

Learning Domain-Invariant Relationship with Instrumental Variable for Domain Generalization

Junkun Yuan, *Student Member, IEEE*, Xu Ma, Kun Kuang, Ruoxuan Xiong, Mingming Gong, and Lanfen Lin, *Member, IEEE*

Abstract—Domain generalization (DG) aims to learn from multiple source domains a model that generalizes well on unseen target domains. Existing methods mainly learn input feature representations with invariant marginal distribution, while the invariance of the conditional distribution is more essential for unknown domain generalization. This paper proposes an instrumental variable-based approach to learn the domain-invariant relationship between input features and labels contained in the conditional distribution. Interestingly, with a causal view on the data generating process, we find that the input features of one domain are valid instrumental variables for other domains. Inspired by this finding, we design a simple yet effective framework to learn the Domain-invariant Relationship with Instrumental Variable (DRIVE) via a two-stage IV method. Specifically, it first learns the conditional distribution of input features of one domain given input features of another domain, and then it estimates the domain-invariant relationship by predicting labels with the learned conditional distribution. Simulation experiments show the proposed method accurately captures the domain-invariant relationship. Extensive experiments on several datasets consistently demonstrate that DRIVE yields state-of-the-art results.

Index Terms—Instrumental variable, causality, visual recognition, deep learning, domain generalization.

I. INTRODUCTION

GENERAL supervised learning extracts statistical patterns by assuming data across training (source) and test (target) sets are independent and identically distributed (i.i.d.). It may lead to poor generalization performance when testing the trained model on the data that is very distinct from the training one, which is known as the *dataset shift* (or domain shift) problem [42]. A prevailing research field for addressing this problem is *domain adaptation* (DA) [3], which adapts the model from source to target by utilizing a given unlabeled target dataset. However, DA methods need to re-collect target data and repeat the model adaptation process for each new target domain, which is time-consuming or even infeasible. *Domain generalization* (DG) [5] is thus proposed to use multiple semantic-related source datasets for generalizable model learning without accessing any target data/information.

(Corresponding author: Kun Kuang.)

J. Yuan, X. Ma, K. Kuang, and L. Lin are with the College of Computer Science and Technology, Zhejiang University, Hangzhou 310027, China (e-mail: yuanjk@zju.edu.cn; maxu@zju.edu.cn; kunkuang@zju.edu.cn; llf@zju.edu.cn.)

R. Xiong is with the Department of Quantitative Theory & Methods, Emory University, Atlanta, GA 30322, USA (e-mail: ruoxuan.xiong@emory.edu).

M. Gong is with the School of Mathematics and Statistics, The University of Melbourne, Melbourne, VIC 3010, Australia (e-mail: mingming.gong@unimelb.edu.au)

Since no information on the target domain is provided in the DG task, it is essential to learn the domain-invariant relationship between input features and labels from the source datasets for generalizable prediction. Numerous DG researches [26], [36], [39], [61] focus on learning discriminative domain-invariant feature representations. Most of them [26], [36], [39] are based on covariate shift assumption that marginal distribution $P(X)$ changes yet conditional distribution $P(Y|X)$ stays invariant across domains. However, it rarely holds in many real scenarios where the relationship $P(Y|X)$ also changes in different environments. To better obtain the stable conditional distribution of the labels, Zhao et al. [61] propose to minimize the dependency between input feature representations and labels by maximizing a conditional entropy regularization term. However, it relies on an adversarial learning strategy that might become unstable and less effective when there are not sufficient source datasets for model adversarial training.

In this paper, we propose to learn the domain-invariant relationship between input features X and labels Y from a causal perspective for domain generalization. We tackle this issue by putting forward a causal view on the data generating process via distinguishing domain-invariant and domain-specific parts as shown in Figure 1(a). In an analyzed *system* of domain m , input features X^m and labels Y^m are (indirectly) determined by domain-invariant factor F^{inv} that contains discriminative semantic information of objects. The domain-specific factor F^m plays the role of a common cause of X^m and Y^m by affecting both of them. X^m and Y^m are also affected by error e_x^m and e_y^m from the *environment* of domain m . In light of this, we build a causal graph of different domains in Figure 1(b). We attribute the changes of the conditional distribution $P(Y|X)$ across domains to the changes of the domain-specific factor F^m . Moreover, we assume that there exists a domain-invariant relationship f between input features and labels, contained in $P(Y|X)$, which we are interested in. With the analysis of Figure 1, we then find that the input features of one domain are valid *instrumental variables* (IVs) [52] (see Section III) of another domain. Inspired by this finding, we propose a two-stage IV process to learn the domain-invariant relationship f , which first learns conditional distribution of input features of one domain given input features of another domain, and then estimates f by predicting labels with the learned conditional distribution. Our proposed framework is simple yet effective, and the presented method is practical that helps the model extract the invariant relationship between input features and labels, effectively improving the out-of-domain generalization performance as verified by experiments.

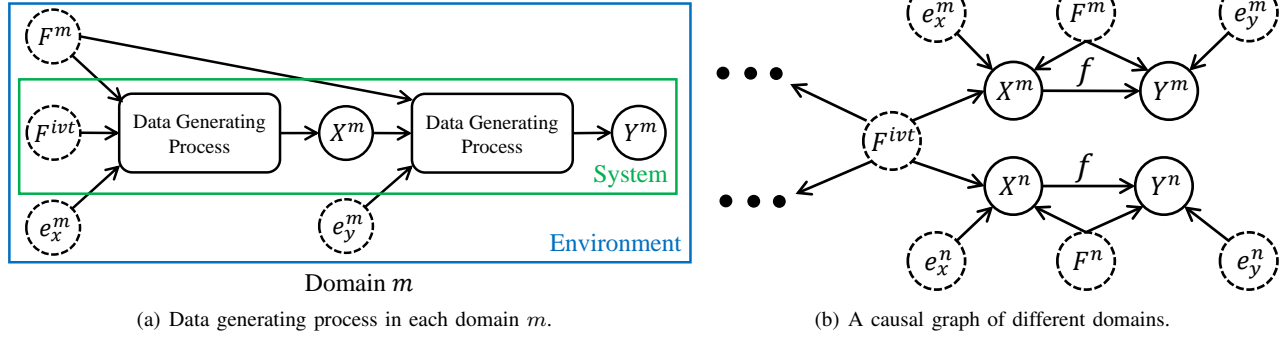


Fig. 1. A causal view on the data generating process (a) and the corresponding causal graph (b). For each domain m , input features X^m and labels Y^m are (indirectly) determined by domain-invariant factor F^{ivt} of the system, and confounded by domain-specific factor F^m of the environment. e_x^m and e_y^m are the error terms. We aim to learn the domain-invariant relationship f between the input features and the labels with instrumental variable-based method for improving model out-of-distribution generalization performance. Note that solid and dashed circles denote observed and latent variables, respectively.

We have the following contributions. (i) We formulate a data generating process by distinguishing domain-invariant and domain-specific parts in data. Based on this, we build a causal graph of different domains to analyze the domain generalization problem from a causal perspective; (ii) We propose a novel method to learn Domain-invariant Relationship with Instrumental Variable (DRIVE), which is simple yet effective by using the input features of one domain as IVs of other domains and implementing an IV-based two-stage generalizable model training process; (iii) We conduct simulation experiments to show that our method accurately estimates the domain-invariant relationship. Extensive experiments on multiple datasets show DRIVE yields state-of-the-art results.

The remainder of this paper is organized as follows. Section II gives a brief review of the related works on domain adaptation/generalization and instrumental variables. The formulation of the investigated domain generalization problem and preliminary of instrumental variables are introduced in Section III. The proposed method DRIVE is presented in Section IV. The experiments on simulated and real-world data and analysis are demonstrated in Section V. Finally, Section VI concludes this paper, with a future research outlook.

II. RELATED WORK

A. Domain Adaptation and Generalization

Domain adaptation (DA) [53], [51], [50], [30], [29], [23], [17], [7] aims to transfer the knowledge from the source domain(s) to the target domain(s). Unsupervised domain adaptation [6], [34], [59], [31] is a prevailing direction to DA that addresses the domain shift problem by minimizing domain gap between a labeled source domain and an unlabeled target domain, via domain adversarial learning [6], [34] or domain distance minimization [59], [31] et al. However, they need to access the data/information of the target domain in advance, which may be expensive or even infeasible in real scenarios.

Domain generalization (DG) is proposed to use multiple labeled source domains for training a generalizable model to unseen target domains. Recent DG methods with a variety of strategies can be included in the following main topics.

The first topic is Domain-invariant representation learning. This line of works learns feature representations that are invariant to domains and discriminative for classification. Some works [12], [22], [41] use auto-encoder structure to obtain invariant representations by performing a data reconstruction task. Li et al. [25] learn invariant class conditional representation with kernel mean embeddings. Piratla et al. [40] decompose networks into common and specific components, making the model rely on the common features. Li et al. [26], [36] introduce adversarial learning to extract effective representations with invariance constraints. Zhao et al. [61] introduce conditional entropy regularization term to learn conditional invariant features. Li et al. [24] model linear dependency in feature space and learn the common information.

Researches under a data augmentation topic are based on the idea that the augmented data generates various new domains, and the model trained on these domains could be more robust. Some approaches [45], [48] use gradient of the model to perturb data and construct new datasets for model training. Some others [8], [49] augment datasets by solving jigsaw puzzles. Zhou et al. [62], [63] employ an adversarial strategy to generate novel domains while keeping semantic information consistent. A recent work [64] mixes instance features to synthesize diverse domains and improves generalization.

Similar to the goal of DG, meta-learning-based methods keeps training the model on a meta-train dataset and improving its performance on a meta-test dataset. Numerous works [2], [19], [10], [21], [27] put forward meta-learning guided training algorithms to improve model out-of-domain generalization. However, it may be difficult to design and implement effective yet efficient meta-learning training algorithms in practice.

Some other methods learn the masks of features [9] or gradient [16] for regularization, or normalize batch and instance [44]. Mahajan et al. [35] provide a causal interpretation of domain generalization, which is similar to ours in letting input features be generated by domain-specific and domain-invariant factors/features. The difference is that they learn causal feature representations that contain both invariant and changing factors, and directly assume the conditional distribution of labels given causal feature representations is invariant. In comparison, we let the conditional distribution $P(Y|X)$

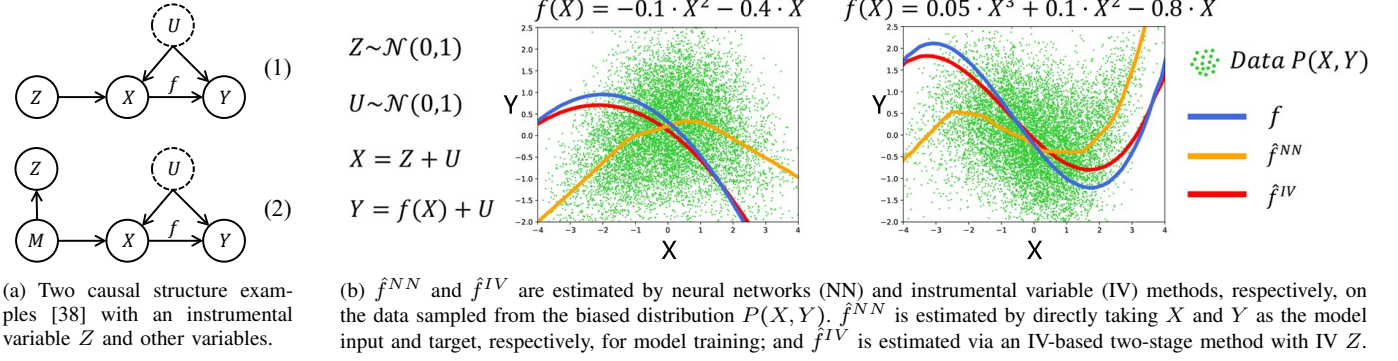


Fig. 2. Causal structures (a) with an instrumental variable Z , and toy experiments (b). We utilize the instrumental variable Z to estimate the invariant relationship f between X and Y by removing the confounding effect of an unobserved confounder U in the biased data distribution $P(X, Y)$.

change across domains and use instrumental variables to estimate the invariant part contained in $P(Y|X)$. As some causality-based domain adaptation works [58], [56] assume the causal direction $Y \rightarrow X$, we argue that our model encodes the same set of distribution change properties as pointed out by a recent work [57]. Our idea is also inspired by recent research on robustness of the neural networks [55], which assumes the domain-specific and domain-invariant causes in the data generating process. While in this paper, we use an IV-based approach to address dataset shift for the DG task.

B. Instrumental Variable

Instrumental variable (IV) method [52] is widely employed to capture causal relationship between variables for counterfactual prediction. Two stage least squares (2SLS) [1] is the most prevailing method in IV-based counterfactual prediction, which learns $\mathbb{E}[\phi(X)|Z]$ with IV Z and linear basis $\phi(\cdot)$, and fits Y by least squares regression with the coefficient $\hat{\phi}(\cdot)$ estimated in the first stage. Some non-parametric researches [37] extend the model basis to more complicated mapping function or regularization, e.g., polynomial basis. DeepIV [14] is proposed to use deep neural networks in the two-stage procedure, it fits a mixture density network $F_\phi(X|Z)$ in the first stage and regresses Y by sampling from the estimated mixture Gaussian distributions of X . KIV [46] is a recent work which maps and learns the relationships among Z , X , and Y in reproducing kernel Hilbert spaces. Another recent progress, DeepGMM [4], extends GMM methods in high-dimensional treatment and IV settings based on variational reformulation of the optimally-weighted GMM. We follow the additive function form used by most of the previous IV-based methods, i.e., $Y = f(X) + e$.

III. PRELIMINARY

In the domain generalization (DG) task, we have Q labeled source datasets $\mathcal{D}_1, \dots, \mathcal{D}_Q$ with different distributions $P_1(X^1, Y^1), \dots, P_Q(X^Q, Y^Q)$ on joint space $\mathcal{X} \times \mathcal{Y}$, where \mathcal{X} and \mathcal{Y} are input feature and label spaces, respectively. In each source domain q , N^q examples are sampled for the dataset \mathcal{D}_q , i.e., $\mathcal{D}_q = \{(x_n^q, y_n^q)\}_{n=1}^{N^q}$. DG aims to train a model with the Q source datasets and improve its generalization on unseen target domains where no data is provided during training.

In causal literature [38], invariant relationship (*response function*) f between X and Y is assumed as shown in Figure 2(a). The unobserved confounder U , which causes changes to both X and Y , introduces bias in data distribution $P(X, Y)$. The estimation of the relationship $P(Y|X)$ by learning from $P(X, Y)$ hence varies across domains with the changes of U . Instrumental variable (IV) [52] Z is a powerful tool for tracking the bias from the unobserved confounder U . A valid IV should satisfy the following conditions [38], [14]: (i) **Relevance**. Z and X should be relevant, i.e., $P(X|Z) \neq P(X)$; (ii) **Exclusion**. Z is correlated to Y only through X , i.e., $Z \perp\!\!\!\perp Y|(X, U)$; and (iii) **Unconfounded instrument**. Z is independent of U , i.e., $Z \perp\!\!\!\perp U$. These conditions make Z a valid IV, which allows us to learn the true relationship f between X and Y by considering the changes of Z .

We compare the functions estimated by direct neural networks (NN) and IV method by conducting toy experiments with 4000 data points sampled for training and test, respectively, as shown in Figure 2(b). We see that NN (orange line) directly learns $P(X, Y)$ that is biased by U , while IV method (red line) uses Z to eliminate the bias from U and estimates f (blue line), i.e. invariant causal relationship, more accurately.

For theoretical analysis, we consider a simple model

$$\mathbf{Y} = \mathbf{X} \cdot \lambda + \mathbf{U},$$

where $\mathbf{Y}, \mathbf{U} \in \mathbb{R}^n$, $\mathbf{X} \in \mathbb{R}^{n \times d_x}$, $\lambda \in \mathbb{R}^{d_x}$, n and d_x are the number of observations and dimension of X , respectively. The invariant relationship f is assumed as a linear mapping vector λ . We estimate λ via a two-stage IV method ($\hat{\lambda}^{IV}$) and an ordinary least squares (OLS) method ($\hat{\lambda}^{OLS}$). We have

$$\begin{aligned} \hat{\lambda}^{OLS} &= ((\mathbf{X})^\top \mathbf{X})^{-1} (\mathbf{X})^\top \mathbf{Y} \\ &= ((\mathbf{X})^\top \mathbf{X})^{-1} (\mathbf{X})^\top (\mathbf{X} \cdot \lambda + \mathbf{U}) \\ &= \lambda + \underbrace{((\mathbf{X})^\top \mathbf{X})^{-1} (\mathbf{X})^\top \mathbf{U}}_{\text{not tend to 0 for } X \text{ is correlated with } U}, \end{aligned}$$

$$\begin{aligned} \hat{\lambda}^{IV} &= ((\mathbf{Z})^\top \mathbf{X})^{-1} (\mathbf{Z})^\top \mathbf{Y} \\ &= ((\mathbf{Z})^\top \mathbf{X})^{-1} (\mathbf{Z})^\top (\mathbf{X} \cdot \lambda + \mathbf{U}) \\ &= \lambda + \underbrace{((\mathbf{Z})^\top \mathbf{X})^{-1} (\mathbf{Z})^\top \mathbf{U}}_{\text{tend to 0 for } Z \text{ is independent of } U}. \end{aligned}$$

The IV method utilizes Z to eliminate the bias of U and the estimator $\hat{\lambda}^{IV}$ converges to λ ; but the OLS estimator is biased. We extend this two-stage IV process in our proposed method to train generalizable model for domain generalization.

IV. LEARNING DOMAIN-INVARIANT RELATIONSHIP WITH INSTRUMENTAL VARIABLE

A. A Causal View on Domain Generalization

Data generating process (DGP). The general supervised learning imposes an i.i.d. assumption, however, changes in the external environment of a new domain will lead to changes in the analyzed system (i.e., variables and their relationships). The general supervised model trained on one domain may overfit the domain-specific information, leading to degradation of generalization ability on a new domain where the external environment changes. Nevertheless, we see that human can easily identify relationship in data no matter how the environment changes, e.g., to recognize images of animals with different backgrounds. We argue that the robust perception of human is based on the ability to distinguish domain-invariant and domain-specific parts in data via causal reasoning [57]. In light of this, it is necessary to analyze the dataset shift problem from a causal view by defining the latent DGP first.

Taking a visual recognition task of animals as an example. As shown in Figure 1(a), X^m and Y^m are images and the corresponding classes sampled from a specific dataset/domain m . There may exist multiple causes in the DGP of X^m and Y^m . Inspired by [55], we argue that X^m is determined by: (i) domain-invariant factor F^{ivt} , which is the key part of the recognized animals like size and limbs; (ii) domain-specific factor F^m that changes with the external environment, like light condition and background when taking pictures; (iii) an error term e_x^m . In the DGP, the factor F^{ivt} is invariant across domains, but the factor F^m plays the role of a common cause of X^m and Y^m , leading to confounding bias and the distribution shift problem across domains. For human, no matter how the images change, the corresponding classes can always be identified. It allows us to argue that there exists a latent domain-invariant relationship between input features X^m and labels Y^m . Therefore, we let Y^m be determined by X^m that contains invariant information of F^{ivt} , and Y^m is also affected by F^m and e_y^m . It should be noted that only the input features and labels are observed, while others are latent.

A causal graph of different domains. Based on the DGP, we build a causal graph of different domains as shown in Figure 1(b). Input X^m/X^n from the domain m/n shares domain-invariant factor F^{ivt} and is affected by different domain-specific factor F^m/F^n and error e_x^m/e_x^n , and the corresponding label Y^m/Y^n is determined by X^m/X^n through the relationship f , and influenced by F^m/F^n and e_y^m/e_y^n .

Assumption 1. Data distributions of different domains satisfy the data generating process and causal graph in Figure 1, where only the factor F^{ivt} and relationship f are invariant.

In each domain m , general supervised learning trains the model to learn the conditional distribution of the label

$$P(Y^m|X^m) = \int P(Y^m|X^m, F^m)P(F^m|X^m)dP(F^m). \quad (1)$$

The domain-specific factor F^m is a common cause of X^m and Y^m , leading to spurious correlation between X^m and Y^m , thus the conditional distribution of changes across domains. Since F^m is latent and can not be controlled, the introduced bias in data may not be removed directly. The model trained by minimizing risk on one domain overfits the bias and may have terrible performance on a new domain where the spurious correlation is different with the changes of the domain-specific factor. Since directly minimizing the risk on target domains is impossible as the data is unknown, instead, we propose to learn the relationship f between the input features and the labels which is invariant across domains. In causal literature [38], utilizing instrumental variable (IV) is an effective way to address the spurious correlation from the unobserved factor. By finding that the input features of one domain are valid IVs for other domains, we propose an IV-based two-stage method to learn the invariant relationship f for stable domain generalization, which is introduced in the following.

B. Learning Invariant Relationship with Instrumental Variable

Under Assumption 1, we give the following conclusions by using d-separation criterion [38].

Proposition 1. For any two domains m and n , if $m \neq n$, then the following conditions hold: (1) $X^n \not\perp\!\!\!\perp X^m$; (2) $X^n \perp\!\!\!\perp Y^m | (X^m, F^m)$; (3) $X^n \perp\!\!\!\perp F^m$; and (4) $X^n \perp\!\!\!\perp e_y^m$.

Based on Proposition 1, we have the following finding.

Theorem 1. For any two domains m and n , if $m \neq n$, then X^n is a valid instrumental variable of domain m .

Theorem 1 can be proved by referring to the conditions of IV in Section III. It indicates that one may adopt the input features of one source dataset as valid IVs to estimate the domain-invariant relationship f with another source dataset via a two-stage IV process (see Section III). That is, we first estimate the conditional distribution $P(X^m|X^n)$, and then predict labels Y^m with $P(X^m|X^n)$ instead of the input features X^m . Since X^n is independent of F^m , the changes of X^n through X^m to Y^m is stable to the changes of F^m . The estimation process can be understood as indirectly learning the changes of X^m with the changes of F^{ivt} , i.e., parent of X^m and X^n . F^{ivt} determines the class of the analyzed system, hence the estimated relationship between X^m and Y^m via this two-stage procedure is discriminative for classification yet insensitive to domain changes. By following [14], [46], [4], without loss of generality, we assume that the label Y^m is structurally determined by the following DGP:

$$Y^m = f(X^m) + F^m + e_y^m, \quad (2)$$

where $f(\cdot)$ is an unknown and non-linear continuous function, and $\mathbb{E}[e_y^m] = \mathbb{E}[F^m] = 0$. Note that we mainly focus on a generic model form, i.e., $Y^m = f(X^m) + F^m + e_y^m$, by

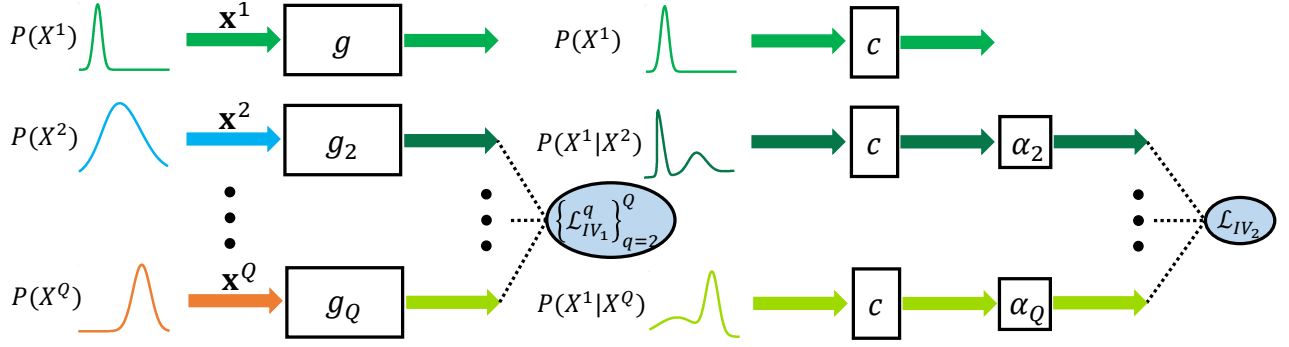


Fig. 3. The proposed DRIVE framework. We learn the invariant relationship f in the classifier c via an IV-based two-stage method, which first learns conditional distributions of X^1 with loss $\{\mathcal{L}_{IV_1^q}\}_{q=2}^Q$, and then use the approximation conditional distribution to predict Y^1 with loss \mathcal{L}_{IV_2} for debiasing c .

Algorithm 1 DRIVE

Input: Datasets $\mathcal{D}_1, \dots, \mathcal{D}_Q$, mixed dataset \mathcal{D}_{mix} , batchsize B , epochs E^{pre}, E^{IV} ;

Output: Well-trained \hat{g} and \hat{c} ;

- 1: **for** $epoch = 1$ to E^{pre} **do** // model pretraining
 - 2: Sample B examples from \mathcal{D}_{mix} and optimize g, c by minimizing \mathcal{L}_{pre} as Eq. (4);
 - 3: **end for**
 - 4: Initialize g_q by $g_q \leftarrow g$ for each $q \in \{2, \dots, Q\}$;
 - 5: **for** $epoch = 1$ to E^{IV} **do** // a two-stage IV method
 - 6: **for** $q = 2$ to Q **do**
 - 7: Sample B examples from \mathcal{D}_1 and \mathcal{D}_q and optimize g_q by minimizing $\mathcal{L}_{IV_1^q}$ as Eq. (5);
 - 8: **end for**
 - 9: Sample B examples from $\mathcal{D}_q, q = 2, \dots, Q$, and optimize c by minimizing \mathcal{L}_{IV_2} as Eq. (6).
 - 10: **end for**
-

following [14], [46], [4], while some other more general forms like $Y^m = c(f(X^m) + F^m + e_y^m)$ can be transformed by $c^{-1}(Y^m) = f(X^m) + F^m + e_y^m$. By taking the expectation of Y^m conditional on X^n , we have:

$$\begin{aligned} \mathbb{E}[Y^m|X^n] &= \mathbb{E}[f(X^m)|X^n] + \underbrace{\mathbb{E}[F^m|X^n] + \mathbb{E}[e_y^m|X^n]}_{=0 \text{ for } X^n \perp\!\!\!\perp F^m \text{ and } X^n \perp\!\!\!\perp e_y^m} \\ &= \int f(X^m) dP(X^m|X^n), \end{aligned} \quad (3)$$

It yields a two-stage strategy of learning the invariant relationship f with the instrumental variable X^n , i.e., estimating $P(X^m|X^n)$ in the first stage and learning f in the second stage by using the approximation of $P(X^m|X^n)$ learned in the first stage. Since f is the distribution function for the observable variables X^m, X^n , and Y^m , it is identified [37].

C. Framework and Algorithm

Following our analysis, we propose our method DRIVE with framework and algorithm as shown in Figure 3 and Algorithm 1, respectively. We adopt the common framework of DG with a feature extractor (backbone) g and a classifier c . g learns feature representations from data, and c outputs the classification results based on the extracted feature representations. Networks g_2, \dots, g_Q are assigned model parameters

from g after a model pretraining process, then they perform an IV-based two-stage method to learn the invariant relationship f in the classifier c by debiasing it. We then introduce the details of our method DRIVE in the following.

Model pretraining. We first pretrain the feature extractor g and the classifier c to initialize their discriminability, and reduce the dimensions of input features for the two-stage IV method. We randomly mix the sources $\{\mathcal{D}_q\}_{q=1}^Q$ to build a mixed dataset \mathcal{D}_{mix} and use it to pretrain g and c with a cross-entropy classification loss \mathcal{L}_{pre} , that is,

$$\mathcal{L}_{pre} = \mathbb{E}_{(x,y) \in \mathcal{D}_{mix}} \ell(c \circ g(x), y), \quad (4)$$

where ℓ is the cross-entropy loss function. Through model pretraining, the feature extractor g learns to extract feature representations of different datasets, and the classifier c is initialized to classify the extracted feature representations. However, c is biased because of the domain-specific information from the source datasets. Therefore, we then perform an IV-based two-stage method to debias c for learning the invariant relationship between the input feature (representations) and the labels.

Stage 1 of the IV method. We first assign the parameters of g to $\{g_q\}_{q=2}^Q$ to initialize them, which is effective for the first stage of the IV method to our empirical experience. The two-stage IV method is conducted by: (ii) learning conditional distributions of X^1 via optimizing the feature extractors $\{g_q\}_{q=2}^Q$; (ii) using the learned conditional distribution to debias the classifier c for learning the domain-invariant relationship f . Specifically, for the first stage, we use g_q to estimate $P(X^1|X^q)$ with the Maximum Mean Discrepancy (MMD) [13], i.e., $d_k^2(v, w) \triangleq \|\mathbb{E}_v[\phi(g_q(x^q))] - \mathbb{E}_w[\phi(g_q(x^1))]\|_{\mathcal{H}_k}^2$. The distributions of the extracted input feature representations $g_q(x^q)$ and $g(x^1)$, i.e., v and w , satisfy $v = w$ iff $d_k^2(v, w) = 0$. A characteristic kernel $k(g_q(x^q), g(x^1)) = \langle \phi(g_q(x^q)), \phi(g(x^1)) \rangle$ is defined as a convex combination of o PSD kernels $\{k_u\}$, i.e., $\mathcal{K} \triangleq \{k = \sum_{u=1}^o \beta_u k_u : \sum_{u=1}^o \beta_u = 1, \beta_u \geq 0, \forall u\}$, where β_u guarantees the characteristic of multi-kernel k [31], [33]. We then estimate $P(X^1|X^q)$ by optimizing g_q , which minimizes the MMD distance between the feature representations of X^1 and X^q with the loss function

$$\mathcal{L}_{IV_1^q} = p_{q,1} d_k^2(g_q(x^q), g(x^1)). \quad (5)$$

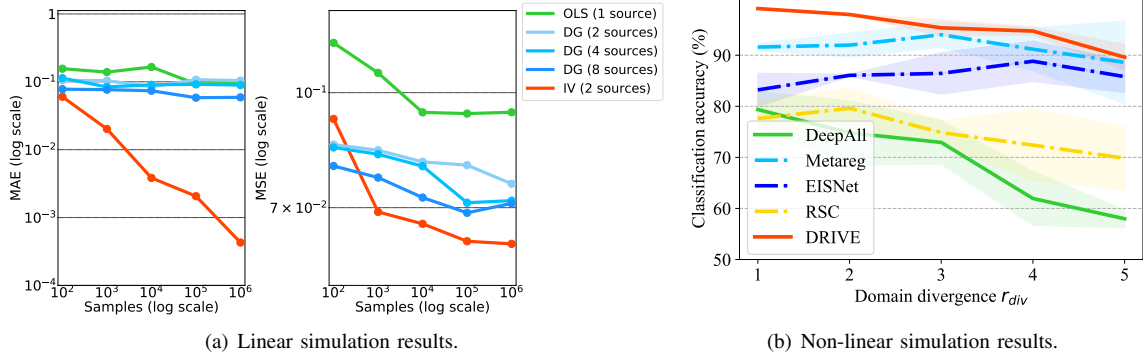


Fig. 4. Simulation results in linear setting (a) for evaluating invariant relationship learning and target regression, and non-linear (b) setting for label prediction.

where $p_{q,1} := \mathbb{1}(y^q = y^1)$, i.e., $p_{q,1} = 1$ when $y^q = y^1$, otherwise $p_{q,1} = 0$. It is used to guarantee only the MMD distance of the input features from the same classes are minimized, which helps g_q to learn a more accurate conditional distribution $P(X^1|X^q)$ for each $q \in \{2, \dots, Q\}$.

Stage 2 of the IV method. In the second stage, we sample points from the conditional distributions estimated in the first stage and use them to predict the labels. We optimize the classifier c with a classification loss of the estimated conditional distribution, that is,

$$\mathcal{L}_{IV_2} = \frac{1}{Q-1} \sum_{q=2}^Q \alpha_q \mathbb{E}_{(x^q, y^q), (x^1, y^1)} [p_{q,1} \ell(c \circ g_q(x^q), y^1)] \quad (6)$$

where α_q is a hyper-parameter. We use $p_{q,1}$ to guarantee the data used for debiasing are in the same classes. By optimizing \mathcal{L}_{IV_2} , the classifier c removes the domain-specific bias of the source datasets introduced in the model pretraining. It allows c to capture the domain-invariant relationship, improving the model out-of-domain generalization ability.

V. EXPERIMENTS

We first conduct simulation experiments to verify the relationship learned by our method DRIVE. Then, we perform experiments on multiple real-world datasets to further testify the model generalization performance achieved by DRIVE.

A. Experiments on Simulated Datasets

Linear simulations. We sample variables for each domain m with $F^{ivt}, F^m \sim \mathcal{N}(\mu_f, 1)$ and $e_x^m, e_y^m \sim \mathcal{N}(\mu_e, 0.1)$. We sample μ_f once from uniform distribution $\text{Unif}(-1, 1)$ and sample μ_e once from $\text{Unif}(-0.1, 0.1)$ for each domain, making the divergence in each domain be random. We first consider linear setting with one-dimensional variables. The DGP of Figure 1(a) is assumed as

$$\begin{aligned} X^m &= \phi_m \cdot F^{ivt} + \alpha_m \cdot F^m + e_x^m, \\ Y^m &= \lambda_{ivt} \cdot X^m + \beta_m F^m + e_y^m, \end{aligned} \quad (7)$$

where λ_{ivt} is the invariant relationship that we are interested in. We sample ϕ_m, λ_{ivt} once from $\text{Unif}(-1, 1)$ and sample α_m, β_m once from $\text{Unif}(-0.5, 0.5)$ for each domain. Note that we let domain-invariant factor and relationship, i.e., F^{ivt} and

λ_{ivt} , be the same in all domains. In each run, we randomly generate 8 source domains for training and a target domain for test with 20,000 points in each domain. We run each method with linear regression, and report the MAE of domain-invariant relationship estimation, i.e., $\mathbb{E}[\hat{\lambda}_{ivt} - \lambda_{ivt}]$, and the MSE of the target domain label Y^t prediction, i.e., $\mathbb{E}[(\hat{Y}^t - Y^t)^2]$. We implement **OLS** method by training the model on one source domain. The general **DG** method is implemented by estimating the coefficient in each domain and average them to get a robust coefficient. **DG (n) sources** is denoted as the coefficient estimated in this way with n sources. **IV** method only needs two sources, i.e., the input features of one is used as IV to estimate the relationship on another source domain. We plot the results in Figure 4(a). Obviously, with the increase of sample size, IV method outperforms others in invariant relationship λ_{ivt} estimation and target label prediction when only using two source datasets. Although more source datasets allow the general DG methods to eliminate the domain-specific bias, they are still fooled by the introduced bias in data.

Non-linear simulations. We then compare DRIVE with other DG methods (introduced in Section II), i.e., Metareg [2], EISNet [49], and RSC [16], in non-linear relationship f estimation. We follow the linear simulations, but sample μ_f from $\text{Unif}(-r_{div}, r_{div})$ for each domain, where a larger value of r_{div} indicates the larger the domain divergence probably be. The invariant relationship λ_{ivt} is replaced with a non-linear abs function. We set dimensions of factor and input features to 1500 and 600, respectively. 5000 points are sampled for two classes respectively. All the methods are implemented by their public code, but their networks are replaced with 4 fully-connected layers with 600, 256, 128, and 64 units, respectively, for fair comparison. We use SGD optimizer with learning rate 0.01, and run 4000 iterations with batchsize 64. The results in Figure 4(b) illustrates that DRIVE with IV-based two-stage method outperforms other state-of-the-art DG methods. It is worth mentioning that DRIVE only utilizes two source domains to train, while other methods have 8 sources. We attribute the significant performance of DRIVE to its domain-invariant relationship learning ability, which makes full use of the two sources to obtain the invariant part contained in the conditional distribution of the labels given the input features. Besides, we find that data augmentation based

TABLE I
RESULTS (%) FOR DOMAIN GENERALIZATION ON PACS DATASET.

Methods	Art	Cartoon	Photo	Sketch	Average
DeepAll [8]	78.96	72.93	96.28	70.59	79.94
JiGen [8]	79.42	75.25	96.03	71.35	80.51
MASF [10]	80.29	77.17	94.99	71.69	81.04
DGER [61]	80.70	76.40	96.65	71.77	81.38
Epi-FCR [21]	82.1	77.0	93.9	73.0	81.5
MMLD [36]	81.28	77.16	96.09	72.29	81.83
EISNet [49]	81.89	76.44	95.93	74.33	82.15
L2A-OT [63]	83.3	78.2	96.2	73.6	82.8
DDAIG [62]	84.2	78.1	95.3	74.7	83.1
DRIVE w/o IV	79.40 \pm 0.10	76.93 \pm 0.09	95.75 \pm 0.10	74.44 \pm 0.07	81.63 \pm 0.03
DRIVE w/o pre	81.95 \pm 0.25	77.55 \pm 0.31	96.64 \pm 0.34	75.65 \pm 0.10	82.95 \pm 0.14
DRIVE	83.36 \pm 0.70	78.76 \pm 0.08	96.87 \pm 0.18	78.68 \pm 0.96	84.42 \pm 0.11

method, i.e., Metareg and EISNet, show the robustness to the domain divergence. It is may because the divergence domains give full play to the advantages of generating novel domains for model training, improving the model generalization.

B. Experiments on Real-World Datasets

Datasets and Implementations. We first conduct experiments on **PACS** [20], which has 7 categories over 4 domains, that is, Art, Cartoon, Sketch, and Photo. Then we have **Office-Home** dataset [47] that consists of 15,500 images of 65 categories over 4 domains, i.e., Art, Clipart, Product, and Real-World. Example images are shown in Figure 5. We follow the training and test split in previous works [47], [20], [61], and perform leave-one-domain-out experiments, i.e., one domain is held out as the target domain for test. We follow [8], [10], [16] by using the pretrained ResNet-18 [15] network. We use SGD optimizer with learning rate 0.01 and batchsize 64. The epochs for the pretraining (E^{pre}) and the IV method (E^{IV}) are both set to 20. As one domain is chosen as the target domain, any of the rest domains can be used as \mathcal{D}^1 , we use a held-out validation set to choose the optimal \mathcal{D}^1 as well as the corresponding hyper-parameters of Eq. (6). We conduct the experiments with CPU Intel i7-8700K \times 1 and GPU Nvidia RTX 3090 \times 1. We run each experiment 3 times with random seed, and cite the results of other methods in their papers.

Since when a domain is used as the target domain, any source could be treated as the \mathcal{D}^1 , and other sources are used to learn the conditional distribution of X^1 . Therefore, we first set all the weights (hyper-parameters) α to 1, and conduct different source combination experiments for PACS (Table III) and Office-Home (Table IV) datasets. After we have the best domain combinations, we then conduct weight combination experiments on PACS (Table V) and Office-Home (Table VI) datasets. Finally, we use the “target- \mathcal{D}^1 ” combinations “Art-Photo”, “Cartoon-Photo”, “Photo-Art”, “Sketch-Photo” with weights $\alpha_1 = 1.25, \alpha_2 = 0.75$ for PACS dataset; and use “Art-Clipart”, “Clipart-Art”, “Product-Art”, “Real-World-Art” with weights $\alpha_1 = 1.5, \alpha_2 = 0.5$ for Office-Home datasets.

Results. Table I and Table II report the results on PACS and Office-Home datasets, respectively. We first find that DRIVE outperforms other methods on both datasets by performing the best on most of the DG sub-tasks and achieving the highest averaged accuracy. We attribute it to that DRIVE learns

TABLE II
RESULTS (%) FOR DOMAIN GENERALIZATION ON OFFICE-HOME DATASET.

Methods	Art	Clipart	Product	Real-World	Average
DeepAll [8]	52.15	45.86	70.86	73.15	60.51
JiGen [8]	53.04	47.51	71.47	72.79	61.20
DSON [44]	59.37	44.70	71.84	74.68	62.90
RSC [16]	58.42	47.90	71.63	74.54	63.12
DRIVE w/o IV	55.53 \pm 0.21	45.92 \pm 0.50	71.64 \pm 0.35	74.49 \pm 0.05	61.90 \pm 0.20
DRIVE w/o pre	59.30 \pm 0.06	47.65 \pm 0.30	72.03 \pm 0.57	75.55 \pm 0.24	63.63 \pm 0.11
DRIVE	60.40 \pm 0.26	47.73 \pm 0.28	72.63 \pm 0.18	76.14 \pm 0.10	64.23 \pm 0.09

TABLE III
RESULTS (%) OF DIFFERENT COMBINATIONS FOR DOMAIN GENERALIZATION ON PACS DATASET.

$\mathcal{D}^1 \setminus \text{Target}$	Art	Cartoon	Photo	Sketch
Art	-	78.10 \pm 0.37	97.17 \pm 0.12	76.91 \pm 0.03
Cartoon	82.21 \pm 0.97	-	96.75 \pm 0.17	76.78 \pm 0.98
Photo	83.77 \pm 0.57	78.34 \pm 0.58	-	77.48 \pm 0.32
Sketch	81.46 \pm 0.10	78.20 \pm 0.76	97.01 \pm 0.27	-

to capture the invariant part (relationship) contained in the conditional distribution for better model generalization. We let DRIVE discard the IV method and pretraining as **w/o IV** and **w/o pre**, respectively in Table I and Table II. It shows that each part is important for DRIVE to yield significant performance, especially the IV method. It is may because the pretraining initialize the discriminability of the feature extractor for better conditional distribution estimation, and the IV method helps model learn domain-invariant relationship by debiasing the classifier. We plot the t-SNE feature visualization in Figure 6. It indicates that IV method helps DRIVE to learn discriminative and domain-invariant feature during the IV-based two-stage process by separating the features of different classes while aggregating the features of different domains.

C. Experiments on Biased Data

We also evaluate DRIVE on the unsupervised domain adaptation (UDA) task where DRIVE uses the input features of the target dataset as IVs to learn the invariant relationship with the given source dataset. We adopt two biased datasets for this task. The first is **Dogs and Cats** [18], where TB1 domain contains bright dogs and dark cats; but TB2 domain contains dark dogs and bright cats. The second is **IMDB face** dataset [18]. Women in a domain EB1 aged 0-29 and in another domain EB2 aged 40+; but men in EB1 aged 40+ and in EB2 aged 0-29. There is clear bias between the domains in the two datasets, which challenges the methods to learn stable relationship between the images and labels. We compare DRIVE with representative DA approaches, DAN [31], DANN [11], JAN [33], MDD [60], CDAN [32], MCD [43]. All the experiments are implemented using the same training setting for fair comparison. Following [32], [54], [28], We employ the pre-trained ResNet-50 [15] as the feature extractor, where the last layer is replaced by one FC layer with 256 units. Classifier is a FC layer put after feature extractor for classification. We train each method through back-propagation by SGD with batch-size 64, learning rate 0.01, momentum 0.9, and weight decay 0.001. Each method



Fig. 5. Example images of the adopted public datasets from left to right: PACS, Office-Home, Dogs and Cats, IMDB face. The former two datasets are used for the domain generalization task; and the latter two datasets are used for the unsupervised domain adaptation task.

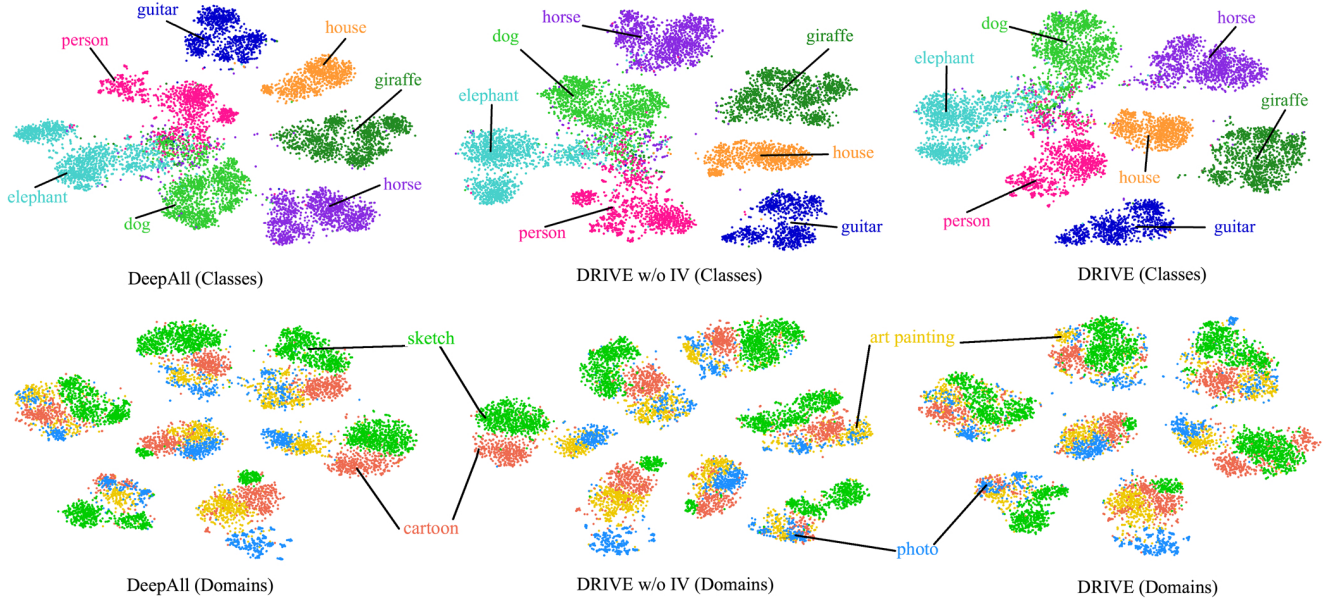


Fig. 6. T-SNE visualization of the learned feature representations of DeepAll, DRIVE w/o IV, and DRIVE, on PACS dataset. Different colors in the above and below sub-figures represent different classes and domains, respectively. The points gather separately for classes and compactly for domains indicates the learned feature representations are more discriminative and domain-agnostic, respectively.

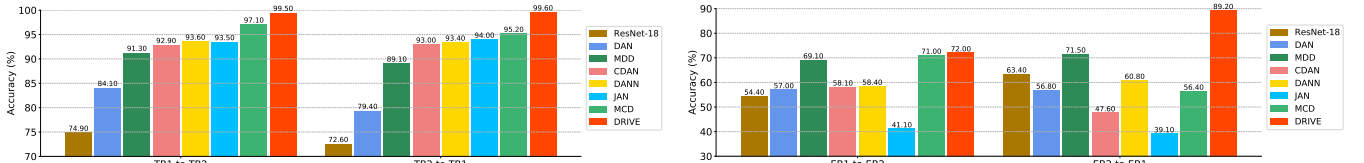


Fig. 7. Results for the unsupervised domain adaptation task on Dogs and Cats (left) and IMDB face (right) biased datasets.

TABLE IV
RESULTS (%) OF DIFFERENT COMBINATIONS FOR DOMAIN GENERALIZATION ON OFFICE-HOME DATASET.

$\mathcal{D}^1 \setminus \text{Target}$	Art	Clipart	Product	Real-World
Art	-	45.71 \pm 0.20	72.31 \pm 0.22	76.88 \pm 0.08
Clipart	60.89 \pm 0.17	-	72.21 \pm 0.05	76.88 \pm 0.12
Product	60.41 \pm 0.20	45.10 \pm 0.53	-	76.83 \pm 0.14
Real-World	60.64 \pm 0.29	45.65 \pm 0.23	72.30 \pm 0.41	-

are run 10 epochs on Dogs and Cats dataset and 5 epochs on IMDB face dataset for fair comparison. Results in Figure 7

show that DRIVE performs much better than others on the two challenging biased datasets. Moreover, we find that DRIVE achieves significant improvement on IMDB face dataset. It is probably because IV method needs sufficient samples to obtain the invariant relationship (see Figure 4(a)), and IMDB face is a large public dataset with 460,723 images.

VI. CONCLUSIONS

In this paper, we first give a causal view on the domain generalization problem, and then propose to learn domain-invariant relationship with instrumental variable via an IV-

TABLE V

RESULTS (%) WITH DIFFERENT WEIGHTS FOR DOMAIN GENERALIZATION ON PACS DATASET.

α_1	α_2	Art	Cartoon	Photo	Sketch	Average
0	2	82.68 \pm 0.23	78.25 \pm 0.15	97.15 \pm 0.18	77.43 \pm 0.29	83.88 \pm 0.88
0.25	1.75	82.33 \pm 0.30	79.15 \pm 0.57	97.07 \pm 0.12	78.16 \pm 0.64	84.18 \pm 0.26
0.5	1.5	82.21 \pm 0.75	78.19 \pm 0.50	97.11 \pm 0.07	77.51 \pm 0.96	83.75 \pm 0.09
0.75	1.25	82.60 \pm 0.88	78.88 \pm 0.89	97.09 \pm 0.03	78.65 \pm 0.71	84.31 \pm 0.62
1	1	83.77 \pm 0.57	78.34 \pm 0.58	97.17 \pm 0.12	77.48 \pm 0.32	84.19 \pm 0.25
1.25	0.75	83.36 \pm 0.70	78.76 \pm 0.08	96.87 \pm 0.18	78.68 \pm 0.96	84.42 \pm 0.11
1.5	0.5	81.89 \pm 0.08	78.60 \pm 0.29	97.35 \pm 0.12	78.20 \pm 1.02	84.01 \pm 0.18
1.75	0.25	81.98 \pm 0.15	79.17 \pm 0.73	96.87 \pm 0.38	78.53 \pm 0.14	84.14 \pm 0.10
2	0	82.14 \pm 0.17	78.22 \pm 0.10	97.05 \pm 0.12	77.46 \pm 1.58	83.72 \pm 0.38

TABLE VI

RESULTS (%) WITH DIFFERENT WEIGHTS FOR DOMAIN GENERALIZATION ON OFFICE-HOME DATASET.

α_1	α_2	Art	Clipart	Product	Real-World	Average
0	2	60.63 \pm 0.25	47.40 \pm 0.08	72.51 \pm 0.08	76.12 \pm 0.59	64.16 \pm 0.10
0.25	1.75	60.71 \pm 0.13	46.48 \pm 0.26	72.59 \pm 0.08	76.93 \pm 0.17	64.18 \pm 0.02
0.5	1.5	60.79 \pm 0.11	46.18 \pm 0.16	72.62 \pm 0.13	76.10 \pm 0.17	63.92 \pm 0.04
0.75	1.25	60.53 \pm 0.09	46.36 \pm 0.36	72.60 \pm 0.14	76.69 \pm 0.11	64.05 \pm 0.08
1	1	60.89 \pm 0.17	45.71 \pm 0.20	72.31 \pm 0.22	76.88 \pm 0.08	63.95 \pm 0.05
1.25	0.75	60.90 \pm 0.39	46.20 \pm 0.35	72.54 \pm 0.23	77.06 \pm 0.25	64.17 \pm 0.10
1.5	0.5	60.40 \pm 0.26	47.73 \pm 0.28	72.63 \pm 0.18	76.14 \pm 0.10	64.23 \pm 0.09
1.75	0.25	60.58 \pm 0.20	46.48 \pm 0.21	72.54 \pm 0.14	76.89 \pm 0.34	64.12 \pm 0.11
2	0	60.95 \pm 0.15	46.28 \pm 0.11	72.44 \pm 0.04	76.88 \pm 0.21	64.14 \pm 0.04

based two-stage method. Extensive experiments show the significant performance of our method. However, the theoretical analysis is based on the additive assumption of the domain-specific factor and error term. We may give the proofs with more moderate assumptions in future work. Our paper benefits the research field of domain generalization and may not have a negative impact of the society to our knowledge.

ACKNOWLEDGMENT

This work was supported in part by National Key Research and Development Program of China (No. 2018AAA0101900), National Natural Science Foundation of China (No. 61625107, No. 62006207), Key R & D Projects of the Ministry of Science and Technology (No. 2020YFC0832500), the Fundamental Research Funds for the Central Universities and Zhejiang Province Natural Science Foundation (No. LQ21F020020).

REFERENCES

- [1] J. D. Angrist and J.-S. Pischke. *Mostly harmless econometrics: An empiricist's companion*. Princeton university press, 2008.
- [2] Y. Balaji, S. Sankaranarayanan, and R. Chellappa. Metareg: Towards domain generalization using meta-regularization. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 998–1008, 2018.
- [3] S. Ben-David, J. Blitzer, K. Crammer, A. Kulesza, F. Pereira, and J. W. Vaughan. A theory of learning from different domains. *Machine learning*, 79(1-2):151–175, 2010.
- [4] A. Bennett, N. Kallus, and T. Schnabel. Deep generalized method of moments for instrumental variable analysis. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 3564–3574, 2019.
- [5] G. Blanchard, G. Lee, and C. Scott. Generalizing from several related classification tasks to a new unlabeled sample. *Advances in neural information processing systems (NeurIPS)*, 24:2178–2186, 2011.
- [6] G. Cai, Y. Wang, L. He, and M. Zhou. Unsupervised domain adaptation with adversarial residual transform networks. *IEEE transactions on neural networks and learning systems (TNNLS)*, 31(8):3073–3086, 2019.
- [7] R. Cai, J. Li, Z. Zhang, X. Yang, and Z. Hao. Dach: Domain adaptation without domain information. *IEEE transactions on neural networks and learning systems (TNNLS)*, 31(12):5055–5067, 2020.
- [8] F. M. Carlucci, A. D’Innocente, S. Bucci, B. Caputo, and T. Tommasi. Domain generalization by solving jigsaw puzzles. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2224–2233, 2019.
- [9] P. Chattopadhyay, Y. Balaji, and J. Hoffman. Learning to balance specificity and invariance for in and out of domain generalization. In *European Conference on Computer Vision (ECCV)*, pages 301–318. Springer, 2020.
- [10] Q. Dou, D. C. de Castro, K. Kamnitsas, and B. Glocker. Domain generalization via model-agnostic learning of semantic features. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 6450–6461, 2019.
- [11] Y. Ganin and V. Lempitsky. Unsupervised domain adaptation by backpropagation. In *International conference on machine learning (ICML)*, pages 1180–1189. PMLR, 2015.
- [12] M. Ghifary, W. Bastiaan Kleijn, M. Zhang, and D. Balduzzi. Domain generalization for object recognition with multi-task autoencoders. In *Proceedings of the IEEE international conference on computer vision (ICCV)*, pages 2551–2559, 2015.
- [13] A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, and A. Smola. A kernel two-sample test. *The Journal of Machine Learning Research (JMLR)*, 13(1):723–773, 2012.
- [14] J. Hartford, G. Lewis, K. Leyton-Brown, and M. Taddy. Deep iv: A flexible approach for counterfactual prediction. In *International Conference on Machine Learning (ICML)*, pages 1414–1423, 2017.
- [15] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.
- [16] Z. Huang, H. Wang, E. P. Xing, and D. Huang. Self-challenging improves cross-domain generalization. In *European Conference on Computer Vision (ECCV)*, pages 124–140, 2020.
- [17] Q. Kang, S. Yao, M. Zhou, K. Zhang, and A. Abusorrah. Effective visual domain adaptation via generative adversarial distribution matching. *IEEE Transactions on Neural Networks and Learning Systems (TNNLS)*, 2020.
- [18] B. Kim, H. Kim, K. Kim, S. Kim, and J. Kim. Learning not to learn: Training deep neural networks with biased data. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9012–9020, 2019.
- [19] D. Li, Y. Yang, Y.-Z. Song, and T. Hospedales. Learning to generalize: Meta-learning for domain generalization. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, volume 32, 2018.
- [20] D. Li, Y. Yang, Y.-Z. Song, and T. M. Hospedales. Deeper, broader and artier domain generalization. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 5542–5550, 2017.
- [21] D. Li, J. Zhang, Y. Yang, C. Liu, Y.-Z. Song, and T. M. Hospedales. Episodic training for domain generalization. *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 1446–1455, 2019.
- [22] H. Li, S. Jialin Pan, S. Wang, and A. C. Kot. Domain generalization with adversarial feature learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5400–5409, 2018.
- [23] H. Li, S. J. Pan, S. Wang, and A. C. Kot. Heterogeneous domain adaptation via nonlinear matrix factorization. *IEEE transactions on neural networks and learning systems (TNNLS)*, 31(3):984–996, 2019.
- [24] H. Li, Y. Wang, R. Wan, S. Wang, T. Li, and A. C. Kot. Domain generalization for medical imaging classification with linear-dependency regularization. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- [25] Y. Li, M. Gong, X. Tian, T. Liu, and D. Tao. Domain generalization via conditional invariant representations. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- [26] Y. Li, X. Tian, M. Gong, Y. Liu, T. Liu, K. Zhang, and D. Tao. Deep domain generalization via conditional invariant adversarial networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 624–639, 2018.
- [27] Y. Li, Y. Yang, W. Zhou, and T. Hospedales. Feature-critic networks for heterogeneous domain generalization. In *International Conference on Machine Learning (ICML)*, pages 3915–3924. PMLR, 2019.
- [28] J. Liang, D. Hu, and J. Feng. Do we really need to access the source data? source hypothesis transfer for unsupervised domain adaptation. In *International Conference on Machine Learning (ICML)*. PMLR, 2020.
- [29] F. Liu, G. Zhang, and J. Lu. Heterogeneous domain adaptation: An unsupervised approach. *IEEE transactions on neural networks and learning systems (TNNLS)*, 31(12):5588–5602, 2020.
- [30] Y. Liu, B. Du, W. Tu, M. Gong, Y. Guo, and D. Tao. Logdet metric-based domain adaptation. *IEEE transactions on neural networks and learning systems (TNNLS)*, 31(11):4673–4687, 2020.

- [31] M. Long, Y. Cao, J. Wang, and M. Jordan. Learning transferable features with deep adaptation networks. In *International conference on machine learning (ICML)*, pages 97–105. PMLR, 2015.
- [32] M. Long, Z. Cao, J. Wang, and M. I. Jordan. Conditional adversarial domain adaptation. In *Advances in neural information processing systems (NeurIPS)*, pages 1640–1650, 2018.
- [33] M. Long, H. Zhu, J. Wang, and M. I. Jordan. Deep transfer learning with joint adaptation networks. In *International conference on machine learning (ICML)*, pages 2208–2217. PMLR, 2017.
- [34] A. Ma, J. Li, K. Lu, L. Zhu, and H. T. Shen. Adversarial entropy optimization for unsupervised domain adaptation. *IEEE Transactions on Neural Networks and Learning Systems (TNNLS)*, 2021.
- [35] D. Mahajan, S. Tople, and A. Sharma. Domain generalization using causal matching. *arXiv preprint arXiv:2006.07500*, 2020.
- [36] T. Matsuura and T. Harada. Domain generalization using a mixture of multiple latent domains. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, 2020.
- [37] W. K. Newey and J. L. Powell. Instrumental variable estimation of nonparametric models. *Econometrica*, pages 1565–1578, 2003.
- [38] J. Pearl. *Causality*. Cambridge university press, 2009.
- [39] X. Peng, Q. Bai, X. Xia, Z. Huang, K. Saenko, and B. Wang. Moment matching for multi-source domain adaptation. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 1406–1415, 2019.
- [40] V. Piratla, P. Netrapalli, and S. Sarawagi. Efficient domain generalization via common-specific low-rank decomposition. In *International Conference on Machine Learning (ICML)*, 2020.
- [41] F. Qiao, L. Zhao, and X. Peng. Learning to learn single domain generalization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 12556–12565, 2020.
- [42] J. Quionero-Candela, M. Sugiyama, A. Schwaighofer, and N. D. Lawrence. *Dataset shift in machine learning*. The MIT Press, 2009.
- [43] K. Saito, K. Watanabe, Y. Ushiku, and T. Harada. Maximum classifier discrepancy for unsupervised domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3723–3732, 2018.
- [44] S. Seo, Y. Suh, D. Kim, J. Han, and B. Han. Learning to optimize domain specific normalization for domain generalization. In *European Conference on Computer Vision (ECCV)*, 2020.
- [45] S. Shankar, V. Piratla, S. Chakrabarti, S. Chaudhuri, P. Jyothi, and S. Sarawagi. Generalizing across domains via cross-gradient training. In *International Conference on Learning Representations (ICLR)*, 2018.
- [46] R. Singh, M. Sahani, and A. Gretton. Kernel instrumental variable regression. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 4593–4605, 2019.
- [47] H. Venkateswara, J. Eusebio, S. Chakraborty, and S. Panchanathan. Deep hashing network for unsupervised domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5018–5027, 2017.
- [48] R. Volpi, H. Namkoong, O. Sener, J. C. Duchi, V. Murino, and S. Savarese. Generalizing to unseen domains via adversarial data augmentation. In *Advances in neural information processing systems (NeurIPS)*, pages 5334–5344, 2018.
- [49] S. Wang, L. Yu, C. Li, C.-W. Fu, and P. Heng. Learning from extrinsic and intrinsic supervisions for domain generalization. In *European conference on computer vision (ECCV)*, 2020.
- [50] W. Wang, H. Li, Z. Ding, F. Nie, J. Chen, X. Dong, and Z. Wang. Rethinking maximum mean discrepancy for visual domain adaptation. *IEEE Transactions on Neural Networks and Learning Systems (TNNLS)*, 2021.
- [51] Z. Wang, B. Du, and Y. Guo. Domain adaptation with neural embedding matching. *IEEE transactions on neural networks and learning systems (TNNLS)*, 31(7):2387–2397, 2019.
- [52] P. G. Wright. *Tariff on animal and vegetable oils*. Macmillan Company, New York, 1928.
- [53] X. Wu, S. Zhang, Q. Zhou, Z. Yang, C. Zhao, and L. J. Latecki. Entropy minimization versus diversity maximization for domain adaptation. *IEEE Transactions on Neural Networks and Learning Systems (TNNLS)*, 2021.
- [54] R. Xu, G. Li, J. Yang, and L. Lin. Larger norm more transferable: An adaptive feature norm approach for unsupervised domain adaptation. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 1426–1435, 2019.
- [55] C. Zhang, K. Zhang, and Y. Li. A causal view on robustness of neural networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- [56] K. Zhang, M. Gong, B. Schölkopf, et al. Multi-source domain adaptation: A causal view. In *AAAI Conference on Artificial Intelligence (AAAI)*, volume 1, pages 3150–3157, 2015.
- [57] K. Zhang, M. Gong, P. Stojanov, B. Huang, Q. Liu, and C. Glymour. Domain adaptation as a problem of inference on graphical models. *Advances in Neural Information Processing Systems (NeurIPS)*, 33, 2020.
- [58] K. Zhang, B. Schölkopf, K. Muandet, and Z. Wang. Domain adaptation under target and conditional shift. In *International Conference on Machine Learning (ICML)*, pages 819–827, 2013.
- [59] L. Zhang, J. Fu, S. Wang, D. Zhang, Z. Dong, and C. P. Chen. Guide subspace learning for unsupervised domain adaptation. *IEEE transactions on neural networks and learning systems (TNNLS)*, 31(9):3374–3388, 2019.
- [60] Y. Zhang, T. Liu, M. Long, and M. I. Jordan. Bridging theory and algorithm for domain adaptation. In *International Conference on Machine Learning (ICML)*, 2019.
- [61] S. Zhao, M. Gong, T. Liu, H. Fu, and D. Tao. Domain generalization via entropy regularization. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- [62] K. Zhou, Y. Yang, T. Hospedales, and T. Xiang. Deep domain-adversarial image generation for domain generalisation. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, volume 34, pages 13025–13032, 2020.
- [63] K. Zhou, Y. Yang, T. Hospedales, and T. Xiang. Learning to generate novel domains for domain generalization. In *European Conference on Computer Vision (ECCV)*, pages 561–578, 2020.
- [64] K. Zhou, Y. Yang, Y. Qiao, and T. Xiang. Domain generalization with mixstyle. In *International Conference on Learning Representations (ICLR)*, 2021.