the incorrect pseudo labels, we predict the pseudo labels for the target samples via exploring the structural information of the target domain and source classifier. Moreover, based on curriculum learning, we propose a strategy to select a part of target samples, whose pseudo labels are more likely to be correct, to participate in the next iteration for learning the transformation matrix. One simple way to predict pseudo labels for target samples is to use the source class centers $\{\mu_c\}_{c=1}^C$ (the prototypes for each class) to classify the target samples. Therefore the conditional probability of a given target sample x_t belonging to class y is defined as:

$$p_s(y|x_t) = \frac{\exp(-P^T \|x_t - \mu_{s,y}\|)}{\sum_{c=1}^C \exp(-P^T \|x_t - \mu_{s,c}\|)} \quad (25)$$

Because there exists distribution discrepancy across domains, only using source prototypes is not enough for pseudo-labeling, which will lead to some incorrect pseudo labels. We further consider the structural information in the target domain, which can be exploited by unsupervised clustering. In this paper, K-Means clustering is used in the target domain. The cluster center $\mu_{t,c}$ is initialized with corresponding class center $\mu_{s,c}$ in the source domain, which ensures one-to-one mapping for each class. Thus, based on target clustering, the conditional probability of a given target sample x_t belonging to class y is defined by:

$$p_t(y|x_t) = \frac{\exp(-P^T \|x_t - \mu_{t,y}\|)}{\sum_{c=1}^C \exp(-P^T \|x_t - \mu_{t,c}\|)} \quad (26)$$

After getting $p_s(y|x_t)$ and $p_t(y|x_t)$, we can obtain two different kinds of pseudo labels \hat{y}_s^t and \hat{y}_t^t for target samples x_t :

$$\hat{y}_s^t = \arg \max_{y \in Y_t} p_s(y|x_t) \quad \hat{y}_t^t = \arg \max_{y \in Y_t} p_t(y|x_t) \quad (27)$$

Based on these two kinds of pseudo labels, a curriculum learning based strategy is proposed to select a part of target samples for training. We firstly select the target samples whose pseudo labels predicted by $p_s(y|x_t)$ and $p_t(y|x_t)$ are the same (i.e., $\hat{y}_s^t = \hat{y}_t^t$). And these samples are considered to satisfy the *label consistency* and are likely to be correct. Then, we progressively select a subset containing top tn_t/T samples with highest prediction probabilities from the samples satisfying the label consistency, where T is the number of total iterations and t is the number of current iteration. Finally, we combine $p_s(y|x_t)$ and $p_t(y|x_t)$ in an iterative weighting method. Formally, the final class conditional probability and the pseudo label for x_t are as follows:

$$p(y|x_t) = (1-t/T) \times p_s(y|x_t) + t/T \times p_t(y|x_t) \quad (28)$$

$$\hat{y}_t = \arg \max_{y \in Y_t} p(y|x_t)$$

To avoid the class imbalance problem when selecting samples, we take the class-wise selection into consideration to ensure that each class will have a certain proportion of samples to be selected, namely,

$$N_{t,c} = \min(n_{t,c} \times t/T, n_{t,c}^{con}) \quad (29)$$

where $N_{t,c}$ is the number of target samples being selected of class c , $n_{t,c}^{con}$ denotes the number of target samples satisfying the label consistency in the class c and t is the current epoch.

Remark: CDEM is composed of two processes: learning transformation matrix P and selecting target samples. We firstly learn the transformation matrix P via solving the optimization problem (24). Then, we select the target samples in the transformed feature space. We perform the two processes in an alternative manner as previous method [39].

5 Experiment

In this section, we evaluate the performance of CDEM by extensive experiments on three widely-used common datasets. The source code of CDEM is available at <https://github.com/yuntaodu/CDEM>.

5.1 Data Preparation

The **Office-Caltech** dataset [11] consists of images from 10 overlapping object classes between Office31 and Caltech-256 [13]. Specifically, we have four domains, **C** (*Caltech-256*), **A** (*Amazon*), **W** (*Webcam*), and **D** (*DSLR*). By randomly selecting two different domains as the source domain and target domain respectively, we construct $3 \times 4 = 12$ cross-domain object tasks, e.g. **C** \rightarrow **A**, **C** \rightarrow **W**, ..., **D** \rightarrow **W**.

The **Office-31** dataset [31] is a popular benchmark for visual domain adaptation. The dataset contains three real-world object domains, *Amazon* (**A**, images downloaded from online merchants), *Webcom* (**W**, low-resolution images by a web camera), and *DSLR* (**D**, high-resolution images by a digital camera). It has 4652 images of 31 classes. We evaluate all methods on six transfer tasks: **A** \rightarrow **W**, **A** \rightarrow **D**, **W** \rightarrow **A**, **W** \rightarrow **D**, **D** \rightarrow **A**, and **D** \rightarrow **W**.

ImageCLEF-DA¹ is a dataset organized by selecting 12 common classes shared by three public datasets, each is considered as a domain: *Caltech-256* (**C**), *ImageNet ILSVRC 2012* (**I**), and *Pascal VOC 2012* (**P**). We evaluate all methods on six transfer tasks: **I** \rightarrow **P**, **P** \rightarrow **I**, **I** \rightarrow **C**, **C** \rightarrow **I**, **C** \rightarrow **P**, and **P** \rightarrow **C**.

5.2 Baseline Methods

We compare the performance of CDEM with several state-of-the-art traditional and deep domain adaptation methods:

- Traditional domain adaptation methods: **INN** [8], **SVM** [10] and **PCA** [15], Transfer Component Analysis (**TCA**) [27], Joint Distribution Alignment (**JDA**) [22], CORrelation Alignment (**CORAL**) [34], Joint Geometrical and Statistical Alignment (**JGSA**) [43], Manifold Embedded Distribution Alignment (**MEDA**) [38], Confidence-Aware Pseudo Label Selection (**CAPLS**) [40] and Selective Pseudo-Labeling (**SPL**) [39].
- Deep domain adaptation methods: Deep Domain Confusion (**DDC**) [36], Deep Adaptation Network (**DAN**) [19], Deep CORAL (**DCORAL**) [35], Residual Transfer Network (**RTN**) [23], Multi Adversarial Domain Adaptation (**MADA**) [29], Conditional Domain Adversarial Network (**CDAN**) [20], Incremental CAN (**iCAN**) [44], Domain Symmetric Networks (**SymNets**) [45], Generate To Adapt (**GTA**) [32] and Joint Domain alignment and Discriminative feature learning (**JDDA**) [4].

¹ <http://imageclef.org/2014/adaptation>

Table 1: Classification Accuracy (%) on Office-Caltech dataset using Decaf6 features.

Method	C→A	C→W	C→D	A→C	A→W	A→D	W→C	W→A	W→D	D→C	D→A	D→W	Average
DDC [36]	91.9	85.4	88.8	85.0	86.1	89.0	78.0	84.9	100.0	81.1	89.5	98.2	88.2
DAN [19]	92.0	90.6	89.3	84.1	91.8	91.7	81.2	92.1	100.0	80.3	90.0	98.5	90.1
DCORAL [35]	92.4	91.1	91.4	84.7	-	-	79.3	-	-	82.8	-	-	-
INN [8]	87.3	72.5	79.6	71.7	68.1	74.5	55.3	62.6	98.1	42.1	50.0	91.5	71.1
SVM [10]	91.6	80.7	86.0	82.2	71.9	80.9	67.9	73.4	100.0	72.8	78.7	98.3	82.0
PCA [15]	88.1	83.4	84.1	79.3	70.9	82.2	70.3	73.5	<u>99.4</u>	71.7	79.2	98.0	81.7
TCA [27]	89.8	78.3	85.4	82.6	74.2	81.5	80.4	84.1	100.0	82.3	89.1	<u>99.7</u>	85.6
JDA [22]	89.6	85.1	89.8	83.6	78.3	80.3	84.8	90.3	100.0	85.5	91.7	<u>99.7</u>	88.2
CORAL [34]	92.0	80.0	84.7	83.2	74.6	84.1	75.5	81.2	100.0	76.8	85.5	99.3	84.7
JGSA [43]	91.4	86.8	93.6	84.9	81.0	88.5	85.0	90.7	100.0	86.2	92.0	<u>99.7</u>	90.0
MEDA [38]	<u>93.4</u>	<u>95.6</u>	91.1	<u>87.4</u>	88.1	88.1	93.2	99.4	<u>99.4</u>	87.5	<u>93.2</u>	97.6	92.8
CAPLS [40]	90.8	85.4	95.5	86.1	87.1	<u>94.9</u>	88.2	92.3	100.0	<u>88.8</u>	93.0	100.0	91.8
SPL [39]	92.7	93.2	98.7	<u>87.4</u>	<u>95.3</u>	89.2	87.0	92.0	100.0	88.6	92.9	98.6	<u>93.0</u>
CDEM (Ours)	93.5	97.0	<u>96.2</u>	88.7	98.0	95.5	<u>89.1</u>	<u>93.5</u>	100.0	90.1	93.4	<u>99.7</u>	94.6

5.3 Experimental Setup

To fairly compare our method with the state-of-the-art methods, we adopt the deep features commonly used in existing unsupervised domain adaption methods. Specifically, DeCaf6 [9] features (activations of the 6th fully connected layer of a convolutional neural network trained on ImageNet, $d = 4096$) are used for Office-Caltech dataset, ResNet50 [14] features ($d = 2048$) are used for Office-31 dataset and ImageCLEF-DA dataset. In this way, we can compare our proposed method with these deep models.

In our experiments, we adopt the PCA algorithm to decrease the dimension of the data before learning to accelerate the proposed method. We set the dimensionality of PCA space $m = 128$ for Office-Caltech dataset and $m = 256$ for Office-31 and ImageCLEF-DA datasets. For the dimensionality of the transformation matrix P , we set $k = 32, 128$ and 64 for Office-Caltech, Office-31 and ImageCLEF-DA respectively. The number of iterations for CDEM to converge is $T = 11$ for all datasets. For regularization parameter δ , we set $\delta = 1$ for Office-Caltech and ImageCLEF-DA datasets and $\delta = 0.1$ for Office-31 dataset. As for the other hyper-parameters, we set β, λ, γ and η by searching through the grid with a range of $\{0.0001, 0.001, 0.01, 0.1, 1, 10\}$. In addition, the coming experiment on parameter sensitivity shows that our method can keep robustness with a wide range of parameter values.

5.4 Results and Analysis

The results on the *Office-Caltech* dataset are reported in Table 1, where the highest accuracy of each cross-domain task is boldfaced. The results of baselines are directly reported from original papers if the protocol is the same. The CDEM method significantly outperforms all the baseline methods on most transfer tasks (7 out of 12) in this dataset. It is desirable that CDEM promotes the classification accuracies significantly on hard transfer tasks, e.g., $\mathbf{A} \rightarrow \mathbf{D}$ and $\mathbf{A} \rightarrow \mathbf{W}$, where the source and target domains are substantially different [31]. Note that CDEM performs better than SPL in most tasks, which only learns domain-invariant features across domains.

The results on *Office-31* dataset are reported in Table 2. The CDEM method outperforms the comparison methods on most transfer tasks. Compared with the best shallow baseline method (CAPLS), the accuracy is improved by 1.7%. Note that the CDEM method

Table 2: Accuracy (%) on Office-31 dataset using either ResNet50 features or ResNet50 based deep models.

Method	A→W	D→W	W→D	A→D	D→A	W→A	Avg
RTN [23]	84.5	96.8	99.4	77.5	66.2	64.8	81.6
MADA [29]	90.0	97.4	99.6	87.8	70.3	66.4	85.2
GTA [32]	89.5	97.9	<u>99.8</u>	87.7	72.8	71.4	86.5
iCAN [44]	92.5	98.8	100.0	90.1	72.1	69.9	87.2
CDAN-E [20]	<u>94.1</u>	98.6	100.0	92.9	71.0	69.3	87.7
JDDA [4]	82.6	95.2	99.7	79.8	57.4	66.7	80.2
SymNets [45]	90.8	98.8	100.0	<u>93.9</u>	74.6	72.5	<u>88.4</u>
TADA [18]	94.3	<u>98.7</u>	<u>99.8</u>	91.6	72.9	73.0	<u>88.4</u>
MEDA [38]	86.2	97.2	99.4	85.3	72.4	74.0	85.7
CAPLS [40]	90.6	98.6	99.6	88.6	<u>75.4</u>	<u>76.3</u>	88.2
CDEM (Ours)	91.1	98.4	99.2	94.0	77.1	79.4	89.9

Table 3: Accuracy (%) on ImageCLEF-DA dataset using either ResNet50 features or ResNet50 based deep models.

Method	I→P	P→I	I→C	C→I	C→P	P→C	Avg
RTN [23]	75.6	86.8	95.3	86.9	72.7	92.2	84.9
MADA [29]	75.0	87.9	96.0	88.8	75.2	92.2	85.8
iCAN [44]	79.5	89.7	94.7	89.9	78.5	92.0	87.4
CDAN-E [20]	77.7	90.7	97.7	91.3	74.2	94.3	87.7
SymNets [45]	<u>80.2</u>	93.6	97.0	93.4	78.7	<u>96.4</u>	89.9
MEDA [38]	79.7	92.5	95.7	92.2	78.5	95.5	89.0
SPL [39]	78.3	<u>94.5</u>	96.7	<u>95.7</u>	<u>80.5</u>	96.3	90.3
CDEM (ours)	80.5	96.0	<u>97.2</u>	96.3	82.1	96.8	91.5

outperforms some deep domain adaptation methods, which implies the performance of CDEM in domain adaptation is better than several deep methods.

The results on *ImageCLEF-DA* dataset are reported in Table 3. The CDEM method substantially outperforms the comparison methods on most transfer tasks, and with more rooms for improvement. An interpretation is that the three domains in *ImageCLEF-DA* are visually dissimilar with each other, and are difficult in each domain with much lower in-domain classification accuracy [22]. MEDA and SPL are the representative shallow domain adaptation methods, which both focus on learning domain-invariant features. Moreover, SPL also uses selective target samples for adaptation. Consequently, the better performance of CDEM implies that minimizing cross-domain errors can further reduce the discrepancy across domains and achieve better adaptation.

5.5 Effectiveness Analysis

Ablation Study We conduct an ablation study to analyse how different components of our method contribute to the final performance. When learning the final classifier, CDEM involves four components: the empirical error minimization (ERM), the distribution alignment (DA), the cross-domain error minimization (CDE) and discriminative feature learning (DFL). We empirically evaluate the importance of each component. To this end, we investigate different combinations of four components and report average classification accuracy on three datasets in Table 4. Note that the result of the first setting (only ERM used) is like

Table 4: Results of ablation study.

Method				Office-Caltech	Office31	ImageCLEF-DA
ERM	DA	CDE	DFL			
✓	✗	✗	✗	90.2	86.6	87.5
✓	✓	✗	✗	91.5	87.2	88.6
✓	✓	✓	✗	94.0	89.2	90.8
✓	✓	✓	✓	94.6	89.9	91.5

the result of the source-only method, where no adaptation is performed across domains. It can be observed that methods with distribution alignment or cross-domain error minimization outperform those without distribution alignment or cross-domain error minimization. Moreover, discriminative learning can further improve performance and CDE achieves the biggest improvement compared with other components. Summarily, using all the terms together achieves the best performance in all tasks.

Evaluation of Selective Target Samples We further perform experiments to show the effectiveness of selective target samples. We compare several variants of the proposed method: **a)** No selection: We use all the target samples for training without any samples removed. **b)** Only label consistency: We only select the samples where the predicted label by $p_s(x_t)$ is the same with $p_t(x_t)$. **c)** Only high probabilities: We only select the target samples with high prediction confidence. **d)** The proposed method. As shown in Fig 2(a), “No selection” leads to a model with the worst performance due to the catastrophic error accumulation. The “Only label consistency” and “Only high probabilities” achieve significantly better results than “No selection”, but are still worse than the proposed method, which verifies that our method of explicitly selecting easier samples can make the model more adaptive and less likely to be affected by the incorrect pseudo labels.

Feature Visualization. In Fig 2(b-d), we visualize the feature representations of task $A \rightarrow D$ (10 classes) by t-SNE [24] as previous methods [39] using JDA and CDEM. Before adaptation, we can see that there is a large discrepancy across domains. After adaptation, JDA learns domain-invariant features which can reduce distribution discrepancy, the source domain and the target domain can become closer. While CDEM further considers the cross-domain errors, achieving a better performance.

6 Conclusion

In this paper, we propose the *Cross-Domain Error Minimization* (CDEM), which not only learns domain-invariant features across domains but also performs cross-domain error minimization. These two goals complement each other and contribute to better domain adaptation. Apart from these two goals, we also integrate the empirical error minimization and discriminative learning into a unified learning process. Moreover, we propose a method to select the target samples to alleviate error accumulation problem caused by incorrect pseudo labels. Through a large number of experiments, it is proved that our method is superior to other strong baseline methods.

7 Acknowledgements

This paper is supported by the National Key Research and Development Program of China (Grant No. 2018YFB1403400), the National Natural Science Foundation of China (Grant No. 61876080), the Collaborative Innovation Center of Novel Software Technology and Industrialization at Nanjing University.

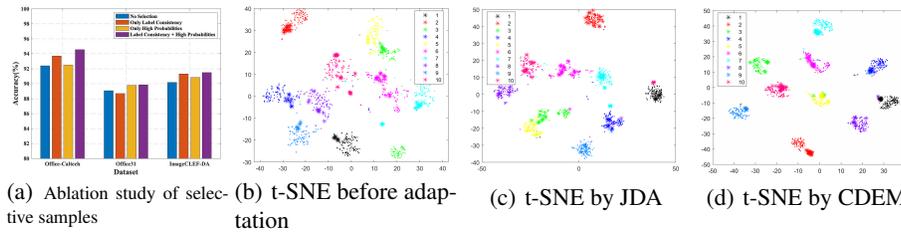


Fig. 2: Ablation study of selective samples, t-SNE visualization and parameter sensitivity

References

- Ben-David, S., Blitzer, J., Crammer, K., Kulesza, A., Pereira, F.C., Vaughan, J.W.: A theory of learning from different domains. *Machine Learning* **79**, 151–175 (2009)
- Bengio, Y., Louradour, J., Collobert, R., Weston, J.: Curriculum learning. In: *ICML '09* (2009)
- Cao, Y., Long, M., Wang, J.: Unsupervised domain adaptation with distribution matching machines. In: *AAAI* (2018)
- Chen, C., Chen, Z., Jiang, B., Jin, X.: Joint domain alignment and discriminative feature learning for unsupervised deep domain adaptation. In: *AAAI* (2019)
- Chen, C., Xie, W., Xu, T., Bing Huang, W., Rong, Y., Ding, X., Huang, Y., Huang, J.: Progressive feature alignment for unsupervised domain adaptation. *CVPR* pp. 627–636 (2018)
- Chen, Q., Du, Y., Tan, Z., Zhang, Y., Wang, C.: Unsupervised domain adaptation with joint domain-adversarial reconstruction networks. In: *ECML/PKDD* (2020)
- Dai, W., Yang, Q., Xue, G.R., Yu, Y.: Boosting for transfer learning. In: *ICML '07* (2007)
- Delany, S.J.: k-nearest neighbour classifiers (2007)
- Donahue, J., Jia, Y., Vinyals, O., Hoffman, J., Zhang, N., Tzeng, E., Darrell, T.: Decaf: A deep convolutional activation feature for generic visual recognition. In: *ICML* (2014)
- Evgeniou, T., Pontil, M.: Support vector machines: Theory and applications. In: *Machine Learning and Its Applications* (2001)
- Gong, B., Shi, Y., Sha, F., Grauman, K.: Geodesic flow kernel for unsupervised domain adaptation. *CVPR* pp. 2066–2073 (2012)
- Graf, A.B.A., Bousquet, O., Rätsch, G., Schölkopf, B.: Prototype classification: insights from machine learning. *Neural computation* **21** **1**, 272–300 (2009)
- Griffin, G., Holub, A., Perona, P.: Caltech-256 object category dataset (2007)
- He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. *CVPR* pp. 770–778 (2016)
- Hotelling, H.: Analysis of a complex of statistical variables into principal components. *Journal of Educational Psychology* **24**, 498–520 (1933)
- Jing, M., Li, J., Zhao, J., Lu, K.: Learning distribution-matched landmarks for unsupervised domain adaptation. In: *DASFAA* (2018)
- Liang, J., He, R., Sun, Z., Tan, T.: Aggregating randomized clustering-promoting invariant projections for domain adaptation. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **41**, 1027–1042 (2019)
- Liu, H., Long, M., Wang, J., Jordan, M.I.: Transferable adversarial training: A general approach to adapting deep classifiers. In: *ICML* (2019)
- Long, M., Cao, Y., Wang, J., Jordan, M.I.: Learning transferable features with deep adaptation networks. In: *ICML* (2015)
- Long, M., Cao, Z., Wang, J., Jordan, M.I.: Conditional adversarial domain adaptation. In: *NeurIPS* (2018)
- Long, M., Wang, J., Ding, G., Pan, S.J., Yu, P.S.: Adaptation regularization: A general framework for transfer learning. *IEEE TKDE* **26**, 1076–1089 (2014)

22. Long, M., Wang, J., Ding, G., Sun, J.G., Yu, P.S.: Transfer feature learning with joint distribution adaptation. *ICCV* pp. 2200–2207 (2013)
23. Long, M., Zhu, H., Wang, J., Jordan, M.I.: Unsupervised domain adaptation with residual transfer networks. In: *NIPS* (2016)
24. Maaten, L.V.D., Hinton, G.E.: Visualizing data using t-sne. *Journal of Machine Learning Research* **9**, 2579–2605 (2008)
25. Mansour, Y., Mohri, M., Rostamizadeh, A.: Domain adaptation: Learning bounds and algorithms. *NeurIPS* (2009)
26. Mika, S., Rätsch, G., Weston, J., Scholkopf, B., Mullers, K.R.: Fisher discriminant analysis with kernels. *Neural Networks for Signal Processing IX* pp. 41–48 (1999)
27. Pan, S.J., Tsang, I.W.H., Kwok, J.T., Yang, Q.: Domain adaptation via transfer component analysis. *IEEE Transactions on Neural Networks* **22**, 199–210 (2011)
28. Pan, S.J., Yang, Q.: A survey on transfer learning. *IEEE TKDE* **22**(10), 1345–1359 (2010)
29. Pei, Z., Cao, Z., Long, M., Wang, J.: Multi-adversarial domain adaptation. In: *AAAI* (2018)
30. Raghu, M., Zhang, C., Kleinberg, J.M., Bengio, S.: Transfusion: Understanding transfer learning for medical imaging. In: *NeurIPS* (2019)
31. Saenko, K., Kulis, B., Fritz, M., Darrell, T.: Adapting visual category models to new domains. In: *ECCV* (2010)
32. Sankaranarayanan, S., Balaji, Y., Castillo, C.D., Chellappa, R.: Generate to adapt: Aligning domains using generative adversarial networks. *CVPR* pp. 8503–8512 (2018)
33. Smith, N., Gales, M.J.F.: Speech recognition using svms. In: *NIPS* (2001)
34. Sun, B., Feng, J., Saenko, K.: Return of frustratingly easy domain adaptation. In: *AAAI* (2015)
35. Sun, B., Saenko, K.: Deep coral: Correlation alignment for deep domain adaptation. In: *ECCV Workshops* (2016)
36. Tzeng, E., Hoffman, J., Zhang, N., Saenko, K., Darrell, T.: Deep domain confusion: Maximizing for domain invariance. *ArXiv abs/1412.3474* (2014)
37. Wang, J., Chen, Y., Hao, S., Feng, W., Shen, Z.: Balanced distribution adaptation for transfer learning. *ICDM* pp. 1129–1134 (2017)
38. Wang, J., Feng, W., Chen, Y., Yu, H., Huang, M., Yu, P.S.: Visual domain adaptation with manifold embedded distribution alignment. In: *MM* 18 (2018)
39. Wang, Q., Breckon, T.P.: Unsupervised domain adaptation via structured prediction based selective pseudo-labeling. *AAAI* (2020)
40. Wang, Q., Bu, P., Breckon, T.P.: Unifying unsupervised domain adaptation and zero-shot visual recognition. *IJCNN* pp. 1–8 (2019)
41. Wyner, A.J., Olson, M., Bleich, J., Mease, D.: Explaining the success of adaboost and random forests as interpolating classifiers. *J. Mach. Learn. Res.* **18**, 48:1–48:33 (2017)
42. Yang, L., Liang, X., Wang, T., Xing, E.P.: Real-to-virtual domain unification for end-to-end autonomous driving. In: *ECCV* (2018)
43. Zhang, J., Li, W., Ogunbona, P.: Joint geometrical and statistical alignment for visual domain adaptation. *CVPR* pp. 5150–5158 (2017)
44. Zhang, W., Ouyang, W., Li, W., Xu, D.: Collaborative and adversarial network for unsupervised domain adaptation. *CVPR* pp. 3801–3809 (2018)
45. Zhang, Y., Tang, H., Jia, K., Tan, M.: Domain-symmetric networks for adversarial domain adaptation. *CVPR* pp. 5026–5035 (2019)
46. Zhang, Y., Liu, T., Long, M., Jordan, M.I.: Bridging theory and algorithm for domain adaptation. In: *ICML* (2019)
47. Zhao, H., des Combes, R.T., Zhang, K., Gordon, G.J.: On learning invariant representation for domain adaptation. *ICML* (2019)