# Towards Out-Of-Distribution Generalization: A Survey

Zheyan Shen*, Jiashuo Liu*, Yue He, Xingxuan Zhang, Renzhe Xu, Han Yu, Peng Cui†, *Senior Member, IEEE*

**Abstract**—Classic machine learning methods are built on the $i.i.d.$ assumption that training and testing data are independent and identically distributed. However, in real scenarios, the $i.i.d.$ assumption can hardly be satisfied, rendering the sharp drop of classic machine learning algorithms' performances under distributional shifts, which indicates the significance of investigating the Out-of-Distribution generalization problem. Out-of-Distribution (OOD) generalization problem addresses the challenging setting where the testing distribution is unknown and different from the training. This paper serves as the first effort to systematically and comprehensively discuss the OOD generalization problem, from the definition, methodology, evaluation to the implications and future directions. Firstly, we provide the formal definition of the OOD generalization problem. Secondly, existing methods are categorized into three parts based on their positions in the whole learning pipeline, namely unsupervised representation learning, supervised model learning and optimization, and typical methods for each category are discussed in detail. We then demonstrate the theoretical connections of different categories, and introduce the commonly used datasets and evaluation metrics. Finally, we summarize the whole literature and raise some future directions for OOD generalization problem. The summary of OOD generalization methods reviewed in this survey can be found at http://out-of-distribution-generalization.com.

**Index Terms**—Out-of-Distribution Generalization, Causal Inference, Invariant Learning, Stable Learning, Representation Learning, Distributionally Robust Optimization

✦

## 1 INTRODUCTION

Modern machine learning techniques have illustrated their excellent capabilities in many areas like computer vision, natural language processing and recommendation, etc. While enjoying the human-surpassing performance in experimental conditions, many researches have revealed the vulnerability of machine learning model when exposed to data with different distributions. Such massive gap is induced by the violation of a fundamental assumption that training and test data are identically and independently distributed (a.k.a. $i.i.d.$ assumption), upon which most of the existing learning models are developed. In many real cases where $i.i.d.$ assumption can hardly be satisfied, especially those high-stake applications such as healthcare, military and autonomous driving, instead of generalization within the training distribution, the ability to generalize under distribution shift is of more critical significance. Therefore, the investigation of out-of-distribution generalization is of great urgency in both academic and industry fields.

Despite the importance of OOD generalization problem, classic supervised learning methods can not be directly deployed to deal with this problem. Theoretically, one of the most fundamental assumptions of classic supervised learning is the $i.i.d.$ assumption, which assumes that the training and testing data are independent and identically distributed. However, distributional shifts are inevitable

---

\* *Equal Contributions*
† *Corresponding Author*
E-mail: *shenzy17@mails.tsinghua.edu.cn, liujiashuo77@gmail.com, heyue18@mails.tsinghua.edu.cn, xingxuanzhang@hotmail.com, xrz199721@gmail.com, h-yu16@hotmail.com, cuip@tsinghua.edu.cn.*
*All authors are with the Department of Computer Science and Technology, Tsinghua University, Beijing, China.*

in OOD generalization problem, which ruins the $i.i.d.$ assumption and renders classic learning theory inapplicable. Empirically, classic supervised learning methods typically are optimized by minimizing their training errors, which greedily absorb all correlations found in data for prediction. Though proved to be effective in $i.i.d$ settings, it would hurt the performance under distributional shifts, since not all correlations will hold in unseen testing distributions. As shown in many literatures [1], [2], [3], [4], [5], when involving strong distributional shifts, models optimized solely with training errors fail dramatically, and sometimes are even worse than random guess, which indicates the urgency to design methods for OOD generalization problems.

In order to deal with the OOD generalization problem, there remain several vital problems to be solved. Firstly, since training and testing data can be drawn from different distributions, how to formally characterize the distributional shifts is still an open problem. In the OOD generalization literature, different branches of methods adopt different ways to model the potential testing distribution. Domain generalization methods [6], [7], [8], [9] mainly focus on real scenarios and utilize data collected from different domains. Causal learning methods [2], [10], [11] formulate training and testing distributions with causal structures and the distributional shifts mainly originate from interventions or confounders. Stable learning methods [4], [12], [13] introduce distributional shifts via selection bias. Secondly, how to design an algorithm with good OOD generalization performance is the research hot spot and there are many branches of methods with different research focuses, including unsupervised representation learning methods, supervised learning models and optimization methods. Thirdly,

the evaluation of the OOD performance of different methods remains challenging, which requires specific datasets and evaluation metrics, as the classic ways for $i.i.d.$ setting are inapplicable under distributional shifts. And this also motivates the generation of different datasets and evaluations.

In this paper, we aim to provide a systematic and comprehensive review of research efforts in terms of a fairly broader sense of OOD generalization, covering the whole life cycle of OOD problem from the definition, methodology, evaluation to the implications and future directions. To the best of our knowledge, we serve as the first effort to discuss out-of-distribution generalization in such a large scope and self-contained form. There exist some previous works discussing the related topics. For example, [14], [15] mainly focus on the discussion of domain generalization; [16] discuss the evaluation benchmarks for OOD generalization. Each of the previous works serves as a piece of puzzle for the whole out-of-distribution generalization problem, and in this work, we organically integrate all the ingredients in a clear and concise way. Specifically, we divide the existing methods into three categories based on their positions in the whole learning pipeline. We also elaborate the theoretical connections between different methods through the lens of causality. To promote the future research on OOD generalization, we provide an exhaustive survey on how to evaluate the learning methods under distribution shifts.

The structure of this paper is organized as follows. We formulate the out-of-generalization problem, discuss its relationship with existing research areas and provide a categorization of methods in Section 2. In the following Section 3, 4, 5, we describe the representative methods of each category in detail respectively. We provide some theoretical connections and insights between different methods in Section 6. We also summarize the applicable benchmarks for OOD generalization and its possible implications in Section 7, 8. Finally, we conclude this paper in Section 9 and raise some promising directions for future research.

## 2 PROBLEM DEFINITION AND CATEGORIZATION OF METHODS

In this section, we first formulate the general out-of-distribution (OOD) generalization problem, illustrate its connections and differences with classic $i.i.d.$ learning problem. Then we categorize the existing methods which aim at addressing OOD generalization into several strands based on their positions through the whole learning pipeline.

### 2.1 Problem Definition

Let $\mathcal{X}$ be the feature space and $\mathcal{Y}$ the label space. A parametric model is defined as $f_\theta : \mathcal{X} \rightarrow \mathcal{Y}$, which serves as a mapping function from original features to the label with learnable paratemeter $\theta$. A loss function $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathcal{R}$, which measures the distance between predicted label and groundtruth. Then we can define the classic supervised learning problem.

**Problem 1** (Supervised Learning). *Given a set of $n$ training samples of the form $\{(x_1, y_1), \ldots, (x_n, y_n)\}$ which are drawn*

*from training distribution $P_{tr}(X, Y)$, a supervised learning problem is to find an optimal model $f_\theta^*$ which can generalize best on data drawn from test distribution $P_{te}(X, Y)$:*

$$f_\theta^* = \arg \min_{f_\theta} \mathbb{E}_{X,Y \sim P_{te}}[\ell(f_\theta(X), Y)]. \tag{1}$$

### 2.1.1 *i.i.d. Learning*

Traditional learning algorithms usually assume that the training samples and test samples are both $i.i.d.$ realizations from a common underlying distribution, which means $P_{tr}(X, Y) = P_{te}(X, Y)$. Based on such hypothesis, Empirical Risk Minimization (ERM) [17] which minimizes the average loss on training samples could obtain an optimal model which can successfully generalize to test distribution [18]. Specifically, ERM minimizes the following objective:

$$\mathcal{L}_{ERM} = \frac{1}{n} \sum_{i=1}^{n} \ell(f_\theta(x_i), y_i). \tag{2}$$

### 2.1.2 *Out-of-Distribution Generalization*

The admirable properties brought by $i.i.d.$ assumption offer a firm ground for the development of a plethora of learning models during the last few decades. However, in real cases, the test distribution on which model been deployed may deviate from training distribution [19], that is $P_{tr}(X, Y) \neq P_{te}(X, Y)$. The distribution shift can come into reality for many reasons such as the temporal/spatial evolution of data or the sample selection bias in data collection process. In any case, it renders Problem 1 more complicated than $i.i.d.$ learning scenario. Moreover, the test distribution we may encounter is usually unknown due to the nature of applications like stream-based online scenario, where test data are generated in the future. To sum up, the general out-of-distribution (OOD) generalization problem can be defined as the instantiation of supervised learning problem where the test distribution $P_{te}(X, Y)$ shifts from the training distribution $P_{tr}(X, Y)$ and remains unknown during the training phase.

In general, the OOD generalization problem is infeasible unless we make some assumptions on how test distribution may change. Among the different distribution shifts, the covariate shift is the most common one and has been studied in depth. Specifically, the covariate shift is defined as $P_{tr}(Y|X) = P_{te}(Y|X)$ and $P_{tr}(X) \neq P_{te}(X)$, which means the marginal distribution of $X$ shifts from training phase to test phase while the label generation mechanism keeps unchanged. In this work, we mainly focus on the covariate shift and decribe various efforts which have been made to handle it [1].

A relative field with OOD generalization is domain adaptation, which assumes the availability of test distribution either labeled ($P_{te}(X, Y)$) or unlabeled ($P_{te}(X)$). In a sense, domain adaptation can be seen as a special instantiation of OOD generalization where we have some prior knowledge on test distribution. Under such mild conditions, domain adaptation could enjoy theoretical guarantees [22] which retain the optimality of trained model on test distribution. The detailed illustration of domain adaptation methods is beyond the scope of this paper, curious readers may

---

1. Other forms of distribution shifts like label shift [20] and concept shift [21] are developed relatively independently and can be found in some well established surveys respectively.

refer to some well established surveys [23], [24], [25], [26]. In this paper, we are more concerned about the general out-of-distribution generalization problem, which aims to improve the efficiency of model when evaluated on unknown test distribution.

## 2.2 The Categorization of OOD Methods

To tackle the challenges brought by unknown distribution shift, tremendous efforts have been made in out-of-distribution generalization, resulting in a rich literature of related methods. The adopted techniques vary greatly ranging from causality to representation learning and from structure-based to optimization-based. However, to the best of our knowledge, little effort has been made to systematically and comprehensively survey these diverse methods in terms of a fairly broader sense of OOD generalization, as well as clarify the differences and connections between these work. In this paper, we try to first fill this gap by reviewing the related methods of OOD generalization.

Generally speaking, the supervised learning problem defined as equation 1 can be divided into three relatively independent components, namely (1) the representation of features $X$ (e.g. $g(X)$); (2) the mapping function $f_\theta(X)$ from features $X$ (or $g(X)$) to the label $Y$, which generally also known as model or the inductive bias; (3) the formulation of optimization objective. Therefore, we categorize the existing methods into three parts based on their positions in the whole learning pipeline accordingly:

- **Unsupervised Representation Learning for OOD Generalization** includes Disentangled Representation Learning and Causal Representation Learning, which exploits the unsupervised representation learning techniques (e.g. variational Bayes) to embed prior knowledge into learning process.
- **Supervised Model Learning for OOD Generalization** includes Causal Learning, Stable Learning and Domain Generalization, which designs various model architectures and learning strategies to achieve OOD generalization.
- **Optimization for OOD Generalization** includes Distributionally Robust Optimization and Invariance-Based Optimization, which directly formulates the objective of OOD generalization and conduct optimization with theoretical guarantees for OOD optimality.

In the following sections, we provide a comprehensive and detailed review of these methods corresponding to the above order and discuss their differences and theoretical connections. An overview of OOD generalization methods is shown in Table 1.

## 3 UNSUPERVISED REPRESENTATION LEARNING

In this section, we review methods focusing on the unsupervised representation learning, which can be mainly divided into two branches, namely disentangled representation learning and causal representation learning. These methods leverage human's prior knowledge to design and restrict the representation learning procedure, which endows the learned representation with certain properties that are potentially helpful for OOD generalization.

### 3.1 Disentangled Representation Learning

Disentangled representation learning aims to learn representations where distinct and informative factors of variations in data are seperated [27], [30], which is considered as one property of good representation and potentially benefits out-of-distribution generalization. Works for fulfilling disentanglement, typified by VAE-based methods [28], [29], emphasize both interpretability and sparsity, where sparsity means small distribution changes usually mainifest themselves in a sparse or local way in the disentangled factorization [36].

$\beta$-VAE [28] introduces an extra hyperparameter $\beta$ into vanilla VAE objective function, making a trade-off between latent bottleneck capacity and independence constraints, thus encouraging the model to learn more efficient representation. The objective function of $\beta$-VAE is as follows:

$$\mathcal{L} = \mathbb{E}_{q(z|x)}[\log p(x|z)] - \beta \mathrm{KL}(q(z|x)\|p(z)) \qquad (3)$$

where $z$ is the latent representation, $x$ observed data, $p(z)$ the prior distribution of latent factors, $p(x|z)$ the decoding distribution, and $q(z|x)$ the encoding posterior distribution, When $\beta$ is set to $1.0$, this formulation degenerates to vanilla VAE, and when $\beta$ is appropriately tuned, $\beta$-VAE can learn disentangled representation from data in an unsupervised way.

FactorVAE [29] adds the term of Total Correlation into the objective function, which is formulated as the KL-divergence between marginal posterior $q(z)$ and its corresponding factorized distribution $\bar{q}(z)$:

$$\mathcal{L} = \mathbb{E}_{q(z|x)}[\log p(x|z)] - \mathrm{KL}(q(z|x)\|p(z)) - \gamma \mathrm{KL}(q(z)\|\bar{q}(z)) \quad (4)$$

where $\bar{q}(z) := \prod_{j=1}^{d} q(z_j)$. This formulation encourages independence for the posterior latent representation. Since the Total Correlation term cannot be computed directly, an extra discriminator is added for density ratio estimation.

Despite of the success of disentangled representation learning, Locatello *et al.* [30] challenge some common assumptions of unsupervised disentangled representation learning (e.g., independence of latent factors). It also questions whether disentanglement can improve downstream task performances, inspiring later works to take downstream tasks into consideration, OOD generalization performance included.

However, whether disentangled representation benefits OOD generalization remains controversial. Leeb *et al.* [31] conduct some quantitative extrapolation experiments, finding that the learned disentangled representation fails to extrapolate to unseen data, while Träuble *et al.* [33] and Dittadi *et al.* [32] empirically verify the ability to generalize under OOD circumstances. The advantage of disentangled representation on OOD tasks still requires further research and discussion.

### 3.2 Causal Representation Learning

Similar to conventional disentangled representation learning, causal representation learning aims to learn variables in the causal graph in an unsupervised or semi-supervised way. Further, causal representation can be viewed as the ultimate goal of disentanglement, which satisfies the informal definition of disentangled representation in terms of interpretability and sparsity. With the learned causal

| Unsupervised Representation Learning (**Section 3**) | Disentangled Representation Learning (**Section 3.1**) | [27], [28], [29], [30], [31], [32], [33](Section 3.1) |
|---|---|---|
| | Causal Representation Learning (**Section 3.2**) | [34], [35], [36](Section 3.2) |
| Supervised Learning Models (**Section 4**) | Domain Generalization (**Section 4.1**) | [7], [8], [9], [37], [38], [39], [40], [41], [42], [43], [44], [45], [46], [47], [48], [49] [37], [39], [50], [51], [51], [52], [53], [53], [54], [54], [55], [56], [57], [58], [59], [60] [6], [61], [62], [63], [63], [64], [65], [66], [67], [68], [69], [70](Section 4.1.1) [52], [71], [72], [73], [74], [75], [76], [77], [78], [79], [80], [81], [82] [83], [84], [85], [86], [87], [88], [89], [90], [91], [92](Section 4.1.2) [93], [94], [95], [96], [97], [98], [99], [100], [101], [102], [103], [104], [105](Section 4.1.3) [14], [15](Other surveys) |
| | Causal Learning & Invariant Learning (**Section 4.2**) | [10], [106], [107], [108], [109], [110](Section 4.2.2) [2], [3], [56], [111], [112], [113], [114], [115], [116], [117](Section 4.2.3) |
| | Stable Learning (**Section 4.3**) | [4], [12], [13], [118], [119], [120], [121], [122], [123], [124](Section 4.3) |
| Optimization Methods (**Section 5**) | Distributionally Robust Optimization (**Section 5.1**) | [1], [125], [126], [127], [128], [129], [130], [131], [132], [133], [134](Section 5.1.1∼ 5.1.3) [135], [136], [137], [138] (Section 5.1.4) |
| | Invariance-Based Optimization (**Section 5.2**) | [5], [139], [140](Section 5.2) |

TABLE 1
An overview of OOD generalization methods.

representation, one can capture the latent data generation process, which can help to resist the distributional shifts induced by interventions.

In real scenarios where observations are made in the form of images or sentences instead of structured data, high-level abstract information needs to be extracted from low-level data [27], and a few existing works [34], [35], [36] propose to recover causal factorization through disentanglement.

CausalVAE [34] combines linear Structural Causal Model (SCM) into VAE model to endow the learned latent representation with causal structure. Specifically, the causal structure is depicted by an adjacency matrix $A$ as:

$$z = A^T z + \epsilon = (I - A^T)^{-1}\epsilon, \quad \epsilon \sim \mathcal{N}(0, I) \quad (5)$$

where $\epsilon$ represents the exogenous factors. In practice, a mild nonlinear function $g_i$ is introduced for stability as $z_i = g_i(A_i \cdot z; \eta_i) + \epsilon_i$. Further, extra labels $u$ of latent causal variables are used in CausalVAE, which gives the objective function as:

$$\mathcal{L} = -\text{ELBO} + \alpha \text{DAG}(A) + \beta l_u + \gamma l_m \quad (6)$$

where ELBO represents the Evidence Lower Bound, $\text{DAG}(A)$ the Directed Acyclic Graph (DAG) constraint, $l_u = \mathbb{E}_{q_\mathcal{X}} ||u - \sigma(A^T u)||_2^2$ measures how well $A$ describes causal relations among labels, and $l_m = \mathbb{E}_{z \sim q_\phi} \Sigma_{i=1}^n ||z_i - g_i(A_i \cdot z; \eta_i)||_2^2$ measures how well $A$ describes causal relations among latent codes. $q_\mathcal{X}$ is the empirical data distribution and $q_\phi$ the approximate posterior distribution.

Move a step on, DEAR [35] incorporates nonlinear SCM with a bidirectional generative model and assumes the known causal graph structure and extra supervised information of latent factors. The objective function is given as:

$$\mathcal{L}(E, G, F) = \mathcal{L}_{gen}(E, G, F) + \mathcal{L}_{sup}(E) \quad (7)$$

where $E, G$ denotes the encoder and generator, respectively. The first part $\mathcal{L}_{gen}(E, G, F) = \text{KL}(q_E(x, z), p_{G,F}(x, z))$ resembles the VAE loss. The difference lies in the prior distribution of $z$. In DEAR, this prior is generated by the nonlinear SCM, while in vanilla VAE, it is simply a factorized

Gaussian. The second part is $\mathcal{L}_{sup}(E) = \mathbb{E}_{x,y}\text{CE}(\bar{E}(x), u)$, where CE is the cross entropy loss function, $\bar{E}$ represents the deterministic part of $E$ and $u$ the extra labels.

# 4 SUPERVISED MODEL LEARNING FOR OOD GENERALIZATION

Apart from learning representations solely unsupervisedly, there are branches of works incorporating supervised information to design various model architectures and corresponding learning strategies. In this section, we review such methods that focus on end-to-end model learning to achieve OOD generalization ability, including domain generalization, causal & invariant learning and stable learning.

## 4.1 Domain Generalization

Incorporating data from several source domains, Domain Generalization (DG) aims to learn models that generalizes well on unseen target domains, which focuses mostly on computer vision related classification problems as predictions are prone to be affected by disturbance on images (e.g., style, background, light, rotation, etc.). Regarding to different methodological focuses, similar to Wang *et al.* [14], we divide DG methods into three branches, namely representation learning, training strategy and data augmentation. Here we briefly introduce the three branches of methods, and one can refer to existing surveys of this field [14], [15] for more details of DG methods.

### 4.1.1 Representation Learning for DG

Representation learning remains a critical part in DG. There are works [6], [44] that theoretically or empirically prove that if the representations remain invariant when the domain varies, the representations are transferable and robust on different domains. Methods that attempts to learn invariant representations among different domains can be mainly divided into three categories, namely domain adversarial

learning, domain alignment and kernel based methods.

**Domain adversarial learning.** Ganin *et al.* [7], [8] propose domain-adversarial neural network (DANN) for domain adaptation. DANN learns representations that are discriminative and invariant to the change of domains by jointly optimizing the underlying features, a label predictor that predicts class labels and is used both in the training and inference phase, and a domain classifier that discriminates between the source and the target domains during training. The representations are trained to confuse the domain classifier so that domain invariant features are learned. Li *et al.* [9] adopt this idea for DG and Gong *et al.* [42] further extend adversarial training to a manifold space. Li *et al.* [40] propose to learn class-specific adversarial networks via a conditional invariant adversarial network (CIAN). Garg *et al.* [48] and Sicilia *et al.* [49] provide theoretically guarantees via generalization bounds for domain adversarial training. Rahman *et al.* [47] introduce a correlation-aware adversarial framework for both Domain Adaptation (DA) and DG that jointly uses the correlation alignment metric and adversarial learning to decrease the domain discrepancy of the source and target data. Shao *et al.* [41], Jia *et al.* [43] and Wang *et al.* [45] propose to use domain adversarial learning for face anti-spoofing and unseen target stance detection. Zhao *et al.* [46] propose an additional entropy-regularization approach that learns invariant representations via minimizing the KL-divergence between the conditional distribution of different source domains.

**Domain alignment.** Apart from domain adversarial learning, some works [9], [37], [39] propose to learn domain invariant representations via the alignment of the features. Motiian *et al.* [38] propose to learn semantic alignment between different domains by minimizing the distance between samples from different domains but the same class, and maximizing the distance between samples from different domains and classes. Some works minimize the divergence of feature distributions by minimizing the maximum mean discrepancy (MMD) distance [37], [39], [50], [51], Wasserstein distance [52], the second order correlation [53], [54], [55], etc., for DA or DG.

**Feature normalization.** There are works [51], [53], [54], [56] adopting feature normalization to constraint the domain discrepancy. Pan *et al.* [57] present IBN-Net, which integrates both Instance Normalization (IN) [58] and Batch Normalization (BN) [59] as building blocks, and learns to capture and eliminate domain variance. Empirically, Pan *et al.* [57] find that IN provides visual and appearance invariance but may harm the discriminative information of representations. Thus they propose to combine IN and BN in shallow layers and only BN in deeper layers. Huang *et al.* [60] propose to enable arbitrary style transfer via an adaptive IN. Nam *et al.* [61] assume that each feature map of a convolutional encoder can be divided into a style-related part and shape-related part, and explicitly combine BN and IN to learn style-invariant representations. Jin *et al.* [62] propose a Style Normalization and Restitution (SNR), which distills task-related features from the residual after the style normalization to ensure the discrimination.

**Kernel methods.** Kernel methods are also extensively applied in DG. Blanchard *et al.* [63], [64] first address the domain generalization problem via kernel methods and propose to learn a domain-invariant kernel with the training data. Muandet *et al.* [6] propose a classic kernel method for DG named Domain-Invariant Component Analysis (DICA), which minimizes distributional variance between samples from source domains. Grubinger *et al.* [65] propose Multi-TCA, an extension of transfer component analysis (TCA) [50] for transfer learning, which enables TCA to deal with multiple source and target domains by learning the shared subspace among source domains. Gan *et al.* [66] extend DICA with attribute regularization. Erfani *et al.* [67] propose Elliptical Summary Randomization (ESRand), which projects domains into a latent space and minimizes the dissimilarity between data distributions with randomized kernel and elliptical data summarization. Ghifary *et al.* [68] raise a unified framework named Scatter Component Analyis (SCA) for both DA and DG. Based on scatter operating on reproducing kernel Hilbert space, SCA learn a domain invariant representation space by trading among minimizing the mismatch between domains, maximizing the discrepancy between classes and maximizing the variance of the whole data. Hu *et al.* [69] present a more fine-grained Multi-domain Discriminant Analysis (MDA) and analyze the bound on excess risk and generalization error for kernel-based domain invariant representation learning methods. Other theoretical analysis for kernel methods for DG are introduced in [63], [70].

### 4.1.2 Training strategy

Given that DG focuses mostly on computer vision related classification problems, some works study training strategies to enhance the generalization ability of deep models on image data, which can be divided into four categories, namely meta learning, ensemble learning, unsupervised/semi-supervised DG and others.

**Meta learning.** Meta learning gains experience over multiple learning episodes via an alternative learning paradigm [141]. Finn *et al.* [71] propose model-agnostic meta-learning (MAML) for DA, which introduces the concept of "episodes" in the training phase and greatly influence the research of meta learning for DG. Li *et al.* [72] first apply meta-learning to DG and some improvements are made in the following works [73], [74], [75], [76], [77], [78], [79], [80], [81]. The main idea of meta learning for DG is to divide the source domains into meta-train and meta-test, where the loss of meta-train and meta-test are optimized simultaneously.

**Model-ensemble learning.** Typically, model-ensemble learning based methods learn ensembles of multiple specific models for different source domains to improve the generalization ability. Some works adopt domain specific subnetworks for different source domains, together with one single classifier [52], [82], [83] or multiple domain-specific classifier heads [84]. Others use domain-specific batch normalization for different domains to learn a better normalization [85], [86].

**Unsupervised/Semi-supervised DG.** Recently, inspired by unsupervised and semi-supervised domain adaptation [142], [143], [144], some works propose to enhance model's generalization ability by unsupervised or semi-supervised learning [87], [88]. Zhang *et al.* [87] propose Domain-Irrelevant Unsupervised Learning (DIUL) to cope with the

distribution shifts between source domains and target domains. They select valid negative samples for any given queue samples according to the similarity between different domains to learn domain-irrelevant representations.

**Others** Inspired by the self-supervised learning method [145], Carlucci *et al.* [89] combine a self-learning puzzles task with the classification task to learn robust representations. Ryu *et al.* [90] propose to sample positive and negative samples via random forest. Li *et al.* [91] alternatively train convolutional layers and the classifier. Huang *et al.* [92] introduce a self-challenging dropout algorithm to prevent the model from overfitting to source domains.

### 4.1.3 Data augmentation

Data augmentation is a common and effective method in deep learning especially in computer vision. The generalization ability of deep models largely depends on the heterogeneity of available data. Thus heterogeneity brought by data augmentation could prevent overfitting and improve the generalization ability. The data augmentation methods for DG can be divided into randomization based augmentation [93], [94], [95], [96], [97], [98], gradient based augmentation [99], [100], [101] and generation based augmentation [102], [103], [104], [105].

## 4.2 Causal & Invariant Learning

Comparing with domain generalization which typically targets on vision tasks, causal learning and invariant learning stems from causal inference literature and addresses the OOD generalization problem from a more principle way, which aims to explore causal variables for prediction and becomes more practical recently. We firstly introduce some preliminaries of causal learning, and then review branches of typical works.

### 4.2.1 Preliminaries

For causal learning methods, it is often assumed that there exist data heterogeneity and causal relationship inside data. Specifically, causal learning assumes that one has access to data from multiple environments, as demonstrated in Definition 1.

**Definition 1** (Heterogeneity of Data). *Consider a dataset $D = \{D^e\}_{e \in \text{supp}(\mathcal{E}_{tr})}$, which is a mixture of data $D^e = \{X^e, Y^e\}$ collected from multiple training environments $e \in \text{supp}(\mathcal{E}_{tr})$, $X^e \subset \mathcal{X}$ and $Y^e \subset \mathcal{Y}$ are the collection of data and label from environment $e$, respectively. $\mathcal{E}_{tr}$ and $\mathcal{E}$ are random variables on indices of training environments and all possible environments, respectively, such that $\text{supp}(\mathcal{E}) \supset \text{supp}(\mathcal{E}_{tr})$. $P^e$ denotes the distribution of data and label in environment $e$. Usually, for all $e \in \text{supp}(\mathcal{E}) \setminus \text{supp}(\mathcal{E}_{tr})$, the data and label distribution $P^e(X, Y)$ can be quite different from that of training environments $\mathcal{E}_{tr}$.*

**Assumption 1** (Causality Assumption [11]). *The structural equation models:*

$$Y^e \leftarrow f_Y(X_{pa(Y)}^e, \epsilon_Y^e), \ \epsilon_Y^e \perp X_{pa(Y)}^e \qquad (8)$$

*remains the same across all environments $e \in supp(\mathcal{E}_{all})$, that is, $\epsilon_Y^e$ has the same distribution as $\epsilon_Y$ for all environments. $pa(Y)$ denotes the direct causes of $Y$.*

And the assumption of causality is given by Assumption 1, which originates from causal inference literature and assumes the causally invariant relationship between the target $Y$ and its direct causes $X_{\text{pa}(Y)}$. Such assumptions indicate that causal variables $X_{\text{pa}(Y)}$ are stable across different environments or data selection biases, which motivates branches of works to achieve OOD generalization via leveraging only causal variables.

### 4.2.2 Causal Inference-Based Methods

Firstly, we review methods related to causal inference, which try to obtain causal variables from heterogeneous data.

It is well known that a gold standard for identifying causal effect of a variable is to conduct randomized experiments like A/B testing, but fully randomized experiments are usually expensive and even infeasible in real applications. Since causal inference or causal structural learning is very ambitious, these inference techniques should be considered as "groundtruth" but not necessarily able to be realized in practice (like typical machine learning settings). Therefore, it is more practical to design such techniques that have a more "causal explanation" than a standard regression or classification framework and could also gain some sort of invariance across environments. Follow such intuition, a strand of methods [10], [106], [107], [108], [109], [110] have been developed by leveraging the heterogeneity inside data (e.g., multiple environments).

**Assumption 2** (Invariance Assumption). *There exists a subset $S^* \subseteq \{1, \ldots, p\}$ of the covariate indices (including the empty set) such that*

$$P(Y^e | X_{S^*}^e) \text{ is the same for all } e \in \mathcal{E}. \qquad (9)$$

*That is, when conditioning on the covariates from $S^*$ (denoted by $X_{S^*}^e$), the conditional distribution is invariant across all environments from $\mathcal{E}$.*

Peters *et al.* [10] first try to investigate the fact that "invariance" could, to some extent, infer the causal structure under necessary conditions and propose Invariant Causal Prediction(ICP). Specifically, they leverage the fact that when considering all direct causes of a target variable, the conditional distribution of the target given the direct causes will not change when interfering all other variables in the model except the target itself. Then they perform a statistical test whether a subset of covariates $S$ satisfies the invariance assumption 2 for the observed environments in $\mathcal{E}$. The null-hypothesis for testing is:

$$H_{0,S}(\mathcal{E}) : \text{ invariance assumption holds.}$$

and all subsets of covariates $S$ which lead to invariance are intersected, that is:

$$\hat{\mathcal{S}}(\mathcal{E}) = \bigcap_S \{S; \ H_{0,S}(\mathcal{E}) \text{ not rejected by test at significance level } \alpha\}.$$

Under the assumption of structural equation model and gaussian residual described in [10], ICP with Chow test [146] could, at least with controllable probability 1-$\alpha$, discover subsets of true causal variables, which reads as:

$$\mathbb{P}[\hat{\mathcal{S}}(\mathcal{E}) \subseteq \text{pa}(Y)] \geq 1 - \alpha, \qquad (10)$$

where $\mathrm{pa}(Y)$ denotes the direct causes of target $Y$ (e.g. the parental variables of $Y$ in the causal graph).

Though being the first attempt to connect the invariance to causality, ICP has several limitations. The most straightforward one is the strict requirements for heterogeneity, since the power of ICP depends highly on the quality of available environments $\mathcal{E}_{tr}$(or perturbations). If the available perturbed subpopulations are not enough, or even single environment, the efficacy of ICP will be lost. As discussed in [106], naively estimating the environments from data and then applying ICP may yield less powerful results, so instead of using static data, Pfister *et al.* [106] propose to leverage the sequential data from non-stationary environment to detect instantaneous causal relations in multivariate linear time series, which relaxes the assumption of known environments. Besides environmental specification, there are other works trying to consolidate the coverage of the so-called invariance-based method. For example, Heinze-Deml *et al.* [108] extend the ICP into non-linear model and continuous environments; Gamella *et al.* [109] apply the ICP into an active learning setting where the interventions (a.k.a. environments) can be proactively chosen during training.

ICP serves as a milestone towards inferring causal structure via invariance property. However, the invariance assumption may be violated in more complicated scenarios. Among which, the most common case is the existence of hidden confounders. Instrument variable(IV) method is one typical method for dealing with hidden confounders, which require the instrument variable $E$ not to directly act on hidden confounding variable $H$ and outcome variable $Y$, as shown in Figure 1a.



(a) SCM of traditional instrument variable model
(b) SCM of anchor regression model

Fig. 1. Comparision of SCM between IV model and Anchor Regression model.

Rothenhäusler *et al.* [107] investigate more relaxed conditions than the standard IV model which allow the direct effect of instrument variable (which they called anchor variabls) on $H$ and $Y$, as shown in Figure 1b. They realize that despite the attractive notion of invariance guarantee against arbitrarily large intervention or perturbation, one seldom encounters such extreme cases, and exact invariance could be too conservative for moderately perturbed data. Specifically, they the following structural equation:

$$Y = X^T \beta + H^T \alpha + A^T \xi + \epsilon_Y, \tag{11}$$

with $X \in \mathbb{R}^p$, $H \in \mathbb{R}^q$ and $A \in \mathbb{R}^r$, and proposed a regularized formulation of ordinary least squares model by the error projection to the space spanned by anchor variables.

$$\hat{\beta}(\gamma) = \mathrm{argmin}_b \left( \|(I - \Pi_{\mathbf{A}})(\mathbf{Y} - \mathbf{X}b)\|_2^2/n + \gamma \|\Pi_{\mathbf{A}}(\mathbf{Y} - \mathbf{X}b)\|_2^2/n \right). \tag{12}$$

where $\Pi_{\mathbf{A}}$ denotes the projection in $\mathbb{R}^n$ onto the column space of $\mathbf{A}$. For $\gamma = 1$, $\hat{\beta}(1)$ equals the ordinary least squares

estimator, for $\gamma \to \infty$ it obtains the two-stage least squares procedure from IV regression.

Sometimes the observation of instrument $A$ is still hard to fulfill, Oberst *et al.* [110] further relax the assumption by introducing the noisy proxy of $A$ and prove the robustness of the method under bounded shifts.

### 4.2.3 Invariant Learning

Deriving from causal inference based methods, invariant learning methods, typified by invariant risk minimization(IRM, [2]), target on latent causal mechanisms and extend ICP to more practical and general settings. Different from causal prediction methods that assumes on raw variable level, IRM makes the invariance assumption as:

**Assumption 3** (IRM's Invariance Assumption). *There exists data representation $\Phi(X)$ such that for all $e, e' \in supp(\mathcal{E}_{tr})$, $\mathbb{E}[Y|\Phi(X^e)] = \mathbb{E}[Y|\Phi(X^{e'})]$, where $\mathcal{E}_{tr}$ denotes the available training environments.*

which generalizes the former invariance assumption to representation level. Then Arjovsky *et al.* [2] propose to find data representation $\Phi(X)$ that can both predict well and elicit an invariant linear predictor $w$ across $\mathcal{E}_{tr}$, which results in the following objective function:

$$\min_{\Phi(X),w} \sum_{e \in \mathrm{supp}(\mathcal{E}_{tr})} \mathcal{L}^e(w \odot \Phi(X), Y) \tag{13}$$

$$\text{s.t. } w \in \arg\min_{\overline{w}} \mathcal{L}^e(\overline{w} \odot \Phi(X)), \text{ for all } e \in \mathrm{supp}(\mathcal{E}_{tr}) \tag{14}$$

In order to achieve invariance across $\mathcal{E}_{all}$ by enforcing low error of equation 13 on $\mathcal{E}_{tr}$, IRM requires sufficient diversity across environments and makes the following assumption.

**Assumption 4** (IRM's Condition, Assumption 8 in [2]). *A set of training environments $\mathcal{E}_{tr}$ lie in linear general position of degree $r$ if $|\mathcal{E}_{tr}| > d - r + \frac{d}{r}$ for some $r \in \mathbb{N}$, and for all non-zero $x \in \mathbb{R}^d$:*

$$\dim\left(\mathrm{span}\left(\{\mathbb{E}_{X^e}[X^e X^{eT}]x - \mathbb{E}_{X^e,\epsilon^e}[X^e \epsilon^e]\}_{e \in \mathcal{E}_{tr}}\right)\right) > d - r \tag{15}$$

With Assumption 4, IRM proves that in linear case, if a representation $\Phi(X)$ of rank $r$ elicits an invariant predictor $w \circ \Phi(X)$ across $\mathcal{E}_{tr}$ and $\mathcal{E}_{tr}$ lies in linear general position of degree $r$, then $w \circ \Phi(X)$ is invariant across $\mathcal{E}_{all}$(Therorem 9 in [2]). The assumptions about linearity, centered noise, and independence between the noise $\epsilon^e$ and the causal variables from the theoretical analysis in IRM (Theorem 9 in [2]) also appear in ICP [10], while IRM does not assume as ICP [10] that the data is Gaussian, the existence of a causal graph, or that the training environments arise from specific types of interventions. And the result extends to latent causal variables while ICP [10] restricts to raw causal feature level.

Based on IRM which is motivated by the representation $\Phi(X)$ such that $\mathbb{E}[Y|\Phi(X)]$ remains invariant, follow-up works have proposed variations on this objective with stronger regularization of the invariance assumption 3, resulting in similar alternatives. Ahuja *et al.* [111] incorporate game theory and replace the linear classifier in IRM with an ensemble of classifiers from different environments. Jin *et al.* [56] replace the regularizer of IRM with a predictive regret and impose stronger constraints on $\Phi(X)$. Krueger

*et al.* [112] propose to penalize the variance of the risks across environments, while Xie *et al.* [113] raise almost the same objective but replace the original penalty with the square root of the variance. Mahajan *et al.* [114] introduce a contrastive regularizer which matches the representation of same objects across environments. And Creager *et al.* targets on the problem of missing environment labels of IRM, and raise Environment Inference for Invariant Learning(EIIL [3]) to learn environments that maximize the IRM's penalty. The whole algorithm is two-stage, which firstly generates environments according to a biased reference model, and then performs invariant learning with learned environments.

While the results in IRM seems promising, Rosenfeld *et al.* [115] point out some problems of IRM on classification tasks. In the linear case, they give simple conditions under which the optimal solution succeeds or, more often, fails to recover the optimal invariant predictor. Specifically, Rosenfeld *et al.* [115] show that there exists a feasible solution which uses *only the environmental features* yet performs better than the optimal invariant predictor on all $e \in \mathcal{E}_{all}$ (Theorem 5.3 in [115]). In nonlinear regime, they show that IRM can fail catastrophically unless the test data are sufficiently similar to the training distribution (Theorem 6.1 in [115]. And Kamath *et al.* [116] also demonstrate that it is possible for IRM to learn a sub-optimal predictor, due to the loss function not being invariant across environments. Ahuja *et al.* [117] compare IRM with ERM from the sample complexity perspective in different shift patterns (Table 1 in [117]), which indicates that under covariate shifts, IRM has no obvious advantage to ERM, and for other distribution shifts such as those involving confounders or anti-causal variables, IRM is guaranteed to be close to the desired OOD solutions in the finite sample regime.

## 4.3 Stable Learning

Compared with domain generalization and causal learning, stable learning inspires another way for incorporating causal inference with machine learning, which significantly relaxes the requirements for multiple environments. The problem setting of stable learning is as following:

**Problem 2** (Settings of Stable Learning)**.** *Given training data* $D^e = (X^e, Y^e)$ *from one environment* $e \in supp(\mathcal{E}_{all})$, *the goal of stable learning is to learn a predictive model with* **uniformly good performance** *on any possible environments in* $supp(\mathcal{E}_{all})$.

To solve such difficult problem, motivated by the literature of variable balancing strategies [147], [148], [149], Shen *et al.* [118] propose to consider all the variables as the treatment and learn a set of global sample weights that could remove the confounding bias for all the potential treatments from data distribution. They derive a global balancing loss which can be easily plugged into standard machine learning tasks as a regularizer, as depicted by equation 16:

$$\sum_{j=1}^{p} \left\| \frac{X_{-j}^T \cdot (W \odot I_j)}{W^T \cdot I_j} - \frac{X_{-j}^T \cdot (W \odot (1 - I_j))}{W^T \cdot (1 - I_j)} \right\|_2^2, \quad (16)$$

where $W$ denotes the sample weights, $\left\| \frac{X_{-j}^T \cdot (W \odot I_j)}{W^T \cdot I_j} - \frac{X_{-j}^T \cdot (W \odot (1-I_j))}{W^T \cdot (1-I_j)} \right\|_2^2$ represents the loss of confounder balancing when setting feature $j$ as treatment variable, and $X_{-j}$ is the remaining features (i.e. confounders) except $j^{th}$

column. The $I_j$ means the $j^{th}$ column of $I$, and $I_{ij}$ refers to the treatment status of unit $i$ when setting feature $j$ as treatment variable. By minimizing the global balancing loss, the confounding bias could be removed in a global scale.

Further, Kuang *et al.* [12] incorporate unsupervised feature representation into the global balancing stage with auto-encoders [150] and modify the original regularizer to a "deep" version as equation 17:

$$\sum_{j=1}^{p} \left\| \frac{\phi(\mathbf{X}_{\cdot, -j})^T \cdot (W \odot \mathbf{X}_{\cdot j})}{W^T \cdot \mathbf{X}_{\cdot, j}} - \frac{\phi(\mathbf{X}_{\cdot, -j})^T \cdot (W \odot (1 - \mathbf{X}_{\cdot, j}))}{W^T \cdot (1 - \mathbf{X}_{\cdot, j})} \right\|_2. \quad (17)$$

Above methods are limited to binary features since the main stream discussions about causal inference are within the scope of binary treatment. When the treatment variable is categorical or continuous, traditional balancing methods are no long feasible as the treatment level could be infinitely large. To mitigate such limitation and address the continuous treatment problem, Kuang *et al.* [13] propose to learn a set of sample weights so that the weighted distribution of treatment and confounder could satisfy the independent condition, which corresponds to the fact that if the treatment and confounder are independent, the treatment effect can be estimated accurately.

Apart from methods addressing confounder bias, Shen *et al.* [4] target on the model mis-specification problem for linear model in stable learning. The main challenge of stable learning in linear model is the inevitable model mis-specification in real scenarios. In other words, besides the linear part, the true generation process actually contains an additional mis-specification term which could be non-linear term or interactions between input variables.

$$y = x^\top \bar{\beta}_{1:p} + \bar{\beta}_0 + b(x) + \epsilon \quad (18)$$

They find that the collinearity between variables play an crucial role on learning a stable model. If a mis-specified model is used at the training time, the existence of collinearity among variables can inflate a small misspecification error to arbitrarily large, thus causes instability of prediction performance across different distributed test data.

In order to alleviate the collinearity among variables, Shen *et al.* [4] propose to learn a set of sample weights which can make the design matrix near orthogonal. Technically, they construct an uncorrelated design matrix $\tilde{X}$ from original $X$ as the 'oracle', and learn the sample weights $w(x)$ by estimating the density ratio $w(x) = p_{\tilde{D}}(x)/p_D(x)$ of underlying uncorrelated distribution $\tilde{D}$ and original distribution $D$.

Further, to mitigate the large variance and shrinkage of the effective sample size brought by sample reweighting, Shen *et al.* [119] propose to leverage the unlabeled data collected from multiple environments and recover the hidden cluster structures among variables. Under several technical assumptions, they [119] prove that, decorrelating the variables between clusters instead of each other would be sufficient to achieve a stable estimation while preventing the variance inflation.

Recently, Zhang *et al.* [120] propose StableNet, as shown in figure 2, which extends former linear frameworks [4], [12], [13] to incorporate deep models. As the complex non-linear dependencies among features obtained in deep model are much more difficult to measure and eliminate than the

Fig. 2. The overall architecture of StableNet from [120]. LSWD refers to *learning sample weighting for decorrelation* and *Final loss* is used to optimized the classification network.

linear case, to remedy for this, StableNet proposes a novel nonlinear feature decorrelation approach based on Random Fourier Features (RFF) [151]. Specifically, StableNet iteratively optimizes sample weights $\mathbf{w}$, representation function $f$, and prediction function $g$ as follows:

$$
\begin{aligned}
f^{(t+1)}, g^{(t+1)} =&\arg\min_{f,g} \sum_{i=1}^n w_i^{(t)} \mathcal{L}(g(f(\mathbf{X}_i)), y_i), \\
\mathbf{w}^{(t+1)} =&\arg\min_{\mathbf{w}\in\Delta_n} \sum_{1\le i<j\le m_Z} \left\| \hat{\Sigma}_{\mathbf{Z}_{:,i}^{(t+1)}\mathbf{Z}_{:,j}^{(t+1)};\mathbf{w}} \right\|_F^2 .
\end{aligned}
\tag{19}
$$

where $\mathbf{Z}^{(t+1)} = f^{(t+1)}(\mathbf{X})$, $\mathcal{L}(\cdot, \cdot)$ represents the cross entropy loss function and $t$ represents the time stamp. The sample reweighting module and the representation learning network are jointly optimized in StableNet and they can effectively partial out the environment related features and leverage truly category related and discriminative features for prediction, leading to more stable performances in the wild non-stationary environments.

Moreover, Kuang *et al.* [124] subsample the data to reduce the confounding effects induced by the distributional shifts. Wang *et al.* [121] incorporate the decorrelation mechanism into clustering and achieve better clustering performance under data selection bias. Zhang *et al.* [122] propose a Deconfounded Visio-Linguistic Bert framework to mitigate the potential data biases. And Yuan *et al.* [123] propose to identify causal features with meta-learning mechanism for OOD generalization.

## 5 OPTIMIZATION FOR OOD GENERALIZATION

To address the Out-of-Distribution (OOD) generalization problem, apart from unsupervised representation learning and end-to-end learning models, optimization methods with theoretical guarantees have recently aroused much attention, which are both model agnostic and data structure agnostic. In this section, we firstly introduce the objective of these OOD optimization methods, and then review the optimization methods, including Distributionally Robust Optimization [1], [126], [135], [152], [153] and Invariant-Based Optimization [5], [139], [140].

In order to address the problem from the optimization perspective, the OOD generalization problem is formulated to controlling the worst-case prediction error among $\mathcal{E}_{all}$, which takes the form of:

$$
\arg\min_f \max_{e\in\text{supp}(\mathcal{E}_{all})} \mathcal{L}(f|e) \tag{20}
$$

where $\mathcal{E}_{all}$ is the random variable on indices of all possible environments, and for all $e \in \text{supp}(\mathcal{E}_{all})$, the data and label distribution $P^e(X, Y)$ can be quite different from that of training distribution $P_{tr}(X, Y)$; $\mathcal{L}(f|e) = \mathbb{E}[l(f(X), Y)|e] = \mathbb{E}^e[l(f(X^e), Y^e)]$ is the risk of predictor $f$ on environment $e$, and $l(\cdot, \cdot) : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}^+$ is the loss function. Intuitively, optimization methods aim to guarantee the worst-case performance under distributional shifts.

### 5.1 Distributionally Robust Optimization

Distributionally robust optimization (DRO), from robust optimization literature, directly solves the OOD generalization problem by optimizing for the worst-case error over an uncertainty distribution set, so as to protect the model against the potential distributional shifts within the uncertainty set, which is often constrained by moment or support conditions [154], [155], $f$-divergence [1], [153], [156] and Wasserstein distance [126], [135], [152]. The objective function of DRO methods can be summarized as:

$$
\arg\min_f \sup_{Q\in\mathcal{P}(P_{tr})} \mathbb{E}_{X,Y\sim Q}[\ell(f(X), Y)] \tag{21}
$$

where $\mathcal{P}(P_{tr})$ is the distribution set lying around the training distribution $P_{tr}$ and $\ell(\cdot; \cdot) : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}^+$ is the loss function. Different DRO methods adopt different kinds of constraints to formulate the distribution set $\mathcal{P}(P_{tr})$ and correspondingly different optimization algorithm. In this paper, we only introduce two typical DRO methods whose distribution sets are formulated by $f$-divergence and Wasserstein distance respectively, and for a more thorough introduction to DRO methods, one can refer to [157].

#### 5.1.1 $f$-Divergence Constraints

The distribution set $\mathcal{P}(P_{tr})$ in $f$-divergence DRO [1] is formulated as:

$$
\mathcal{P}(P_{tr}) = \{Q : D_f(Q\|P_{tr}) \le \rho\} \tag{22}
$$

where $\rho > 0$ controls the extent of the distributional shift, and $D_f(Q\|P_{tr}) = \int f(\frac{dQ}{dP_{tr}})dP_{tr}$ is the $f$-divergence between $Q$ and $P_{tr}$. Intuitively, if the potential testing distribution $P^{e_{test}}(X, Y) \in \mathcal{P}(P_{tr})$, DRO methods can achieve good generalization performance even if $P^{e_{test}}(X, Y) \neq P_{tr}(X, Y)$. As for the optimization, a simplified dual formulation for the Cressie-Read family of $f$-divergence can be obtained.

**Lemma 1** (Optimization of $f$-divergence [1]). *For* $f_k(t) = \frac{t^k - kt + k - 1}{k(k-1)}$ *and* $k \in (1, +\infty)$, $k_* = k/(k-1)$, *and any* $\rho > 0$, *we have for all* $\theta \in \Theta$:

$$\mathcal{R}_k(\theta; P_{tr}) = \inf_{\eta \in \mathbb{R}} \left\{ c_k(\rho) \mathbb{E}_{P_{tr}}[(\ell(f(X), Y) - \eta)_+^{k_*}]^{\frac{1}{k_*}} + \eta \right\} \quad (23)$$

*where* $c_k(\rho) = (k(k-1)\rho + 1)^{\frac{1}{k}}$.

### 5.1.2 Wasserstein Distance Constraints

Since the calculation of $f$-divergence requires the supports of two distributions to be the same while Wasserstein distance does not, the distribution set $\mathcal{P}(P_{tr})$ formulated by Wasserstein distance is more flexible. Wasserstein distance is defined as:

**Definition 2.** *Let* $\mathcal{Z} \subset \mathbb{R}^{m+1}$ *and* $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$, *given a transportation cost function* $c : \mathcal{Z} \times \mathcal{Z} \rightarrow [0, \infty)$, *which is nonnegative, lower semi-continuous and satisfies* $c(z, z) = 0$, *for probability measures* $P$ *and* $Q$ *supported on* $\mathcal{Z}$, *the Wasserstein distance between* $P$ *and* $Q$ *is :*

$$W_c(P, Q) = \inf_{M \in \Pi(P,Q)} \mathbb{E}_{(z,z') \sim M}[c(z, z')] \quad (24)$$

*where* $\Pi(P, Q)$ *denotes the couplings with* $M(A, \mathcal{Z}) = P(A)$ *and* $M(\mathcal{Z}, A) = Q(A)$ *for measures* $M$ *on* $\mathcal{Z} \times \mathcal{Z}$.

Then the distribution set $\mathcal{P}(P_{tr})$ of Wasserstein DRO is formulated as:

$$\mathcal{P}_c(P_{tr}) = \{Q : W_c(Q, P_{tr}) \leq \rho\} \quad (25)$$

where the subscript $c$ denotes the transportation cost function $c(\cdot, \cdot)$. However, Wasserstein DRO is difficult to optimize, and works targeting different models and transportation cost functions have been proposed. Wasserstein DRO for logistic regression was proposed by Abadeh *et al.* [125]. Sinha *et al.* [126] achieved moderate levels of robustness with little computational cost relative to empirical risk minimization with a Lagrangian penalty formulation of WDRL.

### 5.1.3 Robustness Guarantees

Here we briefly review some theoretical results in DRO literatures, including the relationship with regularization and robustness guarantees.

In order to demonstrate how the robust formulation (21) provides distributional robustness, several works establish the relationship between distributional robustness and regularization. For norm-based DRO methods, El Ghaoui *et al.* [127] build the equivalence between the worst-case squared residual within a Frobenius norm-based distribution set and the Tikhonov regularization. Xu *et al.* [128] show the equivalence between robust linear regression with feature perturbations and the Least Absolute Shrinkage and Selection Operator(LASSO), and Yang *et al.* [129] and Bertsimas *et al.* [130] make further extensions. For $f$-divergence-based DRO methods, Duchi *et al.* [133] show that the formulation (21) with distribution set $\mathcal{P}(P_{tr}) = \mathcal{P}_{\rho,n} = \{p \in \mathbb{R}^n : p^T \mathbf{1} = 1, p \geq 0, D_f(p\|\mathbf{1}/n) \leq \rho/n\}$ is a convex approximation to regularizing the empirical risk by variance. For Wasserstein-based DRO methods, Shafieezadeh-Abadeh *et al.* [125] investigate the Wasserstein DRO of logistic regression, and show that the regularized logistic regression is one special

case of it. Chen *et al.* [134] also build the connection between the Wasserstein DRO of linear regression with $\ell_1$ loss function and regularization constraints on the regression coefficients. And Shafieezadeh-Abadeh *et al.* [131] and Gao *et al.* [132] connect the Wasserstein DRO and regularizations in a unified framework.

As for the OOD generalization ability, in fact, the guarantees for OOD generalization of DRO methods naturally derive their formulation (21). Since DRO methods directly optimize for the worst-case risk within the distribution set $\mathcal{P}(P_{tr})$, as long as the potential testing distribution $P_{te} \in \mathcal{P}(P_{tr})$, the OOD generalization ability is guaranteed. Therefore, the remaining work is to provide the finite sample convergence guarantees, which ensure that the population-level objective $\sup_{Q \in \mathcal{P}(P_{tr})} \mathbb{E}_Q[\ell(f(X), Y)]$ can be optimized empirically with finite samples.

**Theorem 1** (Theorem 2 in [1]). *Assume that* $\ell(\theta; x) \in [0, M]$ *for all* $\theta \in \Theta$ *and* $x \in \mathcal{X}$, *and define* $c_k = (k(k-1)\rho + 1)^{1/k}$. *For a fixed* $\theta \in \Theta$ *and* $t > 0$, *with probability at least* $1 - 2e^{-t}$,

$$|\mathcal{R}_k(\theta; \hat{P}_{tr,n}) - \mathcal{R}_k(\theta; P_{tr})| \leq c_k M n^{-\frac{1}{k_* \vee 2}} \left(\frac{2}{k} + \sqrt{t}\right) \quad (26)$$

*where* $\mathcal{R}_k(\theta; P) = \sup_{Q \ll P_{tr}} \{\mathbb{E}_Q[\ell(f_\theta(X), Y)] : D_{f_k}(Q\|P_{tr}) \leq \rho\}$, *and* $k_* = \frac{k}{k-1}$.

Given by Duchi *et al.* [1], Theorem 1 shows that the pointwise concentration of the finite sample objective $\mathcal{R}(\theta; \hat{P}_{tr,n})$ converges to its population counterpart, which also gives that for the empirical minimizer $\hat{\theta}_n$ with probability at least $1 - 2N(\mathcal{F}, \frac{t}{3}, \|\cdot\|_{L^\infty(\mathcal{X})})e^{-t}$

$$\mathcal{R}_k(\hat{\theta}_n; P_{tr}) \leq \inf_{\theta \in \Theta} \mathcal{R}_k(\theta; P_{tr}) + 2c_k M n^{-\frac{1}{k_* \vee 2}} \left(\frac{2}{k} + \sqrt{6}t\right) \quad (27)$$

where $N(V, \epsilon, \|\cdot\|)$ is the covering number of $V$ with respect to $\|\cdot\|$, $\mathcal{F} = \{\ell(\theta, \cdot) : \theta \in \Theta\}$ equipped with sup-norm $\|h\|_{L^\infty(\mathcal{X})} = \sup_{x \in \mathcal{X}} \|h(x)\|$. Sinha *et al.* [126], Chen *et al.* [134] and Liu *et al.* [135] also provide similar generalization bounds for Wasserstein DRO as Theorem 2.

**Theorem 2** (Theorem 3 in [126]). *Assume* $|\ell(f_\theta(X), Y)| \leq M_\ell$ *for all* $\theta \in \Theta$ *and* $X \in \mathcal{X}, Y \in \mathcal{Y}$. *Then for a fixed* $t > 0$ *and numerical constants* $b_2, b_2 > 0$, *with probability at least* $1 - e^{-t}$, *simultaneously for all* $\theta \in \Theta$, $\rho \geq 0$, $\gamma \geq 0$,

$$\sup_{Q : W_c(Q, P_{tr}) \leq \rho} \mathbb{E}_Q[\ell(f_\theta(X), Y)] \leq \gamma\rho + \mathbb{E}_{\hat{P}_{tr,n}}[\phi_\gamma(\theta; X, Y)] + \epsilon_n(t) \quad (28)$$

### 5.1.4 Over-Pessimism Problem

Although DRO methods could theoretically guarantee the out-of-distribution generalization ability when $P^{e_{test}}(X, Y) \in \mathcal{P}(P_{tr})$, there has been work questioning their real effects in practice. Intuitively, in order to achieve good OOD generalization ability, the potential testing distribution should be captured in the built distribution set. However, in real scenarios, to contain the possible true testing distribution, the uncertainty set is often overwhelmingly large, making the learned model make decisions with fairly low confidence, which is also referred to as the low confidence problem. Specifically, Hu *et al.* [136] proved that in classification tasks, DRO ends up being optimal for the training distribution $P_{tr}$, which is due to the over-flexibility of the built distribution set. And Fronger *et al.* [137] also

pointed out the problem of overwhelmingly-large decision set for Wasserstein DRO.

In order to overcome such problem, Blanchet *et al.* propose a data-driven way to select the transportation cost function. Fronger *et al.* [137] propose to further restrict the distribution set with large number of unlabeled data. Liu *et al.* [135] notice that in real scenarios, different covariates may be perturbed in a non- uniform way, and form a more reasonable distribution set according to the stability of covariates across environments. Duchi *et al.* [138] assume that $P(Y|X)$ stays invariant and propose to only perturb the marginal distribution $P(X)$ to deal with covariate shifts.

## 5.2 Invariance-Based Optimization

Different from DRO methods which directly optimize for the worst-case without any additional assumptions, invariance-based optimization methods [5], [139], [140] assume the invariance property inside data and leverages multiple environments to find such invariance for generalizing under distributional shifts. In this section, we introduce two typical methods, namely Maximal Invariant Predictor (MIP, [139], [140]) and Heterogeneous Risk Minimization (HRM, [5]).

Chang *et al.* [139] and Koyama *et al.* [140] formulate the desired invariant representation using information theory, and propose to find the maximal invariant predictor (MIP) across training environments to control the worst-case error in equation 20. The maximal invariant predictor is defined as:

**Definition 3.** *The invariance set $\mathcal{I}$ with respect to $\mathcal{E}$ is defined as:*

$$\begin{aligned} \mathcal{I}_{\mathcal{E}} &= \{\Phi(X) : Y \perp \mathcal{E}|\Phi(X)\} \\ &= \{\Phi(X) : H[Y|\Phi(X)] = H[Y|\Phi(X), \mathcal{E}]\} \end{aligned} \quad (29)$$

*where $H[\cdot]$ is the Shannon entropy of a random variable. The corresponding maximal invariant predictor (MIP) of $\mathcal{I}_{\mathcal{E}}$ is defined as:*

$$S_{\mathcal{E}} = \arg \max_{\Phi \in \mathcal{I}_{\mathcal{E}}} I(Y; \Phi) \quad (30)$$

*where $I(\cdot; \cdot)$ measures Shannon mutual information between two random variables.*

With the invariant predictor $S_{\mathcal{E}_{all}}$, Koyama *et al.* [140] prove that the OOD optimal model is given by $\mathbb{E}[Y|S_{\mathcal{E}_{all}}]$. Further, to obtain the MIP solution from training environments $\mathcal{E}_{tr}$, Koyama *et al.* [140] derive a regularizer as:

$$\text{trace}\left(\text{Var}_{\mathcal{E}_{tr}}(\nabla_\theta \mathcal{L}^e(\theta))\right) \quad (31)$$

where the variance is taken with respect to training environments $\mathcal{E}_{tr}$. Under the controllability condition proposed by Koyama *et al.* [140] which assumes that there exists an environment $e$ such that $X \perp Y|\Phi(X), e$, the optimality of $\mathbb{E}[Y|\Phi(X)]$ can be verified, as shown in Proposition 3.1 in [140] and Theorem 1 in [158].

Move one step on, Liu *et al.* [5] theoretically characterize the role of environments in invariance-based optimization methods. Since $\text{supp}(\mathcal{E}_{tr}) \subset \text{supp}(\mathcal{E}_{all})$, $\mathcal{I}_{\mathcal{E}_{all}} \subseteq \mathcal{I}_{\mathcal{E}_{tr}}$ holds, which shows that the invariance set regularized by training environments is likely to contain undesired variant components and makes the controllability condition for such optimality in [140] is hard to satisfy since in practice. Further, in real scenarios, data are often collected from



Fig. 3. The framework of Heterogeneous Risk Minimization from [5].

different sources without explicit environment labels, which renders the above methods inapplicable. Inspired by this, they propose Heterogeneous Risk Minimization (HRM, [5]), an optimization framework to achieve joint learning of the latent heterogeneity among the data and the invariant predictor. As shown in Figure 3, HRM contains two interactive parts, the frontend $\mathcal{M}_c$ for heterogeneity identification and the backend $\mathcal{M}_p$ for invariant prediction. Given the pooled heterogeneous data, it starts with the heterogeneity identification module $\mathcal{M}_c$ leveraging the learned variant representation $\Psi(X)$ to generate heterogeneous environments $\mathcal{E}_{learn}$. Then the learned environments are used by OOD prediction module $\mathcal{M}_p$ to learn the MIP $\Phi(X)$ as well as the invariant prediction model $f(\Phi(X))$. After that, the better variant part $\Psi(X)$ is derived from the learned MIP $\Phi(X)$ to further boost the process of heterogeneity identification.

# 6 THEORETICAL CONNECTIONS

For branches of methods for OOD generalization, there are some inherent connections among them. In this section, we will demonstrate the connections among causal learning methods, distributionally robust optimization (DRO) methods and stable learning methods, which may benefit the understanding of OOD generalization methods.

## 6.1 DRO and Causality

Recall that DRO methods aim to optimize the worst-case error over a pre-defined distribution set, so as to protect the learned model from potential distributional shifts, which often takes the form of:

$$\arg \min_f \sup_{Q \in \mathcal{P}(P_{tr})} \mathbb{E}_{X,Y \sim Q}[\ell(f(X), Y)] \quad (32)$$

where $\mathcal{P}(P_{tr})$ is the distribution set built around the training distribution $P_{tr}$. Although in DRO literatures, $\mathcal{P}(P_{tr})$ is often characterized by $f$-divergence or Wasserstein distance, different choices of $\mathcal{P}(P_{tr})$ will render DRO equivalent to causal inference in the structural equation model (SEM) context [159], which shows the inherent relationship between causality-based methods and DRO methods. Taking linear equation models for example, suppose we have a directed acyclic graph $G = (V, E)$ with $p$ nodes $V = \{1, \ldots, p\}$ and correspondingly a $p$-dimension random variable $Z$, then the training distribution is determined by the structural causal model (SCM) as:

$$Z = BZ + \epsilon \quad (33)$$

where $Z = (X, Y) \in \mathbb{R}^p$ is the random variable of interest, $B \in \mathbb{R}^{p \times p}$ is the coefficient matrix and $\epsilon \sim P_\epsilon$ the random noise. We will show that finding causal coefficients for predicting $Y$ can be reformulated as performing DRO on interventional distribution set, including do-interventional and shift-interventional distributions.

Do-interventions on variables $S \subseteq V$ can be formulated as:

$$Z_k = (BZ)_k + \epsilon \qquad \text{for } k \notin S \qquad (34)$$
$$Z_K = A_k \qquad \text{for } k \in S \qquad (35)$$

where $A \in \mathbb{R}^p$ and the value of the do-intervention on variable $k \in S$ is $A_k$. Then the error distribution $P_\epsilon$, coefficient matrix $B$, intervention set $S \subseteq V$ and intervention value $A \in \mathbb{R}^p$ induces a distribution for a random variable $Z(A, S)$, denoted as $Z(A, S) \sim P_{A,S}^{(do)}$. And the corresponding do-interventional distribution set can be formulated as $\mathcal{P}^{(do)} = \{P_{A,V/\{p\}}^{(do)} : A \in \mathbb{R}^p\}$. Analogously to do-interventions, the shift-interventions is defined as:

$$Z = BZ + \epsilon + A \qquad (36)$$

where $A \in \mathbb{R}^p$ is the shift direction, and the induced distribution is denoted as $Z(A) \sim P_A^{(shift)}$ and the shift-interventional distribution set can be formulated as $\mathcal{P}^{(shift)} = \{P_A^{(shift)} : A_p = 0\}$.

When performing DRO on $\mathcal{P}^{(do)}$ or $\mathcal{P}^{(shift)}$, causal coefficients can be obtained [159] since

$$\min_\theta \sup_{Q \in \mathcal{P}^{(do)}} \mathbb{E}[\ell(f_\theta(X), Y)] = \begin{cases} \infty, & \text{if } \theta \neq \theta_{causal} \\ \text{Var}(\epsilon_p), & \text{if } \theta = \theta_{causal} \end{cases} \qquad (37)$$

and

$$\min_\theta \sup_{Q \in \mathcal{P}^{(shift)}} \mathbb{E}[\ell(f_\theta(X), Y)] = \begin{cases} \infty, & \text{if } \theta \neq \theta_{causal} \\ \text{Var}(\epsilon_p), & \text{if } \theta = \theta_{causal} \end{cases} \qquad (38)$$

which reveals that causal inference can also be viewed as a special case of distributional robustness.

### 6.2 Stable Learning and Causality

Stable learning algorithms have connections with causality, which can be considered as a feature selection mechanism according to the regression coefficients. These algorithms will eliminate all variant variables, and the invariant variables selected for predicting $Y$ are closely related to Markov blankets and Markov boundaries [160] in local causal discovery literature.

Fig. 4. Three data-generating processes in [13].

Here take the data-generating processes studied in [13] as example, as shown in Figure 4. In all three cases, $V$ are the variant variables to predict $Y$, i.e., they are not correlated with outcome $Y$ (Figure 4a) or influence $Y$ by means of $S$ (Figure 4b and 4c). The algorithm proposed by Shen *et al.* [4] will eliminate $V$, i.e., the coefficients of weighted least squares on $V$ will be zero. For simplicity, assume $S, V$, and $Y$ are all scalars, and in these cases, $Y \perp V \mid S$. Consider any proper weighting function $w(S, V)$ such that $S \perp V$ in the weighted distribution $\tilde{P}(S, V)$. The existence of $w$ can be fulfilled under proper assumptions. Then coefficients of weighted least squares under infinite samples are given by:

$$\begin{pmatrix} \beta_S \\ \beta_V \end{pmatrix} = \begin{pmatrix} \text{Var}_{\tilde{P}}(S) & \text{Cov}_{\tilde{P}}(S, V) \\ \text{Cov}_{\tilde{P}}(S, V) & \text{Var}_{\tilde{P}}(V) \end{pmatrix}^{-1} \begin{pmatrix} \text{Cov}_{\tilde{P}}(S, Y) \\ \text{Cov}_{\tilde{P}}(V, Y) \end{pmatrix}. \qquad (39)$$

In addition, the conditional distribution $P(Y|S, V)$ will stay the same if weighting function $w$ is defined on $(S, V)$ only. Hence, conditional independence $Y \perp V \mid S$ is also satisfied in the weighted distribution $\tilde{P}(S, V, Y)$. As a result, $Y \perp V$ in $\tilde{P}$, which makes $\text{Cov}_{\tilde{P}}(V, Y) = 0$. Hence,

$$\beta_V = (\text{Var}_{\tilde{P}}(V))^{-1} \text{Cov}_{\tilde{P}}(V, Y) = 0. \qquad (40)$$

Actually, the assumption $Y \perp V \mid S$ is closely related to Markov blankets and Markov boundaries and $S$ is a Markov blanket with respect to $Y$ by definition. One straightforward corollary is that stable learning algorithms [4] will eliminate the variables that do not belong to the Markov boundary with respect to outcome $Y$.

## 7 DATASETS AND EVALUATIONS METRICS

To promote the development of OOD generalization research, it is of vital importance to evaluate the OOD generalization performances of different algorithms reasonably and accurately. The evaluation of an algorithm typically consists of two disentangled components, datasets and evaluation metrics. In this section, we summarize the datasets and evaluation metrics that are commonly used as OOD benchmarks in existing publications.

### 7.1 Datasets

Datasets can be classified according to different criterions, for example, synthetic data and real-world data, raw feature data and complicated data, and different research fields utilize different kinds of datasets, for example, traditional statistical learning often uses synthetic data; computer vision often uses real-world image data. As for OOD generalization problem, different from traditional machine learning tasks which are based on *i.i.d.* assumption, it is necessary to involve or generate the distributional shifts to simulate the unknown testing distribution, so as to test whether an algorithm can generalize to unseen distributions. Therefore, in line with recent works on OOD generalization, we find it necessary to use both the synthetic & simple data and real-world & complicated data to verify the effectiveness of an OOD generalization method.

#### 7.1.1 Synthetic Data

Synthetic data are important tools to simulate explainable and controllable distributional shifts. As Aubin *et al.* [161] find that recent OOD generalization methods perform poorly on some simple low-dimensional linear problems, it is necessary to test OOD generalization methods on such simple but challenging data, which can reflect whether and to what extent an algorithm can resist certain kind of distributional shifts.

Generally speaking, there are three mechanisms to simulate distributional shifts across environments, named by selection bias, confounding bias and anti-causal effect, with which one can simulate certain kind of distributional shifts to various degrees and clearly justify an algorithm's real effect. In these mechanisms, the covariates $X$ are divided into two groups as $X = [S, V]^T$, corresponding to the invariant and variant parts inside data. And it is assumed that $P(Y|S)$ remains invariant across environments, and $P(Y|V)$ is perturbed with different mechanisms, which brings distributional shifts.

**(a) Selection Bias**

Kuang *et al.* [13] propose a selection bias mechanism to introduce distributional shifts, and similar settings are also adopted in [5], [135]. In this setting, $P(Y|V)$ is perturbed with selection bias. The data generation process is as following:

$$Y = f(S) + \epsilon = \theta_S^T S + \beta S_1 \cdot S_2 \cdot S_3 + \epsilon \qquad (41)$$

and the selection probability $\hat{P}(x, y)$ of certain data point $(x, y)$ is calculated as:

$$\hat{P} = \prod_{v_i \in V} |r|^{-5*|f(S) - \text{sign}(r) * v_i|} \qquad (42)$$

where $|r| > 1$ is the selection factor and controls the distributional shifts. Intuitively, $r$ eventually controls the strengths and direction of the spurious correlation between $V$ and $Y$(i.e. if $r > 0$, a data point whose $V$ is close to its $Y$ is more probably to be selected.). The larger value of $|r|$ means the stronger spurious correlation between $V$ and $Y$, and $r \geq 0$ means positive correlation and vice versa. Since the distributional shift is absolutely controlled by $r$, one can adopt different $r$ to simulate different extents of distributional shifts, as done in [5], [13], [135].

**(b) Confounding Bias**

Confounding bias is also one of the most common means of distributional shifts [2], [162], [163]. Compared to selection bias, in this setting, the variant feature $V$ is related to target $Y$ owing to the unobserved confounder $C$. For example, the data generation process proposed by Subbaswamy *et al.* [163] is as following:

$$V = W_V^e C + \epsilon_V \qquad (43)$$
$$Y = W_S^T S + W_c C + \epsilon_Y \qquad (44)$$

where coefficient $W_V^e$ controls the relationship between $V$ and $Y$, and one can change $W_V^e$ across environments to simulate the distributional shifts.

**(c) Anti-Causal Effect**

Besides the above mechanisms, Arjovsky *et al.* [2] and Liu *et al.* [5] introduce an anti-causal mechanism to change $P(Y|V)$. In this setting, the data generation process is as following:

$$Y = W_S^T S + \epsilon_Y \qquad (45)$$
$$V = W_V^e Y + \epsilon_V^e \qquad (46)$$

where coefficient $W_V^e$ and experimental $\epsilon_y, \epsilon_V^e$ controls the relationship between $V$ and $Y$. Intuitively, larger $\epsilon_Y$ and smaller $\epsilon_V^e$ will make model easier to utilize $V$ for prediction, making OoD generalization more challenging.

### 7.1.2 Real-World Data

Synthetic data are easy to generate distributional shifts to varying degrees, and they can sufficiently test whether an algorithm can generalize to unseen distributions. However, synthetic data are rather simple and it is difficult to generate complicated data (e.g., images). Further, whether an algorithm can solve real-world distributional shift problems is also an important criterion to evaluate an OOD generalization method. Therefore, it is necessary to involve real-world datasets. Here, we describe the most popular real-world datasets used in OOD literature, including image datasets and other forms of dataset (e.g., tabular

data, language data). A summary of these datasets can be found in Table 2.

**Image Datasets**

With the rapid development of computer vision, there are a number of image datasets been raised, and we classify them into three categories with respect to the flexibility of simulating distributional shifts, namely synthetic transformation data, fixed wild data and controllable wild data.

(a) SYNTHETIC TRANSFORMATION DATA.

Although most image datasets are not produced for OOD generalization, people modify them with some synthetic transformations to simulate distributional shifts. The most typical ones, including ImageNet [164] variants (e.g. ImageNet-A [165], ImageNet-C [166], ImageNet-R [167]) adopt special data selection policy or perturbations to generate testing data with distributional shifts. Others, typified by MNIST [168] variants (e.g. Colored MNIST [2], MNIST-R [169]), simulate different environments by coloring or rotating original images. Being well-designed, these datasets make it available for preliminary study and effectiveness verification of OOD generalization algorithms.

(b) FIXED WILD DATA.

There are a few datasets built to support OOD generalization validation, which mainly utilize real-world backgrounds or environments. Widely used in domain generalization, PACS [170] and Office-Home [171] adopt the image style (e.g., art, cartoon) to differentiate domains/distributions, and VLCS [172] takes data collected independently from four sources as environments. Besides, Camelyon17 [173] contains tissue slides sampled and post-processed in different hospitals and DomainNet [174] extends PACS to a far larger scale, consisting of more domains and categories. Waterbirds [153] contains birds' images with either water or land backgrounds. iWildCam [175] collect images from different locations and produce realistic distributional shifts. Recently, Koh *et al.* collect several datasets together and produce Wilds [176] as a benchmark for OOD generalization.

(c) CONTROLLABLE WILD DATA.

Recently, there are datasets enabling more flexible and controllable ways to simulate distributional shifts, typified by NICO [177]. NICO [177] elaborately selects visual contexts with various types, including background, attribute, view and etc., and produce various environments. With diverse contexts, NICO could simulate different types of realistic distributional shifts, and with balanced sample size in each context, different degrees of distributional shifts could be easily achieved. These two properties make NICO enable flexible setups of distributional shifts. Besides, FMoW [178] collects satellite images of building or land use token at different times and regions, and PovertyMap [179] contains images of an urban or rural area from disjoint sets of countries.

**Other Datasets**

Besides image data, there are still many other datasets

related to OOD literature. A house sales price dataset on Kaggle[2] is used in [4], [5], [135], where environments are generated according to the built year of houses. OGB-MolPCBA [180] collects molecular graphs in over 100,000 scaffolds and formulate a molecular property prediction task across different scaffolds. CivilComments [181] and Amazon [182] gather the individual comments of different users and distinctive groups (e.g. male and female). Towards auto-engineering, Py150 [183] contains codes from 8,421 git repositories for code completion generalization.

## 7.2 Evaluation Metric

Besides datasets, evaluation metrics are also important for evaluating an OOD generalization algorithm. Compared to $i.i.d.$ problem where only one testing distribution is considered, there are often multiple testing distributions in OOD generalization problems, so as to better capture the unseen distributions. Further, Ye $et$ $al.$ [184] empirically find that the test accuracy in single environment would mislead the judgement to algorithms in OOD scenario. Therefore, to sufficiently evaluate the OOD generalization algorithms, more statistics of accuracy of multiple test environments should be taken into consideration. Here, we discuss three evaluation metrics, including the average accuracy, worst-case accuracy and the standard deviation of accuracies. For convenience, we assume the model's accuracies of $K$ testing environments are $\{acc_1, ..., acc_K\}$, respectively.

### 7.2.1 Average Accuracy

The average accuracy $\overline{Acc}$ over test distributions is the most straightforward way to assess the effectiveness of OOD algorithms, which is commonly used in OOD generalization literatures [4], [5], [13], which is calculated as:

$$\overline{Acc} = \frac{1}{K} \sum_{k=1}^{K} acc_k \qquad (47)$$

The average accuracy measures the overall performances among testing edistributions, but it cannot describe the volatility of an algorithm's performance. Further, the average accuracy equally treat all testing distributions without considering the property of each, which may be misled by distributions that frequently occurr. We think it possible to weigh the accuracy according to the discrepancy of each testing distribution with training, which is still an open problem.

### 7.2.2 Worst-Case Accuracy

Widely-used in DRO literatures [1], [137], [152], [185], the worst-case accuracy $Acc_{worst}$ is defined as the worst-case accuracy across testing distributions:

$$Acc_{worst} = \min_{k \in [K]} acc_k \qquad (48)$$

which corresponds to the objective functions [1], [137], [152], [185]. The worst-case accuracy reflects the reliability of an algorithm and is crucial in high-stake applications such as medical diagnoses, criminal justices and financial security.

2. https://www.kaggle.com/c/house-prices-advanced-regression-techniques/data

### 7.2.3 Standard Deviation (STD)

The STD of accuracies over testing distributions $Acc_{std}$ measures the variation of performance across different distributions, which is defined as:

$$Acc_{std} = \sqrt{\frac{1}{K-1} \sum_{k=1}^{K} (acc_k - \overline{Acc})^2} \qquad (49)$$

This metric measures the sensitivity of an algorithm and reflects the robustness and stability of algorithm, which is importance for an OOD generalization algorithm.

## 8 IMPLICATION FOR FAIRNESS AND EXPLAINABILITY

### 8.1 Fairness

Nowadays, fairness issues have raised great concerns in decision-making systems such as loan applications [186], hiring processes [187], and criminal justice [188]. Poorly designed algorithms tend to amplify data bias, resulting in discriminations against specific subgroups of individuals based on their inherent characteristics, which are often named sensitive attributes in fairness problems. Many works define their fairness and propose corresponding fair algorithms, from which the definition of fairness can be divided into three types: individual fairness [189], [190], group fairness [191], [192], [193], and causality-based fairness notions [194], [195], [196]. However, different fairness notions are in conflict [197]. Methods that mitigate unfairness in the algorithms fall under three categories: pre-processing [198], [199], [200], in-processing [201], [202], [203], and post-processing [191] algorithms.

Fairness has recently been linked to OOD issues, according to [3]. Generally speaking, subgroups split by sensitive attributes in fairness literature correspond to environments in OOD literature. Following that, both areas need to specify learning objectives with respect to the subgroups/environments. In fairness literature, the learning objectives represent context-specific fairness notions, while in OOD literature, the learning objectives should be designed according to invariance assumptions. Similar learning objectives could be adopted in both areas. For example, objectives similar to fairness criterion equalized odds [191] are adopted in OOD literature [40], [204] to deal with simplicity bias [205]. The learning objective of IRM [2] is also similar to calibration in fairness literature [206]. Meanwhile, classical approaches from OOD literature could be applied to address fairness issues. Fair representation learning methods [193], [207], [208], [209] originated from domain adaptation (DA) methods [8], [22]. When sensitive attributes are unknown, DRO and adversarially learning were introduced in fairness literature [210], [211], [212] to obtain a distributionally robust predictor and ensure the worst subgroup performance. [192], [213], [214] also adopt adversarially learning methods to ensure all computationally identifiable subgroups are treated equally. As a result, Pursuing OOD could be considered as pursuing fairness concerning the subgroups/environments if the invariance assumption adopted for OOD could be viewed as a fairness notion.

TABLE 2
Commonly used image datasets for OOD generalization. Shift type denotes the type of distributional shifts, and the mixed type in image type means that there are both real and unreal images.

| Image Data Set | ImageNet-Variant [165], [166], [167] | Colored MNIST [2] | MNIST-R [169] | Waterbirds [153] | Camelyon17 [173] | VLCS [172] | PACS [170] |
|---|---|---|---|---|---|---|---|
| # Domains | - | 3 | 6 | 2 | 5 | 4 | 4 |
| # Categories | - | 2 | 10 | 2 | 2 | 5 | 7 |
| # Examples | - | - | 6,000 | 4,800 | 45w | 2,800 | 10w |
| Shift Type | Adversarial Policy | Color | Angle | Background | Hospital | Data Source | Style |
| Image Type | Mixed Type | Digits | Digits | Birds | Tissue Slides | Real Objects | Mixed Type |

| Image Data Set | Office-Home [171] | DomainNet [174] | iWildCam [175] | FMoW [178] | PovertyMap [179] | NICO [177] | |
|---|---|---|---|---|---|---|---|
| # Domains | 4 | 6 | 323 | $16 \times 5$ | $23 \times 2$ | 188 | |
| # Categories | 65 | 345 | 182 | 62 | Real Value | 19 | |
| # Examples | 15w | 50w | 20w | 50w | 2w | 2.5w | |
| Shift Type | Style | Style | Location | Time, Location | Country, Urban/Rural | Background, Attribute, Action, View and Co-occurring Object | |
| Image Type | Mixed Type | Mixed Type | Real Animals | Satellite | Satellite | Real Objects | |

In addition to considering subgroups as environments, [215] investigate another scenario in which the environment is a separable variable. They studied fair classifiers that are robust to perturbations in the training distribution and devised a DRO-like method to reach their goal. Fair and robust learning is also applied in [216]. These works differ from works listed in the last paragraph in that fairness and robustness are two objectives here whereas the aforementioned works consider them the same.

## 8.2 Explainability

Explanation methods can be generally divided into post hoc analyses and model-based methods [217]. There exist several works in both directions. Post hoc analyses usually explain a black-box model by calculating feature importance [218]. Typical methods include gradient-based [219], [220], [221], [222], influence function [223], and Shapley values [224]. Model-based explanation methods often adopt simpler hypothesis such as linear regression [225], LASSO [226], generalized additive models [227], decision trees [225], and rule-based methods [228], [229].

Causality [162] has recently been introduced to model explanation, especially in deep learning methods. Traditional deep learning algorithms are rarely used in high-stakes applications due to their lack of explainability. Causality could provide a way to shed light on the explainability of deep learning. For example, several works adopt causality to explain deep models in textual and visual explanation [230], [231], [232]. Futhermore, the Causal And-Or Graph was proposed in robotics [233] and object tracking [234] to build explainable algorithms with the knowledge of causality. [235] also applied a causal filtering step in self-driving automobile problems.

$$\text{OOD} \longleftarrow \text{Causality} \longrightarrow \text{Explainability}. \qquad (50)$$

Actually, causality is the crux for both OOD and explainability as shown in equation 50. The models will have good OOD performance and explainability simultaneously if they utilize the causal relationship between the features and the outcomes. Hence, explainability would be a side product when pursuing OOD with causality.

## 9 CONCLUSION AND FUTURE DIRECTIONS

Out-of-Distribution generalization problem has aroused much research attention recently, and is critical for the deployment of machine learning algorithms. In this paper, we systematically and comprehensively review the definition, the main branches of methods, theoretical connections among different methods and the datasets and evaluation metrics of OOD generalization problem. Based on our analysis, we come up with several potential challenges that could be the directions of future research. And we hope this paper could inspire the future research of OOD generalization problem.

(A) THEORETICAL CHARACTERIZATION
Although growing popular recently, the theoretical characterization a learnable OOD generalization problem remains vague in recent literatures. Characterizing the learnability of a problem is a basic question in machine learning. Though previous research efforts have been made in $i.i.d.$ setting, the learnability is difficult to define and analyze under distributional shifts, since it is impossible to enable models to generalize to arbitrary and unknown distributions. Therefore, in OOD generalization problem, figuring out what kind of distributional shifts should be taken into consideration is critical for the analysis of learnability. There is very few exploration [236] on this and more research efforts need be paid on this.

(B) DEMANDS FOR ENVIRONMENTS
The majority of OOD generalization methods require multiple training environments. However, modern datasets are often assembled by merging data from multiple sources without maintaining source labels, which greatly restricts the deployment of OOD generalization methods in real scenarios. Therefore, it is more practical and realistic that we only have access to one training environment with latent heterogeneity. Recently, while there are some works [3], [5] try to leverage the latent heterogeneity and relax the demands for environments, how to explore and make good use of the latent heterogeneity inside data is critical for the deployment of OOD generalization methods and is a promising future direction.

## (C) REASONABLE EVALUATIONS

Although the evaluation criterions for classic machine learning algorithms under *i.i.d.* assumption are well-developed, including testing data, model selection mechanisms and so on, they cannot directly be deployed to OOD scenarios. Since the testing distribution is both unknown and different from the training, how to design reasonable and realistic experimental settings remains a challenging problem. Further, the model selection mechanism also matters, since the choice of validation data is non-trivial in OOD scenarios, and Gulrajani *et al.* [237] also demonstrate that domain generalization algorithms without a model selection strategy are incomplete. Also, Gulrajani *et al.* [237] notice that the real effects of many domain generalization methods are weak, which indicates that existing evaluation criterions are inadequate to validate an OOD generalization algorithm. Therefore, it is critical for the community to develop more reasonable evaluation criterions for OOD generalization.

## (D) INCORPORATE SELF-LEARNING

Recently, there is a rapid development of large-scale self-learning (or pre-trained models), such as BERT [238], GPT-3 [239], SimCLR [240], which propose to firstly pre-train models on large-scale datasets and then finetune them on downstream tasks. Since the distributional shifts between downstream tasks and pre-training datasets are inevitable, how to design efficient self-learning methods with good OOD generalization ability or incorporate self-learning methods to achieve better OOD generalization performance remains a promising direction for future.

## ACKNOWLEDGMENTS

## REFERENCES

[1] J. Duchi and H. Namkoong, "Learning models with uniform performance via distributionally robust optimization," *arXiv preprint arXiv:1810.08750*, 2018.

[2] M. Arjovsky, L. Bottou, I. Gulrajani, and D. Lopez-Paz, "Invariant risk minimization," *arXiv preprint arXiv:1907.02893*, 2019.

[3] E. Creager, J.-H. Jacobsen, and R. Zemel, "Environment inference for invariant learning," in *International Conference on Machine Learning*. PMLR, 2021, pp. 2189–2200.

[4] Z. Shen, P. Cui, T. Zhang, and K. Kunag, "Stable learning via sample reweighting," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 4, 2020, pp. 5692–5699.

[5] J. Liu, Z. Hu, P. Cui, B. Li, and Z. Shen, "Heterogeneous risk minimization," in *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, Virtual Event*, ser. Proceedings of Machine Learning Research. PMLR, 2021.

[6] K. Muandet, D. Balduzzi, and B. Schölkopf, "Domain generalization via invariant feature representation," in *International Conference on Machine Learning*. PMLR, 2013, pp. 10–18.

[7] Y. Ganin and V. Lempitsky, "Unsupervised domain adaptation by backpropagation," in *International conference on machine learning*. PMLR, 2015, pp. 1180–1189.

[8] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky, "Domain-adversarial training of neural networks," *The journal of machine learning research*, vol. 17, no. 1, pp. 2096–2030, 2016.

[9] H. Li, S. J. 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*, 2018, pp. 5400–5409.

[10] J. Peters, P. Bühlmann, and N. Meinshausen, "Causal inference by using invariant prediction: identification and confidence intervals," *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, vol. 78, no. 5, pp. 947–1012, 2016.

[11] P. Bühlmann, "Invariance, causality and robustness," *arXiv preprint arXiv:1812.08233*, 2018.

[12] K. Kuang, P. Cui, S. Athey, R. Xiong, and B. Li, "Stable prediction across unknown environments," pp. 1617–1626, 2018.

[13] K. Kuang, R. Xiong, P. Cui, S. Athey, and B. Li, "Stable prediction with model misspecification and agnostic distribution shift," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 4, 2020, pp. 4485–4492.

[14] J. Wang, C. Lan, C. Liu, Y. Ouyang, W. Zeng, and T. Qin, "Generalizing to unseen domains: A survey on domain generalization," *arXiv preprint arXiv:2103.03097*, 2021.

[15] K. Zhou, Z. Liu, Y. Qiao, T. Xiang, and C. C. Loy, "Domain generalization: A survey," *arXiv preprint arXiv:2103.02503*, 2021.

[16] N. Ye, K. Li, L. Hong, H. Bai, Y. Chen, F. Zhou, and Z. Li, "OoD-Bench: Benchmarking and Understanding Out-of-Distribution Generalization Datasets and Algorithms," *CoRR*, 2021.

[17] V. Vapnik, "Principles of risk minimization for learning theory," in *Advances in Neural Information Processing Systems 4*, vol. 4, 1991, pp. 831–838.

[18] ——, "An overview of statistical learning theory," *IEEE Transactions on Neural Networks*, vol. 10, no. 5, pp. 988–999, 1999.

[19] J. Quiñonero-Candela, M. Sugiyama, N. D. Lawrence, and A. Schwaighofer, *Dataset shift in machine learning*. Mit Press, 2009.

[20] S. Garg, Y. Wu, S. Balakrishnan, and Z. Lipton, "A unified view of label shift estimation," in *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, Eds., vol. 33. Curran Associates, Inc., 2020, pp. 3290–3300. [Online]. Available: https://proceedings.neurips.cc/paper/2020/file/219e052492f4008818b8adb6366c7ed6-Paper.pdf

[21] J. Gama, I. Žliobaitė, A. Bifet, M. Pechenizkiy, and A. Bouchachia, "A survey on concept drift adaptation," *ACM computing surveys (CSUR)*, vol. 46, no. 4, pp. 1–37, 2014.

[22] S. Ben-David, J. Blitzer, K. Crammer, A. Kulesza, F. Pereira, and J. W. Vaughan, "A theory of learning from different domains," *Machine learning*, vol. 79, no. 1-2, pp. 151–175, 2010.

[23] S. J. Pan and Q. Yang, "A survey on transfer learning," *IEEE Transactions on knowledge and data engineering*, vol. 22, no. 10, pp. 1345–1359, 2009.

[24] S. Sun, H. Shi, and Y. Wu, "A Survey of Multi-source Domain Adaptation," *Information Fusion*, 2015.

[25] F. Zhuang, Z. Qi, K. Duan, D. Xi, Y. Zhu, H. Zhu, H. Xiong, and Q. He, "A comprehensive survey on transfer learning," *Proceedings of the IEEE*, vol. 109, no. 1, pp. 43–76, 2021.

[26] G. Csurka, "Domain adaptation for visual applications: A comprehensive survey," *arXiv preprint arXiv:1702.05374*, 2017.

[27] Y. Bengio, A. Courville, and P. Vincent, "Representation learning: A review and new perspectives," *IEEE transactions on pattern analysis and machine intelligence*, vol. 35, no. 8, pp. 1798–1828, 2013.

[28] I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, and A. Lerchner, "beta-vae: Learning basic visual concepts with a constrained variational framework," 2016.

[29] H. Kim and A. Mnih, "Disentangling by factorising," in *International Conference on Machine Learning*. PMLR, 2018, pp. 2649–2658.

[30] F. Locatello, S. Bauer, M. Lucic, G. Raetsch, S. Gelly, B. Schölkopf, and O. Bachem, "Challenging common assumptions in the unsupervised learning of disentangled representations," in *international conference on machine learning*. PMLR, 2019, pp. 4114–4124.

[31] F. Leeb, G. Lanzillotta, Y. Annadani, M. Besserve, S. Bauer, and B. Schölkopf, "Structure by Architecture: Disentangled Representations without Regularization," *arXiv e-prints*, p. arXiv:2006.07796, Jun. 2020.

[32] A. Dittadi, F. Träuble, F. Locatello, M. Wüthrich, V. Agrawal, O. Winther, S. Bauer, and B. Schölkopf, "On the transfer of disentangled representations in realistic settings," *arXiv preprint arXiv:2010.14407*, 2020.

[33] F. Träuble, E. Creager, N. Kilbertus, F. Locatello, A. Dittadi, A. Goyal, B. Schölkopf, and S. Bauer, "On disentangled representations learned from correlated data," in *International Conference on Machine Learning*. PMLR, 2021, pp. 10401–10412.

[34] M. Yang, F. Liu, Z. Chen, X. Shen, J. Hao, and J. Wang, "Causal-vae: disentangled representation learning via neural structural causal models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 9593–9602.

[35] X. Shen, F. Liu, H. Dong, Q. Lian, Z. Chen, and T. Zhang, "Disentangled generative causal representation learning," *arXiv preprint arXiv:2010.02637*, 2020.

[36] B. Schölkopf, F. Locatello, S. Bauer, N. R. Ke, N. Kalchbrenner, A. Goyal, and Y. Bengio, "Toward causal representation learning," *Proceedings of the IEEE*, vol. 109, no. 5, pp. 612–634, 2021.

[37] E. Tzeng, J. Hoffman, N. Zhang, K. Saenko, and T. Darrell, "Deep domain confusion: Maximizing for domain invariance," *arXiv preprint arXiv:1412.3474*, 2014.

[38] S. Motiian, M. Piccirilli, D. A. Adjeroh, and G. Doretto, "Unified deep supervised domain adaptation and generalization," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 5715–5725.

[39] J. Wang, W. Feng, Y. Chen, H. Yu, M. Huang, and P. S. Yu, "Visual domain adaptation with manifold embedded distribution alignment," in *Proceedings of the 26th ACM international conference on Multimedia*, 2018, pp. 402–410.

[40] 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)*, 2018, pp. 624–639.

[41] R. Shao, X. Lan, J. Li, and P. C. Yuen, "Multi-adversarial discriminative deep domain generalization for face presentation attack detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 10 023–10 031.

[42] R. Gong, W. Li, Y. Chen, and L. V. Gool, "Dlow: Domain flow for adaptation and generalization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2477–2486.

[43] Y. Jia, J. Zhang, S. Shan, and X. Chen, "Single-side domain generalization for face anti-spoofing," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 8484–8493.

[44] I. Albuquerque, J. Monteiro, M. Darvishi, T. H. Falk, and I. Mitliagkas, "Adversarial target-invariant representation learning for domain generalization," 2020.

[45] Z. Wang, Q. Wang, C. Lv, X. Cao, and G. Fu, "Unseen target stance detection with adversarial domain generalization," in *2020 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2020, pp. 1–8.

[46] S. Zhao, M. Gong, T. Liu, H. Fu, and D. Tao, "Domain generalization via entropy regularization," *Advances in Neural Information Processing Systems*, vol. 33, 2020.

[47] M. M. Rahman, C. Fookes, M. Baktashmotlagh, and S. Sridharan, "Correlation-aware adversarial domain adaptation and generalization," *Pattern Recognition*, vol. 100, p. 107124, 2020.

[48] V. Garg, A. T. Kalai, K. Ligett, and S. Wu, "Learn to expect the unexpected: Probably approximately correct domain generalization," in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2021, pp. 3574–3582.

[49] A. Sicilia, X. Zhao, and S. J. Hwang, "Domain adversarial neural networks for domain generalization: When it works and how to improve," *arXiv preprint arXiv:2102.03924*, 2021.

[50] S. J. Pan, I. W. Tsang, J. T. Kwok, and Q. Yang, "Domain adaptation via transfer component analysis," *IEEE transactions on neural networks*, vol. 22, no. 2, pp. 199–210, 2010.

[51] J. Wang, Y. Chen, W. Feng, H. Yu, M. Huang, and Q. Yang, "Transfer learning with dynamic distribution adaptation," *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 11, no. 1, pp. 1–25, 2020.

[52] F. Zhou, Z. Jiang, C. Shui, B. Wang, and B. Chaib-draa, "Domain generalization with optimal transport and metric learning," *arXiv preprint arXiv:2007.10573*, 2020.

[53] B. Sun and K. Saenko, "Deep coral: Correlation alignment for deep domain adaptation," in *European conference on computer vision*. Springer, 2016, pp. 443–450.

[54] B. Sun, J. Feng, and K. Saenko, "Return of frustratingly easy domain adaptation," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 30, no. 1, 2016.

[55] X. Peng and K. Saenko, "Synthetic to real adaptation with generative correlation alignment networks," in *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 2018, pp. 1982–1991.

[56] W. Jin, R. Barzilay, and T. Jaakkola, "Domain extrapolation via regret minimization," *arXiv preprint arXiv:2006.03908*, 2020.

[57] X. Pan, P. Luo, J. Shi, and X. Tang, "Two at once: Enhancing learning and generalization capacities via ibn-net," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 464–479.

[58] D. Ulyanov, A. Vedaldi, and V. Lempitsky, "Improved texture networks: Maximizing quality and diversity in feed-forward stylization and texture synthesis," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 6924–6932.

[59] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *International conference on machine learning*. PMLR, 2015, pp. 448–456.

[60] X. Huang and S. Belongie, "Arbitrary style transfer in real-time with adaptive instance normalization," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 1501–1510.

[61] H. Nam and H.-E. Kim, "Batch-instance normalization for adaptively style-invariant neural networks," *arXiv preprint arXiv:1805.07925*, 2018.

[62] X. Jin, C. Lan, W. Zeng, and Z. Chen, "Style normalization and restitution for domaingeneralization and adaptation," *arXiv preprint arXiv:2101.00588*, 2021.

[63] G. Blanchard, G. Lee, and C. Scott, "Generalizing from several related classification tasks to a new unlabeled sample," *Advances in neural information processing systems*, vol. 24, pp. 2178–2186, 2011.

[64] G. Blanchard, A. A. Deshmukh, U. Dogan, G. Lee, and C. Scott, "Domain generalization by marginal transfer learning," *arXiv preprint arXiv:1711.07910*, 2017.

[65] T. Grubinger, A. Birlutiu, H. Schöner, T. Natschläger, and T. Heskes, "Domain generalization based on transfer component analysis," in *International Work-Conference on Artificial Neural Networks*. Springer, 2015, pp. 325–334.

[66] C. Gan, T. Yang, and B. Gong, "Learning attributes equals multi-source domain generalization," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 87–97.

[67] S. Erfani, M. Baktashmotlagh, M. Moshtaghi, X. Nguyen, C. Leckie, J. Bailey, and R. Kotagiri, "Robust domain generalisation by enforcing distribution invariance," in *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence (IJCAI-16)*. AAAI Press, 2016, pp. 1455–1461.

[68] M. Ghifary, D. Balduzzi, W. B. Kleijn, and M. Zhang, "Scatter component analysis: A unified framework for domain adaptation and domain generalization," *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 7, pp. 1414–1430, 2016.

[69] S. Hu, K. Zhang, Z. Chen, and L. Chan, "Domain generalization via multidomain discriminant analysis," in *Uncertainty in Artificial Intelligence*. PMLR, 2020, pp. 292–302.

[70] A. A. Deshmukh, Y. Lei, S. Sharma, U. Dogan, J. W. Cutler, and C. Scott, "A generalization error bound for multi-class domain generalization," *arXiv preprint arXiv:1905.10392*, 2019.

[71] C. Finn, P. Abbeel, and S. Levine, "Model-agnostic meta-learning for fast adaptation of deep networks," in *International Conference on Machine Learning*. PMLR, 2017, pp. 1126–1135.

[72] D. Li, Y. Yang, Y.-Z. Song, and T. M. Hospedales, "Learning to generalize: Meta-learning for domain generalization," in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.

[73] Y. Balaji, S. Sankaranarayanan, and R. Chellappa, "Metareg: Towards domain generalization using meta-regularization," *Advances in Neural Information Processing Systems*, vol. 31, pp. 998–1008, 2018.

[74] S. Otálora, M. Atzori, V. Andrearczyk, A. Khan, and H. Müller, "Staining invariant features for improving generalization of deep convolutional neural networks in computational pathology," *Frontiers in bioengineering and biotechnology*, vol. 7, p. 198, 2019.

[75] Y. Li, Y. Yang, W. Zhou, and T. Hospedales, "Feature-critic networks for heterogeneous domain generalization," in *International Conference on Machine Learning*. PMLR, 2019, pp. 3915–3924.

[76] Q. Dou, D. Coelho de Castro, K. Kamnitsas, and B. Glocker, "Domain generalization via model-agnostic learning of semantic features," *Advances in Neural Information Processing Systems*, vol. 32, pp. 6450–6461, 2019.

[77] Q. Liu, Q. Dou, and P.-A. Heng, "Shape-aware meta-learning for generalizing prostate mri segmentation to unseen domains," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2020, pp. 475–485.

[78] Y. Du, J. Xu, H. Xiong, Q. Qiu, X. Zhen, C. G. Snoek, and L. Shao, "Learning to learn with variational information bottleneck for domain generalization," in *European Conference on Computer Vision*. Springer, 2020, pp. 200–216.

[79] Y. Du, X. Zhen, L. Shao, and C. G. Snoek, "Metanorm: Learning to normalize few-shot batches across domains," in *International Conference on Learning Representations*, 2020.

[80] Y. Zhao, Z. Zhong, F. Yang, Z. Luo, Y. Lin, S. Li, and N. Sebe, "Learning to generalize unseen domains via memory-based multi-source meta-learning for person re-identification," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 6277–6286.

[81] S. Choi, T. Kim, M. Jeong, H. Park, and C. Kim, "Meta batch-instance normalization for generalizable person re-identification," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 3425–3435.

[82] A. D'Innocente and B. Caputo, "Domain generalization with domain-specific aggregation modules," in *German Conference on Pattern Recognition*. Springer, 2018, pp. 187–198.

[83] Z. Ding and Y. Fu, "Deep domain generalization with structured low-rank constraint," *IEEE Transactions on Image Processing*, vol. 27, no. 1, pp. 304–313, 2017.

[84] S. Wang, L. Yu, K. Li, X. Yang, C.-W. Fu, and P.-A. Heng, "Dofe: Domain-oriented feature embedding for generalizable fundus image segmentation on unseen datasets," *IEEE Transactions on Medical Imaging*, vol. 39, no. 12, pp. 4237–4248, 2020.

[85] M. Segu, A. Tonioni, and F. Tombari, "Batch normalization embeddings for deep domain generalization," *arXiv preprint arXiv:2011.12672*, 2020.

[86] M. Mancini, S. R. Bulo, B. Caputo, and E. Ricci, "Best sources forward: domain generalization through source-specific nets," in *2018 25th IEEE international conference on image processing (ICIP)*. IEEE, 2018, pp. 1353–1357.

[87] X. Zhang, L. Zhou, R. Xu, P. Cui, Z. Shen, and H. Liu, "Domain-irrelevant representation learning for unsupervised domain generalization," *arXiv preprint arXiv:2107.06219*, 2021.

[88] Y. Liao, R. Huang, J. Li, Z. Chen, and W. Li, "Deep semisupervised domain generalization network for rotary machinery fault diagnosis under variable speed," *IEEE Transactions on Instrumentation and Measurement*, vol. 69, no. 10, pp. 8064–8075, 2020.

[89] F. M. Carlucci, A. D'Innocente, S. Bucci, B. Caputo, and T. Tommasi, "Domain generalization by solving jigsaw puzzles," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2229–2238.

[90] J. Ryu, G. Kwon, M.-H. Yang, and J. Lim, "Generalized convolutional forest networks for domain generalization and visual recognition," in *International Conference on Learning Representations*, 2019.

[91] D. Li, J. Zhang, Y. Yang, C. Liu, Y.-Z. Song, and T. M. Hospedales, "Episodic training for domain generalization," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 1446–1455.

[92] Z. Huang, H. Wang, E. P. Xing, and D. Huang, "Self-challenging improves cross-domain generalization," in *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*. Springer, 2020, pp. 124–140.

[93] X. Yue, Y. Zhang, S. Zhao, A. Sangiovanni-Vincentelli, K. Keutzer, and B. Gong, "Domain randomization and pyramid consistency: Simulation-to-real generalization without accessing target domain data," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 2100–2110.

[94] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel, "Domain randomization for transferring deep neural networks from simulation to the real world," in *2017 IEEE/RSJ international conference on intelligent robots and systems (IROS)*. IEEE, 2017, pp. 23–30.

[95] X. B. Peng, M. Andrychowicz, W. Zaremba, and P. Abbeel, "Sim-to-real transfer of robotic control with dynamics randomization," in *2018 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2018, pp. 3803–3810.

[96] R. Khirodkar, D. Yoo, and K. Kitani, "Domain randomization for scene-specific car detection and pose estimation," in *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 2019, pp. 1932–1940.

[97] J. Tremblay, A. Prakash, D. Acuna, M. Brophy, V. Jampani, C. Anil, T. To, E. Cameracci, S. Boochoon, and S. Birchfield, "Training deep networks with synthetic data: Bridging the reality gap by domain randomization," in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2018, pp. 969–977.

[98] A. Prakash, S. Boochoon, M. Brophy, D. Acuna, E. Cameracci, G. State, O. Shapira, and S. Birchfield, "Structured domain randomization: Bridging the reality gap by context-aware synthetic data," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 7249–7255.

[99] S. Shankar, V. Piratla, S. Chakrabarti, S. Chaudhuri, P. Jyothi, and S. Sarawagi, "Generalizing across domains via cross-gradient training," *arXiv preprint arXiv:1804.10745*, 2018.

[100] R. Volpi, H. Namkoong, O. Sener, J. Duchi, V. Murino, and S. Savarese, "Generalizing to unseen domains via adversarial data augmentation," *arXiv preprint arXiv:1805.12018*, 2018.

[101] 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*, vol. 34, no. 07, 2020, pp. 13 025–13 032.

[102] F. Qiao, L. Zhao, and X. Peng, "Learning to learn single domain generalization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 12 556–12 565.

[103] M. M. Rahman, C. Fookes, M. Baktashmotlagh, and S. Sridharan, "Multi-component image translation for deep domain generalization," in *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 2019, pp. 579–588.

[104] K. Zhou, Y. Yang, T. Hospedales, and T. Xiang, "Learning to generate novel domains for domain generalization," in *European Conference on Computer Vision*. Springer, 2020, pp. 561–578.

[105] N. Somavarapu, C.-Y. Ma, and Z. Kira, "Frustratingly simple domain generalization via image stylization," *arXiv preprint arXiv:2006.11207*, 2020.

[106] N. Pfister, P. Bühlmann, and J. Peters, "Invariant Causal Prediction for Sequential Data," *Journal of the American Statistical Association*, 2018.

[107] D. Rothenhäusler, N. Meinshausen, P. Bühlmann, and J. Peters, "Anchor regression: heterogeneous data meets causality," *arXiv preprint arXiv:1801.06229*, 2018.

[108] C. Heinze-Deml, J. Peters, and N. Meinshausen, "Invariant Causal Prediction for Nonlinear Models," *Journal of Causal Inference*, 2018.

[109] J. L. Gamella and C. Heinze-Deml, "Active Invariant Causal Prediction: Experiment Selection through Stability," in *NIPS*, 2020.

[110] M. Oberst, N. Thams, J. Peters, and D. Sontag, "Regularizing towards Causal Invariance: Linear Models with Proxies," in *ICML*, 2021.

[111] K. Ahuja, K. Shanmugam, K. Varshney, and A. Dhurandhar, "Invariant risk minimization games," in *International Conference on Machine Learning*. PMLR, 2020, pp. 145–155.

[112] D. Krueger, E. Caballero, J.-H. Jacobsen, A. Zhang, J. Binas, D. Zhang, R. Le Priol, and A. Courville, "Out-of-distribution generalization via risk extrapolation (rex)," in *International Conference on Machine Learning*. PMLR, 2021, pp. 5815–5826.

[113] C. Xie, F. Chen, Y. Liu, and Z. Li, "Risk variance penalization: From distributional robustness to causality," *arXiv e-prints*, pp. arXiv–2006, 2020.

[114] D. Mahajan, S. Tople, and A. Sharma, "Domain generalization using causal matching," in *Proceedings of the 38th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, M. Meila and T. Zhang, Eds., vol. 139. PMLR, 18–24 Jul 2021, pp. 7313–7324. [Online]. Available: http://proceedings.mlr.press/v139/mahajan21b.html

[115] E. Rosenfeld, P. Ravikumar, and A. Risteski, "The risks of invariant risk minimization," *arXiv preprint arXiv:2010.05761*, 2020.

[116] P. Kamath, A. Tangella, D. Sutherland, and N. Srebro, "Does invariant risk minimization capture invariance?" in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2021, pp. 4069–4077.

[117] K. Ahuja, J. Wang, A. Dhurandhar, K. Shanmugam, and K. R. Varshney, "Empirical or invariant risk minimization? a sample complexity perspective," *arXiv preprint arXiv:2010.16412*, 2020.

[118] Z. Shen, P. Cui, K. Kuang, B. Li, and P. Chen, "Causally regularized learning with agnostic data selection bias." in *ACM Multimedia*, 2018, pp. 411–419.

[119] Z. Shen, P. Cui, J. Liu, T. Zhang, B. Li, and Z. Chen, "Stable learning via differentiated variable decorrelation," in *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery &amp; Data Mining*, 2020, pp. 2185–2193.

[120] X. Zhang, P. Cui, R. Xu, L. Zhou, Y. He, and Z. Shen, "Deep stable learning for out-of-distribution generalization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 5372–5382.

[121] X. Wang, S. Fan, K. Kuang, C. Shi, J. Liu, and B. Wang, "Decorrelated clustering with data selection bias," in *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI 2020*, C. Bessiere, Ed. ijcai.org, 2020, pp. 2177–2183. [Online]. Available: https://doi.org/10.24963/ijcai.2020/301

[122] S. Zhang, T. Jiang, T. Wang, K. Kuang, Z. Zhao, J. Zhu, J. Yu, H. Yang, and F. Wu, "Devlbert: Learning deconfounded visio-linguistic representations," in *MM '20: The 28th ACM International Conference on Multimedia, Virtual Event / Seattle, WA, USA, October 12-16, 2020*. ACM, 2020, pp. 4373–4382. [Online]. Available: https://doi.org/10.1145/3394171.3413518

[123] Z. Yuan, X. Peng, X. Wu, B.-k. Bao, and C. Xu, "Meta-learning causal feature selection for stable prediction," in *2021 IEEE International Conference on Multimedia and Expo (ICME)*, 2021, pp. 1–6.

[124] K. Kuang, H. Zhang, F. Wu, Y. Zhuang, and A. Zhang, "Balance-subsampled stable prediction," *CoRR*, vol. abs/2006.04381, 2020. [Online]. Available: https://arxiv.org/abs/2006.04381

[125] S. Shafieezadeh-Abadeh, P. M. Esfahani, and D. Kuhn, "Distributionally robust logistic regression," 2015.

[126] A. Sinha, H. Namkoong, and J. Duchi, "Certifying some distributional robustness with principled adversarial training," *International Conference on Learning Representations*, 2018.

[127] L. El Ghaoui and H. Lebret, "Robust solutions to least-squares problems with uncertain data," *SIAM Journal on Matrix Analysis and Applications*, vol. 18, no. 4, pp. 1035–1064, 1997. [Online]. Available: https://doi.org/10.1137/S0895479896298130

[128] H. Xu, C. Caramanis, and S. Mannor, "Robust regression and lasso," *CoRR*, vol. abs/0811.1790, 2008. [Online]. Available: http://arxiv.org/abs/0811.1790

[129] W. Yang and H. Xu, "A unified robust regression model for lasso-like algorithms," in *Proceedings of the 30th International Conference on Machine Learning, ICML 2013, Atlanta, GA, USA, 16-21 June 2013*, ser. JMLR Workshop and Conference Proceedings, vol. 28. JMLR.org, 2013, pp. 585–593. [Online]. Available: http://proceedings.mlr.press/v28/yang13e.html

[130] D. Bertsimas and M. S. Copenhaver, "Characterization of the equivalence of robustification and regularization in linear and matrix regression," 2017.

[131] S. Shafieezadeh-Abadeh, D. Kuhn, and P. M. Esfahani, "Regularization via mass transportation," 2019.

[132] R. Gao, X. Chen, and A. J. Kleywegt, "Wasserstein distributional robustness and regularization in statistical learning," *arXiv e-prints*, pp. arXiv–1712, 2017.

[133] J. C. Duchi, P. W. Glynn, and H. Namkoong, "Statistics of robust optimization: A generalized empirical likelihood approach," *Mathematics of Operations Research*, 2021.

[134] R. Chen and I. C. Paschalidis, "A robust learning approach for regression models based on distributionally robust optimization," *The Journal of Machine Learning Research*, vol. 19, no. 1, pp. 517–564, 2018.

[135] J. Liu, Z. Shen, P. Cui, L. Zhou, K. Kuang, B. Li, and Y. Lin, "Stable adversarial learning under distributional shifts," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 10, 2021, pp. 8662–8670.

[136] W. Hu, G. Niu, I. Sato, and M. Sugiyama, "Does distributionally robust supervised learning give robust classifiers?" *Proceedings of the 35th International Conference on Machine Learning, PMLR 80:2029-2037*, 2016.

[137] C. Frogner, S. Claici, E. Chien, and J. Solomon, "Incorporating unlabeled data into distributionally robust learning," *arXiv preprint arXiv:1912.07729*, 2019.

[138] J. C. Duchi, T. Hashimoto, and H. Namkoong, "Distributionally robust losses against mixture covariate shifts," *Under review*, 2019.

[139] S. Chang, Y. Zhang, M. Yu, and T. S. Jaakkola, "Invariant rationalization," in *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, ser. Proceedings of Machine Learning Research, vol. 119. PMLR, 2020, pp. 1448–1458. [Online]. Available: http://proceedings.mlr.press/v119/chang20c.html

[140] M. Koyama and S. Yamaguchi, "Out-of-distribution generalization with maximal invariant predictor," *CoRR*,

vol. abs/2008.01883, 2020. [Online]. Available: https://arxiv.org/abs/2008.01883

[141] T. Hospedales, A. Antoniou, P. Micaelli, and A. Storkey, "Meta-learning in neural networks: A survey," *arXiv preprint arXiv:2004.05439*, 2020.

[142] Y. Sun, E. Tzeng, T. Darrell, and A. A. Efros, "Unsupervised domain adaptation through self-supervision," *arXiv preprint arXiv:1909.11825*, 2019.

[143] H. Daumé III, A. Kumar, and A. Saha, "Frustratingly easy semi-supervised domain adaptation," in *Proceedings of the 2010 Workshop on Domain Adaptation for Natural Language Processing*, 2010, pp. 53–59.

[144] M. Xiao and Y. Guo, "Feature space independent semi-supervised domain adaptation via kernel matching," *IEEE transactions on pattern analysis and machine intelligence*, vol. 37, no. 1, pp. 54–66, 2014.

[145] M. Noroozi and P. Favaro, "Unsupervised learning of visual representations by solving jigsaw puzzles," in *European conference on computer vision*. Springer, 2016, pp. 69–84.

[146] G. C. Chow, "Tests of equality between sets of coefficients in two linear regressions," *Econometrica*, vol. 28, no. 3, pp. 591–605, 1960.

[147] S. Athey, G. W. Imbens, and S. Wager, "Approximate residual balancing: De-biased inference of average treatment effects in high dimensions," *arXiv preprint arXiv:1604.07125*, 2016.

[148] J. R. Zubizarreta, "Stable weights that balance covariates for estimation with incomplete outcome data," *Journal of the American Statistical Association*, vol. 110, no. 511, pp. 910–922, 2015.

[149] J. Hainmueller, "Entropy balancing for causal effects: A multivariate reweighting method to produce balanced samples in observational studies," *Political Analysis*, p. mpr025, 2011.

[150] B. Schölkopf, J. Platt, and T. Hofmann, "Greedy layer-wise training of deep networks," *Advances in Neural Information Processing Systems*, vol. 19, pp. 153–160, 2007.

[151] A. Rahimi, B. Recht *et al.*, "Random features for large-scale kernel machines." in *NIPS*, vol. 3, no. 4. Citeseer, 2007, p. 5.

[152] P. M. Esfahani and D. Kuhn, "Data-driven distributionally robust optimization using the wasserstein metric: Performance guarantees and tractable reformulations," *Mathematical Programming*, vol. 171, no. 1-2, pp. 115–166, 2018.

[153] S. Sagawa, P. W. Koh, T. B. Hashimoto, and P. Liang, "Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization," *arXiv preprint arXiv:1911.08731*, 2019.

[154] E. Delage and Y. Ye, "Distributionally robust optimization under moment uncertainty with application to data-driven problems," *Oper. Res.*, vol. 58, no. 3, p. 595–612, May 2010.

[155] D. Bertsimas, V. Gupta, and N. Kallus, "Data-driven robust optimization," *Mathematical Programming*, vol. 167, no. 2, pp. 235–292, 2018.

[156] H. Namkoong and J. C. Duchi, "Stochastic gradient methods for distributionally robust optimization with f-divergences," *Neural Information Processing Systems (NIPS)*, pp. 2208–2216, 2016.

[157] H. Rahimian and S. Mehrotra, "Distributionally robust optimization: A review," *arXiv preprint arXiv:1908.05659*, 2019.

[158] M. Rojas-Carulla, B. Schölkopf, R. Turner, and J. Peters, "Invariant models for causal transfer learning," *The Journal of Machine Learning Research*, vol. 19, no. 1, pp. 1309–1342, 2018.

[159] N. Meinshausen, "Causality from a distributional robustness point of view," in *2018 IEEE Data Science Workshop (DSW)*, 2018, pp. 6–10.

[160] J. Pearl, *Probabilistic reasoning in intelligent systems: networks of plausible inference*. Elsevier, 2014.

[161] B. Aubin, A. Słowik, M. Arjovsky, L. Bottou, and D. Lopez-Paz, "Linear unit-tests for invariance discovery," *arXiv preprint arXiv:2102.10867*, 2021.

[162] J. Pearl, *Causality*. Cambridge university press, 2009.

[163] A. Subbaswamy and S. Saria, "Counterfactual normalization: Proactively addressing dataset shift and improving reliability using causal mechanisms," *arXiv preprint arXiv:1808.03253*, 2018.

[164] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.

[165] D. Hendrycks, K. Zhao, S. Basart, J. Steinhardt, and D. Song, "Natural adversarial examples," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 15 262–15 271.

[166] D. Hendrycks and T. Dietterich, "Benchmarking neural network robustness to common corruptions and perturbations," *arXiv preprint arXiv:1903.12261*, 2019.

[167] D. Hendrycks, S. Basart, N. Mu, S. Kadavath, F. Wang, E. Dorundo, R. Desai, T. Zhu, S. Parajuli, M. Guo *et al.*, "The many faces of robustness: A critical analysis of out-of-distribution generalization," *arXiv preprint arXiv:2006.16241*, 2020.

[168] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.

[169] M. Ghifary, W. B. 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*, 2015, pp. 2551–2559.

[170] 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*, 2017, pp. 5542–5550.

[171] 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*, 2017, pp. 5018–5027.

[172] C. Fang, Y. Xu, and D. N. Rockmore, "Unbiased metric learning: On the utilization of multiple datasets and web images for softening bias," in *Proceedings of the IEEE International Conference on Computer Vision*, 2013, pp. 1657–1664.

[173] P. Bandi, O. Geessink, Q. Manson, M. Van Dijk, M. Balkenhol, M. Hermsen, B. E. Bejnordi, B. Lee, K. Paeng, A. Zhong *et al.*, "From detection of individual metastases to classification of lymph node status at the patient level: the camelyon17 challenge," *IEEE transactions on medical imaging*, vol. 38, no. 2, pp. 550–560, 2018.

[174] S. Zhao, B. Li, X. Yue, Y. Gu, P. Xu, R. Hu, H. Chai, and K. Keutzer, "Multi-source domain adaptation for semantic segmentation," *arXiv preprint arXiv:1910.12181*, 2019.

[175] S. Beery, A. Agarwal, E. Cole, and V. Birodkar, "The iwildcam 2021 competition dataset," *arXiv preprint arXiv:2105.03494*, 2021.

[176] P. W. Koh, S. Sagawa, S. M. Xie, M. Zhang, A. Balsubramani, W. Hu, M. Yasunaga, R. L. Phillips, I. Gao, T. Lee *et al.*, "Wilds: A benchmark of in-the-wild distribution shifts," in *International Conference on Machine Learning*. PMLR, 2021, pp. 5637–5664.

[177] Y. He, Z. Shen, and P. Cui, "Towards non-iid image classification: A dataset and baselines," *Pattern Recognition*, vol. 110, p. 107383, 2021.

[178] G. Christie, N. Fendley, J. Wilson, and R. Mukherjee, "Functional map of the world," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 6172–6180.

[179] C. Yeh, A. Perez, A. Driscoll, G. Azzari, Z. Tang, D. Lobell, S. Ermon, and M. Burke, "Using publicly available satellite imagery and deep learning to understand economic well-being in africa," *Nature communications*, vol. 11, no. 1, pp. 1–11, 2020.

[180] M. Z. Y. D. H. Ren, B. L. M. C. J. Leskovec, W. Hu, and M. Fey, "Open graph benchmark: Datasets for machine learning on graphs," *arXiv preprint arXiv:2005.00687*, 2020.

[181] D. Borkan, L. Dixon, J. Sorensen, N. Thain, and L. Vasserman, "Nuanced metrics for measuring unintended bias with real data for text classification," in *Companion proceedings of the 2019 world wide web conference*, 2019, pp. 491–500.

[182] J. Ni, J. Li, and J. McAuley, "Justifying recommendations using distantly-labeled reviews and fine-grained aspects," in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 2019, pp. 188–197.

[183] V. Raychev, P. Bielik, and M. Vechev, "Probabilistic model for code with decision trees," *ACM SIGPLAN Notices*, vol. 51, no. 10, pp. 731–747, 2016.

[184] H. Ye, C. Xie, Y. Liu, and Z. Li, "Out-of-distribution generalization analysis via influence function," *arXiv preprint arXiv:2101.08521*, 2021.

[185] A. Sinha, H. Namkoong, R. Volpi, and J. Duchi, "Certifying some distributional robustness with principled adversarial training," *arXiv preprint arXiv:1710.10571*, 2017.

[186] A. Mukerjee, R. Biswas, K. Deb, and A. P. Mathur, "Multi–objective evolutionary algorithms for the risk–return trade–off in bank loan management," *International Transactions in operational research*, vol. 9, no. 5, pp. 583–597, 2002.

[187] L. A. Rivera, "Hiring as cultural matching: The case of elite professional service firms," *American sociological review*, vol. 77, no. 6, pp. 999–1022, 2012.

[188] J. Larson, S. Mattu, L. Kirchner, and J. Angwin, "How we analyzed the compas recidivism algorithm," *ProPublica (5 2016)*, vol. 9, 2016.

[189] C. Dwork, M. Hardt, T. Pitassi, O. Reingold, and R. Zemel, "Fairness through awareness," in *Proceedings of the 3rd innovations in theoretical computer science conference*, 2012, pp. 214–226.

[190] M. Yurochkin and Y. Sun, "Sensei: Sensitive set invariance for enforcing individual fairness," in *International Conference on Learning Representations*, 2021.

[191] M. Hardt, E. Price, and N. Srebro, "Equality of opportunity in supervised learning," *Advances in neural information processing systems*, vol. 29, pp. 3315–3323, 2016.

[192] M. Kearns, S. Neel, A. Roth, and Z. S. Wu, "Preventing fairness gerrymandering: Auditing and learning for subgroup fairness," in *International Conference on Machine Learning*. PMLR, 2018, pp. 2564–2572.

[193] R. Xu, P. Cui, K. Kuang, B. Li, L. Zhou, Z. Shen, and W. Cui, "Algorithmic decision making with conditional fairness," in *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2020, pp. 2125–2135.

[194] N. Kilbertus, M. Rojas-Carulla, G. Parascandolo, M. Hardt, D. Janzing, and B. Schölkopf, "Avoiding discrimination through causal reasoning," in *Advances in Neural Information Processing Systems 30*, 2017.

[195] M. J. Kusner, J. Loftus, C. Russell, and R. Silva, "Counterfactual fairness," in *Advances in Neural Information Processing Systems 30*, 2017.

[196] S. Chiappa, "Path-specific counterfactual fairness," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 7801–7808.

[197] J. Kleinberg, S. Mullainathan, and M. Raghavan, "Inherent trade-offs in the fair determination of risk scores," *arXiv preprint arXiv:1609.05807*, 2016.

[198] H. Wang, B. Ustun, and F. Calmon, "Repairing without retraining: Avoiding disparate impact with counterfactual distributions," in *International Conference on Machine Learning*, 2019, pp. 6618–6627.

[199] M. Feldman, S. A. Friedler, J. Moeller, C. Scheidegger, and S. Venkatasubramanian, "Certifying and removing disparate impact," in *proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*, 2015, pp. 259–268.

[200] F. Kamiran and T. Calders, "Data preprocessing techniques for classification without discrimination," *Knowledge and Information Systems*, vol. 33, no. 1, pp. 1–33, 2012.

[201] M. B. Zafar, I. Valera, M. Gomez Rodriguez, and K. P. Gummadi, "Fairness beyond disparate treatment & disparate impact: Learning classification without disparate mistreatment," in *Proceedings of the 26th international conference on world wide web*, 2017, pp. 1171–1180.

[202] M. B. Zafar, I. Valera, M. G. Rogriguez, and K. P. Gummadi, "Fairness constraints: Mechanisms for fair classification," in *Artificial Intelligence and Statistics*. PMLR, 2017, pp. 962–970.

[203] A. Agarwal, A. Beygelzimer, M. Dudík, J. Langford, and H. Wallach, "A reductions approach to fair classification," in *International Conference on Machine Learning*. PMLR, 2018, pp. 60–69.

[204] F. Ahmed, Y. Bengio, H. van Seijen, and A. Courville, "Systematic generalisation with group invariant predictions," in *International Conference on Learning Representations*, 2021.

[205] H. Shah, K. Tamuly, A. Raghunathan, P. Jain, and P. Netrapalli, "The pitfalls of simplicity bias in neural networks," in *Advances in Neural Information Processing Systems 33*, 2020.

[206] A. Chouldechova, "Fair prediction with disparate impact: A study of bias in recidivism prediction instruments," *Big data*, vol. 5, no. 2, pp. 153–163, 2017.

[207] H. Edwards and A. Storkey, "Censoring representations with an adversary," in *International Conference on Learning Representations*, 2016.

[208] D. Madras, E. Creager, T. Pitassi, and R. Zemel, "Learning adversarially fair and transferable representations," in *International Conference on Machine Learning*. PMLR, 2018, pp. 3384–3393.

[209] H. Zhao, A. Coston, T. Adel, and G. J. Gordon, "Conditional learning of fair representations," in *International Conference on Learning Representations*, 2020.

[210] T. Hashimoto, M. Srivastava, H. Namkoong, and P. Liang, "Fairness without demographics in repeated loss minimization," in *International Conference on Machine Learning*. PMLR, 2018, pp. 1929–1938.

[211] P. Lahoti, A. Beutel, J. Chen, K. Lee, F. Prost, N. Thain, X. Wang, and E. Chi, "Fairness without demographics through adversarially reweighted learning," in *Advances in Neural Information Processing Systems 33*, 2020.

[212] J. Duchi, T. Hashimoto, and H. Namkoong, "Distributionally robust losses for latent covariate mixtures," *arXiv preprint arXiv:2007.13982*, 2020.

[213] U. Hébert-Johnson, M. Kim, O. Reingold, and G. Rothblum, "Multicalibration: Calibration for the (computationally-identifiable) masses," in *International Conference on Machine Learning*. PMLR, 2018, pp. 1939–1948.

[214] M. P. Kim, A. Ghorbani, and J. Zou, "Multiaccuracy: Black-box post-processing for fairness in classification," in *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, 2019, pp. 247–254.

[215] D. Mandal, S. Deng, S. Jana, J. M. Wing, and D. Hsu, "Ensuring fairness beyond the training data," in *Advances in Neural Information Processing Systems 33*, 2020.

[216] Y. Roh, K. Lee, S. Whang, and C. Suh, "Fr-train: A mutual information-based approach to fair and robust training," in *International Conference on Machine Learning*. PMLR, 2020, pp. 8147–8157.

[217] W. J. Murdoch, C. Singh, K. Kumbier, R. Abbasi-Asl, and B. Yu, "Definitions, methods, and applications in interpretable machine learning," *Proceedings of the National Academy of Sciences*, vol. 116, no. 44, pp. 22 071–22 080, 2019.

[218] J. Adebayo, J. Gilmer, M. Muelly, I. Goodfellow, M. Hardt, and B. Kim, "Sanity checks for saliency maps," in *Advances in Neural Information Processing Systems*, 2018, pp. 9505–9515.

[219] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-cam: Visual explanations from deep networks via gradient-based localization," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 618–626.

[220] M. Sundararajan, A. Taly, and Q. Yan, "Axiomatic attribution for deep networks," in *International Conference on Machine Learning*, 2017, pp. 3319–3328.

[221] V. Ramanishka, A. Das, J. Zhang, and K. Saenko, "Top-down visual saliency guided by captions," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 7206–7215.

[222] P. Dabkowski and Y. Gal, "Real time image saliency for black box classifiers," in *Advances in Neural Information Processing Systems*, 2017, pp. 6967–6976.

[223] P. W. Koh and P. Liang, "Understanding black-box predictions via influence functions," in *International Conference on Machine Learning*, 2017, pp. 1885–1894.

[224] S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," in *Advances in neural information processing systems*, 2017, pp. 4765–4774.

[225] J. Friedman, T. Hastie, R. Tibshirani *et al.*, *The elements of statistical learning*. Springer series in statistics New York, 2001, vol. 1, no. 10.

[226] R. Tibshirani, "Regression shrinkage and selection via the lasso," *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 58, no. 1, pp. 267–288, 1996.

[227] T. J. Hastie and R. J. Tibshirani, *Generalized additive models*. CRC press, 1990, vol. 43.

[228] J. H. Friedman, B. E. Popescu *et al.*, "Predictive learning via rule ensembles," *The Annals of Applied Statistics*, vol. 2, no. 3, pp. 916–954, 2008.

[229] B. Letham, C. Rudin, T. H. McCormick, D. Madigan *et al.*, "Interpretable classifiers using rules and bayesian analysis: Building a better stroke prediction model," *The Annals of Applied Statistics*, vol. 9, no. 3, pp. 1350–1371, 2015.

[230] D. Alvarez-Melis and T. Jaakkola, "A causal framework for explaining the predictions of black-box sequence-to-sequence models," in *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, 2017, pp. 412–421.

[231] L. Anne Hendricks, R. Hu, T. Darrell, and Z. Akata, "Grounding visual explanations," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 264–279.

[232] Y. Goyal, Z. Wu, J. Ernst, D. Batra, D. Parikh, and S. Lee, "Counterfactual visual explanations," in *International Conference on Machine Learning*, 2019, pp. 2376–2384.

[233] C. Xiong, N. Shukla, W. Xiong, and S.-C. Zhu, "Robot learning with a spatial, temporal, and causal and-or graph," in *2016 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2016, pp. 2144–2151.

[234] Y. Xu, L. Qin, X. Liu, J. Xie, and S.-C. Zhu, "A causal and-or graph model for visibility fluent reasoning in tracking interacting objects," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2178–2187.

[235] J. Kim and J. Canny, "Interpretable learning for self-driving cars by visualizing causal attention," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2942–2950.

[236] H. Ye, C. Xie, T. Cai, R. Li, Z. Li, and L. Wang, "Towards a Theoretical Framework of Out-of-Distribution Generalization," *CoRR*, 2021.

[237] I. Gulrajani and D. Lopez-Paz, "In search of lost domain generalization," in *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. [Online]. Available: https://openreview.net/forum?id=lQdXeXDoWtI

[238] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.

[239] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, "Language models are few-shot learners," in *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, Eds., 2020. [Online]. Available: https://proceedings.neurips.cc/paper/2020/hash/1457c0d6bfcb4967418bfb8ac142f64a-Abstract.html

[240] T. Chen, S. Kornblith, M. Norouzi, and G. E. Hinton, "A simple framework for contrastive learning of visual representations," in *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, ser. Proceedings of Machine Learning Research, vol. 119. PMLR, 2020, pp. 1597–1607. [Online]. Available: http://proceedings.mlr.press/v119/chen20j.html

**Zheyan Shen** is a Ph.D candidate in Department of Computer Science and Technology, Tsinghua University. He received his B.S. from the Department of Computer Science and Technology, Tsinghua University in 2017. His research interests include causal inference, stable prediction under selection bias and interpretability of machine learning.

**Jiashuo Liu** received his BE degree from the Department of Computer Science and Technology, Tsinghua University in 2020. He is currently pursuing a Ph.D. Degree in the Department of Computer Science and Technology at Tsinghua University. His research interests focus on invariant learning and distributionally robust learning, especially in developing algorithms with stable performance under distributional shifts.

**Yue He** is a Ph.D candidate in Department of Computer Science and Technology, Tsinghua University. He received his B.S. from the School of Computer Science and Technology, Beihang University in 2018. His research interests include Causal Discovery, Out-of-Distribution Generalizaion, Network Computing and Recommender System.

**Xingxuan Zhang** is a Ph.D student in Department of Computer Science and Technology, Tsinghua University. He received his M.S and B.S from Shanghai Jiao Tong University in 2017 and 2013, respectively. His research interests include out-of-distribution generalization, domain generalization and cross-domain object detection.

**Renzhe Xu** is a third year PhD student in Computer Science and Technology at Tsinghua University, where he obtained a bachelor's degree in Computer Science and Technology in 2019. His research interests include fair machine learning, causal inference, and data mining.

**Han Yu** is a Ph.D candidate in Department of Computer Science and Technology, Tsinghua University. He received his B.S. from the Department of Computer Science and Technology, Tsinghua University in 2021. His research interests include deep learning theory, out-of-distribution generalization and stable learning.

**Peng Cui** is an Associate Professor with tenure in Tsinghua University. He got his PhD degree from Tsinghua University in 2010. His research interests include causally-regularized machine learning, network representation learning, and social dynamics modeling. He has published more than 100 papers in prestigious conferences and journals in data mining and multimedia. His recent research won the IEEE Multimedia Best Department Paper Award, SIGKDD 2016 Best Paper Finalist, ICDM 2015 Best Student Paper Award, SIGKDD 2014 Best Paper Finalist, IEEE ICME 2014 Best Paper Award, ACM MM12 Grand Challenge Multimodal Award, and MMM13 Best Paper Award. He is PC co-chair of CIKM2019 and MMM2020, SPC or area chair of WWW, ACM Multimedia, IJCAI, AAAI, etc., and Associate Editors of IEEE TKDE, IEEE TBD, ACM TIST, and ACM TOMM etc. He received ACM China Rising Star Award in 2015, and CCF-IEEE CS Young Scientist Award in 2018. He is now a Distinguished Member of ACM and CCF, and a Senior Member of IEEE.