

Discriminative Noise Robust Sparse Orthogonal Label Regression-based Domain Adaptation

Lingkun Luo, Liming Chen, *Senior Member, IEEE*, and Shiqiang Hu,

Abstract—Domain adaptation (DA) aims to enable a learning model trained from a source domain to generalize well on a target domain, despite the mismatch of data distributions between the two domains. State-of-the-art DA methods have so far focused on the search of a latent shared feature space where source and target domain data can be aligned either statistically and/or geometrically. In this paper, we propose a novel unsupervised DA method, namely *Discriminative Noise Robust Sparse Orthogonal Label Regression-based Domain Adaptation (DOLL-DA)*. The proposed DOLL-DA derives from a novel integrated model which searches a shared feature subspace where source and target domain data are, through optimization of some *repulse force* terms, discriminatively aligned statistically, while at same time regresses orthogonally data labels thereof using a label embedding trick. Furthermore, in minimizing a novel *Noise Robust Sparse Orthogonal Label Regression (NRS_OLR)* term, the proposed model explicitly accounts for data outliers to avoid negative transfer and introduces the property of sparsity when regressing data labels. We carry out comprehensive experiments in comparison with 32 state of the art DA methods using 8 standard DA benchmarks and 49 cross-domain image classification tasks. The proposed DA method demonstrates its effectiveness and consistently outperforms the state-of-the-art DA methods with a margin which reaches 17 points on the CMU PIE dataset. To gain insight into the proposed DOLL-DA, we also derive three additional DA methods based on three partial models from the full model, namely OLR, CDDA+, and JOLR-DA, highlighting the added value of 1) discriminative statistical data alignment; 2) Noise Robust Sparse Orthogonal Label Regression; and 3) their joint optimization through the full DA model. In addition, we also perform time complexity and an in-depth empiric analysis of the proposed DA method in terms of its sensitivity *w.r.t.* hyper-parameters, convergence speed, impact of the base classifier and random label initialization as well as performance stability *w.r.t.* target domain data being used in training.

Index Terms—Domain adaptation, Transfer Learning, Visual classification, Noise Robust Sparse Orthogonal Label Regression

I. INTRODUCTION

Traditional machine learning tasks assume that both training and testing data are drawn from a same data distribution [43], [45]. However, in many real-life applications, due to different factors as diverse as sensor difference, lighting changes, view-point variations, *etc.*, data from a target domain may have a different data distribution with respect to the labeled data in a source domain where a predictor can not be reliably learned. On the other hand, manually labeling enough target domain

data for the purpose of training an effective predictor can be very expensive, tedious and thus prohibitive.

Domain adaptation (DA) [43], [45], [37] aims to leverage possibly abundant labeled data from a *source* domain to learn an effective predictor for data in an unseen domain, namely *target* domain, despite the data distribution discrepancy between the source and target. While DA can be *semi-supervised* by assuming that a certain amount of labeled data is available in the target domain, in this paper we are interested in *unsupervised* DA where we assume that no labels are available for target domain data.

While there exists an increasing number of deep learning-based unsupervised DA methods, we focus in this paper on *shallow* DA methods as they are easier to train and can provide insights into the design decisions of deep DA methods. The relationships between *shallow* and *deep* DA methods will be discussed in depth in Sect. II on related works. State of the art *shallow* DA methods can be categorized into *instance*-based [43], [10], *feature*-based [44], [33], [63], or *classifier*-based. Classifier-based DA is widely applied in semi-supervised DA as it aims to fit a classifier trained on source domain data to target domain data through adaptation of its parameters, and thereby require some labels in the target domain [56]. The instance-based approach generally assumes that 1) the conditional distributions of source and target domain are identical [64], and 2) certain portion of the data in the source domain can be reused [43] for learning in the target domain through re-weighting. Feature-based adaptation [54], [21], [14], [57], [39] relaxes such a strict assumption and only requires that there exists a mapping from the input data space to a latent shared feature representation space. This latent shared feature space captures the information necessary for training classifiers for source and target tasks. In this paper, we propose a novel *hybrid* DA method using both feature adaptation and classifier optimization.

A common method to approach feature adaptation is to seek a shared latent subspace between the source and target domain [45], [44] via optimization. State-of-the-art features three main lines of approaches, namely, data geometric structure alignment-based (**DGSA**), data distribution centered (**DDC**) or their hybridization. **DGSA** based approaches [54], [51], [63] seek a subspace where source and target data can be well aligned and interlaced in preserving inherent hidden geometric data structure via low rank constraint and/or sparse representation. **DDC** methods [38], [35], [30], [36] aim to search a latent subspace where the discrepancy between the source and target data distributions is minimized, via various distances, *e.g.*, Bregman divergence [53], Geodesic distance [16], Wasserstein

K. Luo, S. Qiang are with School of Aeronautics and Astronautics, Shanghai Jiao Tong University, 800 Dongchuan Road, Shanghai, China e-mail: lolinkun@gmail.com, sqhu@sjtu.edu.cn.

L. Chen is with LIRIS, CNRS UMR 5205, Ecole Centrale de Lyon, 36 avenue Guy de Collongue, Ecully, France e-mail: (liming.chen)@ec-lyon.fr.

Manuscript received May 15, 2018.

distance[6], [7] or Maximum Mean Discrepancy[18] (MMD). The most popular distance is MMD due to its simplicity and solid theoretical foundation.

Our previously proposed DA method, namely **DGA-DA**[39], is a hybridization of **DDC** and **DGSA** approaches. **DGA-DA** leverages the advantages of both **DDC** and **DGSA** methods and demonstrates a state of the art performance on a number of DA benchmarks. **DGA-DA** relies upon the analysis of a cornerstone theoretical result in DA [2], [1], [26], which estimates an error bound of a learned hypothesis h on a target domain as follows:

$$e_{\mathcal{T}}(h) \leq e_{\mathcal{S}}(h) + d_{\mathcal{H}}(\mathcal{D}_{\mathcal{S}}, \mathcal{D}_{\mathcal{T}}) + \min \{ \mathcal{E}_{\mathcal{D}_{\mathcal{S}}} [\|f_{\mathcal{S}}(\mathbf{x}) - f_{\mathcal{T}}(\mathbf{x})\|], \mathcal{E}_{\mathcal{D}_{\mathcal{T}}} [\|f_{\mathcal{S}}(\mathbf{x}) - f_{\mathcal{T}}(\mathbf{x})\|] \} \quad (1)$$

Our previously proposed **DGA-DA** provides a unified framework which jointly optimizes term 2 in Eq.(1) as **DDC-DA** methods when aligning data distributions, and term 3 in Eq.(1) as **DGSA-DA** approaches when performing label inference through the underlying data geometric structure. It further introduces a *repulsive force*(**RF**) term using both source and target domain data when seeking the latent feature space and makes the proposed DA method discriminative. In this paper, we go one step further and propose a novel DA method, namely **Discriminative Noise Robust Sparse Orthogonal Label Regression-based Domain Adaptation (DOLL-DA)**, which optimizes at the same time the three terms of the right-hand of Eq.(1), including in particular the first term on classification error in Eq.(1), when seeking a discriminative latent feature space.

Specifically, the proposed **DOLL-DA** derives from an integrated DA model which: **1)** searches a shared feature subspace where the source and target data distributions are discriminatively aligned using an improved *repulsive force* (**RF**) term added to the **MMD** constraints, thereby optimizing the second term in Eq.(1) and indirectly improving its first term; **2)** makes use an embedding trick to immerse data labels into the shared feature space and projects each data sample within the vicinity of its label vector orthogonal to other ones, thereby further regularizing the improved **MMD** constraints and avoiding potential contradictions among sub-domains when *repulsive force* is applied in the search of the shared feature subspace; **3)** linearly regresses data labels in the shared feature subspace, thereby further explicitly optimizes the first term of the right-hand in Eq.(1). Moreover, data outliers are accounted for and the property of sparsity in label regression is introduced to circumvent negative transfer and over-fitting, leading to a novel *Noise Robust Sparse Orthogonal Label Regression (NRS_OLR)* term in our DA model; **4)** leverages the true labels available in the source domain and ensures a *label consistency* between the source and target domain through an iterative integrated linear label regression, thereby minimizing the third term of Eq.(1); Fig.1 depicts the general framework of the proposed **DOLL-DA** method.

To sum up, the contributions of this paper are as follows:

- Improved *repulsive force* (**RF**)-based MMD constraints are introduced to enable discriminative alignment of data distributions between the source and the target domain.
- Orthogonal label subspace is proposed through a label embedding trick to further regularize the improved **RF**-based **MMD** constraints using a novel *Orthogonal Label Regression (ORL)* constraint, thereby circumventing potential conflicts which could arise when optimizing the improved **RF**-based MMD constraints and further enhancing the discriminative power of the proposed DA method.
- The hypothesis for classifying the target domain data is learned simultaneously through a single feature projection matrix **A** when aligning discriminatively the source and target domain data in the shared feature subspace. Furthermore, a property of sparsity is introduced through a $l_{2,1}$ -norm constraint on **A** when regressing data labels and data outliers are also accounted for within the model, leading to a *Noise Robust Sparse Orthogonal Label Regression(NRS_OLR)* term to ensure the proposed DA model to avoid negative transfer as well as overfitting.
- A novel generalized power iteration method is introduced to solve the optimization problem of the full integrated DA model, resulting in the proposed **DOLL-DA**. Furthermore, we perform time complexity analysis and also derive three additional DA methods based on three partial models from the full integrated DA in order to highlight the individual contribution of **RF** and **NRS_OLR** term, respectively, as well as their added value when they are jointly optimized.
- We perform extensive experiments on 49 image classification DA tasks using 8 popular **DA** benchmarks and demonstrate the effectiveness of the proposed **DOLL-DA** which consistently outperforms thirty state-of-the-art **DA** algorithms with a margin which can reach 17 points. Moreover, we also carry out in-depth analysis of the proposed **DA** method, in particular *w.r.t.* their hyper-parameters, convergence speed, the choice of the base classifier, random label initialization and impact of the quantity of target domain data for performance stability.

The article is organized as follows. Sect.II discusses the related work. Sect.III presents the method. Sect.IV benchmarks the proposed DA method and provides in-depth analysis. Sect.V draws the conclusion.

II. RELATED WORK

Last years have seen **DA** techniques applied to multiple computer vision applications. State-of-the-art has so far featured two main research streams: 1) Shallow **DA**; 2) Deep **DA**. They are overviewed in Sect.II-A and sect.II-B, respectively, and discussed in comparison with the proposed **DOLL-DA** in sect.II-C.

A. Shallow Domain Adaptation

1) *Feature-based DA*: The rationale of the *feature*-based domain adaptation is to assume a shared latent feature space between the source and target domain which is searched in narrowing the existing distribution discrepancies across the domains. A popular strategy for searching such a shared latent feature space is to embrace the dimensionality reduction and

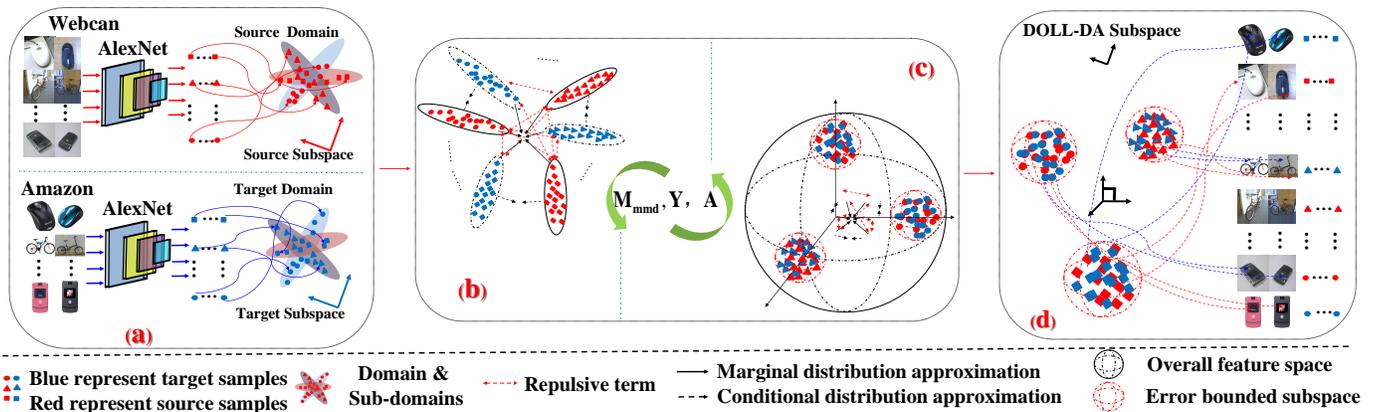


Fig. 1: Illustration of the proposed **DOLL-DA** method. Fig.1 (a): source domain data and target domain data, *e.g.*, mouse, bike, smartphone images, with different distributions and inherent hidden data geometric structures between the source in red and the target in blue. Samples of different class labels are represented by different geometrical shapes, *e.g.*, round, triangle and square; Fig.1 (b) illustrates **DOLL-DA** which aligns data distributions closely yet discriminatively through the use of the nonparametric distance, *i.e.*, Maximum Mean Discrepancy (MMD). Fig.1 (c): accounts for well regularization of the designed **MMD** distances, which intends to enable the different sub-domains, *w.r.t.*, its orthogonal subspace meanwhile cares about the noise data; In **DOLL-DA**, the optimized common subspace **A** and label matrix **Y** are updated iteratively within the processes in Fig.1 (b-c); Fig.1 (d): the achieved latent joint subspace where both marginal and class conditional data distributions are aligned discriminatively through the well proposed orthogonal regularization; Furthermore, noise data are well accounted as well as the hidden manifold structure, the formal one intends to avoid the negative transfer while the latter one improve the transferability of the learned model based on the source domain.

propose to explicitly minimize some predefined distance measure to reduce the mismatch between source and target in terms of marginal distribution [52] [42], or conditional distribution [33]. These methods can be further distinguished based on whether they incorporate some form of data discriminativeness or not in the search of such a shared latent feature space.

Nondiscriminative distribution alignment (NDA): **NDA** strategies propose to align the marginal and conditional distributions across the source and target domains in reducing different data distribution distance measurements, *e.g.*, Bregman Divergence [52], Wasserstein distance [6], [7], *Maximum Mean Discrepancy* (MMD) to explicitly shrink the cross-domain divergence of marginal data distributions [42], and both the marginal and conditional data distributions [33]. The obvious drawback of **NDA** is that it ignores the discriminative knowledge among different labeled source sub-domains, thereby increasing the burden of the required classifier.

Discriminative distribution alignment (DDA): **DDA** methods improve **NDA** ones by explicitly leveraging the discriminative information in source domain data labeled into different sub-domains. **ILS**[20] learns a discriminative latent space using Mahalanobis metric and makes use of Riemannian optimization strategy to match statistical properties across different domains. [36] adapts linear discriminant analysis (**LDA**) and leverages the discriminative information from the target domain to estimate the common feature space. **OBTL**[25] proposes bayesian transfer learning based domain adaptation, which explicitly discusses the relatedness across different sub-domains. **SCA**[15] achieves discriminativeness in optimizing the interplay of the between and within-class scatters. Our proposed **DGA-DA**[39] also introduces a specific *repulsive*

force term to capture the data discriminativeness. However, our **DGA-DA** further cares about the underlying data manifold structure when performing label inference. However, despite improvement over **NDA** methods, **DDA** methods rely upon temporary and unreliable pseudo labels of target domain data to capture the repulsive force in the target domain, and thereby can mislead the search of an optimized shared discriminative latent space.

2) *Subspace alignment-based DA:* In line with [12], an increasing number of DA methods, *e.g.*, [40], [51], [63], [54], [9], emphasize the importance of aligning the underlying data subspace and manifold structures between the source and the target domain for effective DA. In these methods, low-rank and sparse constraints are introduced into DA to extract a low-dimension feature subspace where target samples can be sparsely reconstructed from source samples [51], or interleaved by source samples [63], thereby aligning the geometric structures of the underlying data manifolds. A few recent DA methods, *e.g.*, **RSA-CDDA**[40], **JGSA**[64], further propose unified frameworks to reduce the shift between domains both statistically and geometrically. **HCA**[31] improves **JGSA** using a homologous constraint on the two transformations for the source and target domains, respectively, to make the transformed domains related and hence alleviate negative domain adaptation.

However, in light of the upper error bound as defined in Eq.(1), we can see that subspace alignment based DA methods account for the underlying data geometric structure and expect but without theoretical guarantee the alignment of discriminative data distributions. Our proposed **DOLL-DA** improves subspace alignment-based DA by jointly aligning the

data distributions discriminatively and enhancing the hidden manifold structure through label regression to regress different sub-domains *w.r.t* its orthogonal label subspace.

B. Deep Domain Adaptation

Recently, **DA** has been intensively investigated under the paradigm of deep learning (**DL**), and has featured the following two main approaches.

1) *Statistic matching-based DA*: These methods aim to reduce the divergence across domains using statistic measurements incorporated into **DL** frameworks. **DAN**[32] reduces the marginal distribution divergence in incorporating the multi-kernel MMD loss on the fully connected layers of AlexNet. **JAN**[35] improves **DAN** by jointly decreasing the divergence of both the marginal and conditional distributions. **D-CORAL**[55] further introduces the second-order statistics into the AlexNet[27] framework for more effective **DA** strategy.

2) *Adversarial loss-based DA*: These methods make use of **GAN**[17] and propose to align data distributions across domains in making sample features indistinguishable *w.r.t* the domain labels through an adversarial loss on a domain classifier [14], [57], [46]. **DANN**[14] and **ADDA**[57] learn a domain-invariant feature subspace in reducing the marginal distribution divergence. **MADA** [46] additionally make use of multiple domain discriminators, thereby aligning conditional data distributions. Different from the previous approaches, **DSN**[4] achieves domain-invariant representations in explicitly separating the similarities and dissimilarities in the source and target domains. **MADAN**[65] explores knowledge from different multi-source domains to fulfill **DA** tasks. **CyCADA**[22] addresses the distribution divergence using a bi-directional **GAN** based training framework.

The main advantage of these **DL** based **DA** methods is that they jointly shrink the divergence of data distributions across domains and achieve a discriminative feature representation of data through a single unified end-to-end learning framework. However, they also present the drawback that the discriminative force is merely extracted from the labelled source domain while it could rely upon simultaneously the source and target domains in leveraging the underlying data geometric structures. Furthermore, these **DL** based **DA** approaches mostly work as a black-box and suffer from interpretability, thereby falling short to provide deep insights for further improving **DA** methods.

C. Discussion

Fig.2 compares our proposed **DOLL-DA** with our previously proposed **DA** method, **DGA-DA**[39] as well as **MEDA**[62], and highlights their similarities and differences according to the following 6 properties:

- **Dis**(Discriminateness): both **DGA-DA** and **DOLL-DA** introduce a (**RF**) term to achieve discriminativeness across domains, while **MEDA** ignores this merit. **DGA-DA** takes into account **RF** term across domain, while the proposed **DOLL-DA** improves it in extending it with a **RF** term within the source domain.

Method	Dis	Reg	Cons	Joint	Error	Manifold
DGA-DA	✓	X	✓	X	X	Laplace Graph
MEDA	X	X	X	✓	X	Laplace Graph
DOLL-DA	✓	✓	✓	✓	✓	Regression

Dis=Discriminateness; Reg=Regularization of discriminativeness; Cons=Constraint ($A^T X H X^T A = I$); Joint= Joint optimization;

Fig. 2: Model comparison

- **Reg**(Regularization): Different from **DGA-DA** and **MEDA**, **DOLL-DA** proposes orthogonal label regression which further regularizes the effectiveness of the improved **RF**-based **MMD** constraints, thereby circumventing potential contradictions when reinforcing these **MMD** constraints.
- **Cons**(Constraint): Both **DGA-DA** and **DOLL-DA** are optimized using the constraint $A^T X H X^T A = I$, which removes an arbitrary scaling factor in the embedding and prevents the optimization from collapsing into a subspace of dimension less than the required dimensions. **MEDA** does not have such a constraint to obtain an analytical solution.
- **Joint**(Joint optimization): Both **MEDA** and **DOLL-DA** achieve data distribution alignment and classifier optimization through a unified joint optimization model, whereas the optimization in **DGA-DA** for data distribution alignment and classifier optimization is carried out separately.
- **Error**: **DOLL-DA** designs an error tolerated subspace to care about the noise in data, therefore decreasing the risk of potential negative transfer.
- **Manifold**: both **MEDA** and **DGA-DA** capture the underlying data manifold structures for label inference. While effective, they also suffer from the computational burden due to the singular value decomposition. **DOLL-DA** also cares about data manifold structures, but through computation effective iterative integrated linear label regression.

III. THE PROPOSED METHOD

Sect.III-A defines the notations and states the **DA** problem. Sect.III-B formulates our **DA** model while Sect.III-C presents the generalized power iteration method to solve the proposed **DA** model and derives the algorithm of **DOLL-DA**. Sect.III-D extends **DOLL-DA** to non-linear problem through kernel mapping. Sect.III-E performs time complexity analysis of the proposed **DOLL-DA**.

A. Notations and Problem Statement

Matrices are written as boldface uppercase letters. Vectors are written as boldface lowercase letters. For matrix $\mathbf{M} = (m_{ij})$, its i -th row is denoted as \mathbf{m}^i , and its j -th column is denoted by \mathbf{m}_j . We define the Frobenius norm $\|\cdot\|_F$ and $l_{2,1}$ norm as: $\|\mathbf{M}\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^l m_{ij}^2}$ and $\|\mathbf{M}\|_{2,1} = \sum_{i=1}^n \sqrt{\sum_{j=1}^l m_{ij}^2}$. A domain D is defined as an l -dimensional feature space χ and a marginal probability distribution $P(x)$, i.e., $\mathcal{D} = \{\chi, P(x)\}$ with $x \in \chi$. Given a

specific domain D , a task T is composed of a C-cardinality label set \mathcal{Y} and a classifier $f(x)$, i.e., $T = \{\mathcal{Y}, f(x)\}$, where $f(x) = \mathcal{Q}(y|x)$ can be interpreted as the class conditional probability distribution for each input sample x .

In unsupervised domain adaptation, we are given a source domain $\mathcal{D}_S = \{x_i^s, y_i^s\}_{i=1}^{n_s}$ with n_s labeled samples $\mathbf{X}_S = [x_1^s \dots x_{n_s}^s]$, which are associated with their class labels $\mathbf{Y}_S = \{y_1, \dots, y_{n_s}\}^T \in \mathbb{R}^{n_s \times C}$, and an unlabeled target domain $\mathcal{D}_T = \{x_j^t\}_{j=1}^{n_t}$ with n_t unlabeled samples $\mathbf{X}_T = [x_1^t \dots x_{n_t}^t]$, whose labels $\mathbf{Y}_T = \{y_{n_s+1}, \dots, y_{n_s+n_t}\}^T \in \mathbb{R}^{n_t \times C}$ are unknown. Here, $y_i \in \mathbb{R}^C$ ($1 \leq i \leq n_s + n_t$) is a one-vs-all label hot vector in which $y_i^j = 1$ if x_i belongs to the j -th class, and 0 otherwise. We define the data matrix $\mathbf{X} = [\mathbf{X}_S, \mathbf{X}_T] \in \mathbb{R}^{l \times n}$ ($l = \text{feature dimension}$; $n = n_s + n_t$) in packing both the source and target data. The source domain \mathcal{D}_S and target domain \mathcal{D}_T are assumed to be different, i.e., $\chi_S = \chi_T$, $\mathcal{Y}_S = \mathcal{Y}_T$, $\mathcal{P}(\chi_S) \neq \mathcal{P}(\chi_T)$, $\mathcal{Q}(\mathcal{Y}_S|\chi_S) \neq \mathcal{Q}(\mathcal{Y}_T|\chi_T)$. We also define the notion of *sub-domain*, i.e., class, denoted as $\mathcal{D}_S^{(c)}$, representing the set of samples in \mathcal{D}_S with the class label c . It is worth noting that, the definition of sub-domains in the target domain, namely $\mathcal{D}_T^{(c)}$, requires a base classifier, e.g., Nearest Neighbor (NN), to attribute pseudo labels for samples in \mathcal{D}_T .

The maximum mean discrepancy (MMD) is an effective non-parametric distance-measure that compares the distributions of two sets of data by mapping the data into Reproducing Kernel Hilbert Space[3] (RKHS). Given two distributions \mathcal{P} and \mathcal{Q} , the MMD between \mathcal{P} and \mathcal{Q} is defined as:

$$Dist(P, Q) = \left\| \frac{1}{n_1} \sum_{i=1}^{n_1} \phi(p_i) - \frac{1}{n_2} \sum_{i=1}^{n_2} \phi(q_i) \right\|_{\mathcal{H}} \quad (2)$$

where $P = \{p_1, \dots, p_{n_1}\}$ and $Q = \{q_1, \dots, q_{n_2}\}$ are two random variable sets from distributions \mathcal{P} and \mathcal{Q} , respectively, and \mathcal{H} is a universal RKHS with the reproducing kernel mapping $\phi: f(x) = \langle \phi(x), f \rangle$, $\phi: \mathcal{X} \rightarrow \mathcal{H}$.

The aim of the proposed DOLL-DA is to search jointly a transformation matrix $\mathbf{A} \in \mathbb{R}^{l \times k}$ projecting discriminatively both the source and target domain data of dimension l into a latent shared orthogonal feature subspace of dimension k as well as a label regressor while minimizing simultaneously the three terms of the upper error bound in Eq.(1).

B. Formulation

Our final model **DOLL-DA** (sect.III-B5) starts from **JDA** (sect.III-B1), which is improved in Sect.III-B2 and Sect.III-B3 for discriminative data distribution alignment (**DDA**) by leveraging the discriminative knowledge from the source and target domains. Fig.3 summarizes these steps which aim to minimize $d_{\mathcal{H}}(\mathcal{D}_S, \mathcal{D}_T)$ (term.2 of Eq.(1)). Sect.III-B4 further introduces an orthogonal label regressor using an embedding trick and accounts for noisy data as well as sparsity in label regression to derive Noise Robust Sparse Orthogonal Label Regression (**NRS_OLR**). Sect.III-B5 integrates **DDA** and **NRS_OLR** to achieve our final model and thereby optimizes at the same time the three terms of the right-hand of Eq.(1).

1) *Matching Marginal and Conditional Distributions*: As shown in Fig.3.a and Fig.3.b, our model starts from **JDA**, which makes use of MMD in RKHS to measure the distances between the expectations of the source domain/sub-domain and target domain/sub-domain. Specifically, **1)** The empirical distance of the source and target domains is defined as $Dist^m$; **2)** The conditional distance $Dist^c$ is defined as the sum of the empirical distances between sub-domains in \mathcal{D}_S and \mathcal{D}_T with a same label; **3)** $Dist_{Clo}$ is defined as the sum of $Dist^m$ and $Dist^c$.

$$\begin{aligned} Dist_{Clo} &= Dist^m(\mathcal{D}_S, \mathcal{D}_T) + Dist^c \sum_{c=1}^C (D_S^c, D_T^c) \\ &= \left\| \frac{1}{n_s} \sum_{i=1}^{n_s} \mathbf{A}^T x_i - \frac{1}{n_t} \sum_{j=n_s+1}^{n_s+n_t} \mathbf{A}^T x_j \right\|^2 \\ &\quad + \left\| \frac{1}{n_s^{(c)}} \sum_{x_i \in \mathcal{D}_S^{(c)}} \mathbf{A}^T x_i - \frac{1}{n_t^{(c)}} \sum_{x_j \in \mathcal{D}_T^{(c)}} \mathbf{A}^T x_j \right\|^2 \\ &= tr(\mathbf{A}^T \mathbf{X} (\mathbf{M}_0 + \sum_{c=1}^C \mathbf{M}_c) \mathbf{X}^T \mathbf{A}) \end{aligned} \quad (3)$$

- $Dist^m(\mathcal{D}_S, \mathcal{D}_T)$: where \mathbf{M}_0 is the MMD matrix between \mathcal{D}_S and \mathcal{D}_T with $(\mathbf{M}_0)_{ij} = \frac{1}{n_s n_s}$ if $(\mathbf{x}_i, \mathbf{x}_j \in \mathcal{D}_S)$, $(\mathbf{M}_0)_{ij} = \frac{1}{n_t n_t}$ if $(\mathbf{x}_i, \mathbf{x}_j \in \mathcal{D}_T)$ and $(\mathbf{M}_0)_{ij} = 0$ otherwise. Thus, the difference between the marginal distributions $\mathcal{P}(\chi_S)$ and $\mathcal{P}(\chi_T)$ is reduced when minimizing $Dist^m(\mathcal{D}_S, \mathcal{D}_T)$.
- $Dist^c(\mathcal{D}_S, \mathcal{D}_T)$: where C is the number of classes, $\mathcal{D}_S^{(c)} = \{\mathbf{x}_i : \mathbf{x}_i \in \mathcal{D}_S \wedge y(\mathbf{x}_i) = c\}$ represents the c^{th} sub-domain in the source domain, in which $n_s^{(c)} = \left\| \mathcal{D}_S^{(c)} \right\|_0$ is the number of samples in the c^{th} source sub-domain. $\mathcal{D}_T^{(c)}$ and $n_t^{(c)}$ are defined similarly for the target domain but using pseudo-labels. Finally, \mathbf{M}_c denotes as the MMD matrix between the sub-domains with labels c in \mathcal{D}_S and \mathcal{D}_T with $(\mathbf{M}_c)_{ij} = \frac{1}{n_s^{(c)} n_s^{(c)}}$ if $(\mathbf{x}_i, \mathbf{x}_j \in \mathcal{D}_S^{(c)})$, $(\mathbf{M}_c)_{ij} = \frac{1}{n_t^{(c)} n_t^{(c)}}$ if $(\mathbf{x}_i, \mathbf{x}_j \in \mathcal{D}_T^{(c)})$, $(\mathbf{M}_c)_{ij} = \frac{-1}{n_s^{(c)} n_t^{(c)}}$ if $(\mathbf{x}_i \in \mathcal{D}_S^{(c)}, \mathbf{x}_j \in \mathcal{D}_T^{(c)})$ or $(\mathbf{x}_i \in \mathcal{D}_T^{(c)}, \mathbf{x}_j \in \mathcal{D}_S^{(c)})$ and $(\mathbf{M}_c)_{ij} = 0$ otherwise. As a consequence, the mismatch of conditional distributions between \mathcal{D}_S^c and \mathcal{D}_T^c is reduced in minimizing $Dist^c$.

2) *Across domain Repulsive force term*: As shown in Fig.3.(a,b), sect.III-B1 merely cares about shrinking the MMD distances in order to align data marginal and conditional distributions between the source and target domain, and ignores discriminative knowledge within data. Here, a *repulsive force*(**RF**) term $Dist_{S \rightarrow T}^{re} + Dist_{T \rightarrow S}^{re}$ is introduced to enable discriminative **DA** as shown in Fig.3.c. Specifically, we denote $S \rightarrow T$ and $T \rightarrow S$ to index the distances computed from \mathcal{D}_S to \mathcal{D}_T and \mathcal{D}_T to \mathcal{D}_S , respectively, and $Dist_{S \rightarrow T}^{re}$ as the sum of the distances between each source sub-domain $\mathcal{D}_S^{(c)}$ and all the target sub-domains $\mathcal{D}_T^{(r)}$; $r \in \{1 \dots C\} - \{c\}$ excluding the c -th target sub-domain. Symmetrically, $Dist_{T \rightarrow S}^{re}$ is defined in a similar way as $Dist_{S \rightarrow T}^{re}$. These two distances are computed

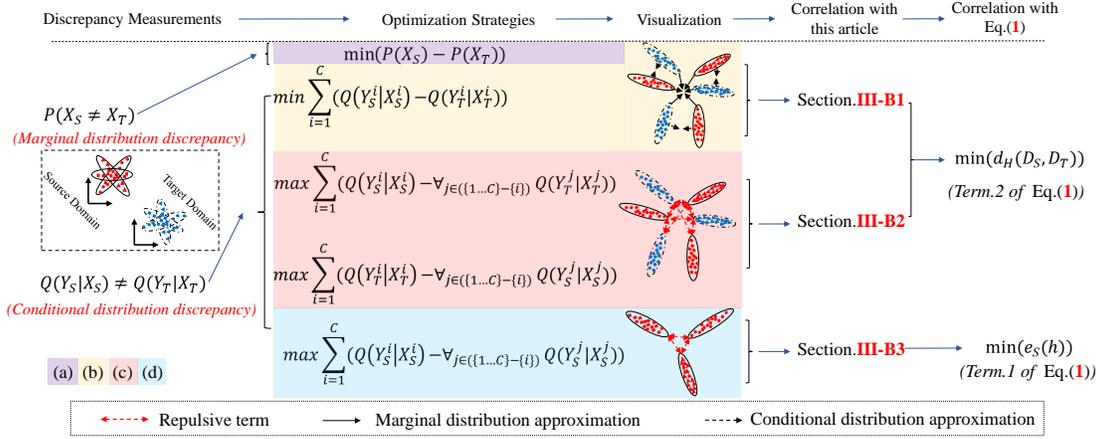


Fig. 3: Fig.3 (a): marginal distribution matching; Fig.3 (b): conditional distribution matching; Fig.3 (c): *Repulsive force* term proposed on the source domain; Fig.3 (d): *Repulsive force* term proposed across the source and target domains. (purple, yellow, blue, and red parts represent Fig.3 (a), Fig.3 (b), Fig.3 (c), and Fig.3 (d) respectively.)

as:

$$\begin{aligned} Dist_{S \rightarrow T}^{re} + Dist_{T \rightarrow S}^{re} &= Dist^c \sum_{c=1}^C (D_S^c, D_T^{r \in \{(1...C) - \{c\}\}}) \\ &+ Dist^c \sum_{c=1}^C (D_T^c, D_S^{r \in \{(1...C) - \{c\}\}}) \quad (4) \\ &= \sum_{c=1}^C tr(\mathbf{A}^T \mathbf{X} (\mathbf{M}_{S \rightarrow T} + \mathbf{M}_{T \rightarrow S}) \mathbf{X}^T \mathbf{A}) \end{aligned}$$

Where:

- $\mathbf{M}_{S \rightarrow T}$ is defined as: $(\mathbf{M}_{S \rightarrow T})_{ij} = \frac{1}{n_s^{(c)} n_s^{(c)}}$ if $(\mathbf{x}_i, \mathbf{x}_j \in D_S^{(c)})$, $\frac{1}{n_t^{(r)} n_t^{(r)}}$ if $(\mathbf{x}_i, \mathbf{x}_j \in D_T^{(r)})$, $\frac{-1}{n_s^{(c)} n_t^{(r)}}$ if $(\mathbf{x}_i \in D_S^{(c)}, \mathbf{x}_j \in D_T^{(r)} \text{ or } \mathbf{x}_i \in D_T^{(r)}, \mathbf{x}_j \in D_S^{(c)})$ and 0 otherwise.
- $\mathbf{M}_{T \rightarrow S}$ is defined as: $(\mathbf{M}_{T \rightarrow S})_{ij} = \frac{1}{n_t^{(c)} n_t^{(c)}}$ if $(\mathbf{x}_i, \mathbf{x}_j \in D_T^{(c)})$, $\frac{1}{n_s^{(r)} n_s^{(r)}}$ if $(\mathbf{x}_i, \mathbf{x}_j \in D_S^{(r)})$, $\frac{-1}{n_t^{(c)} n_s^{(r)}}$ if $(\mathbf{x}_i \in D_T^{(c)}, \mathbf{x}_j \in D_S^{(r)} \text{ or } \mathbf{x}_i \in D_S^{(r)}, \mathbf{x}_j \in D_T^{(c)})$ and 0 otherwise.

Therefore, maximizing Eq.(4) increases the distances of each sub-domain with the other remaining sub-domains across domain, *i.e.*, the between-class distances across domain, and thereby facilitates a discriminative DA. This across domain **RF** term was introduced in our previously proposed **DGA-DA** [39] and has already shown its effectiveness.

3) *Repulsive force term within the source domain*: While Sect.III-B1 and sect.III-B2 have so far endeavored to minimize the second term of the right-hand in Eq.(1), we turn our attention here to optimize the first term of Eq.(1) as shown in Fig.3.(d). Specifically, we introduce a *repulsive force* term $Dist_{S \rightarrow S}^{re}$ (Fig.3.(d)), so as to increase the discriminative power on the labeled source domain, thereby making it possible for a better predictive model on the source domain. Using $S \rightarrow S$ to index the distances computed from D_S to D_S , we can compute, similarly as in eq.(4), $Dist_{S \rightarrow S}^{re}$ as the sum of the distances from each source sub-domain $D_S^{(c)}$ to all the other source sub-domains $D_S^{(r)}$: $r \in \{(1...C) - \{c\}\}$, excluding the c -th source sub-domain:

$$\begin{aligned} Dist_{S \rightarrow S}^{re} &= Dist^c \sum_{c=1}^C (D_S^c, D_S^{r \in \{(1...C) - \{c\}\}}) \quad (5) \\ &= \sum_{c=1}^C tr(\mathbf{A}^T \mathbf{X} (\mathbf{M}_{S \rightarrow S}) \mathbf{X}^T \mathbf{A}) \end{aligned}$$

where $\mathbf{M}_{S \rightarrow S}$ is defined as: $(\mathbf{M}_{S \rightarrow S})_{ij} = \frac{1}{n_s^{(c)} n_s^{(c)}}$ if $(x_i, x_j \in D_S^{(c)})$, $\frac{1}{n_s^{(r)} n_s^{(r)}}$ if $(x_i, x_j \in D_S^{(r)})$, $\frac{-1}{n_s^{(c)} n_s^{(r)}}$ if $(x_i \in D_S^{(c)}, x_j \in D_S^{(r)} \text{ or } x_i \in D_S^{(r)}, x_j \in D_S^{(c)})$ and 0 otherwise.

Maximizing Eq.(5) increases the between-class distances in the source domain and thereby optimizing the first term of the right-hand in Eq.(1) on classification errors on the source domain. In our model, the **RF** term within the source domain as defined by Eq.(5) added to the across domain **RF** term as defined by Eq.(4) is named improved *repulsive force* (**RF**) and optimized simultaneously.

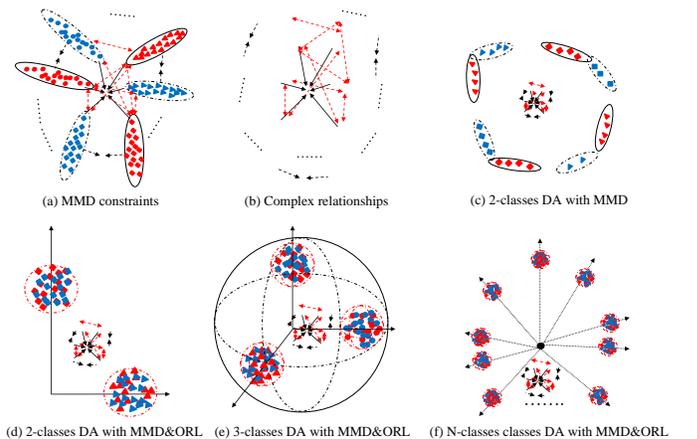


Fig. 4: Illustration of MMD constraints and orthogonal projection based label consistent regression.

4) *Noise Robust Sparse Orthogonal Label Regression*: The improved *repulse force* term formulated in sect.III-B2 and sect.III-B3 aims to increase the between-class distances, across domain through Eq.(4) and within the source domain through

Eq.(5), and enable discriminative alignment of marginal and conditional data distributions (eq.(3)) in sect.III-B1), but does not seek to decrease the intra-class distances. As a result, situations as shown in Fig.4.(a,b,c), where the instances of a class are pushed away from other class instances but don't get closer each other within the class, could happen. More importantly, while eq.(4) and eq.(5) make it possible to decrease the classification errors of a hypothesis on the source domain, the first term of the error bound in eq.(1) is not explicitly optimized. Here, we propose to solve the previous two issues through orthogonal label regression as shown in Fig.4.(d, e, f) where the instances of each class (sub-domain) across domain get close to their one-vs-all hot label vector which is orthogonal to those of other classes.

To this end, we introduce a novel *orthogonal label regression (OLR)* constraint $\Phi(\mathbf{A}, \mathbf{Y}_S, \mathbf{X}_S)$, where \mathbf{A} is the transformation matrix projecting both the source and target data onto a latent shared feature subspace of dimension k . Specifically, we first introduce an embedding trick which consists of immersing a C -dimensional one-vs-all hot label vector into the k -dimensional latent shared feature subspace simply by adding $(k - C + 1)$ times 0, e.g., a 3-dimensional one-vs-all label vector $(0, 1, 0)^t$ is represented by its corresponding one-vs-all hot vector $(0, 1, 0, 0, 0)^t$ in the 5-dimensional feature subspace using the embedding trick. We can then perform least square regression (**LSR**): $\min \|\mathbf{X}^T \mathbf{A} - \mathbf{Y}\|_F^2$ st. $\mathbf{Y} \geq \mathbf{0}, \mathbf{Y}\mathbf{1} = \mathbf{1}$, with \mathbf{Y} the class label matrix as defined in sect.III-A and extended into a $n \times k$ matrix by embedding each one-vs-all hot label vector into a k -dimensional one using the embedding trick. Minimization of **LSR** thus simply requires that each labeled sample be projected within the vicinity of its corresponding label hot vector in the k -dimensional feature space.

It is worth noting that the proposed **OLR** enjoys the following three properties:

- **Orthogonality**, i.e., $(\mathbf{Y}_{i=c} \bullet \mathbf{Y}_{i \neq c}) = 0$ with \bullet denoting the dot product. This constraint simply expresses that one-vs-all hot label vectors in the shared latent feature subspace are orthogonal each other. As a result, projecting each data sample into the vicinity of its corresponding label vector keeps the samples of each sub-domain far away from those of other sub-domains and thereby improving data discriminativeness and optimizing term.1 of Eq.(1);
- **Label Embedding Constraint**, i.e., $Q(\mathbf{Y}_{i \notin \{1, \dots, C\}} | \chi_S \cup \chi_T) = 0$. This property simply denotes the fact that we have made use of the embedding trick for immersing each C dimensional one-vs-all hot label vector into the k dimensional shared feature space and there are no label vectors for classes ranging from $C + 1$ to k ;
- **Sharing of the feature projection and label regression matrix** through \mathbf{A} . Thanks to the label embedding constraint, the projection matrix \mathbf{A} through eq.(3), eq.(4) and eq.(5) for the search of a shared latent feature space aligning discriminatively marginal and conditional data distributions between the source and target domain can be shared with the one used for the orthogonal label regression (**OLR**) constraint, thereby jointly optimizing

term.1 and term.2 of Eq.(1) within a single unified feature and label subspace. Furthermore, as shown in Fig.4.(d), (e)), data of each sub-domain, i.e. data with a same label, from the source and target domain, are projected within the vicinity of their corresponding one-vs-all hot label vector, thereby also decreasing Term.3 of Eq.(1), i.e., the errors of their respective labelling functions on the source and target domain.

However, data from the source and target domain can be noisy. We account for data noise through an error matrix \mathbf{E} and the **OLR** constraint can therefore be reformulated as: $\min \|\mathbf{X}^T \mathbf{A} - \mathbf{Y} + \mathbf{E}\|_F^2$ st. $\mathbf{Y} \geq \mathbf{0}, \mathbf{Y}\mathbf{1} = \mathbf{1}$. The error matrix \mathbf{E} makes possible a certain tolerance of errors when projecting data into the vicinity of its corresponding label vector in the latent shared feature space, thereby enabling to account for outliers and alleviating the influence of negative transfer. Additionally, given the fact that, in real-life applications, e.g. visual object recognition, data of a given class generally lie within a manifold of much lower dimension in comparison with the original data space, e.g., pixel number of images, we further introduce a $l_{2,1}$ -norm constraint so as to fulfill the property that the class label of a data sample should be regressed from a sparse combination of features in the latent shared feature subspace. This constraint introduces a regularization term on \mathbf{A} for discriminative subspace projection, which also optimizes Term.3 of Eq.(1). Putting all these together, the initial *Orthogonal Label Regression (OLR)* constraint becomes *Noise Robust Sparse Orthogonal Label Regression (NRS_OLR)* and is finally formulated as:

$$\min \|\mathbf{X}^T \mathbf{A} - \mathbf{Y} + \mathbf{1e}^T\|_F^2 + \beta \|\mathbf{A}\|_{2,1}^2 \quad (6)$$

$$\text{st. } \mathbf{Y} \geq \mathbf{0}, \mathbf{Y}\mathbf{1} = \mathbf{1}, (\mathbf{Y}_{i=c} \bullet \mathbf{Y}_{i \neq c}) = 0$$

5) *The final model*: By integrating all the properties introduced in sect.III-B1 through sect.III-B4, we obtain our final DA model, formulated as Eq.(19)

$$\min_{\mathbf{A}, \mathbf{e}, \mathbf{Y}_{\mathbf{U}} \geq \mathbf{0}, \mathbf{Y}_{\mathbf{U}} \mathbf{1} = \mathbf{1}} (tr(\mathbf{A}^T \mathbf{X} \mathbf{M}^* \mathbf{X}^T \mathbf{A}) + \alpha \|\mathbf{A}\|_F^2 + \beta \|\mathbf{A}\|_{2,1}^2) + \|\mathbf{X}^T \mathbf{A} + \mathbf{1e}^T - \mathbf{Y}\|_F^2 \quad (7)$$

$$\text{st. } \mathbf{M}^* = \mathbf{M}_0 + \sum_{c=1}^C (\mathbf{M}_c) - \mathbf{M}_{REP}, \mathbf{Y} \geq \mathbf{0}, \mathbf{Y}\mathbf{1} = \mathbf{1},$$

$$\mathbf{A}^T \mathbf{X} \mathbf{H} \mathbf{X}^T \mathbf{A} = \mathbf{I}, (\mathbf{Y}_{i=c} \bullet \mathbf{Y}_{i \neq c}) = 0$$

where $\mathbf{M}_{REP} = \mathbf{M}_{S \rightarrow T} + \mathbf{M}_{T \rightarrow S} + \mathbf{M}_{S \rightarrow S}$ is the improved overall *repulsive force* constraint matrix, $\mathbf{H} = \mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}^T$ is the centering matrix, the constraint $\mathbf{Y} \geq \mathbf{0}$ and $\mathbf{Y}\mathbf{1} = \mathbf{1}$ simply expresses the fact that each data sample has a label vector whose class probability sums to 1, whereas the constraint $\mathbf{A}^T \mathbf{X} \mathbf{H} \mathbf{X}^T \mathbf{A} = \mathbf{I}$ derives from Principal Component Analysis (**PCA**) to preserve the intrinsic data covariance of both domains and avoid the trivial solution for \mathbf{A} .

Through iterative optimization of Eq.(19), our DA method searches jointly a well regularized discriminative latent feature subspace shared between the source and target domain and a noise robust sparse label regression model, thereby optimizing at the same time the three terms of the right-hand of Eq.(1).

C. Solving the model

Eq.(19) is not convex, therefore we propose an effective method which optimizes the key variables, e.g., $\mathbf{A}, \mathbf{M}, \mathbf{Y}$, in

a coordinate descent manner. Main steps for solving Eq.(19) are as follows. All the key steps have closed form solution:

Step.1 (Initialization of \mathbf{M}^*) \mathbf{M}^* can be initialized by calculating \mathbf{M}_0 since there is no labels or pseudo labels on the target domain initially. We obtain $\mathbf{M}^* = \mathbf{M}_0$ where \mathbf{M}_0 is the MMD matrix as defined in Eq.(3).

Step.2 (Initialization of \mathbf{A}) Similar to **JDA**, \mathbf{A} can be initialized to reduce the marginal and conditional distributions between $\mathcal{P}(\mathcal{X}_S)$ and $\mathcal{P}(\mathcal{X}_T)$ through an adaptive feature subspace via the Rayleigh quotient algorithm, in solving Eq.(8):

$$(\mathbf{X}(\mathbf{M}_0 + \sum_{c=1}^{c=C} \mathbf{M}_c)\mathbf{X}^T + \alpha\mathbf{I})\mathbf{A} = \mathbf{X}\mathbf{H}\mathbf{X}^T\mathbf{A}\Phi \quad (8)$$

where $\mathbf{X} = \mathbf{X}_S \cup \mathbf{X}_T$ and $\Phi = \text{diag}(\phi_1, \dots, \phi_k) \in R^{k \times k}$ are Lagrange multipliers, $\sum_{c=1}^{c=C} \mathbf{M}_c$ is obtained through labeled source domain data $\mathbf{A}^T\mathbf{X}_S$ and pseudo labels inferred on the target domain data $\mathbf{A}^T\mathbf{X}_T$. \mathbf{A} is then initialized as the k smallest eigenvectors of Eq.(8), with k defining the dimension of the latent shared feature subspace between the source and target domain.

Step.3 (Update of \mathbf{e}) \mathbf{e} is updated in solving Eq.(19) with other variables held fixed. To update \mathbf{e} , one should solve Eq.(9)

$$\mathbf{e}' = \arg \min_{\mathbf{e}} \|\mathbf{X}^T\mathbf{A} + \mathbf{1}\mathbf{e}^T - \mathbf{Y}\|_F^2 \quad (9)$$

In setting to 0 the partial derivative of Eq.(9) with respect to \mathbf{e} , we achieve the optimal solution of \mathbf{e} as

$$\mathbf{e} = \frac{1}{n}(\mathbf{Y}^T\mathbf{1} - \mathbf{A}^T\mathbf{X}\mathbf{1}) \quad (10)$$

Step.4 (Update of \mathbf{A}) \mathbf{A} is updated by solving the optimization problem in Eq.(19) with other variables held fixed. To ensure that Eq.(19) is differentiable, we regularize $\|\mathbf{A}\|_{2,1}$ as $(\sum_{j=1}^k \sqrt{\|\mathbf{a}^j\|_2^2 + \varepsilon})$ to avoid $\|\mathbf{A}\|_{2,1} = 0$. As a result, Eq.(19) becomes Eq.(11)

$$\mathbf{A}' = \arg \min_{\mathbf{A}, \mathbf{A}^T\mathbf{X}\mathbf{H}\mathbf{X}^T\mathbf{A}=\mathbf{I}} ((\text{tr}(\mathbf{A}^T\mathbf{X}(\mathbf{M}^*)\mathbf{X}^T\mathbf{A}) + \alpha\|\mathbf{A}\|_F^2 + \|\mathbf{X}^T\mathbf{A} + \mathbf{1}\mathbf{e}^T - \mathbf{Y}\|_F^2 + \beta(\sum_{j=1}^k \sqrt{\|\mathbf{a}^j\|_2^2 + \varepsilon})) \quad (11)$$

ε is infinitely close to zero, thereby making Eq.(11) closely equivalent to Eq.(19). Solving directly Eq.(11) is non-trivial, we introduce a new variable $\mathbf{G} \in R^{k \times k}$ which is a diagonal matrix with $g_{jj} = (\sum_{i=1}^k \sqrt{\|\mathbf{a}^i\|_2^2 + \varepsilon}) \div (\sqrt{\|\mathbf{a}^j\|_2^2 + \varepsilon})$. \mathbf{G} and \mathbf{A} can be optimized iteratively. With \mathbf{G} held fixed and \mathbf{e} computed as in Eq.(10), we can reformulate Eq.(11) as Eq.(12)

$$\mathbf{A}' = \arg \min_{\mathbf{A}, \mathbf{A}^T\mathbf{X}\mathbf{H}\mathbf{X}^T\mathbf{A}=\mathbf{I}} (\text{tr}(\mathbf{A}^T\mathbf{X}(\mathbf{M}^*)\mathbf{X}^T\mathbf{A}) + \alpha\|\mathbf{A}\|_F^2 + \|\mathbf{H}\mathbf{X}^T\mathbf{A} - \mathbf{H}\mathbf{Y}\|_F^2 + \beta\text{tr}(\mathbf{A}^T\mathbf{G}\mathbf{A})) \quad (12)$$

Eq.(12) is a least square problem on the Stiefel manifold, which is a non-convex optimization problem. Therefore, it cannot be directly solved via the Lagrangian method or an analytical solution. Inspired by previous research on solving quadratic problem on the Stiefel manifold[13], [41], [19], we propose a novel generalized power iteration (**GPI**) method to optimize the projection matrix \mathbf{A} that rotates the factor matrix to best fit the hypothesis subspace.

Step.i We propose Cholesky factorization of \mathbf{H} , which aims to obtain a lower triangular matrix \mathbf{h} , so that $\mathbf{h}\mathbf{h}^T = \mathbf{H}$.

Step.ii Eq.(12) is reformulated as:

$$\arg \min_{\mathbf{A}, \mathbf{A}^T\mathbf{X}\mathbf{H}\mathbf{X}^T\mathbf{A}=\mathbf{I}} \text{tr}(\mathbf{A}^T(\mathbf{X}\mathbf{H}^T\mathbf{H}\mathbf{X}^T + \beta\mathbf{G} + \alpha\mathbf{I} + (\mathbf{X}(\mathbf{M}^*)\mathbf{X}^T))\mathbf{A}) - 2\text{tr}(\mathbf{A}^T\mathbf{X}\mathbf{H}^T\mathbf{H}\mathbf{Y}) \quad (13)$$

We set:

$$\begin{aligned} \mathbf{W}^T &= \mathbf{A}^T\mathbf{X}\mathbf{h} \\ \mathbf{C} &= \mathbf{h}^{-1}\mathbf{H}^T\mathbf{H}\mathbf{Y} \\ \mathbf{B} &= (\mathbf{X}\mathbf{h})^{-1}((\mathbf{X}\mathbf{H}^T\mathbf{H}\mathbf{X}^T + \beta\mathbf{G} + \alpha\mathbf{I} + (\mathbf{X}(\mathbf{M}^*)\mathbf{X}^T))(\mathbf{h}^T\mathbf{X}^T)^{-1}) \end{aligned} \quad (14)$$

Using Eq.(14), Eq.(13) can be written for short as:

$$\arg \min_{\mathbf{W}, \mathbf{W}^T\mathbf{W}=\mathbf{I}} \text{tr}(\mathbf{W}^T\mathbf{B}\mathbf{W} - 2\mathbf{W}^T\mathbf{C}) \quad (15)$$

where \mathbf{B} is a symmetric matrix.

Step.iii Initialize $\mathbf{W}^T = \mathbf{A}^T\mathbf{X}\mathbf{h}$ to satisfy $\mathbf{W}^T\mathbf{W} = \mathbf{I}$, and set $\mathbf{B}' = \mu\mathbf{I} - \mathbf{B}$, which ensures \mathbf{B}' is a positive definite matrix.

Step.iv Set $\mathbf{Z} = 2\mathbf{B}'\mathbf{W} + 2\mathbf{C}$. Then, optimize $\mathbf{U}\mathbf{S}\mathbf{V}^T = \mathbf{Z}$ via singular value decomposition method on \mathbf{Z} .

Step.v Update $\mathbf{W}^T = \mathbf{U}\mathbf{V}^T$ and $\mathbf{A}^T = \mathbf{W}^T(\mathbf{X}\mathbf{h})^{-1}$.

Eventually, the final optimization of Eq.(11) is $\mathbf{A}' = (\mathbf{W}^T(\mathbf{X}\mathbf{h})^{-1})^T$. Algorithm 1 details the whole process to update \mathbf{A}' .

Algorithm 1: Power iteration method for solving Eq.(11)

Input: Data \mathbf{X} , Source domain label \mathbf{Y}_S , MMD matrix \mathbf{M}^* , fixed matrix \mathbf{G} , \mathbf{H} , regularization parameters β and α

- 1 **1:** Initialize \mathbf{W} and \mathbf{B}' as introduced in **Step.iii**;
- 2 **2:** Update \mathbf{Z} , \mathbf{A}^T , \mathbf{W}^T ;
- 3 **if** *Non convergence* **then**
- 4 (i) Do **Step.iv**;
- 5 (ii) Do **Step.v**;
- 6 **else**
- 7 **break**;
- 8 **3:** Update \mathbf{A}' via solving $\mathbf{A}' = (\mathbf{W}(\mathbf{X}\mathbf{h})^{-1})^T$

Output: $\mathbf{A} \leftarrow \mathbf{A}'$

Step.5 (Update of \mathbf{Y}) The label matrix \mathbf{Y} contains two parts: true labels $\mathbf{Y}_S = \{y_1, \dots, y_{n_s}\}^T \in \mathbb{R}^{n_s \times C}$, and pseudo labels $\mathbf{Y}_T = \{y_{n_s+1}, \dots, y_{n_s+n_t}\}^T \in \mathbb{R}^{n_t \times C}$. Our aim is to iteratively refine the latter ones. Given fixed \mathbf{A} , \mathbf{e} and \mathbf{M}^* , each $y_i \in \mathbf{Y}_T$ can be updated by solving the following problem:

$$y_i' = \arg \min_{y_i \geq 0, y_i^T \mathbf{1} = 1} \|\mathbf{X}^T\mathbf{A} - y_i + \mathbf{e}\|_F^2 \quad (16)$$

Using Lagrangian multipliers method, the final optimal solution of y_i is

$$y_i' = (\mathbf{A}^T\mathbf{x}_i + \mathbf{e} + \partial) \quad (17)$$

where ∂ is coefficient of Lagrangian constraint $y_i^T \mathbf{1} - 1 = 0$, which can be obtained by solving $y_i^T \mathbf{1} = 1$.

Step.6 (Update of \mathbf{M}^*) With the labeled source domain data $\mathbf{A}^T\mathbf{X}_S$ and the labels inferred on the target domain data $\mathbf{A}^T\mathbf{X}_T$ as in **Step.5**, we can update \mathbf{M}^* as

$$\mathbf{M}^* = \mathbf{M}_0 + \sum_{c=1}^C (\mathbf{M}_c) - \mathbf{M}_{REP} \quad (18)$$

where \mathbf{M}_c and \mathbf{M}_{REP} are defined in Eq.(19).

The complete learning algorithm is summarized in Algorithm 2 - **DOLL-DA**.

Algorithm 2: Discriminative Label Consistent Domain Adaptation (DOLL-DA)

Input: Data \mathbf{X} , Source domain label \mathbf{Y}_S , subspace dimension k , iterations T , regularization parameters β and α

- 1 **1:** Initialize $\mathbf{M}^* = \mathbf{M}_0$ as defined in Eq.(3) ;
- 2 **2:** Initialize \mathbf{A} by solving Eq.(8); ($t := 0$)
- 3 **while** $\sim isempty(\mathbf{X}, \mathbf{Y}_S)$ and $t < T$ **do**
- 4 **3:** Update \mathbf{M}^* by solving Eq.(18)
- 5 **4:** Update \mathbf{e} by solving Eq.(10)
- 6 **5:** Update \mathbf{A} ; ($t_1 := 0$.)
- 7 **if** $t_1 < T$ **then**
- 8 (i) Initialize \mathbf{G} as an identity matrix;
- 9 (ii) Update \mathbf{A} by solving **Algorithm1**;
- 10 (iii) Update \mathbf{G} by calculating
- 11
$$g_{jj} = (\sum_{i=1}^k \sqrt{\|\mathbf{a}^i\|_2^2 + \varepsilon}) \div (\sqrt{\|\mathbf{a}^j\|_2^2 + \varepsilon});$$
- 12 (iv) $t_1 = t_1 + 1$;
- 13 **else**
- 14 **break**;
- 15 **6:** Update \mathbf{Y} by solving Eq.(17)
- 16 **7:** Update pseudo target labels
- $$\mathbf{Y}_{\mathcal{T}}^{(T)} = \mathbf{Y}[:, (n_s + 1) : (n_s + n_t)];$$
- 8:** $t = (t + 1)$;

Output: \mathbf{A} , $\mathbf{Z} = \mathbf{A}^T \mathbf{X}$, \mathbf{Y}

D. Kernelization

The proposed **DOLL-DA** method is extended to nonlinear problems in a Reproducing Kernel Hilbert Space via the kernel mapping $\phi : x \rightarrow \phi(x)$, or $\phi(\mathbf{X}) : [\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n)]$, and the kernel matrix $\mathbf{K} = \phi(\mathbf{X})^T \phi(\mathbf{X}) \in \mathbb{R}^{n \times n}$. We utilize the representer theorem to formulate the Kernel **DOLL-DA** as:

$$\begin{aligned} \min_{\mathbf{A}, \mathbf{e}, \mathbf{Y}_{\mathcal{U}} \geq 0, \mathbf{Y}_{\mathcal{U}} \mathbf{1} = 1} & (tr(\mathbf{A}^T \mathbf{K} \mathbf{M}^* \mathbf{K}^T \mathbf{A}) + \alpha \|\mathbf{A}\|_F^2 + \beta \|\mathbf{A}\|_{2,1}^2) \\ & + \|\mathbf{K}^T \mathbf{A} + \mathbf{1e}^T - \mathbf{Y}\|_F^2 \\ \text{st. } \mathbf{M}^* &= \mathbf{M}_0 + \sum_{c=1}^C (\mathbf{M}_c) - \mathbf{M}_{REP}, \mathbf{Y} \geq 0, \mathbf{Y} \mathbf{1} = 1, \\ \mathbf{A}^T \mathbf{K} \mathbf{H} \mathbf{K}^T \mathbf{A} &= \mathbf{I}, (\mathbf{Y}_{i=c} \bullet \mathbf{Y}_{i \neq c}) = 0 \end{aligned} \tag{19}$$

E. Time Complexity Analysis

Given n the number of data samples including both the source and target domain and l the feature dimension, we denote by t, t_1 the number of iterations with $t, t_1 \prec \min(l, n)$. The major computational burden of the proposed Algorithm.2 lies in Step 2, 3 and 5 as sketched in Sect.III-C. In Step 2, the singular value decomposition(SVD) is computed on a $n * n$ matrix, its computational complexity is $O(n^3)$. Step 3. constructs the \mathbf{M}_{cvd} matrix, whose computational complexity is $O(4Cn^2)$ with C the number of classes, *i.e.*, the class cardinality. Step 4.(ii) makes use of SVD to solve optimization on a $n * k$ size matrix as introduced in step.(iv) of Algorithm.1, its computational complexity is thus $O(n^2k + nk^2 + k^3)$. In Step.5 and Step.6, $O(nk)$ operations are required for all other lines. Therefore, the overall computational complexity of the proposed Algorithm - **DOLL-DA** is $O(n^3 + 4tCn^2 + tnk + t_1(n^2k + nk^2 + k^3))$.

IV. EXPERIMENTS

A. Benchmarks and Features

As illustrated in Fig.5, USPS[23]+MNIST[28], COIL20[33], PIE[33], office+Caltech[33], Office-Home[60] and SVHN-MNIST[4] are standard benchmarks for the purpose of evaluation and comparison with state-of-the-art in DA. In this paper, we follow the data preparation as most previous works[15], [8], [40], [4], [30] do. We construct 49 datasets for different image classification tasks.

Office+Caltech consists of 2533 images of 10 categories (8 to 151 images per category per domain)[15]. These images come from four domains: (A) AMAZON, (D) DSLR, (W) WEBCAM, and (C) CALTECH. AMAZON images were acquired in a controlled environment with studio lighting. DSLR consists of high resolution images captured by a digital SLR camera in a home environment under natural lighting. WEBCAM images were acquired in a similar environment to DSLR, but with a low-resolution webcam. CALTECH images were collected from Google Images.

We use two types of image features extracted from these datasets, *i.e.*, **SURF** and **DeCAF6**, that are publicly available. The **SURF**[16] features are *shallow* features extracted and quantized into an 800-bin histogram using a codebook computed with K-means on a subset of images from Amazon. The resultant histograms are further standardized by z-score. The **Deep Convolutional Activation Features (DeCAF6)**[11] are *deep* features computed as in **AELM**[59] which makes use of VLFeat MatConvNet library with different pretrained CNN models, including in particular the Caffe implementation of **AlexNet**[27] trained on the ImageNet dataset. The outputs from the 6th layer are used as *deep* features, leading to 4096 dimensional **DeCAF6** features. In this experiment, we denote the dataset **Amazon**, **Webcam**, **DSLR**, and **Caltech-256** as **A**, **W**, **D**, and **C**, respectively. The arrow “ \rightarrow ” is proposed to denote the direction from “source” to “target”. For example, “**W** \rightarrow **D**” means the Webcam image dataset is considered as the labeled *source* domain whereas the DSLR image dataset the unlabeled *target* domain.

USPS+MNIST shares 10 common digit categories from two subsets, namely USPS and MNIST, but with very different data distributions (see Fig.5). We construct a first DA task **USPS vs MNIST** by randomly sampling first 1,800 images in USPS to form the source data, and then 2,000 images in MNIST to form the target data. Then, we switch the source/target pair to get another DA task, *i.e.*, **MNIST vs USPS**. We uniformly rescale all images to size 16×16 , and represent each one by a feature vector encoding the gray-scale pixel values. We also extract deep feature from softmax layer[47] of LeNet[28] architecture, leading to a 10 dimensional feature. Thus the source and target domain data share the same feature space. As a result, we have defined two cross-domain DA tasks, namely **USPS** \rightarrow **MNIST** and **MNIST** \rightarrow **USPS**.

COIL20 contains 20 objects with 1440 images (Fig.5). The images of each object were taken in varying its pose about 5 degrees, resulting in 72 poses per object. Each image has a resolution of 32×32 pixels and 256 gray levels per pixel. In this experiment, we partition the dataset into two subsets,

namely COIL 1 and COIL 2 [63]. COIL 1 contains all images taken within the directions in $[0^{\circ}, 85^{\circ}] \cup [180^{\circ}, 265^{\circ}]$ (quadrants 1 and 3), resulting in 720 images. COIL 2 contains all images taken in the directions within $[90^{\circ}, 175^{\circ}] \cup [270^{\circ}, 355^{\circ}]$ (quadrants 2 and 4) and thus the number of images is also 720. In this way, we construct two subsets with relatively different distributions. In this experiment, the COIL20 dataset with 20 classes is split into two DA tasks, *i.e.*, $COIL1 \rightarrow COIL2$ and $COIL2 \rightarrow COIL1$.

PIE face database consists of 68 subjects with each under 21 various illumination conditions [8], [33]. We adopt five pose subsets: C05, C07, C09, C27, C29, which provide a rich basis for domain adaptation, that is, we can choose one pose as the source and any remaining one as the target. Therefore, we obtain $5 \times 4 = 20$ different source/target combinations. Finally, we combine all five poses together to form a single dataset for large-scale transfer learning experiment. We crop all images to 32×32 and only adopt the pixel values as the input. Finally, with different face poses, of which five subsets are selected, denoted as PIE1, PIE2, *etc.*, resulting in $5 \times 4 = 20$ DA tasks, *i.e.*, $PIE1$ vs $PIE2 \dots PIE5$ vs $PIE4$, respectively.

Office-Home dataset as shown in Fig.5 is a novel DA dataset recently introduced in [60]. This dataset contains 4 domains. Each domain contains 65 categories. This dataset is used in a similar manner as the **Office+Caltech** dataset. From the 4 domains, *i.e.*, the Art (Ar), Clipart (Cl), Product (Pr) and Real-World (Rw), we generate 12 DA tasks, namely $Ar \rightarrow Cl \dots Pr \rightarrow Rw$, respectively. **DeCAF6** features are extracted to evaluate the performance of the proposed DA algorithms on this dataset.

SVHN-MNIST contains the MNIST dataset as introduced in **USPS-MNIST** but also the Street View House Numbers (SVHN) which is a collection of house numbers collected from Google street view images (see Fig.5). SVHN is quite distinct from the dataset of handwriting digits, *i.e.*, digits in MNIST. Moreover, all the 2 domains are quite large, each having at least 60k samples over 10 classes. We propose to make use of the LeNet architecture [28] and domain classifier as introduced in [47] to extract features for our DA tasks.

B. Baseline Methods

The proposed **DOLL-DA** method is compared with **thirty-two** methods of the literature, including deep learning-based approaches for unsupervised domain adaption. They are:

- **Shallow methods:** (1) 1-Nearest Neighbor Classifier(NN); (2) Principal Component Analysis (PCA); (3) **GFK** [16]; (4) **TCA** [42]; (5) **TSL** [52]; (6) **JDA** [33]; (7) **ELM** [59]; (8) **AELM** [59]; (9) **SA** [12]; (10) **mSDA** [5]; (11) **TJM** [34]; (12) **RTML** [8]; (13) **SCA** [15]; (14) **CDML** [61]; (15) **LTSL** [51]; (16) **LRSR** [63]; (17) **KPCA** [49]; (18) **JGSA** [64]; (19) **CORAL** [54]; (20) **RVDLR** [24]; (21) **LPJT** [29]; (22) **DGA-DA**[39].
- **Deep methods:** (23) **AlexNet** [27]; (24) **DAH** [60]; (25) **DANN** [14]; (26) **ADDA** [57]; (27) **LTRU** [50]; (28) **ATU** [48]; (29) **BSWD** [47]; (30) **DSN** [4]; (31) **DDC** [58]; (32) **DAN** [32].

Direct comparison of the proposed **DOLL-DA** using shallow features against these **DL**-based DA approaches could be

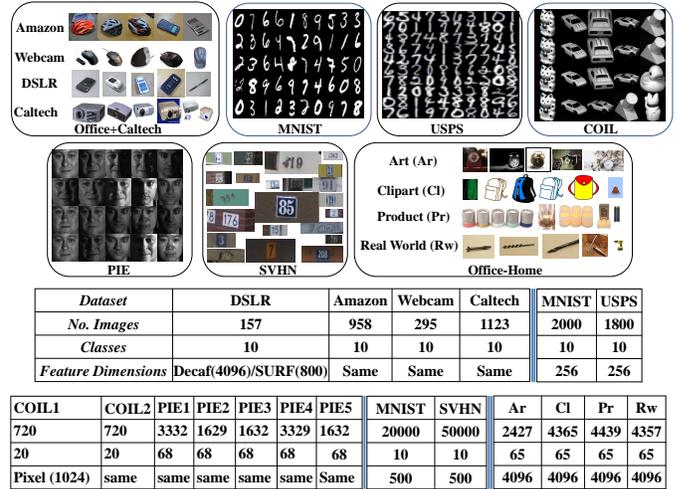


Fig. 5: Sample images from eight datasets used in our experiments. Each dataset represents a different domain. The Office dataset in Office+Caltech contains three sub-datasets, namely DSLR, Amazon and Webcam.

unfair. However, in order to give an idea on the performance gap between shallow and deep DA methods, we still compare their results with those of our shallow **DOLL-DA**. For the purpose of fair comparison, we follow the experiment settings of **DGA-DA**, **JGSA** and **BSWD**, and apply DeCAF6 as the input features for some methods to be evaluated. Whenever possible, the reported performance scores of the **thirty-two** methods of the literature are directly collected from their original papers or previous research [57], [59], [29], [15], [47], [64], [39]. They are assumed to be their *best* performance.

C. Experimental Setup

For the problem of domain adaptation, it is not possible to tune a set of optimal hyper-parameters, given the fact that the target domain has no labeled data. Following the setting of previous research [39], [33], [63], we also evaluate the proposed **DOLL-DA** by empirically searching in the parameter space for the *optimal* settings. Specifically, the proposed **DOLL-DA** method has three hyper-parameters, *i.e.*, the subspace dimension k , regularization parameters β and α . In our experiments, we set $k = 300$ and 1) $\beta = 0.1$, and $\alpha = 1$ for **USPS**, **MNIST**, **COIL20** and **PIE**, 2) $\beta = 1$, $\alpha = 1$ for **Office+Caltech**, **SVHN-MNIST** and **Office-Home**.

In our experiment, *accuracy* on the test dataset as defined by Eq.(20) is the performance measurement. It is widely used in literature, *e.g.*, [32], [40], [33], [63], *etc.*

$$Accuracy = \frac{|x:x \in D_T \wedge \hat{y}(x) = y(x)|}{|x:x \in D_T|} \quad (20)$$

where D_T is the target domain treated as test data, $\hat{y}(x)$ is the predicted label and $y(x)$ is the ground truth label for the test data x .

The core model of the proposed **DOLL-DA** method is built on **JDA**, but adds up two optimization terms, namely discriminative data distribution alignment (**DDA**) term as defined in Eq.(4, 5), and Noise Robust Sparse Orthogonal

Label Regression (**NRS_OLR**) term as defined in Eq.(6). For further insight into the proposed **DOLL-DA** and the rational *w.r.t* the **DDA** and **NRS_OLR** term, respectively, we derive from Eq.(19) three additional partial models, namely **OLR**, **CDDA+** and **JOLR-DA**:

- **OLR**: In this setting, the partial model only makes use of the **NRS_OLR** term, *i.e.*, noise robust sparse orthogonal label regression term as in Eq.(6) and ignores the rest of our final model, *i.e.*, data distribution alignment term as defined in Eq.(3) as well as *Discriminative MMD constraint* terms as defined in Eq.(4, 5). This partial model amounts to make use of a particular classifier, *i.e.*, noise robust sparse orthogonal label regression, and enables to quantify the importance of data distribution alignment in DA when contrasted with the baseline **JDA**;
- **CDDA+**: In this setting, the regularization term of **NRS_OLR** (Eq.(6)) is simply replaced by the Nearest Neighbor (NN) predictor. This correspond to our final DA model as defined in Eq.(8) and Eq.(18) which only make use of the *Discriminative MMD constraint* terms but without the **NRS_OLR** term as defined by Eq.(6). In comparison with **CDDA**, *i.e.*, Close yet Discriminative DA, the partial model already studied in our former **DGA-DA** method [39] which demonstrated the effectiveness of *discriminative force*(**RF**) term in **DA**, **CDDA+** includes the improved **RF** term which adds the **RF** within the source domain as defined by Eq.(5) to the across domain **RF** defined in Eq.(4) as in **CDDA**. This partial model makes it possible to emphasize the interest of the joint optimization of the improved **RF** and **NRS_OLR** terms in contrasting **DOLL-DA** with **CDDA+**;
- **JOLR-DA**: In this setting, the partial model of the proposed method only cares about data distribution alignment as in **JDA** as well as the newly introduced **NRS_OLR** term but ignores the *discriminative force* as defined in Sect.III-B2 and Sect.III-B3. In studying **JOLR-DA**, we aim to highlight: 1) the contribution of **NRS_OLR** in regularizing the MMD constraints in contrasting **DOLL-DA** with **JOLR-DA**; 2) the effectiveness of **NRS_OLR** and the proposed joint optimization strategy, in confronting **JOLR-DA** with **JDA**.
- **DOLL-DA**: This setting correspond to our full final model as defined in Eq.(19). It thus contains both **CDDA+** as defined in Eq.(3, 4, 5) and the **NRS_OLR** term as defined by Eq.(6) in sect. III-B4.

D. Experimental Results and Discussion

1) *Experiments on the CMU PIE Data Set*: The CMU PIE dataset is a large face dataset featuring both illumination and pose variations. Fig.6 synthesizes the experimental results for **DA** using this dataset, where top results are highlighted in red color. As expected, without data distribution alignment,**OLR**, with 55.88% average accuracy, performs worse than the base line **JDA** with 60.24% average accuracy. In accounting for the discriminative force in data distribution alignment, **CDDA+** improves over **JDA** by 3 points and achieves 63.22% average accuracy. In adding noise robust sparse orthogonal

label regression to **JDA**, **JOLR-DA** achieves 69.96% average accuracy and improves over **JDA** by 9 points, thereby demonstrating the effectiveness of the **NRS_OLR** term. Now our final model, **DOLL-DA**, with 82.50% average accuracy, achieves the state of the art performance on this dataset and improves over the baseline **JDA** by 22 points and the former state of the art DA method, *i.e.*, **DGA-DA**, by a large margin of 17 points, thereby demonstrates with force the interest of joint optimization of discriminative data distribution terms and the **NRS_OLR** term.

2) *Experiments on the COIL 20 Dataset*: The **COIL** dataset (see fig.5) features the challenge of pose variations between the source and target domain. Fig.7 reports the experimental results on the **COIL** dataset and displays similar patterns as those on the **PIE** dataset. **OLR** performs worse than **JDA** which is improved by **CDDA+** and **JOLR-DA**. Finally, **DOLL-DA** achieves 96.84% and further improves **CDDA+** and **JOLR-DA** by 4 and 3 points, respectively. It is interesting to note that **DGA-DA** achieves 100% average accuracy on this dataset and thereby outperforms **DOLL-DA** by 3.16 points. However, **DGA-DA**'s outstanding performance is mainly related to the two particularities of the **COIL** 20 dataset. Indeed, the **COIL** dataset synthesizes 20 objects as foreground in varying their pose while the background is *pure black*, therefore each sub-domain contains a single object and is naturally distributed within a specific manifold. Furthermore, the **COIL** 20 dataset merely contains 20 classes in contrast with 68 classes in the **PIE** dataset, thereby its classes are much more separated than those in **PIE**. As a matter of fact, the most simple baseline **NN**, *i.e.*, Nearest Neighbor, already achieves a high average accuracy of 83.20% on **COIL** 20 while it only displays 34.76% average accuracy on **PIE**. As a result, in explicitly modeling the hidden data manifold structure through a Laplace graph, **DGA-DA** performs better label inference than **DOLL-DA**.

3) *Experiments on the Office-Home Dataset*: As introduced in **DAH**[60], **Office-Home** is a novel challenging benchmark for the DA task. It contains 4 very different domains with 65 object categories, thereby generating 12 different DA tasks. Fig.8 synthesizes the performance of the proposed DA methods with DeCAF6 features in comparison with the state of the art methods. Both **DAH** and **DAN** are deep DA methods and make use of multi-kernel **MMD** in aligning the source and target domain distributions. With 43.46% and 45.54% average accuracy, respectively, they surpass **JDA** with a margin up to 8 points. In extending **JDA** with repulsive force terms, **CDDA+** slightly improves **JDA** by 0.19 point. However, thanks to the **NRS_OLR** term, **JOLR-DA** displays 44.48% average accuracy and improves **JDA** by 7 points whereas the proposed full model, **DOLL-DA**, achieves the novel state of the art performance with 48.23% average accuracy and improves **CDDA+** by 11.07 points, **JOLR-DA** by roughly 4 points, and the former state of the art performance achieved by **DAH** with a margin of 2.69 points.

4) *Experiments on the USPS+MNIST Data Set*: The **USPS+MNIST** dataset displays different writing styles between source and target. In Fig.9, the left columns of the red vertical bar report the experimental results using shallow

	PIE 1 5→7	PIE 2 5→9	PIE 3 5→27	PIE 4 5→29	PIE 5 7→5	PIE 6 7→9	PIE 7 7→27	PIE 8 7→29	PIE 9 9→5	PIE 10 9→7	PIE 11 9→27	PIE 12 9→29	PIE 13 27→5	PIE 14 27→7	PIE 15 27→9	PIE 16 27→29	PIE 17 29→5	PIE 18 29→7	PIE 19 29→9	PIE 20 29→27	Average
■ NN	26.09	26.59	30.67	16.67	24.49	46.63	54.07	26.53	21.37	41.01	46.53	26.23	32.95	62.68	73.22	37.19	18.49	24.19	28.31	31.24	34.76
■ PCA	24.80	25.18	29.26	16.30	24.22	45.53	53.35	25.43	20.95	40.45	46.14	25.31	31.96	60.96	72.18	35.11	18.85	23.39	27.21	30.34	33.85
■ GFK	26.15	27.27	31.15	17.59	25.24	47.37	54.25	27.08	21.82	43.16	46.41	26.78	34.24	62.92	73.35	37.38	20.35	24.62	28.49	31.33	35.35
■ CDML	53.22	53.12	80.12	48.23	52.39	54.23	68.36	37.34	43.54	54.87	62.76	38.21	75.12	80.53	83.72	52.78	27.34	30.82	36.34	40.61	53.69
■ RTML	60.12	55.21	85.19	52.98	58.13	63.92	76.16	40.38	53.12	58.67	69.81	42.13	81.12	83.92	89.51	56.26	29.11	33.28	39.85	47.13	58.80
■ LTSL	22.96	20.65	31.81	12.07	18.25	16.05	45.15	17.52	22.36	20.26	57.34	24.57	51.20	70.10	72.00	48.28	13.06	21.61	17.03	29.59	31.59
■ mSDA	28.35	26.91	30.39	21.76	28.27	44.19	55.39	28.08	24.83	42.59	50.25	27.83	32.89	63.01	74.70	34.81	25.85	26.33	28.63	32.98	36.41
■ RDALR	40.76	41.79	59.63	29.35	41.81	51.47	64.73	33.70	34.69	47.70	56.23	33.15	55.64	67.83	75.86	40.26	26.98	29.90	29.90	33.64	44.75
■ LRSR	65.87	64.09	82.03	54.90	45.04	53.49	71.43	47.97	52.49	55.56	77.50	54.11	81.54	85.39	82.23	72.61	52.19	49.41	58.45	64.31	63.53
■ TSL	44.08	47.49	62.78	36.15	46.28	57.60	71.43	35.66	36.94	47.02	59.45	36.34	63.66	72.68	83.52	44.79	33.28	34.13	36.58	38.75	49.43
■ TCA	40.76	41.79	59.63	29.35	41.81	51.47	64.73	33.70	34.69	47.70	56.23	33.15	55.64	67.83	75.86	40.26	26.98	29.90	29.90	33.64	44.75
■ JGSA	55.13	53.19	75.01	50.49	64.83	60.91	78.49	51.59	61.10	62.31	77.80	59.87	77.97	79.80	77.08	64.52	58.82	52.92	60.23	67.26	64.47
■ LDADA	52.30	51.90	78.82	41.61	52.67	60.85	70.56	39.28	52.26	51.44	61.58	49.71	74.19	78.08	84.74	45.71	46.01	37.08	39.71	49.14	55.88
■ DGA-DA	65.32	62.81	83.54	56.07	63.69	61.27	82.37	46.63	56.72	61.26	77.83	44.24	81.84	85.27	90.95	53.80	57.44	53.84	55.27	61.82	65.10
■ OLR	43.03	44.55	70.62	43.26	50.87	49.33	71.88	40.50	53.96	44.69	68.58	48.04	80.94	70.84	69.49	60.60	54.47	38.49	50.86	62.63	55.88
■ JDA	58.81	54.23	84.50	49.75	57.62	62.93	75.82	39.89	50.96	57.95	68.45	39.95	80.58	82.63	87.25	54.66	46.46	42.05	53.31	57.01	60.24
■ JOLR-DA	60.96	64.03	87.47	50.24	70.83	65.42	78.96	59.07	69.59	62.19	82.13	67.15	86.25	83.06	78.62	67.22	62.59	63.72	70.28	69.37	69.96
■ CDDA+	60.22	59.80	83.48	53.14	62.33	64.64	79.90	44.00	58.46	59.73	78.48	47.24	83.10	82.26	86.64	58.33	48.02	46.62	52.02	55.99	63.22
■ DOLL-DA	79.56	78.59	91.53	75.94	82.14	75.61	91.86	72.90	78.69	78.59	86.39	72.12	93.50	90.61	84.56	86.97	80.52	86.32	79.60	83.90	82.50

Fig. 6: Accuracy% on the PIE Images Dataset.

	NN	PCA	GFK	TSL	LTSL	LRSR	TCA	DGA-DA	OLR	JDA	JOLR-DA	CDDA+	DOLL-DA
■ C1 → C2	83.61	84.72	72.50	88.06	75.69	88.61	88.47	100	85.2	89.31	92.46	91.58	96.36
■ C2 → C1	82.78	84.03	74.17	87.92	72.22	89.17	85.83	100	80.1	88.47	94.24	93.92	97.32
■ Average	83.20	84.38	73.34	87.99	73.96	88.89	87.15	100	82.65	88.89	93.35	92.75	96.84

Fig. 7: Accuracy% on the COIL Images Dataset.

	Ar→Cl	Ar→Pr	Ar→Rw	Cl→Ar	Cl→Pr	Cl→Rw	Pr→Ar	Pr→Cl	Pr→Rw	Rw→Ar	Rw→Cl	Rw→Pr	Average
— GFK	21.60	37.72	38.83	21.63	34.94	34.20	24.52	25.73	42.92	32.88	28.96	50.89	32.40
— TCA	19.93	32.08	35.71	19.00	31.36	31.74	21.92	23.64	42.12	30.74	27.15	48.68	30.34
— CORAL	27.10	36.16	44.32	26.08	40.03	40.33	27.77	30.54	50.61	38.48	36.36	57.11	37.91
— DAN	30.66	42.17	54.13	32.83	47.59	49.78	29.07	34.05	56.70	43.58	38.25	62.73	43.46
— DANN	33.33	42.96	54.42	32.26	49.13	49.76	30.49	38.14	56.76	44.71	42.66	64.65	44.94
— DAH	31.64	40.75	51.73	34.69	51.93	52.79	29.91	39.63	60.71	44.99	45.13	62.54	45.54
— DGA-DA	33.27	40.02	49.05	32.75	50.07	46.50	30.79	36.33	59.83	47.26	44.30	62.91	44.42
— OLR	20.32	30.77	29.73	25.36	35.69	30.29	32.21	25.46	49.99	34.08	23.02	45.67	31.88
— JDA	25.34	35.98	42.94	24.52	40.19	40.09	25.96	32.72	49.25	35.10	35.35	55.35	36.97
— JOLR-DA	35.40	46.74	50.98	34.24	52.38	49.83	28.22	35.68	50.13	46.39	45.32	58.46	44.48
— CDDA+	25.48	36.69	41.96	24.35	42.26	40.56	24.59	31.25	49.21	34.28	35.68	59.71	37.16
— DOLL-DA	87.11	49.07	52.47	34.78	54.63	56.98	31.97	40.37	60.23	48.59	48.24	64.34	48.23

Fig. 8: Accuracy% on the Office-Home Images Dataset.

features, whereas the right columns of the red bar display the results of the methods using deep features. Once more, the proposed **DOLL-DA** along with its partial models display the same behavior as in the previous experiences. Using shallow features, **CDDA+** and **JOLR-DA** demonstrate the effectiveness of the discriminative force terms and **NRS_OLR** term. With 69.14% and 71.59% average accuracy, **CDDA+** and **JOLR-DA** improve the baseline **JDA** by 5 and 8 points, respectively. When jointly optimizing the discriminative force terms and **NRS_OLR** term, the proposed final model, **DOLL-DA**, further boosts the performance of the baseline **JDA** with a margin as high as 14 points and set a novel state of the art performance with 77.82% average accuracy when shallow features are used.

In right part of the red vertical bar in Fig.9, we compare our approach with deep network-based methods, *i.e.*, **ADDA** and **DANN**, which search a common latent subspace for minimizing the source and target representation divergence through the popular adversarial learning, and improve the performance of traditional **DA** methods. As can be seen in Fig.9, the proposed **DOLL-DA** using shallow features already surpasses **DANN** by 2 points. When using deep features, *i.e.*, those from LeNet [28] as in [47], **DOLL-DA** demonstrates its

effectiveness once more and display a novel state of the art performance of 92.05% average accuracy.

5) *Large-Scale Experiments on the SVHN-MNIST Data Set:* Different with the other datasets, SVHN-MNIST is composed of 50k and 20k digit images from two very different domains, respectively, thereby generating a large scale **DA** benchmark. Following the experiment setting of previous research[39], Fig.10 shows the experimental results. Our proposed **DOLL-DA** along with its partial models display the same patterns on this large scale benchmark as in the previous experiments. In ignoring the domain shifts, **OLR** achieves a poor performance. Both **CDDA+** and **JOLR-DA** improve **JDA** with a large margin. In accounting jointly for the discriminative force and **NRS_OLR**, **DOLL-DA** boosts the baseline **JDA** by 22 points and surpasses **DGA-DA**, the former state of the art DA method on this dataset, by 4.5 points.

6) *Experiments on the Office+Caltech-256 Data Sets:* Fig.11 and Fig.12 synthesize the experimental results in comparison with the state of the art when deep features (*i.e.*, DeCAF6 features) and classic shallow features (*i.e.*, SURF features) are used, respectively.

- As can be seen in Fig.12, using shallow features, **CDDA+** and **JOLR-DA** improve the baseline **JDA** by 2 and 1 points, respectively, thanks to the discriminative force term and **NRS_OLR** term introduced into the baseline model. Taking them together, our final model **DOLL-DA** further improves **JDA** by 4 points and achieves 51.180 average accuracy which is in par with the previous state the art method, namely **JGSA**, with its 50.58% average accuracy. It is worth noting that, **JGSA** also suggests aligning data both statistically, discriminatively and geometrically, and corroborates data geometry aware **DA** approach, and achieves very good performance on this dataset.
- Fig.11 compares the proposed **DA** method using deep features *w.r.t.* the state of the art, in particular end-to-end deep learning-based **DA** methods. As can be seen in Fig.11, the use of deep features has enabled impressive accuracy improvement over shallow features. Simple baseline methods, *e.g.*, **NN**, **PCA**, see their accuracy

	NN	PCA	GFK	TSL	KPC A	SCA	TJM	ELM	SA	mSD A	AEL M	JGSA	TCA	DGA-DA	OLR	JDA	JOLR-DA	CDD A+	DOL L-DA	DAN N	ADD A	DOL L-DA
■ USPS→MNIST	44.70	44.95	46.45	53.75	42.55	48.00	52.52	57.70	40.15	43.20	57.77	68.15	51.05	70.75	50.80	59.65	67.52	62.05	71.50	73.00	90.10	86.90
■ MNIST→USPS	65.94	66.22	67.22	66.06	62.61	65.11	63.28	61.11	48.22	66.94	62.33	80.44	56.28	82.33	63.50	67.28	75.65	76.22	83.05	77.10	89.40	97.20
■ Average	55.32	55.59	56.84	59.91	52.58	56.56	57.90	59.41	44.19	55.07	60.05	74.30	53.67	76.54	57.15	63.47	71.59	69.14	77.82	75.05	89.75	92.05

Fig. 9: Accuracy% on the USPS+MNIST Images Dataset.

	ADD A	DAN N	BSW D	SA	LTR U	ATU	DSN	DGA-DA	OLR	JDA	JOLR-DA	CDD A+	DOL L-DA
■ SHVN → MNIST	76.00	80.70	82.80	59.32	78.8	86.20	82.7	83.3	59.9	65.8	78.59	76.2	87.8

Fig. 10: Accuracy% on the SVHN-MNIST Images Dataset.

	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
— PCA	85.60	66.10	74.50	70.30	57.20	64.90	60.30	62.50	98.70	52.00	62.70	89.10	70.40
— NN	87.05	72.20	80.89	78.54	77.31	80.25	68.21	73.07	100.0	70.08	75.89	97.97	80.12
— ELM	89.07	70.51	78.98	79.61	74.58	80.25	70.61	75.37	100.0	68.21	80.79	98.31	80.52
— GFK	87.27	75.93	83.44	80.32	76.95	80.89	67.76	74.32	100.0	69.10	75.78	98.64	80.87
— SA	87.06	75.59	80.25	79.61	78.31	81.53	68.83	75.16	100.0	69.99	73.49	98.98	80.73
— mSDA	89.67	68.47	82.17	78.81	78.98	79.62	69.46	76.62	100.0	73.29	81.32	98.64	81.42
— TJM	88.10	72.20	74.52	77.65	75.25	82.80	71.42	80.27	100.0	72.57	78.60	98.31	80.97
— AELM	89.46	79.32	81.53	79.96	77.63	85.35	71.24	76.83	100.0	75.60	83.19	98.98	83.25
— RTML	90.20	83.80	88.70	83.10	79.50	83.80	82.90	90.80	100.0	81.60	90.60	98.60	87.80
— SCA	89.46	85.42	87.90	78.81	75.93	85.35	74.80	86.12	100.0	78.09	89.98	98.64	85.88
— JGSA	91.44	86.78	93.63	84.86	81.02	88.54	84.95	90.71	100.0	86.20	91.96	99.66	89.98
— AlexNet	91.90	83.70	87.10	83.00	79.50	87.40	73.00	83.80	100.0	79.00	87.10	97.70	86.10
— DAN	92.00	90.60	89.30	84.10	91.80	91.70	81.20	92.10	100.0	80.30	90.00	98.50	90.10
— DDC	91.90	85.40	88.80	85.00	86.10	89.00	78.00	84.90	100.0	81.10	89.50	98.20	88.20
— MEDA	93.40	95.60	91.10	87.40	88.10	88.10	93.20	99.40	100.0	87.50	93.20	97.60	92.80
--- DGA-DA	91.25	93.56	91.72	85.20	80.98	89.81	86.46	90.81	100.0	86.20	93.11	100.0	90.76
--- OLR	92.07	77.97	84.08	82.55	71.19	75.80	61.80	70.15	100.0	61.18	67.01	98.98	78.57
— JDA	89.70	83.70	86.60	82.20	78.60	80.20	80.50	88.10	100.0	80.10	89.40	98.90	86.50
— JOLR-DA	88.96	78.31	88.54	81.23	88.14	91.08	81.30	90.71	98.09	82.10	91.02	98.86	88.20
— CDDA+	89.46	88.47	89.56	82.66	89.15	89.17	82.55	90.92	100.0	81.92	91.19	100.0	89.59
— DOLL-DA	93.86	93.98	91.71	89.87	91.86	91.72	86.46	93.11	100.0	86.78	95.90	100.0	92.94

Fig. 11: Accuracy% on the Office+Caltech Images with DeCAF6 Features.

	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
■ NN	23.70	25.76	25.48	26.00	29.83	25.48	19.86	22.96	59.24	26.27	28.50	63.39	31.37
■ PCA	36.95	32.54	38.22	34.73	35.59	27.39	26.36	31.00	77.07	29.65	32.05	75.93	39.79
■ GFK	41.02	40.68	38.85	40.25	38.98	36.31	30.72	29.75	80.89	30.28	32.05	75.59	42.95
■ KPCA	40.40	31.53	40.76	37.04	31.86	33.76	27.60	29.44	89.81	27.78	31.00	84.41	42.12
■ SCA	43.74	33.56	39.49	38.29	33.90	34.21	30.63	30.48	92.36	32.32	33.72	88.81	44.29
■ LTSL	25.26	19.32	21.02	16.92	14.58	21.02	34.64	39.56	72.61	35.08	39.67	74.92	34.55
■ LRSR	51.25	38.64	47.13	43.37	36.61	38.85	29.83	34.13	82.80	31.61	33.19	77.29	45.39
■ CDML	47.70	35.60	42.50	40.70	37.30	35.30	31.60	32.40	77.90	32.20	29.40	79.40	43.50
■ mSDA	45.92	37.96	46.49	40.96	40.33	36.30	31.96	33.61	87.26	30.89	35.59	87.45	46.23
■ ELM	49.37	37.79	45.22	40.07	33.56	34.31	31.17	33.85	88.54	28.23	28.50	73.22	43.65
■ AELM	53.13	49.49	50.96	41.14	35.25	36.94	34.11	38.93	89.81	33.83	33.09	80.33	48.08
■ SA	41.02	40.34	47.13	40.16	39.66	35.03	31.17	33.82	85.99	31.26	35.80	84.75	45.51
■ TJM	46.76	38.98	44.59	39.45	42.03	45.22	30.19	29.96	89.17	31.43	32.78	85.42	46.33
■ TSL	44.47	34.24	43.31	37.58	33.90	26.11	29.83	30.27	87.26	28.50	27.56	85.42	42.37
■ TCA	38.20	38.64	41.40	37.76	37.63	33.12	29.30	30.06	87.26	31.70	32.15	86.10	43.61
■ JGSA	53.13	48.47	48.41	41.50	45.08	45.22	33.57	40.81	88.54	30.28	38.73	93.22	50.58
■ DGA-DA	52.09	47.12	45.86	41.32	38.31	38.22	33.30	41.75	89.81	33.66	33.61	93.22	49.02
■ OLR	38.52	26.1	26.52	31.7	30.95	24.2	16.66	20.86	50.89	32.56	30.05	72.54	33.46
■ JDA	44.78	41.69	45.22	39.36	37.97	39.49	31.17	32.78	89.17	31.52	33.09	89.49	46.31
■ JOLR-DA	50.52	48.14	46.50	40.52	45.08	46.50	34.55	39.46	76.43	32.15	32.67	78.31	47.57
■ CDDA+	48.33	44.75	48.41	42.12	41.89	37.58	31.97	37.27	88.08	34.64	33.51	90.51	48.26
■ DOLL-DA	54.18	51.19	47.13	44.88	45.08	46.50	38.29	39.46	86.62	32.41	38.20	90.22	51.18

Fig. 12: Accuracy% on the Office+Caltech Images with SURF-Bow Features.

soared by roughly 40 points, demonstrating the power of deep learning paradigm. Our proposed DA method also takes advantage of this jump and sees its accuracy soared from 48.26 to 89.59 for **CDDA+**, from 47.57 to 88.20 for

JOLR-DA, and from 50.78 to 91.65 for **DOLL-DA**. As for shallow features, **CDDA+** and **JOLR-DA** improve **JDA** by roughly 3 and 2 points, respectively, while the final model, **DOLL-DA**, displays the best average accuracy of 92.94% in par with 92.80% displayed by **MEDA**.

E. Empirical Analysis

Despite the proposed **DOLL-DA** displays state of the art performance over 49 DA tasks through 8 datasets except for the **COIL** dataset where it achieves the second best performance, an important question is how fast the proposed method converges (sect.IV-E2) as well as its sensitivity *w.r.t.* its hyper-parameters (Sect.IV-E1). Additionally, we are curious about how well **DOLL-DA** performs in changing the base classifier (sect.IV-E3) which is required for the formulation of the *repulsive force* terms to enhance data discriminativeness, as well as the **DOLL-DA** leveraging random initialization (sect.IV-E4) instead of using the base classifier for optimization. Furthermore, in analyzing the generalization capacity (sect.IV-E5) of **DOLL-DA**, we evaluate the performance of **DOLL-DA** with unseen target data for detail exploration.

1) *Sensitivity of the proposed DOLL-DA w.r.t. to hyper-parameters*: Three hyper-parameters, namely k , β and α , are introduced in the proposed methods.

Dimensionality analysis: k is the dimension of the searched shared latent feature subspace between the source and target domain. In Fig.13, we plot the classification accuracies of the proposed DA method *w.r.t.* different values of k on the **COIL** and **PIE** datasets. As shown in Fig.13, the subspace dimensionality k varies with $k \in \{20, 40, 60, 80, 100, 150, 200, 300, 400, 450\}$, yet the proposed 3 DA variants, namely, **CDDA+**, **JOLR-DA** and **DOLL-DA**, remain stable *w.r.t.* a wide range of with $k \in \{40 \leq k \leq 400\}$. It can be seen that both **JOLR-DA** and **DOLL-DA** display better robustness than **CDDA+** *w.r.t.* the variation of k , thereby suggesting the effectiveness of the **NRS_OLR** term in the search of the global minimization. Obviously, the larger is k the better the shared subspace can afford complex data distributions, but at the cost of increased computation complexity as highlighted in sect.III-E on time complexity analysis. In our experiments, we set $k = 300$ to balance the efficiency and accuracy.

Sensitivity of α and β : α and β as defined in Eq.(19) are the major hyper-parameters of the proposed **DOLL-DA**. While α aims to regularize the projection matrix A to avoid over-fitting the chosen shared feature subspace *w.r.t.* both source and target domain data, β as expressed in Eq.(6) controls the dimensionality of class dependent data manifold in the searched shared feature subspace, or in other words the sparsity level of the linear combination of the projected

features to regress the class label. We study the sensitivity of the proposed **DOLL-DA** method with a wide range of parameter values, *i.e.*, $\alpha = (0.001, 0.01, 0.1, 1, 10, 20, 50)$ and $\beta = (0.05, 0.1, 1, 5, 10, 100, 200)$. We plot in Fig.14 the results on $D \rightarrow W$, $C \rightarrow D$ and $PIE-27 \rightarrow PIE-5$ datasets on the proposed **DOLL-DA** with k held fixed at 300. As can be seen from Fig.14, the proposed **DOLL-DA** displays its stability as the resultant classification accuracies remain roughly the same despite a wide range of α and β values.

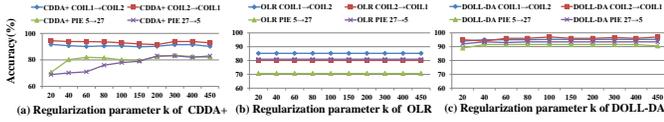


Fig. 13: Sensitivity analysis of the proposed methods: (a) accuracy *w.r.t.* subspace dimension k of **CDDA+**; (b) accuracy *w.r.t.* subspace dimension k of **JOLR-DA**; (c) accuracy *w.r.t.* subspace dimension k of **DOLL-DA**. Three datasets are used, *i.e.*, COIL1, COIL2 and PIE.

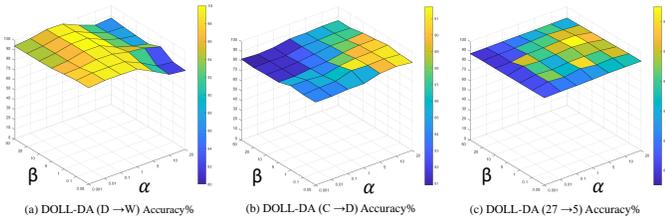


Fig. 14: The classification accuracies of the proposed **DOLL-DA** method vs. the parameters α and β on the selected three cross domains data sets, with k held fixed at 300.

2) **Convergence analysis:** In Fig.15, we further perform convergence analysis of the proposed **DOLL-DA** along with its partial models, *i.e.*, **CDDA+** and **JOLR**, using the **DeCAF6** features on the **Office+Caltech** datasets and pixel value features on the **PIE** dataset. We aim to disclose how fast the proposed methods achieve their best performance *w.r.t.* the number of iterations T . Fig.15 reports 6 cross DA experiments ($C \rightarrow A$, $D \rightarrow W \dots PIE-27 \rightarrow PIE-5$) with the number of iterations $T = (1, 2, 3, 4, 5, 6, 7, 8, 9, 10)$. As shown in Fig.15, **CDDA+**, **JOLR-DA** and **DOLL-DA** converge within 3~5 iterations during optimization, but **JOLR-DA** and **DOLL-DA** seem to converge even faster with a better accuracy, thanks to the **NRS_OLR** term introduced in our DA model.

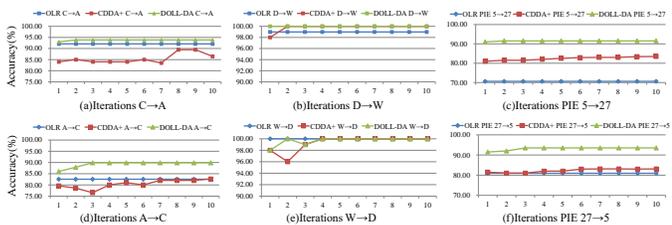


Fig. 15: Convergence analysis using 6 cross-domain image classification tasks on Office+Caltech256 and PIE datasets. (accuracy *w.r.t.* #iterations)

3) **Impact of the base classifier:** The *repulsive force* term across domain as formulated in sect.III-B2 requires using pseudo labels for the target domain data. Therefore, the quality of these pseudo labels could have much impact on the effectiveness of data discriminativeness. In sect.III-C, our model is initialized using pseudo labels in solving Eq.(8), which boils down to **JDA**. The inference of the pseudo labels on the target domain data requires a base classifier trained using the labeled source domain data. We test the sensitivity of the proposed **DOLL-DA** *w.r.t.* the base classifier using two popular classifiers, *i.e.*, NN and SVM. Fig.16 shows that **DOLL-DA-NN**(92.94%) and **DOLL-DA-SVM**(91.51%) achieve almost the same performance. This result suggests that our *repulsive force* term in the proposed **DOLL-DA** displays a certain level of robustness *w.r.t.* the choice of the base classifier.

4) **Random Label Initialization:** Going one step further *w.r.t.* to the experiment in sect.IV-E3, an interesting question is how **DOLL-DA** behaves with randomly initialized pseudo labels at its first iteration as well as its convergence efficiency in this specific experiment setting.

In this setting, **DOLL-DA** makes use of **randomly initialized** labels for the target domain data instead of solving Eq.(8) at its first iteration. As shown in Fig.16, **DOLL-DA** still achieves 90.01% accuracy on the **Office-Caltech256** dataset, thus only slightly below **DOLL-DA-NN**, and converges on average at 3.08 average iterations. This result further supports the robustness of the designed *repulsive force* term regularized by the **NRS_OLR** term in the search of the optimized shared features subspace.

	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
DOLL-DA-NN	93.86	93.98	91.71	89.87	91.86	91.72	86.46	93.11	100.00	86.78	95.90	100.00	92.94
DOLL-DA-SVM	93.86	91.92	87.96	86.59	91.86	90.21	86.46	93.11	98.45	86.21	91.54	100.00	91.51
—Rand Init	93.22	88.47	85.99	88.07	90.85	87.89	83.44	91.89	100.00	82.46	89.04	98.85	90.01
—Iter.(ts)	2.00	2.00	2.00	3.00	2.00	3.00	4.00	3.00	2.00	3.00	8.00	3.00	3.08

Fig. 16: Convergence analysis using 12 cross-domain image classification tasks on Office+Caltech256 datasets with **DeCAF6** Features. (accuracy *w.r.t.* #iterations)

5) **Stability w.r.t. target domain data:** We benchmark the stability of the proposed **DOLL-DA** *w.r.t.* the quantity of target domain data used for training. Specifically, we carry out two additional experiments for 2 DA tasks using the Office-Caltech256 dataset, *i.e.*, $W \rightarrow A$, $D \rightarrow W$, keeping 5% (10%,30%,50%,70%, resp.) of target domain data from being used as auxiliary data in training. Results are reported in Fig.17. As can be seen there, when all target domain data are used in training the proposed **DOLL-DA**, *i.e.*, column 0.00%, **DOLL-DA** displays 93.11% and 100% accuracy for $W \rightarrow A$ and $D \rightarrow W$ DA tasks, respectively. Now when more and more target domain data are kept from being used in training, passing from 5% through 70%, **DOLL-DA** proves quite stable on the $D \rightarrow W$ task but decreases constantly to reach 79.96% on the $W \rightarrow A$ task. However, this result still proves the usefulness of the proposed DA method when only 30% target domain data are used as auxiliary data in training, given the fact that the baseline NN only displays 73.07% accuracy as shown in Fig.11. The performance difference between these

two DA tasks can be explained by their inherent difficulties. The DA task $D \rightarrow W$ is to generalize a classifier trained from the source domain, *i.e.*, labeled images in the DSLR domain, thus with much background, to the target domain, *i.e.*, images in the Webcam domain, much simpler because devoid of the background, and the simple baseline NN already achieves 97.97% accuracy as shown in Fig. 11. On the other side, the DA task $W \rightarrow A$ does exactly the contrary and needs to generalize a classifier trained from the source domain, *i.e.*, labeled images in the Webcam domain, thus without background, to a much more complicated target domain, *i.e.*, images in the Amazon domain with arbitrary background, and the simple baseline NN only achieves 73.07% accuracy as shown in Fig. 11.

$W \rightarrow A$	Unseen	0.00%	5.00%	10.00%	30.00%	50.00%	70.00%
	Accuracy	93.11	92.59	86.07	85.85	84.32	79.96
$D \rightarrow W$	Unseen	0.00%	5.00%	10.00%	30.00%	50.00%	70.00%
	Accuracy	100	100	98.75	98.75	98.75	98.75

Fig. 17: Unseen data sensitivity using 2 cross-domain image classification tasks on Office+Caltech256 datasets with DeCAF6 Features.

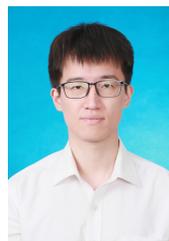
V. CONCLUSION

We have proposed in this paper a novel unsupervised DA method, namely **Discriminative Noise Robust Sparse Orthogonal Label Regression-based Domain Adaptation (DOLL-DA)**, which simultaneously optimizes the three terms of the upper error bound of a learned classifier on the target domain in aligning discriminatively data distributions through a *repulse force* term while orthogonally regressing data labels within the shared feature subspace. Furthermore, the proposed model explicitly accounts for data outliers to avoid negative transfer and introduces the property of sparsity when regressing data labels. Comprehensive experiments using the standard benchmark in DA show the effectiveness of the proposed method which consistently outperform state of the art DA methods. Future work includes embedding of the proposed **DOLL-DA** into the paradigm of deep learning and considers the setting of online learning where target domain data only arrives sequentially one after another.

REFERENCES

- [1] Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman Vaughan. A theory of learning from different domains. *Machine learning*, 79(1):151–175, 2010. 2
- [2] Shai Ben-David, John Blitzer, Koby Crammer, and Fernando Pereira. Analysis of representations for domain adaptation. In *Advances in neural information processing systems*, pages 137–144, 2007. 2
- [3] Karsten M Borgwardt, Arthur Gretton, Malte J Rasch, Hans-Peter Kriegel, Bernhard Schölkopf, and Alex J Smola. Integrating structured biological data by kernel maximum mean discrepancy. *Bioinformatics*, 22(14):e49–e57, 2006. 5
- [4] Konstantinos Bousmalis, George Trigeorgis, Nathan Silberman, Dilip Krishnan, and Dumitru Erhan. Domain separation networks. In *Advances in Neural Information Processing Systems*, pages 343–351, 2016. 4, 9, 10
- [5] Minmin Chen, Zhixiang Eddie Xu, Kilian Q. Weinberger, and Fei Sha. Marginalized denoising autoencoders for domain adaptation. *CoRR*, abs/1206.4683, 2012. 10
- [6] Nicolas Courty, Rémi Flamary, Amaury Habrard, and Alain Rakotomamonjy. Joint distribution optimal transportation for domain adaptation. In *Advances in Neural Information Processing Systems*, pages 3733–3742, 2017. 2, 3
- [7] Nicolas Courty, Rémi Flamary, Devis Tuia, and Alain Rakotomamonjy. Optimal transport for domain adaptation. *IEEE transactions on pattern analysis and machine intelligence*, 39(9):1853–1865, 2017. 2, 3
- [8] Zhengming Ding and Yun Fu. Robust transfer metric learning for image classification. *IEEE Trans. Image Processing*, 26(2):660–670, 2017. 9, 10
- [9] Zhengming Ding and Yun Fu. Robust multiview data analysis through collective low-rank subspace. *IEEE Trans. Neural Netw. Learning Syst.*, 29(5):1986–1997, 2018. 3
- [10] Jeff Donahue, Judy Hoffman, Erik Rodner, Kate Saenko, and Trevor Darrell. Semi-supervised domain adaptation with instance constraints. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 668–675, 2013. 1
- [11] Jeff Donahue, Yangqing Jia, Oriol Vinyals, Judy Hoffman, Ning Zhang, Eric Tzeng, and Trevor Darrell. Decaf: A deep convolutional activation feature for generic visual recognition. In *Proceedings of the 31th International Conference on Machine Learning, ICML 2014, Beijing, China, 21-26 June 2014*, pages 647–655, 2014. 9
- [12] Basura Fernando, Amaury Habrard, Marc Sebban, and Tinne Tuytelaars. Unsupervised visual domain adaptation using subspace alignment. In *IEEE International Conference on Computer Vision, ICCV 2013, Sydney, Australia, December 1-8, 2013*, pages 2960–2967, 2013. 3, 10
- [13] Simone Fiori. Formulation and integration of learning differential equations on the stiefel manifold. *IEEE transactions on neural networks*, 16(6):1697–1701, 2005. 8
- [14] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *The Journal of Machine Learning Research*, 17(1):2096–2030, 2016. 1, 4, 10
- [15] Muhammad Ghifary, David Balduzzi, W. Bastiaan Kleijn, and Mengjie Zhang. Scatter component analysis: A unified framework for domain adaptation and domain generalization. *IEEE Trans. Pattern Anal. Mach. Intell.*, 39(7):1414–1430, 2017. 3, 9, 10
- [16] Boqing Gong, Yuan Shi, Fei Sha, and Kristen Grauman. Geodesic flow kernel for unsupervised domain adaptation. In *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, pages 2066–2073. IEEE, 2012. 1, 9, 10
- [17] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014. 4
- [18] Arthur Gretton, Karsten M Borgwardt, Malte Rasch, Bernhard Schölkopf, and Alex J Smola. A kernel method for the two-sample problem. In *Advances in neural information processing systems*, pages 513–520, 2007. 2
- [19] Mehrtash Harandi, Mathieu Salzmann, and Richard Hartley. Dimensionality reduction on spd manifolds: The emergence of geometry-aware methods. *IEEE transactions on pattern analysis and machine intelligence*, 2017. 8
- [20] Samitha Herath, Mehrtash Harandi, and Fatih Porikli. Learning an invariant hilbert space for domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3845–3854, 2017. 3
- [21] Samitha Herath, Mehrtash Tafazzoli Harandi, and Fatih Porikli. Learning an invariant hilbert space for domain adaptation. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 3956–3965, 2017. 1
- [22] Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu, Phillip Isola, Kate Saenko, Alexei Efros, and Trevor Darrell. CyCADA: Cycle-consistent adversarial domain adaptation. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 1989–1998, Stockholm, Sweden, Stockholm Sweden, 10–15 Jul 2018. PMLR. 4
- [23] Jonathan J. Hull. A database for handwritten text recognition research. *IEEE Trans. Pattern Anal. Mach. Intell.*, 16(5):550–554, 1994. 9
- [24] I-Hong Jhuo, Dong Liu, DT Lee, and Shih-Fu Chang. Robust visual domain adaptation with low-rank reconstruction. In *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, pages 2168–2175. IEEE, 2012. 10
- [25] Alireza Karbalayghareh, Xiaoning Qian, and Edward R Dougherty. Optimal bayesian transfer learning. *IEEE Transactions on Signal Processing*, 66(14):3724–3739, 2018. 3

- [26] Daniel Kifer, Shai Ben-David, and Johannes Gehrke. Detecting change in data streams. In *Proceedings of the Thirtieth international conference on Very large data bases-Volume 30*, pages 180–191. VLDB Endowment, 2004. 2
- [27] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012. 4, 9, 10
- [28] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. 9, 10, 12
- [29] Jingjing Li, Mengmeng Jing, Ke Lu, Lei Zhu, and Heng Tao Shen. Locality preserving joint transfer for domain adaptation. *IEEE Transactions on Image Processing*, 28(12):6103–6115, 2019. 10
- [30] Jian Liang, Ran He, Zhenan Sun, and Tieniu Tan. Aggregating randomized clustering-promoting invariant projections for domain adaptation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018. 1, 9
- [31] Youfa Liu, Weiping Tu, Bo Du, Lefei Zhang, and Dacheng Tao. Homologous component analysis for domain adaptation. *IEEE Transactions on Image Processing*, 29:1074–1089, 2019. 3
- [32] Mingsheng Long, Yue Cao, Jianmin Wang, and Michael I Jordan. Learning transferable features with deep adaptation networks. In *ICML*, pages 97–105, 2015. 4, 10
- [33] Mingsheng Long, Jianmin Wang, Guiguang Ding, Jianguang Sun, and Philip S Yu. Transfer feature learning with joint distribution adaptation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2200–2207, 2013. 1, 3, 9, 10
- [34] Mingsheng Long, Jianmin Wang, Guiguang Ding, Jianguang Sun, and Philip S. Yu. Transfer joint matching for unsupervised domain adaptation. In *2014 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2014, Columbus, OH, USA, June 23-28, 2014*, pages 1410–1417, 2014. 10
- [35] Mingsheng Long, Han Zhu, Jianmin Wang, and Michael I. Jordan. Deep transfer learning with joint adaptation networks. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, pages 2208–2217, 2017. 1, 4
- [36] Hao Lu, Chunhua Shen, Zhiguo Cao, Yang Xiao, and Anton van den Hengel. An embarrassingly simple approach to visual domain adaptation. *IEEE Transactions on Image Processing*, 2018. 1, 3
- [37] Ying Lu, Lingkun Luo, Di Huang, Yunhong Wang, and Liming Chen. Knowledge transfer in vision recognition: A survey. *ACM Comput. Surv.*, 53(2), April 2020. 1
- [38] Lingkun Luo, Liming Chen, Shiqiang Hu, Ying Lu, and Xiaofang Wang. Discriminative and geometry aware unsupervised domain adaptation. *CoRR*, abs/1712.10042, 2017. 1
- [39] Lingkun Luo, Liming Chen, Shiqiang Hu, Ying Lu, and Xiaofang Wang. Discriminative and geometry-aware unsupervised domain adaptation. *IEEE Transactions on Cybernetics*, 2020. 1, 2, 3, 4, 6, 10, 11, 12
- [40] Lingkun Luo, Xiaofang Wang, Shiqiang Hu, and Liming Chen. Robust data geometric structure aligned close yet discriminative domain adaptation. *CoRR*, abs/1705.08620, 2017. 3, 9, 10
- [41] Feiping Nie, Jianjun Yuan, and Heng Huang. Optimal mean robust principal component analysis. In *International conference on machine learning*, pages 1062–1070, 2014. 8
- [42] Sinno Jialin Pan, Ivor W Tsang, James T Kwok, and Qiang Yang. Domain adaptation via transfer component analysis. *IEEE Transactions on Neural Networks*, 22(2):199–210, 2011. 3, 10
- [43] Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*, 22(10):1345–1359, 2010. 1
- [44] Pau Panareda Busto and Juergen Gall. Open set domain adaptation. In *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017. 1
- [45] V. M. Patel, R. Gopalan, R. Li, and R. Chellappa. Visual domain adaptation: A survey of recent advances. *IEEE Signal Processing Magazine*, 32(3):53–69, May 2015. 1
- [46] Zhongyi Pei, Zhangjie Cao, Mingsheng Long, and Jianmin Wang. Multi-adversarial domain adaptation. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018. 4
- [47] Artem Rozantsev, Mathieu Salzmann, and Pascal Fua. Beyond sharing weights for deep domain adaptation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018. 9, 10, 12
- [48] Kuniaki Saito, Yoshitaka Ushiku, and Tatsuya Harada. Asymmetric tri-training for unsupervised domain adaptation. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, pages 2988–2997, 2017. 10
- [49] Bernhard Schölkopf, Alexander J. Smola, and Klaus-Robert Müller. Nonlinear component analysis as a kernel eigenvalue problem. *Neural Computation*, 10(5):1299–1319, 1998. 10
- [50] Ozan Sener, Hyun Oh Song, Ashutosh Saxena, and Silvio Savarese. Learning transferrable representations for unsupervised domain adaptation. In *Advances in Neural Information Processing Systems*, pages 2110–2118, 2016. 10
- [51] Ming Shao, Dmitry Kit, and Yun Fu. Generalized transfer subspace learning through low-rank constraint. *International Journal of Computer Vision*, 109(1-2):74–93, 2014. 1, 3, 10
- [52] S. Si, D. Tao, and B. Geng. Bregman divergence-based regularization for transfer subspace learning. *IEEE Transactions on Knowledge and Data Engineering*, 22(7):929–942, July 2010. 3, 10
- [53] Si Si, Dacheng Tao, and Bo Geng. Bregman divergence-based regularization for transfer subspace learning. *IEEE Transactions on Knowledge and Data Engineering*, 22(7):929–942, 2010. 1
- [54] Baochen Sun, Jiashi Feng, and Kate Saenko. Return of frustratingly easy domain adaptation. In *AAAI*, volume 6, page 8, 2016. 1, 3, 10
- [55] Baochen Sun and Kate Saenko. Deep coral: Correlation alignment for deep domain adaptation. In *European Conference on Computer Vision*, pages 443–450. Springer, 2016. 4
- [56] Yuxing Tang, Josiah Wang, Xiaofang Wang, Boyang Gao, Emmanuel Dellandrea, Robert Gaizauskas, and Liming Chen. Visual and semantic knowledge transfer for large scale semi-supervised object detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017. 1
- [57] Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell. Adversarial discriminative domain adaptation. In *Computer Vision and Pattern Recognition (CVPR)*, volume 1, page 4, 2017. 1, 4, 10
- [58] Eric Tzeng, Judy Hoffman, Ning Zhang, Kate Saenko, and Trevor Darrell. Deep domain confusion: Maximizing for domain invariance. *CoRR*, abs/1412.3474, 2014. 10
- [59] Muhammad Uzair and Ajmal S. Mian. Blind domain adaptation with augmented extreme learning machine features. *IEEE Trans. Cybernetics*, 47(3):651–660, 2017. 9, 10
- [60] Hemanth Venkateswara, Jose Eusebio, Shayok Chakraborty, and Sethuraman Panchanathan. Deep hashing network for unsupervised domain adaptation. *arXiv preprint arXiv:1706.07522*, 2017. 9, 10, 11
- [61] Hao Wang, Wei Wang, Chen Zhang, and Fanjiang Xu. Cross-domain metric learning based on information theory. In *Proceedings of the Twenty-Eighth AAAI Conference on Artificial Intelligence, July 27 -31, 2014, Québec City, Québec, Canada.*, pages 2099–2105, 2014. 10
- [62] Jindong Wang, Wenjie Feng, Yiqiang Chen, Han Yu, Meiyu Huang, and Philip S Yu. Visual domain adaptation with manifold embedded distribution alignment. In *2018 ACM Multimedia Conference on Multimedia Conference*, pages 402–410. ACM, 2018. 4
- [63] Yong Xu, Xiaozhao Fang, Jian Wu, Xuelong Li, and David Zhang. Discriminative transfer subspace learning via low-rank and sparse representation. *IEEE Trans. Image Processing*, 25(2):850–863, 2016. 1, 3, 10
- [64] Jing Zhang, Wanqing Li, and Philip Ogunbona. Joint geometrical and statistical alignment for visual domain adaptation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017. 1, 3, 10
- [65] Sicheng Zhao, Bo Li, Xiangyu Yue, Yang Gu, Pengfei Xu, Runbo Hu, Hua Chai, and Kurt Keutzer. Multi-source domain adaptation for semantic segmentation. In *Advances in Neural Information Processing Systems*, pages 7285–7298, 2019. 4



Lingkun Luo received his Ph.D degree in Shanghai Jiao Tong University. He had served as research assistant and postdoc in Ecole Centrale de Lyon, Department of Mathematics and Computer Science, and a member of LIRIS laboratory. Now, he is a postdoc in Shanghai Jiao Tong University. He has authored more than 20 research articles. His research interests include machine learning, pattern recognition and computer vision.



Liming Chen received the joint B.Sc. degree in mathematics and computer science from the University of Nantes, Nantes, France in 1984, and the M.Sc. and Ph.D. degrees in computer science from the University of Paris 6, Paris, France, in 1986 and 1989, respectively.

He first served as an Associate Professor with the Université de Technologie de Compiègne, before joining École Centrale de Lyon, Écully, France, as a Professor in 1998, where he leads an advanced research team on Computer Vision, Machine Learning and Multimedia. From 2001 to 2003, he also served as Chief Scientific Officer in a Paris-based company, Avivias, specializing in media asset management. In 2005, he served as Scientific Multimedia Expert for France Telecom R&D China, Beijing, China. He was the Head of the Department of Mathematics and Computer Science, École Centrale de Lyon from 2007 through 2016. His current research interests include computer vision, machine learning, image and multimedia with a particular focus on robot vision and learning since 2016. Liming has over 300 publications and successfully supervised over 40 PhD students. He has been a grant holder for a number of research grants from EU FP program, French research funding bodies and local government departments. Liming has so far guest-edited 5 journal special issues. He is an associate editor for Eurasip Journal on Image and Video Processing and a senior IEEE member.



Shiqiang Hu received his PhD degree at Beijing Institute of Technology. He has over 150 publications and successfully supervised over 15 PhD students. Now, he is a full professor in Shanghai Jiao Tong University. His research interests include data fusion technology, image understanding, and nonlinear filter.