

---

# ZIN: When and How to Learn Invariance Without Environment Partition?

---

**Yong Lin**

HKUST

ylindf@connect.ust.hk

**Shengyu Zhu\***

Huawei Noah's Ark Lab

zhushyu@outlook.com

**Lu Tan**

Tsinghua University

tanl21@mails.tsinghua.edu.cn

**Peng Cui**

Tsinghua University

cuip@tsinghua.edu.cn

## Abstract

It is commonplace to encounter heterogeneous data, of which some aspects of the data distribution may vary but the underlying causal mechanisms remain constant. When data are divided into distinct environments according to the heterogeneity, recent invariant learning methods have proposed to learn robust and invariant models using this environment partition. It is hence tempting to utilize the inherent heterogeneity even when environment partition is not provided. Unfortunately, in this work, we show that learning invariant features under this circumstance is fundamentally impossible without further inductive biases or additional information. Then, we propose a framework to jointly learn environment partition and invariant representation, assisted by additional auxiliary information. We derive sufficient and necessary conditions for our framework to provably identify invariant features under a fairly general setting. Experimental results on both synthetic and real world datasets validate our analysis and demonstrate an improved performance of the proposed framework. Our findings also raise the need of making the role of inductive biases more explicit when learning invariant models without environment partition in future works. Codes are available at [https://github.com/linyongver/ZIN\\_official](https://github.com/linyongver/ZIN_official).

## 1 Introduction

Machine learning algorithms with empirical risk minimization (ERM) generally assume independent and identically distributed (i.i.d.) data in training and test sets. Due to changing circumstances, selection bias, or time shifts, data distributions are often heterogeneous across different environments in practical applications. When the distributions of training and test data are indeed different, model performance can severely degrade [47, 17, 38], as ERM based algorithms may exploit the spurious correlations that are particularly useful to ERM in the training data but do not exist in the test data.

To tackle the distribution-shift or out-of-distribution generalization problem, a recent line of methods propose to utilize the causally invariant mechanisms (rather than the spurious correlations in the training data) that are supposed to be stable across different environments. For example, Peters et al. [36] proposes to exploit the invariance principle to learn a linear model that merely relies on the direct causes of the target. Such a model is shown to be robust to potential interventions on any variable

---

\*Corresponding author.

except the target itself. Arjovsky et al. [4] propose invariant risk minimization (IRM) to capture invariant correlations by learning representations that elicit an optimal invariant predictor across multiple training environments. Ahuja et al. [2], Jin et al. [21], Krueger et al. [24], Xie et al. [52], Chen et al. [8, 9] further develop variants of IRM by introducing game theory, regret minimization, variance penalization, multi-objective optimization, etc., and Xu and Jaakkola [54], Chang et al. [7], Lin et al. [26] try to learn invariant features by coupled adversarial neural networks. Some recent works report that IRM is less effective when applied to large neural networks [26, 15]. Lin et al. [27] find that this can be largely attributed to the overfitting problem, and Zhou et al. [57, 56] propose to alleviate this issue through imposing sparsity constraint and sample reweighting, respectively.

Noticeably, the aforementioned invariant learning methods require datasets to be explicitly partitioned into environments (or domains) according to the heterogeneity underlying the data. However, such an environment partition may be unavailable or hard to obtain in practice [29, 10]. Recently, a line of works try to learn invariance without environment indexes where the dataset is assembled by merging data from multiple environments. For instance, Creager et al. [10] propose environmental inference for invariant learning (EUIL), a two-stage method by firstly inferring the environments with an ERM based trained biased model and then performing invariant learning on inferred environments. Liu et al. [29] devise an interactive mechanism, called heterogeneous risk minimization (HRM), with environment inference and invariant learning on raw feature level, and Liu et al. [30] extend it to representation level with kernelized trick.

However, as we will show, learning invariant models under this circumstance (with only input-label pairs) can be *risky*, both theoretically and empirically. As our first contribution, in Section 4, we prove that learning invariant features from heterogeneous data without environment indexes is in theory impossible. Specifically, we provide a counter-example and a general theoretic result (Theorem 1), to show that the causally invariant features are unidentifiable if no environment information is provided. This impossibility result, similar to the identifiability issue in causal discovery [45, 37], independent component analysis (ICA) [20], and unsupervised learning of disentangled representations [32], motivates us to consider when and how to learn invariant features.

In this work, we turn to additional auxiliary variables that encode some information about the latent heterogeneity. This strategy can be treated as a case between the two existing directions: 1) it is less restrictive than requiring ideal environment partition according to distribution heterogeneity; and 2) it also obtains theoretic guarantee which is missing in the case with no environment information at all. Notably, such auxiliary information is often cheaply available for every input in practice [53, 51]. Examples include time index of the data in time series forecasting tasks [6, 33], locations (longitude and latitude) of collected satellite data in remote sensing [14, 40], and meta information associated with an instance [31]. Based on the additionally observed variables, we proceed to propose a framework to jointly learn environment partition and invariant representation in Section 5. Under a fairly general setting, we derive both sufficient and necessary conditions for our framework to identify invariant features. Extensions to other settings are discussed in Appendix A.

Finally, Section 6 presents experimental analysis on both synthetic and real world datasets. The proposed framework achieves an improved performance over existing methods like EUIL and HRM, and has a comparable performance to IRM with ground-truth environment partition.

## 2 Related Work

Invariant learning methods can be interpreted as robust to certain interventions or heterogeneity from a causal perspective [36]. On the other hand, heterogeneous data have also placed challenges and potential benefits in causal inference tasks. For example, Huang et al. [19] consider heterogeneous datasets to identify variables with changing local mechanism and further the causal directions. Here confounders are assumed as functions of the domain index or time. Tillman and Spirtes [48], Danks et al. [11], Huang et al. [18] study the multiple dataset setting with non-identical sets of variables for causal discovery, where the datasets may have different exogenous noise distributions and only part of the variables is present in each dataset.

Besides the invariant learning methods described in Section 1, another large class of methods for generalizing beyond training data is distributionally robust optimization (DRO) [5, 43, 25, 13, 12, 55]. DRO methods propose to optimize the worst-case risk over a set of distributions close to the training distribution. A notable instance of DRO methods is group DRO, which optimizes the worst-case loss over groups in the training data [41]. As with most invariant learning methods, group DRO generally requires a partition of groups obtained according to group annotation and label of each data sample.

Obtaining group annotations, however, can be costly or even impossible in practice. We may not know *a priori* the inherent spurious features associated with data. Thus, a number of works have proposed to infer group partition or directly identify the minority group by looking at features produced by biased models, e.g., [44, 1, 10, 28]. Notice that sometimes a small number of labelled annotations are still needed to find a proper trained model, which is different from the setting considered in this work. For example, Liu et al. [28] upweights samples with high losses from the initial ERM model and relies on a small validation set of annotated data to tune parameters. Additionally, Nam et al. [34], Sanh et al. [42] try to obtain a more robust model by boosting from the wrongly specified samples, based on the observation that models with limited capacity tend to learn spurious correlations or shortcuts. Unfortunately, as we will show in Section 4, it is in theory impossible to distinguish spurious or invariant features from only the observed data consisting of input-label pairs.

### 3 Preliminaries

Throughout this paper, upper-cased letters such as  $\mathbf{X} \in \mathbb{R}^d$  and  $Y \in \mathbb{R}$  denote random variables, and lower-cased letters such as  $x$  and  $y$  denote deterministic instances. Suppose there exist invariant features  $\mathbf{X}_v \in \mathbb{R}^{d_v}$  and spurious (non-invariant) features  $\mathbf{X}_s \in \mathbb{R}^{d_s}$ . We observe a scrambled version  $\mathbf{X} = q(\mathbf{X}_v, \mathbf{X}_s) \in \mathbb{R}^d$  with  $q$  being an injective function. Consider the dataset  $\mathcal{D} := \{(\mathbf{x}_i, y_i)\}_{i=1}^n$  where data are collected from multiple environments  $e \in \mathcal{E}_{supp}$ . We use superscript  $e$  to indicate the environment index with a variable, e.g.,  $\mathbf{X}^e$  and  $Y^e$ , when it is useful to make the index explicit. Notice that the true environment index  $e$  is not provided in training, unless otherwise stated.

In this work, we assume an underlying structural causal model (SCM) governing the data generation process [35, 37]. An SCM considers a set of variables associated with vertices of a directed acyclic graph, where directed edges represent direct causation. Each variable is obtained as a result of an assignment of a deterministic function depending on the parental variables in the graph and an exogenous random variable. As an example, we present the following data generation process that is also considered in [3, 4]:

$$Y^e = g_v(\mathbf{X}_v^e, \epsilon_v), \quad \mathbf{X}_v^e \perp \epsilon_v; \quad \mathbf{X}^e = q(\mathbf{X}_v^e, \mathbf{X}_s^e), \quad (1)$$

where  $g_v$  denote a non-degenerate deterministic function and  $\epsilon_v$  is an independent noise variable. We observe data  $\{\mathbf{X}^e, Y^e\}_{e=1}^E$  from multiple environments, each with probability  $P(\mathbf{X}^e, Y^e)$ . Then the marginal distribution of the mixed data can be represented as  $P(\mathbf{X}, Y) = \sum_{e=1}^E \alpha_e P(\mathbf{X}^e, Y^e)$  for some  $\alpha_e > 0$  and  $\sum_{e=1}^E \alpha_e = 1$ .

We consider the model as the composition of a feature extractor  $\Phi$  and a classifier (or predictor)  $f_\omega$  that is parameterized by  $\omega$ . To seek for out-of-distribution generalization ability, we hope that  $\Phi$  merely encodes the information of invariant features. Let  $\ell(\cdot, \cdot)$  denote a loss function such as cross-entropy loss and squared error. Our goal is to learn a robust and invariant model that minimizes the following objective over all the considered environments  $e \in \mathcal{E}_{supp}$ :

$$\sup_{e \in \mathcal{E}_{supp}} \mathbb{E}_{\mathbf{X}^e, Y^e} [\ell(f_\omega(\Phi(\mathbf{X}^e)), Y^e)].$$

While there is no single, common definition of *invariant feature* in the literature, we treat the direct causal parents of  $Y$  as invariant features, as in [36, 4]. This is because the conditional probability of the outcome  $Y$  given its parental variables or direct causes remains unchanged with any intervention (not on  $Y$ ). In contrast, if a non-parental variable is included in the conditioned variables, then the conditional probability could change after some intervention. As such, the underlying SCM governing data generation determines the desired invariance.

## 4 Impossibility Result

The first question we ask is whether it is possible to learn invariant models from the heterogeneous data of multiple environments with unknown environmental indexes. The following example immediately gives a negative answer.

**Example.** Consider that there are two environments  $e \in \{1, 2\}$  with  $\alpha_1 = \alpha_2 = 0.5$ , and we observe binary-valued features  $X_1, X_2$  and label  $Y$ . Suppose the joint distribution of mixed environments is given in Eq. (2). A learning algorithm that tries to learn invariant feature based on this dataset would result in a model that (deterministically) depends on either  $X_1, X_2$ , or both to predict  $Y$ . However, as discussed in last section, whether the learned features are invariant is determined by the underlying data generation process. We now present two possible data generation processes that generate the same distribution yet have different invariant features:

$$\begin{cases} Y = 0, & w.p. 0.5, \\ Y = 1, & w.p. 0.5, \\ X_2 = X_1 = Y, & w.p. 0.6375, \\ X_2 \neq X_1 = Y, & w.p. 0.1125, \\ X_1 \neq X_2 = Y, & w.p. 0.2125, \\ X_2 = X_1 \neq Y, & w.p. 0.0375. \end{cases} \quad (2)$$

- $X_1$  is the invariant feature while  $X_2$  is spurious, where  $p_s^{e=1} = 0.8$  and  $p_s^{e=2} = 0.9$ :

$$X_1 \sim \text{Bernoulli}(0.5), \quad Y = \begin{cases} X_1, & w.p. 0.75, \\ 1 - X_1, & w.p. 0.25, \end{cases} \quad X_2^e = \begin{cases} Y, & w.p. p_s^e, \\ 1 - Y, & w.p. 1 - p_s^e. \end{cases} \quad (3)$$

- $X_2$  is the invariant feature while  $X_1$  is spurious where  $p_s^{e=1} = 0.8$  and  $p_s^{e=2} = 0.7$ :

$$X_2 \sim \text{Bernoulli}(0.5), \quad Y = \begin{cases} X_2, & w.p. 0.85, \\ 1 - X_2, & w.p. 0.15, \end{cases} \quad X_1^e = \begin{cases} Y, & w.p. p_s^e, \\ 1 - Y, & w.p. 1 - p_s^e. \end{cases} \quad (4)$$

The joint distribution of the mixed environments for each data generation process is consistent with Eq. (2). Thus, from the joint distribution  $P(X_1, X_2, Y)$ , a learned model only depends on either  $X_1, X_2$  or both, and fails to generalize for at least one of the two data generation processes. On the other hand, when the partition is given, one can verify that IRM would correctly identify  $X_1$  (resp.  $X_2$ ) as the invariant feature for the first (resp. second) scenario.

The above toy example is inspired by commonly-used classification tasks in the IRM literature, e.g., CMNIST [4] and CifarMnist [26]. For instance, in CMNIST, invariant feature  $X_1$  denotes the semantic feature of the shape of hand-written digits ‘0’ and ‘1’, and spurious feature  $X_2$  represents the color, which is either red or green. Label  $Y \in \{0, 1\}$  corresponds to the digit shape and is also binary. Indeed, the construction procedure of CMNIST in [4] can be exactly described by the data generation process in Eq. (3). To further demonstrate the impossibility result, we next conduct an empirical validation on a new dataset in accordance to the data generation process in Eq. (4)

**Empirical Validation.** We construct a variant of CMNIST according to the second data generation process. In this new dataset, color is the invariant feature and digit shape is spurious. We use the same notations where  $X_1$  stands for digit shape,  $X_2$  is the color, and  $Y$  denotes label. As described in Eq. (4), the label corresponds to the digit color and the digit shape is spuriously correlated with the label. We name this variant of CMNIST as MCOLOR (short for MNIST-COLOR). In the test domain, we set  $p_s^{e=3} = 0.1$  to simulate the distributional shift, same as in CMNIST. We can then compare ERM, EIIL, IRM with environment partition, and LfF (learning from failures) [34] that tries to learn a robust model by boosting from wrongly specified samples of shallow neural networks. The empirical results are reported in Table 1. For ERM (oracle), we train the model only on the invariant feature, i.e., digit shape in CMNIST and color in MCOLOR.

We can see that the EIIL method performs poorly on the MCOLOR dataset. This is due to the inductive bias of EIIL and it indeed relies on the digit shape as the invariant feature. Since we have no prior knowledge of the data generation process, the true invariant feature can be the color, e.g., in the MCOLOR dataset. Similarly, ERM and LfF all rely on either color or shape as the invariant feature and would fail on at least one of CMNIST and MCOLOR. On the other hand, if environment partition is available, IRM can still learn the desired invariant feature.

**Table 1:** Experimental results on CMNIST and MCOLOR.

Method	Env Partition	CMNIST		MCOLOR	
		Train Acc	Test Acc	Train Acc	Test Acc
ERM (oracle)	No	75.2±0.2	72.1±0.1	85.0±0.0	85.0±0.0
ERM	No	86.4 ±0.0	14.5 ±0.1	86.3±0.1	80.1±0.6
IRM	Yes	71.4 ±0.3	66.4 ±0.3	84.9±0.1	84.7±0.3
EIIL	No	72.5 ±0.7	67.2 ±3.3	74.0±0.4	17.8±0.4
LfF	No	76.7 ±0.3	21.2 ±0.4	76.6±0.0	74.2±0.0

Moreover, such examples are not rare. For any data distribution  $P(\mathbf{X}, Y)$  generated by some data generation process like (1) (which can be replaced by other forms of SCM), we can find a different data generation process that induces the same distribution yet has different invariant/spurious features. Similarly, it is impossible to identify the spurious features only based on  $P(\mathbf{X}, Y)$ . This result is formally summarized in Theorem 1.

**Theorem 1.** *Let  $\mathbf{X}_v$  and  $\mathbf{X}_s$  be respectively the invariant and spurious features, with label  $Y$ . For the joint distribution  $P(\mathbf{X}, Y)$  consisting of data from multiple environments with  $P(\mathbf{X}_s, Y) \neq 0$ , we can find  $\mathbf{X}'_v$  and  $\mathbf{X}'_s$  as invariant and spurious features so that:*

- $\mathbf{X}'_v \neq \mathbf{X}_v$  and  $\mathbf{X}'_v \cap \mathbf{X}_s \neq \emptyset$ ;
- $(\mathbf{X}'_v, \mathbf{X}'_s)$  together with some noise variables generate the same distribution  $P(\mathbf{X}, Y)$ .

A proof is provided in Appendix B.1. Under this circumstance, we have to introduce additional assumptions/conditions or certain “inductive bias” to identify the invariant features [45, 37, 20, 32, 50]. The latter may come implicitly from the proposed algorithm, model, and/or the data. We believe that this is the reason that makes existing methods achieve improved empirical performance in their considered scenarios. For example, in [29, 30], the discrepancy of spurious features among clusters are expected to be larger than that of causal features, and in [10, 34, 42, 28], the ERM model of the first stage should heavily or fully rely on the spurious feature. However, these inductive biases may not be always guaranteed. Thus, as suggested by Theorem 1, the role of inductive biases shall be discussed more explicitly when learning invariant models without environment partition.

Alternatively, in this work, we consider that there exists additionally observed variable  $\mathbf{Z}$  with the data, which has been often considered in the literature [19, 20, 22, 53, 51]. The variable  $\mathbf{Z}$  could be, for example, the time index in a time series, some kind of class labels, concurrently observed variables, or human annotations on potential unmeasured variable [46]. It is worth noting that Xie et al. [53] also take advantage of auxiliary information to help improve OOD performance in a semi-supervised manner. However, they consider a different setting with access to unlabeled (out-of-domain) test data together with additional unlabeled training data. In the next section, we show how to utilize  $\mathbf{Z}$  to learn invariant models from the heterogeneous dataset.

## 5 ZIN: Learning Invariance with Additional Auxiliary Information

We consider that there exists additional auxiliary information  $\mathbf{Z} \in \mathbb{R}^{d_z}$  in company with the data  $(\mathbf{X}, Y)$ . In this section, we propose ZIN, *auxiliary information Z for environmental INference*, for invariant learning from the heterogeneous dataset  $\mathcal{D}$  without environment partition. We will derive conditions for invariance identification and show these conditions are both sufficient and necessary in a fairly general setting.

### 5.1 Method

We aim to learn a function  $\rho(\cdot) : \mathbb{R}^{d_z} \rightarrow \mathbb{R}^K$  that softly assigns a sample to  $K$  environments. Here  $K$  is a pre-specified number (a hyper-parameter) and its choice will be empirically investigated in our experiments. Let  $\rho^{(k)}(\cdot)$  denote the  $k$ -th entry of  $\rho(\cdot)$ , with  $\rho(\mathbf{Z}) \in [0, 1]^K$  and  $\sum_k \rho^{(k)}(\mathbf{Z}) = 1$ . Denote the ERM loss as  $\mathcal{R}(\omega, \Phi) = \frac{1}{n} \sum_{i=1}^n \ell(f_\omega(\Phi(\mathbf{x}_i)), y_i)$  and the loss in the  $k$ -th inferred environment as  $\mathcal{R}_{\rho^{(k)}}(\omega, \Phi) = \frac{1}{n} \sum_{i=1}^n \rho^{(k)}(\mathbf{z}_i) \ell(f_\omega(\Phi(\mathbf{x}_i)), y_i)$ .

Recall that IRM [4] learns an invariant representation  $\Phi$ , upon which there is a classifier  $f_\omega$  that is simultaneously optimal in all environments. Suppose that the environments have been given according to a fixed  $\rho(\cdot)$ . To measure the optimality of  $f_\omega$  in the  $k$ -th environment, we can fit an environment-dependent classifier  $f_{\omega_k}$  on the data from that environment. If  $f_{\omega_k}$  achieves a smaller loss, then we know that  $f_\omega$  is not optimal in this environment. We can further train a set of environment-dependent classifiers  $\{f_{\omega_k}\}_{k=1}^K$ , one for each environment, to measure whether  $f_\omega$  is simultaneously optimal in all environments. Thus, when  $\rho(\cdot)$  is provided, our formulation to learn invariance is as follows:

$$\min_{\omega, \Phi} \max_{\{\omega_k\}} \mathcal{L}(\Phi, \omega, \omega_1, \dots, \omega_K, \rho) := \mathcal{R}(\omega, \Phi) + \underbrace{\lambda \sum_{k=1}^K [\mathcal{R}_{\rho^{(k)}}(\omega, \Phi) - \mathcal{R}_{\rho^{(k)}}(\omega_k, \Phi)]}_{\text{invariance penalty}}. \quad (5)$$

If  $\Phi$  extracts spurious features that are unstable in the inferred environments,  $\mathcal{R}_{\rho^{(k)}}(\omega, \Phi)$  will be larger than  $\mathcal{R}_{\rho^{(k)}}(\omega_k, \Phi)$ , resulting in a non-zero invariance penalty.

Next, we consider how to learn the partition function  $\rho(\cdot)$ . A *good* partition function should generate environments among which the spurious features exhibit instability, so that there is a large penalty if  $\Phi$  extracts spurious features. Thus, we seek for an environment partition that maximizes the invariance penalty. The overall framework is provided below:

$$\min_{\omega, \Phi} \max_{\rho, \{\omega_1, \dots, \omega_K\}} \mathcal{L}(\Phi, \omega, \omega_1, \dots, \omega_K, \rho). \quad (6)$$

The idea of our method can be summarized as inferring an environment partition and then learning invariant models based on the inferred environments. Creager et al. [10] has used a similar intuition and adopted a pre-trained biased model to estimate the environments. However, their two-stage method cannot be jointly optimized, and the environment partition relies on the given model and lacks a theoretical guarantee. In contrast, we will first present sufficient conditions on  $\mathcal{Z}$  so that the proposed framework can provably identify the invariant features. We further show that these conditions are also necessary, and that violation of these conditions will lead to failure of invariance identification. Assigning independent weights to each data sample, which is adopted in EIL [10], is unfortunately an instance of such violations and may fail to identify invariant features.

## 5.2 Sufficient Conditions for Identifiability

In this section, we try to understand ZIN from a theoretical perspective. We start with a simple yet general setting:  $\mathbf{X} = [\mathbf{X}_v; \mathbf{X}_s]$  (i.e., no scramble on the observation),  $\Phi \in \{0, 1\}^d$  is an element-wise feature selection mask, and  $f_\omega$  is a general non-linear function. We also focus on classification tasks with  $\ell(\cdot, \cdot)$  being the cross-entropy loss, as there is an interesting connection to (conditional) Shannon entropy. Extensions to other loss functions and linear feature transformations are left to Appendix A.1 and A.2, respectively.

In this setting, our goal is equivalent to learning the optimal feature mask that merely selects invariant features, i.e.,  $\Phi_v = [\mathbf{1}^{d_v}; \mathbf{0}^{d_s}]$ . For a given feature mask  $\Phi$ , the objective function is equal to  $\hat{\mathcal{L}}(\Phi) = \min_{\omega} \max_{\rho, \{\omega_k\}} \mathcal{L}(\Phi, \omega, \omega_1, \dots, \omega_k, \rho)$ . Then ZIN can correctly identify the invariant features, or solution to Problem 6 is equivalent to  $\Phi_v$ , if and only if  $\hat{\mathcal{L}}(\Phi_v) < \hat{\mathcal{L}}(\Phi)$  for all  $\Phi \neq \Phi_v$ . This observation will be used to establish our main theoretic result.

With a little abuse of notation, we use  $H(Y|\mathbf{X}')$  to denote expected loss of an optimal classifier over some  $\mathbf{X}'$  and  $Y$ , and similarly  $H(Y|\rho(\mathcal{Z}), \mathbf{X}')$  to denote the minimum risk  $\sum_{k=1}^K \mathcal{R}_{\rho^{(k)}}(\omega_k, \Phi)$  for a given  $\rho(\mathcal{Z})$ . With cross-entropy loss and when  $\rho(z_i)$  gives exactly one environment, i.e.,  $\rho(\mathcal{Z})$  is one-hot, it can be verified that the optimal loss  $H(\cdot|\cdot)$  coincides with conditional entropy. In the following we first state the assumptions for our identifiability result.

**Assumption 1.** *For a given feature mask  $\Phi$  and any constant  $\epsilon > 0$ , there exists  $f \in \mathcal{F}$  such that  $\mathbb{E}[\ell(f(\Phi(\mathbf{X})), Y)] \leq H(Y|\Phi(\mathbf{X})) + \epsilon$ .*

**Assumption 2.** *If a feature violates the invariance constraint, adding another feature would not make the penalty vanish, i.e., there exists a constant  $\delta > 0$  so that for spurious feature  $\mathbf{X}_1 \subset \mathbf{X}_s$  and any feature  $\mathbf{X}_2 \subset \mathbf{X}$ ,  $H(Y|\mathbf{X}_1, \mathbf{X}_2) - H(Y|\rho(\mathcal{Z}), \mathbf{X}_1, \mathbf{X}_2) \geq \delta (H(Y|\mathbf{X}_1) - H(Y|\rho(\mathcal{Z}), \mathbf{X}_1))$ .*

**Assumption 3.** For any distinct features  $\mathbf{X}_1, \mathbf{X}_2$ ,  $H(Y|\mathbf{X}_1, \mathbf{X}_2) \leq H(Y|\mathbf{X}_1) - \gamma$  with fixed  $\gamma > 0$ .

Assumption 1 is a common assumption that requires the function space  $\mathcal{F}$  be rich enough such that, given  $\Phi$ , there exists  $f \in \mathcal{F}$  that can fit  $P(Y|\Phi(\mathbf{X}))$  well. Assumption 2 aims to ensure a sufficient positive penalty if a spurious feature is included (see Appendix A.3 for further discussion). Assumption 3 indicates that any feature contains some useful information w.r.t.  $Y$ , which cannot be explained by other features. Otherwise, we can simply remove such a feature (e.g., by variable selection methods [49]), as it does not affect prediction. We next present our sufficient conditions for ZIN to identify invariant features.

**Condition 1** (Invariance Preserving Condition). *Given invariant feature  $\mathbf{X}_v$  and any function  $\rho(\cdot)$ , it holds that  $H(Y|\mathbf{X}_v, \rho(\mathbf{Z})) = H(Y|\mathbf{X}_v)$ .*

**Condition 2** (Non-invariance Distinguishing Condition). *For any feature  $\mathbf{X}_s^k \in \mathbf{X}_s$ , there exists a function  $\rho(\cdot)$  and a constant  $C > 0$  such that  $H(Y|\mathbf{X}_s^k) - H(Y|\mathbf{X}_s^k, \rho(\mathbf{Z})) \geq C$ .*

We remark that Condition 1 can be met if  $H(Y|\mathbf{X}_v, \mathbf{Z}) = H(Y|\mathbf{X}_v)$  (a proof is provided in Appendix B.6). Condition 1 requires that invariant features should remain invariant w.r.t any environment partition induced by  $\rho(\mathbf{Z})$ . Otherwise, if there exists a partition where an invariant feature becomes non-invariant, then this feature would induce a positive penalty. Condition 2 implies that for each spurious feature, there exists at least one partition so that this feature is non-invariant in the split environments. If a spurious feature does not incur any invariance penalty in all possible environment partitions, we can never distinguish it from true invariant features. With these conditions, our main result follows and a proof is given in Appendix B.2.

**Theorem 2** (Identifiability of Invariant Features). *With Assumptions 1-3 and Conditions 1-2, if  $\epsilon < \frac{C\gamma\delta}{4\gamma+2C\delta H(Y)}$  and  $\lambda \in [\frac{H(Y)+1/2\delta C}{\delta C-4\epsilon} - \frac{1}{2}, \frac{\gamma}{4\epsilon} - \frac{1}{2}]$ , then we have  $\hat{\mathcal{L}}(\Phi_v) < \hat{\mathcal{L}}(\Phi)$  for all  $\Phi \neq \Phi_v$ , where  $H(Y)$  denotes the entropy of  $Y$ . Thus, the solution to Problem 6 identifies invariant features.*

### 5.3 Necessary Conditions for Identifiability

Conditions 1 and 2 may appear rather strong. In this section, we prove that they are also necessary conditions for ZIN to identify invariant features. Specifically, Proposition 1 shows that if Condition 1 is violated, then some invariant features will be excluded in the solution of Problem 6; and in Proposition 2, violation of Condition 2 renders some spurious features included in the solution.

**Proposition 1.** *With Assumptions 1-3, if Condition 1 is violated, i.e., there exists  $\rho(\cdot)$  so that  $H(y|\mathbf{X}_v) - H(y|\mathbf{X}_v, \rho(\mathbf{Z})) \geq C' > 0$ , then there exists a mask  $\Phi' \neq \Phi_v$  with  $\hat{\mathcal{L}}(\Phi_v) > \hat{\mathcal{L}}(\Phi')$ .*

A proof is provided in Appendix B.3, which is similar to the first step of the proof of Theorem 2. A question is how Condition 1 could be violated. On the other hand, while we introduce  $\mathbf{Z}$  as additional variables, it is also interesting to know whether we can obtain a valid partition based on only the training data  $(\mathbf{X}, Y)$ . Below we provide examples regarding these questions, with proofs given in Appendix B.4. As  $h(\text{Index}(\mathbf{X}, Y))$  can be treated as learning independent weights for each sample based on  $\{(x_i, y_i)\}_i$ , this case includes EIL [10] as an instance.

**Corollary 1.** *Condition 1 is violated for the following cases: there exists a function  $\rho(\cdot)$  and an injective function  $h(\cdot)$  so that (a)  $\rho(\mathbf{Z}) = h(Y)$ , (b)  $\rho(\mathbf{Z}) = h(\mathbf{X}, Y)$ , or (c)  $\rho(\mathbf{Z}) = h(\text{Index}(\mathbf{X}, Y))$ .*

We now discuss the necessity of Condition 2, with an additional assumption.

**Assumption 4.** *If two features are invariant w.r.t. an environment partition, then the concatenated features are also invariant. That is, for  $\mathbf{X}_1, \mathbf{X}_2 \subset \mathbf{X}$  and  $\rho(\cdot)$ , if  $H(Y|\mathbf{X}_1) - H(Y|\rho(\mathbf{Z}), \mathbf{X}_1) = 0$  and  $H(Y|\mathbf{X}_2) - H(Y|\rho(\mathbf{Z}), \mathbf{X}_2) = 0$ , we have  $H(Y|\mathbf{X}_1, \mathbf{X}_2) - H(Y|\rho(\mathbf{Z}), \mathbf{X}_1, \mathbf{X}_2) = 0$ .*

**Proposition 2.** *With Assumptions 1, 3 and 4, if Condition 2 is violated, i.e., there exists a spurious feature  $\mathbf{X}_s^k \in \mathbf{X}_s$  such that  $H(Y|\mathbf{X}_s^k) - H(Y|\mathbf{X}_s^k, \rho(\mathbf{Z})) = 0$  for any  $\rho(\cdot)$ , then there exists a feature mask  $\Phi' \neq \Phi_v$  with  $\hat{\mathcal{L}}(\Phi_v) > \hat{\mathcal{L}}(\Phi')$ .*

Since there exists a spurious feature  $\mathbf{X}_s^k$  that is ‘‘invariant’’ in all possible environment partition, adding this feature to  $\Phi_v$  does not induce any invariance penalty but increases the prediction power.

Thus, such a feature mask can achieve a smaller loss than  $\hat{\mathcal{L}}(\Phi_v)$ . For completeness, we provide a proof in Appendix B.5. Proposition 2 again indicates in theory that  $\mathbf{Z}$  should be sufficiently diverse and informative so that each spurious feature can be recognized.

#### 5.4 Choice of the Auxiliary Information

Conditions 1 and 2 present theoretical requirements of  $\mathbf{Z}$  to successfully learn invariant features. We now discuss how to find such auxiliary information given certain prior knowledge.

Take Fig. 1 for example. When we collect a data point (e.g., taking a photo), there often exists some meta information (e.g., time slot, coordinate, and temperature). Suppose that the image is generated according to the causal graph on the right panel of Fig. 1. In this case, the invariant feature  $\mathbf{X}_v$  consists of  $X_1$  and  $X_2$ , and the meta information can be used as  $\mathbf{Z}$  which is colored in green in the causal graph. This is because  $H(Y|X_1, X_2) = H(Y|X_1, X_2, \mathbf{Z})$  or equivalently  $Y \perp \mathbf{Z} \mid X_1, X_2$ , so that Condition 1 holds. This example shows some evidence of choosing the auxiliary information:  $\mathbf{Z}$  and  $Y$  should be d-separated by  $\mathbf{X}_v$ . There are other examples satisfying Condition 1, which are given Fig. 5 in Appendix E. We also illustrate some failure choices in Fig. 6 in Appendix E. Specifically, if  $\mathbf{Z}$  is a child or a direct cause of  $Y$ , then it contains additional information of  $Y$  conditional on the invariant feature, violating Condition 1. This kind of information cannot be used in our framework.

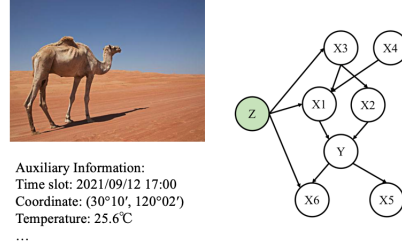


Figure 1: An example of  $\mathbf{Z}$  satisfying Condition 1.

To summarize, with some prior knowledge on the path between  $\mathbf{Z}$  and  $Y$ , we suggest to choose as much auxiliary information as possible under the following condition: *A feasible choice of  $\mathbf{Z}$  should satisfy Condition 1 and should be correlated with certain variables in the causal graph of the SCM. The path between  $\mathbf{Z}$  and  $Y$  should be d-separated by  $\mathbf{X}_v$  (e.g., in Figs. 1 and 5) or there is no path between  $\mathbf{Z}$  and  $Y$ . Notably,  $\mathbf{Z}$  cannot be the parent or child of  $Y$  as shown in Fig. 6.*

## 6 Experiments

This section empirically verifies our theoretic analysis and the effectiveness of ZIN on both synthetic and real world datasets. We compare ZIN with several existing methods: ERM, IRM, [4], group DRO [41], EIIL [10], HRM [29], and LfF [34]. We provide the ground-truth partition to IRM and group DRO. Notice that LfF only works for classification tasks.

We implement both  $\rho$  and  $f$  in ZIN as a two-layer MLP with 32 hidden units. We adopt the first order approximation of Eq. (5), as described in Appendix C. The number of inferred environments  $K$  is set to be 2 as default. We implement  $\Phi$  as a two-layer MLP for the synthetic dataset and house price prediction, ResNet-18 [16] for CelebA, and 1D CNN for Landcover. More details are provided in Appendix C.

### 6.1 Synthetic Dataset

**Temporal Heterogeneity.** We consider temporal heterogeneity with distributional shift w.r.t. time. Let  $t \in [0, 1]$  be time index and  $X_v(t) \in \mathbb{R}$  the invariant feature. The data generation process is

$$X_v(t) \sim \begin{cases} \mathcal{N}(1, \sigma^2), & w.p. 0.5, \\ \mathcal{N}(-1, \sigma^2), & w.p. 0.5, \end{cases} \quad Y(t) \sim \begin{cases} \text{sign}(X_v(t)), & w.p. p_v, \\ -\text{sign}(X_v(t)), & w.p. p'_v, \end{cases} \quad X_s(t) \sim \begin{cases} \mathcal{N}(Y(t), \sigma^2), & w.p. p_s(t), \\ \mathcal{N}(-Y(t), \sigma^2), & w.p. p'_s(t), \end{cases}$$

where  $p_v$  is a constant w.r.t.  $t$ , indicating a stable correlation between  $Y(t)$  and  $X_v(t)$ ,  $p'_v = 1 - p_v$ , and  $p'_s(t) = 1 - p_s(t)$ . Notice here  $p_s(t)$  varies with time  $t$ . A similar setting with two-dimensional spatial variable is considered in Appendix D.1. Our goal is to learn a model that purely relies on  $X_v$ . We simulate two heterogeneous environments along time, namely,  $\{[0, 0.5), [0.5, 1]\}$ , and  $p_s(t)$  will be set differently. We use tuple of  $p_s(t)$  in the two environments to denote a simulated case. For example,  $(0.999, 0.7)$  stands for  $p_s(t) = 0.999, t \in [0, 0.5)$  and  $p_s(t) = 0.7, t \in [0.5, 1]$ . We



**Table 2:** Test Mean and Worst accuracy (%) on four temporal heterogeneity synthetic datasets.

Env Partition	$p_s(t)$	0.999, 0.7				0.999, 0.8				0.999, 0.9			
	$p_v$	0.9		0.8		0.9		0.8		0.9		0.8	
	Test Acc	Mean	Worst	Mean	Worst	Mean	Worst	Mean	Worst	Mean	Worst	Mean	Worst
No	ERM	75.37	57.31	59.65	25.81	68.72	41.97	55.90	15.07	60.61	23.39	52.85	7.57
	EIIL	38.41	16.80	64.89	49.15	50.77	46.67	68.36	56.35	61.99	53.81	70.10	59.36
	HRM	50.00	49.99	49.98	49.93	50.00	49.98	50.01	49.99	50.00	49.98	49.99	49.97
	ZIN	<b>87.50</b>	<b>85.36</b>	<b>77.85</b>	<b>75.39</b>	<b>86.35</b>	<b>82.91</b>	<b>76.79</b>	<b>72.77</b>	<b>83.71</b>	<b>75.89</b>	<b>73.55</b>	<b>64.69</b>
Yes	IRM	87.57	85.47	77.99	75.65	86.57	83.25	77.00	73.39	83.99	76.48	73.84	65.33

evaluate the performance on four distinct test environments with  $p_s \in \{0.999, 0.8, 0.2, 0.1\}$  and  $p_v$  being constant. We use  $t$  as the auxiliary information. More details are provided in Appendix C

**Results.** Table 2 reports the test accuracy. In all simulated settings, the worst accuracy of ERM is much lower than the mean accuracy, indicating that ERM tends to rely on spurious feature  $X_s$ . EIIL can improve the worst accuracy in some cases, e.g., when  $p_s(t) = (0.999, 0.8)$  and  $p_s(t) = (0.999, 0.9)$ . However, its performance is even worse than ERM for some other settings. This may be attributed to the first stage of EIIL, where the trained biased model is not guaranteed and may learn both spurious and invariant features. The proposed method ZIN improves the worst test accuracy significantly. For instance, when  $p_s(t) = (0.999, 0.7)$  and  $p_v = 0.9$ , ZIN outperforms ERM and EIIL by over 28% and 65%, respectively. Moreover, ZIN is very close to IRM that knows the ground-truth environment partition, showing that ZIN can infer the environments effectively. Finally, it seems that HRM does not learn a useful model for this classification task.

**Ablation Study.** In Appendix D.2, we conduct ablation study to verify our theoretical results in Section 5 by choosing different auxiliary variables. We also empirically study the choice of hyper-parameter  $K$ . A notable observation is that  $K$  has a relatively small impact, as shown in Fig. 3.

## 6.2 Real World Datasets

**House Price Prediction.** This experiment considers a real world regression dataset of house sales prices from Kaggle (<https://www.kaggle.com/c/house-prices-advanced-regression-techniques>). The target variable is house price and each sample contains 17 predictive variables like the built year of the house, number of bedrooms, etc. The dataset is split according to the built year, with training dataset in period [1900, 1950] and test dataset in period (1950, 2000]. We normalize the prices of the houses with the same built year, and our target is to predict the normalized price. We choose the built year as auxiliary information for ZIN. As there is no well-defined ground-truth partition for IRM, we manually split the training dataset equally into 5 segments with 10-year range in each segment.

Table 3 reports the mean squared error (MSE), and ZIN again outperforms other methods on test dataset. We also investigate whether this type of additional information  $Z$ , together with  $X$ , helps EIIL and HRM. The corresponding results are marked with (+ $Z$ ) in Table 3. We observe that the results of EIIL and HRM with additional information  $Z$  are worse than those without using  $Z$ . This implies that such information may not be directly utilized by or even can harm the two methods.

**CelebA.** This task is to predict *Smiling* based on the image from CelebA [31], and by construction the target is spuriously correlated with *Gender*. We have access to the meta annotations of CelebA images and they serve as  $Z$  in our framework. In particular, we pick seven features as potential auxiliary variables: {*Young, Blond Hair, Eyeglasses, High Cheekbones, Big Nose, Bags Under Eyes, Chubby*}. Notice that the annotation of *Gender* is not provided to ZIN, EIIL, HRM and ERM, but is made available to IRM and group DRO to create an oracle environment partition.

Empirical results are given in Table 4, where ZIN( $\#$ ) represents ZIN using  $\#$  of the seven features as  $Z$ . ERM achieves the highest training accuracy, while only has 47.58% worst test accuracy. The worst test performances of ERM, EIIL and LfF indicate that they may not learn the causal features. Finally, we observe that ZIN performs better when more auxiliary information is provided; specifically, ZIN(7) achieves the best test performance based on the mixed dataset. This observation validates our discussion on choosing the auxiliary information in Section 5.4.

**Table 3: House price prediction (MSE).**

Method	Env Index	Train	Test Mean	Test Worst
IRM	Yes	0.1327	0.4456	0.6821
group DRO	Yes	0.1213	0.6887	1.0050
ERM	No	<b>0.1141</b>	0.4764	0.6703
EIIL	No	0.6841	0.9625	1.3909
EIIL(+Z)	No	0.6912	0.9701	1.4201
HRM	No	0.3466	0.4621	0.5721
HRM(+Z)	No	0.3190	0.4221	0.5873
ZIN	No	0.2275	<b>0.3339</b>	<b>0.4815</b>

**Table 4: Accuracy (%) on CelebA task.**

Method	Env Index	Train	Test Mean	Test Worst
IRM	Yes	81.30±1.53	78.44±0.48	75.03±1.29
group DRO	Yes	89.32±0.67	74.28±0.22	58.11±0.61
ERM	No	<b>90.97±0.66</b>	70.76±0.26	47.58±0.46
LfF	No	59.89±0.72	52.97±0.56	44.38±2.01
EIIL	No	90.01±0.73	71.45±0.29	50.48±1.98
ZIN(1)	No	90.62±0.78	70.79±0.61	47.62±0.98
ZIN(4)	No	83.57±1.40	75.20±0.71	63.47±1.41
ZIN(7)	No	83.06±1.28	<b>76.29±0.60</b>	<b>67.27±1.15</b>

In Appendix D, Fig. 3 reports the results of different choices of hyper-parameter  $K$  and Fig. 4 visualizes the inferred environments in this experiment. Again, we find that  $K$  has a relatively small impact on the proposed algorithm.

**Landcover.** Our final task is land cover prediction that classifies the land cover type (e.g., grasslands) from satellite data [14, 40]. We take the same setup from [53]: the time series data input dimension is  $46 \times 8$ ; the target  $Y$  is one of six land cover classes; six climate related variables are the auxiliary variables. We also consider non-African locations as training data and Africa as test data. In addition, we take location (latitude and longitude) as possible additional information. For ZIN, we average climate variables over the time dimension and use those means to predict environments. In Table 5, we see that ZIN with location achieves the best OOD performance, while ZIN with climate variables is slightly worse. This is because climate variables indeed contain some information for predicting the target, violating our conditions on  $Z$ . EIIL also has a good performance in this task, and we conjecture that its first stage has learned the spurious features as desired. However, together with previous empirical results, it can be risky as the first stage may not be guaranteed.

**Table 5: Test Accuracy (%) on Landcover task.**

Method	IID Test	OOD Test
ERM [52]	75.92	58.31
ERM (+climate) [52]	76.58	54.78
LfF	66.24	61.69
EIIL	72.61	64.79
ZIN (climate)	72.56	62.50
ZIN (location)	72.18	<b>66.06</b>

### 6.3 Discussion on Auxiliary Information for Related Methods

The proposed framework ZIN relies on additional auxiliary variables. It is interesting to ask whether this type of information is also useful to related methods such as EIIL and HRM. In many scenarios, the additional variable  $Z$  may have no or little information about the target, e.g., 1) in time series tasks, the time index is rarely used to predict the label; and 2) in the CelebA classification task, features like *Young* and *Eyeglasses* are likely to provide little information about predicting *Smiling*. Then  $Z$  may not be useful to ERM and EIIL. In Section 6.2, we also provide  $Z$  to EIIL and HRM, and the experimental result shows that it does not help or even harms the test performance.

## 7 Concluding Remarks

This paper investigates when and how to learn invariance from heterogeneous data without explicit environment indexes. We first show that learning invariant models in this case is generally impossible. Then we propose ZIN to jointly learn environment partition and invariant representation using some additional auxiliary variables. We provide theoretic guarantees in both feature selection and linear feature learning scenarios. Experimental results verify our analysis and demonstrate an improved performance of ZIN over existing methods. Our results also raise the need of future works to make the role of inductive biases more explicit, when learning invariance from heterogeneous data without environment indexes. A limitation of the present work is the lack of theoretic guarantee with general nonlinear feature extractors, which is also an open problem for IRM even with environment partition.

## References

- [1] Faruk Ahmed, Yoshua Bengio, Harm van Seijen, and Aaron C. Courville. Systematic generalization with group invariant predictions. In *ICLR*, 2021.
- [2] Kartik Ahuja, Karthikeyan Shanmugam, Kush Varshney, and Amit Dhurandhar. Invariant risk minimization games. In *International Conference on Machine Learning*, 2020.
- [3] Kartik Ahuja, Jun Wang, Amit Dhurandhar, Karthikeyan Shanmugam, and Kush R. Varshney. Empirical or invariant risk minimization? A sample complexity perspective. *International Conference on Learning Representations*, 2021.
- [4] Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization. *arXiv preprint arXiv:1907.02893*, 2019.
- [5] Aharon Ben-Tal, Dick Den Hertog, Anja De Waegenare, Bertrand Melenberg, and Gijs Rennen. Robust solutions of optimization problems affected by uncertain probabilities. *Management Science*, 59(2):341–357, 2013.
- [6] Joos-Hendrik Böse, Valentin Flunkert, Jan Gasthaus, Tim Januschowski, Dustin Lange, David Salinas, Sebastian Schelter, Matthias Seeger, and Yuyang Wang. Probabilistic demand forecasting at scale. *Proceedings of the VLDB Endowment*, 10(12):1694–1705, 2017.
- [7] Shiyu Chang, Yang Zhang, Mo Yu, and Tommi Jaakkola. Invariant rationalization. In *International Conference on Machine Learning*, 2020.
- [8] Yongqiang Chen, Yonggang Zhang, Han Yang, Kaili Ma, Binghui Xie, Tongliang Liu, Bo Han, and James Cheng. Invariance principle meets out-of-distribution generalization on graphs. *arXiv preprint arXiv:2202.05441*, 2022.
- [9] Yongqiang Chen, Kaiwen Zhou, Yatao Bian, Binghui Xie, Kaili Ma, Yonggang Zhang, Han Yang, Bo Han, and James Cheng. Pareto invariant risk minimization. *arXiv preprint arXiv:2206.07766*, 2022.
- [10] Elliot Creager, Jörn-Henrik Jacobsen, and Richard Zemel. Environment inference for invariant learning. In *International Conference on Machine Learning*, 2021.
- [11] David Danks, Clark Glymour, and Robert Tillman. Integrating locally learned causal structures with overlapping variables. In *Advances in Neural Information Processing Systems*, 2009.
- [12] John C Duchi and Hongseok Namkoong. Learning models with uniform performance via distributionally robust optimization. *The Annals of Statistics*, 49(3):1378–1406, 2021.
- [13] Rui Gao, Xi Chen, and Anton J Kleywegt. Wasserstein distributionally robust optimization and variation regularization. *arXiv preprint arXiv:1712.06050*, 2020.
- [14] Pall Oskar Gislason, Jon Atli Benediktsson, and Johannes R. Sveinsson. Random forests for land cover classification. *Pattern Recognition Letters*, 27(4):294–300, 2006.
- [15] Ishaan Gulrajani and David Lopez-Paz. In search of lost domain generalization. *arXiv preprint arXiv:2007.01434*, 2020.
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [17] Dan Hendrycks and Thomas G. Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. In *International Conference on Learning Representations*, 2019.

- [18] Biwei Huang, Kun Zhang, Mingming Gong, and Clark Glymour. Causal discovery from multiple data sets with non-identical variable sets. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020.
- [19] Biwei Huang, Kun Zhang, Jiji Zhang, Joseph Ramsey, Clark Glymour, and Bernhard Schölpf. Causal discovery from heterogeneous/nonstationary data. *Journal of Machine Learning Research*, 21, May 2020.
- [20] Aapo Hyvarinen, Hiroaki Sasaki, and Richard Turner. Nonlinear ICA using auxiliary variables and generalized contrastive learning. In *International Conference on Artificial Intelligence and Statistics*, 2019.
- [21] Wengong Jin, Regina Barzilay, and Tommi Jaakkola. Domain extrapolation via regret minimization. *arXiv preprint arXiv:2006.03908*, 2020.
- [22] Ilyes Khemakhem, Diederik Kingma, Ricardo Monti, and Aapo Hyvarinen. Variational autoencoders and nonlinear ICA: A unifying framework. In *International Conference on Artificial Intelligence and Statistics*, 2020.
- [23] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2014.
- [24] David Krueger, Ethan Caballero, Joern-Henrik Jacobsen, Amy Zhang, Jonathan Binas, Dinghui Zhang, Remi Le Priol, and Aaron Courville. Out-of-distribution generalization via risk extrapolation (rex). In *International Conference on Machine Learning*, 2021.
- [25] Jaeho Lee and Maxim Raginsky. Minimax statistical learning with wasserstein distances. In *Advances in Neural Information Processing Systems*, 2018.
- [26] Yong Lin, Lian Qing, and Tong Zhang. An empirical study of invariant risk minimization on deep models. 2021.
- [27] Yong Lin, Hanze Dong, Hao Wang, and Tong Zhang. Bayesian invariant risk minimization. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022.
- [28] Evan Zheran Liu, Behzad Haghgoo, Annie S. Chen, Aditi Raghunathan, Pang Wei Koh, Shiori Sagawa, Percy Liang, and Chelsea Finn. Just train twice: Improving group robustness without training group information. In *ICML*, 2021.
- [29] Jiashuo Liu, Zheyuan Hu, Peng Cui, Bo Li, and Zheyuan Shen. Heterogeneous risk minimization. In *International Conference on Machine Learning*, 2021.
- [30] Jiashuo Liu, Zheyuan Hu, Peng Cui, Bo Li, and Zheyuan Shen. Kernelized heterogeneous risk minimization. In *NeurIPS*, 2021.
- [31] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, 2015.
- [32] Francesco Locatello, Stefan Bauer, Mario Lucic, Gunnar Raetsch, Sylvain Gelly, Bernhard Schölkopf, and Olivier Bachem. Challenging common assumptions in the unsupervised learning of disentangled representations. In *International Conference on Machine Learning*, 2019.
- [33] Manfred Mudelsee. Trend analysis of climate time series: A review of methods. *Earth-Science Reviews*, 190:310–322, 2019.
- [34] Jun Hyun Nam, Hyuntak Cha, Sungsoo Ahn, Jaeho Lee, and Jinwoo Shin. Learning from failure: Training debiased classifier from biased classifier. In *NeurIPS*, 2020.
- [35] Judea Pearl. *Causality: Models, Reasoning and Inference*. Cambridge University Press, USA, 2nd edition, 2009.

- [36] Jonas Peters, Peter Bühlmann, and Nicolai Meinshausen. Causal inference by using invariant prediction: Identification and confidence intervals. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, pages 947–1012, 2016.
- [37] Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of Causal Inference - Foundations and Learning Algorithms*. MIT Press, Cambridge, MA, USA, 2017.
- [38] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishal Shankar. Do ImageNet classifiers generalize to ImageNet? In *International Conference on Machine Learning*, 2019.
- [39] Elan Rosenfeld, Pradeep Ravikumar, and Andrej Risteski. The risks of invariant risk minimization. *International Conference on Learning Representations*, 2021.
- [40] Marc Russwurm, Sherrie Wang, Marco Korner, and David Lobell. Meta-learning for few-shot land cover classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2020.
- [41] Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. In *International Conference on Learning Representations*, 2020.
- [42] Victor Sanh, Thomas Wolf, Yonatan Belinkov, and Alexander M. Rush. Learning from others’ mistakes: Avoiding dataset biases without modeling them. In *ICML*, 2021.
- [43] Alexander Shapiro. Distributionally robust stochastic programming. *SIAM Journal on Optimization*, 27(4):2258–2275, 2017.
- [44] Nimit Sharad Sohoni, Jared A. Dunnmon, Geoffrey Angus, Albert Gu, and Christopher Ré. No subclass left behind: Fine-grained robustness in coarse-grained classification problems. In *NeurIPS*, 2020.
- [45] Peter Spirtes, Clark N Glymour, and Richard Scheines. *Causation, Prediction, and Search*. MIT Press, second edition, 2000.
- [46] Megha Srivastava, Tatsunori Hashimoto, and Percy Liang. Robustness to spurious correlations via human annotations. In *International Conference on Machine Learning*, 2020.
- [47] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *International Conference on Learning Representations*, 2014.
- [48] Robert Tillman and Peter Spirtes. Learning equivalence classes of acyclic models with latent and selection variables from multiple datasets with overlapping variables. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, 2011.
- [49] Jorge Vergara and Pablo Estevez. A review of feature selection methods based on mutual information. *Neural Computing and Applications*, 24, 01 2014.
- [50] Ruoyu Wang, Mingyang Yi, Zhitang Chen, and Shengyu Zhu. Out-of-distribution generalization with causal invariant transformations. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [51] Sherrie Wang, William Chen, Sang Michael Xie, George Azzari, and David B. Lobell. Weakly supervised deep learning for segmentation of remote sensing imagery. *Remote Sensing*, 12(2), 2020.
- [52] Chuanlong Xie, Fei Chen, Yue Liu, and Zhenguo Li. Risk variance penalization: From distributional robustness to causality. *arXiv preprint arXiv:2006.07544*, 2020.

- [53] Sang Michael Xie, Ananya Kumar, Robbie Jones, Fereshte Khani, Tengyu Ma, and Percy Liang. In-N-Out: Pre-training and self-training using auxiliary information for out-of-distribution robustness. In *International Conference on Learning Representations*, 2021.
- [54] Yilun Xu and Tommi Jaakkola. Learning representations that support robust transfer of predictors. *arXiv preprint arXiv:2110.09940*, 2021.
- [55] Mingyang Yi, Lu Hou, Jiacheng Sun, Lifeng Shang, Xin Jiang, Qun Liu, and Zhi-Ming Ma. Improved OOD generalization via adversarial training and pre-training. In *International Conference on Machine Learning*, 2021.
- [56] Xiao Zhou, Yong Lin, Renjie Pi, Weizhong Zhang, Renzhe Xu, Peng Cui, and Tong Zhang. Model agnostic sample reweighting for out-of-distribution learning. In *International Conference on Machine Learning*, 2022.
- [57] Xiao Zhou, Yong Lin, Weizhong Zhang, and Tong Zhang. Sparse invariant risk minimization. In *International Conference on Machine Learning*, 2022.

## Checklist

1. For all authors...
  - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
  - (b) Did you describe the limitations of your work? [Yes] In the conclusion section.
  - (c) Did you discuss any potential negative societal impacts of your work? [Yes] In Appendix F.
  - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...
  - (a) Did you state the full set of assumptions of all theoretical results? [Yes]
  - (b) Did you include complete proofs of all theoretical results? [Yes]
3. If you ran experiments...
  - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes]
  - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes]
  - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes]
  - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes]
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
  - (a) If your work uses existing assets, did you cite the creators? [Yes]
  - (b) Did you mention the license of the assets? [N/A]
  - (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
  - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
  - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
5. If you used crowdsourcing or conducted research with human subjects...
  - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
  - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
  - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

## A More Theoretical Results

### A.1 Extension to Other Loss Functions

Sections 5.2 and 5.3 focus on classification task with cross-entropy loss, to establish sufficient and necessary conditions for invariance identification. Similar results can be shown for other loss functions or tasks in a straightforward manner.

To see this, notice that  $H(\cdot|\cdot)$  in Assumptions 1-3 and Conditions 1-2 is used to represent the optimal expected risk that coincides with the conditional entropy, with cross-entropy loss and when  $\rho(z_i)$  gives exactly one environment. For other loss functions  $l(\cdot, \cdot)$  like squared error, we can use  $L(\cdot|\cdot)$  to represent the optimal expected risks, and  $L(\cdot)$  denotes the optimal expected risk with no predictor variables, e.g., variance for the squared error loss. Replacing  $H(\cdot|\cdot)$  with  $L(\cdot|\cdot)$ , we can follow the same proof procedure to obtain similar sufficient and necessary conditions. This is summarized in Theorem 3 for completeness.

**Theorem 3.** *Consider Assumptions 1-3 and Conditions 1-2 where  $H(\cdot|\cdot)$  is replaced with  $L(\cdot|\cdot)$  accordingly. If  $\epsilon < \frac{C\gamma\delta}{4\gamma+2C\delta L(Y)}$  and  $\lambda \in [\frac{L(Y)+1/2\delta C}{\delta C-4\epsilon} - \frac{1}{2}, \frac{\gamma}{4\epsilon} - \frac{1}{2}]$ , we have  $\hat{\mathcal{L}}(\Phi_v) < \hat{\mathcal{L}}(\Phi)$  for all  $\Phi \neq \Phi_v$ , where we assume  $L(Y) < \infty$ . If Condition 1 or Condition 2 is violated, then there exists a feature mask  $\Phi' \neq \Phi_v$  so that  $\hat{\mathcal{L}}(\Phi_v) > \hat{\mathcal{L}}(\Phi')$ .*

### A.2 Beyond Feature Selection: Linear Feature Transformations

In Section 5, both the invariant and spurious features can be directly observed and we focus on the ability of ZIN to identify the invariant features. In this section, we extend feature selection to feature learning, and show that ZIN is able to learn the invariant features given a scrambled observation in a linear form following [4, 39]. Specifically, we consider the same data generation process as [4]:

$$Y^e = \mathbf{X}_v^e \cdot \beta + \epsilon_v, \quad \mathbf{X}_v^e \perp \epsilon_v, \quad \mathbb{E}[\epsilon_v] = 0; \quad \mathbf{X}^e = \mathbf{W} \cdot [\mathbf{X}_v^e; \mathbf{X}_s^e], \quad (7)$$

where  $\beta \in \mathbb{R}^{d_v}$  and  $\mathbf{W} \in \mathbb{R}^{d \times (d_v + d_s)}$ . We assume that there exists  $\tilde{\mathbf{W}} \in \mathbb{R}^{d_v \times d}$  so that  $\tilde{\mathbf{W}}(\mathbf{W}[\mathbf{x}_v; \mathbf{x}_s]) = \mathbf{x}_v$  for all  $\mathbf{x}_v$  and  $\mathbf{x}_s$ . Both the feature extractor and predictor take a linear form, i.e.,  $\Phi$  takes values in  $\mathbb{R}^{d \times d}$  and  $\omega \in \mathbb{R}^d$ . The prediction for  $\mathbf{X}$  is  $\omega \circ \Phi(\mathbf{X}) = (\Phi \mathbf{X})^T \omega$ .

A major difficulty in this setting lies in how to characterize the effect of an invariant/spurious feature in a *quantitative* way: the feature extractor  $\Phi$  may extract an arbitrarily small portion of spurious information. Following [4, 39], we consider a constrained form of Problem 6 for theoretic analysis:

$$\min_{\omega, \Phi} \mathcal{R}(\omega, \Phi), \quad \text{subject to} \quad \max_{\rho, \{\omega_k\}} \sum_{k=1}^K [\mathcal{R}_{\rho^{(k)}}(\omega, \Phi) - \mathcal{R}_{\rho^{(k)}}(\omega_k, \Phi)] = 0. \quad (8)$$

As in Conditions 1 and 2, the auxiliary information should also be sufficiently informative so that the inferred environments can be diverse enough but also maintain the underlying invariance. This is in accordance to existing conditions for identifiability in the linear case [4, 39]. In this paper, we utilize such a condition, *linear general position* condition, from [4]. For our analysis, we take squared error as our loss function and consider that  $\rho(\cdot)$  partitions the environments in a hard manner, i.e., each data sample would be assigned to exactly one environment. As in Appendix A.1, we use  $L(\cdot|\cdot)$  to represent the optimal expected risks. We also assume that the environments are non-degenerate, i.e., each inferred environment contains some data samples; otherwise, we can simply remove such an environment. Our identifiability result for the linear case then follows.

**Proposition 3.** *Assume Condition 1 where  $H(\cdot|\cdot)$  is replaced with  $L(\cdot|\cdot)$  accordingly. Suppose that there exists  $\rho(\cdot)$  such that the generated environments, denoted as  $\{\mathbf{X}^k\}_{k=1}^K$ , lie in linear general position of degree  $r$ , i.e.,  $K > d - r + d/r$  for some  $r \in \mathbb{N}$  and for all non-zero  $\mathbf{x} \in \mathbb{R}^d$ :  $\dim(\text{span}(\{\mathbb{E}_{\mathbf{X}^k}[\mathbf{X}^k \mathbf{X}^{kT}] \mathbf{x} - \mathbb{E}_{\mathbf{X}^k, \epsilon_v}[\mathbf{X}^k \epsilon_v]\}_{k=1}^K)) > d - r$ . If  $\Phi \in \mathbb{R}^{d \times d}$  has rank  $r > 0$ , then Problem 8 results in the desired invariant predictor.*

*Proof. Step 1:* No spurious feature will be learned. Given a partition  $\{\mathbf{X}^k\}_{k=1}^K$  that lies in linear general position of degree  $r$ , Arjovsky et al. [4, Theorem 9] shows that  $\Phi$  and  $\omega$  satisfies the normal equations  $\Phi \mathbb{E}_{\mathbf{X}^k}[\mathbf{X}^k \mathbf{X}^{kT}] \Phi^T \omega = \Phi \mathbb{E}_{\mathbf{X}^k, Y^k}[\mathbf{X}^k Y^k]$  for all  $k$  if and only if  $\Phi$  elicits the desired



invariant predictor  $\Phi^T \omega = \tilde{\mathbf{W}}^T \beta$ . Thus, we only need to show our solution meets the same normal equations. Let  $\Phi'$  and  $\omega'$  denote a solution to Problem 8. According to the constraint and the general linear position condition, we know there exists a partition  $\{\mathbf{X}_k\}_{k=1}^K$  lies in general linear position of degree  $r$ , and we have  $\mathcal{R}^k(\omega', \Phi') = \mathcal{R}^k(\omega_k, \Phi')$ . Notice that  $\omega_k$  minimizes the mean squared error on only the  $k$ -th environment, hence it must satisfy the normal equation  $\Phi' \mathbb{E}_{\mathbf{X}^k} [\mathbf{X}^k \mathbf{X}^{kT}] \Phi'^T \omega_k = \Phi' \mathbb{E}_{\mathbf{X}^k, Y^k} [\mathbf{X}^k Y^k]$ . As  $\omega'$  achieves the same minimum mean squared error on the  $k$ -th environment,  $\omega'$  must satisfy the normal equation, too. Thus, no spurious information will be included and  $\mathbf{X} \Phi = [S \mathbf{X}_v; \mathbf{0}_{d_s}]$  where  $S \in \mathbb{R}^{d_v \times d_v}$  is an invertible matrix and  $\mathbf{0}_{d_s}$  denotes a  $d_s$ -dim vector of all zeros.

**Step 2:** No invariant information will be discarded. Under Condition 1 where  $H(\cdot|\cdot)$  is replaced with  $L(\cdot|\cdot)$ , we have  $L(Y|X_v) - L(Y|X_v, \rho(Z)) = 0$ . Then  $\mathbf{X} \Phi$  will satisfy the constraint of Problem 8. Notice that  $\mathbf{X} \Phi$  achieves the smallest loss when only using invariant feature information.

Combining these two steps completes the proof.  $\square$

### A.3 More Discussions on Assumption 2

This assumption aims to ensure that the invariance penalty cannot be arbitrarily small if a spurious feature, together with other features, is selected by a feature mask. For example, in the extreme case where  $\mathbf{X}_2$  is the only invariant feature (i.e.,  $\mathbf{X}_v$  consists of only  $\mathbf{X}_2$ ) and can perfectly predict  $Y$ , we would have  $H(Y|\mathbf{X}_2) = 0$ . Then for a spurious feature  $\mathbf{X}_1$ ,  $H(Y|\mathbf{X}_1, \mathbf{X}_2) = 0$  and  $H(Y|\mathbf{X}_1, \mathbf{X}_2, \rho(Z)) = 0$  for any  $\rho(\cdot)$ , and we cannot identify  $\mathbf{X}_1$  as the spurious feature. Nevertheless, since there are generally exogenous noise variables in the SCM and  $H(Y|\mathbf{X}_v)$  is positive, we believe that this assumption holds in most cases.

## B Proofs

### B.1 Proof of Theorem 1

*Proof.* First, we assume  $Y$  and  $\mathbf{X}_v$  are univariate variables, i.e.,  $Y, \mathbf{X}_v \in \mathbb{R}$ . Let  $\eta_1 \sim \text{Uniform}(0, 1)$  independent of  $\mathbf{X}_s$  and  $Y$ , and set  $\mathbf{X}'_v = \mathbf{X}_s$ . Define the conditional cumulative distribution function and its inverse as:

$$\begin{aligned} F_{Y|\mathbf{X}_s=\mathbf{x}_s}(y) &= P(Y \leq y | \mathbf{X}_s = \mathbf{x}_s), \\ Y' = f'_1(\mathbf{X}'_v, \eta_1) &= F_{Y|\mathbf{X}_s}^{-1}(\eta_1) = \inf\{y \in \mathbb{R} : F_{Y|\mathbf{X}_s}(y) \geq \eta_1\} \text{ with } \mathbf{X}_s = \mathbf{X}'_v. \end{aligned}$$

By definition, we would have

$$\begin{aligned} P(Y' \leq y | \mathbf{X}'_v) &= P(f'_1(\mathbf{X}'_v, \eta_1) \leq y) \\ &= P(F_{Y|\mathbf{X}_s}^{-1}(\eta_1) \leq y) \\ &= P(F_{Y|\mathbf{X}_s} \circ F_{Y|\mathbf{X}_s}^{-1}(\eta_1) \leq F_{Y|\mathbf{X}_s}(y)) \\ &= P(\eta_1 \leq F_{Y|\mathbf{X}_s}(y)) \\ &= P(Y \leq y | \mathbf{X}_s). \end{aligned} \tag{9}$$

Similarly, we can construct  $\mathbf{X}'_s = f'_2(\mathbf{X}'_v, Y', \eta_2) = F_{\mathbf{X}_v|Y, \mathbf{X}_s}^{-1}(\eta_2)$  so that  $P(\mathbf{X}'_s | Y', \mathbf{X}'_v) = P(\mathbf{X}_v | Y, \mathbf{X}_s)$  with  $\eta_2 \sim \text{Uniform}(0, 1)$ . Thus, we have

$$\begin{aligned} P(\mathbf{X}'_v, \mathbf{X}'_s, Y) &= P(\mathbf{X}'_s | Y', \mathbf{X}'_v) P(Y' | \mathbf{X}'_v) P(\mathbf{X}'_v) \\ &= P(\mathbf{X}_v | Y, \mathbf{X}_s) P(Y | \mathbf{X}_s) P(\mathbf{X}_s) \\ &= P(\mathbf{X}_v, \mathbf{X}_s, Y). \end{aligned} \tag{10}$$

Next, we can easily find a function  $q'(\cdot)$  so that  $\mathbf{X}' = q'(\mathbf{X}'_v, \mathbf{X}'_s) = q(\mathbf{X}_v, \mathbf{X}_s)$ , where we have chosen  $\mathbf{X}'_v = \mathbf{X}_s$  and  $\mathbf{X}'_s = \mathbf{X}_v$ . Together with Equation (10), we conclude that  $P(\mathbf{X}', Y') = P(\mathbf{X}, Y)$ .

Second, we consider that  $\mathbf{X}_v$  has a multi-dimension. We may leave  $Y$  to be a univariate variable as the label in many ML problems is scalar-valued. For this case, we can pick an entry of  $\mathbf{X}_v$ , say,  $\mathbf{X}_v^{(j)}$  for some  $j$ . Then we set  $\mathbf{X}'_v = (\mathbf{X}_v^{(-j)}, \mathbf{X}_s)$  where  $\mathbf{X}_v^{(-j)}$  denotes the rest entries of  $\mathbf{X}_v$  except the  $j$ -th. Then we can similarly use the inverse conditional cumulative distribution function and an independent uniformly distributed noise variable  $\eta_1$  to construct  $Y' = f'(\mathbf{X}'_v, \eta_1)$  so that  $P(Y' \leq y | \mathbf{X}'_v) = P(Y \leq y | \mathbf{X}_v^{(-j)}, \mathbf{X}_s)$ . Similarly, we set  $\mathbf{X}'_s = \mathbf{X}_v^{(j)}$  and we can construct  $\mathbf{X}'_s = f'_2(\mathbf{X}'_v, Y', \eta_2)$  so that  $P(\mathbf{X}'_s | Y', \mathbf{X}'_v) = P(\mathbf{X}_v^{(j)} | Y, \mathbf{X}_s)$ . We can then get the same conclusion following the previous proof.  $\square$

## B.2 Proof of Theorem 2

*Proof.* Our proof proceeds by two steps. First, we show that any feature mask that selects at least one spurious feature would induce a penalty. With sufficiently large  $\lambda$ , the penalty will dominate the expected risk and then exceed  $\hat{\mathcal{L}}(\Phi_v)$ . Second, we show that any proper subset of the invariant features induces a loss larger than  $\hat{\mathcal{L}}(\Phi_v)$ .

**Step 1** Suppose that the feature mask contains at least one spurious features. Denote the selected features as  $\mathbf{X}_{+s}$  and the corresponding feature mask as  $\Phi_{+s}$ . We aim to show that

$$\hat{\mathcal{L}}(\Phi_{+s}) > \hat{\mathcal{L}}(\Phi_v).$$

By Assumption 1 with a given  $\epsilon > 0$ , we have

$$\begin{aligned} \hat{\mathcal{L}}(\Phi_v) &\leq (1 + 2\lambda)\epsilon + H(Y|\mathbf{X}_v) + \lambda(H(Y|\mathbf{X}_v) - H(Y|\mathbf{X}_v, \rho(Z))) \\ &= (1 + 2\lambda)\epsilon + H(Y|\mathbf{X}_v) \\ &\leq (1 + 2\lambda)\epsilon + H(Y). \end{aligned} \quad (11)$$

One the other hand, we have

$$\begin{aligned} \mathcal{L}(\Phi_{+s}) &\geq -(1 + 2\lambda)\epsilon + H(Y|\mathbf{X}_{+s}) + \lambda(H(Y|\mathbf{X}_{+s}) - H(Y|\mathbf{X}_{+s}, \rho(Z))) \\ &\geq -(1 + 2\lambda)\epsilon + \lambda(H(Y|\mathbf{X}_{+s}) - H(Y|\mathbf{X}_{+s}, \rho(\mathbf{Z}))) \\ &\geq -(1 + 2\lambda)\epsilon + \lambda\delta C, \end{aligned} \quad (12)$$

where the last inequality is due to Assumption 2 and Condition 2. Thus, if we choose  $\epsilon < \delta C/4$  and  $\lambda > \frac{H(Y)+2\epsilon}{\delta C-4\epsilon}$ , we can get

$$\mathcal{L}(\Phi_v) < \mathcal{L}(\Phi_{+s}).$$

**Step 2** Denote a proper subset of invariant features as  $\mathbf{X}_{-v} \subsetneq \mathbf{X}_v$ , and similarly the feature mask as  $\Phi_{-v}$ .

In Step 1, we have shown that

$$\hat{\mathcal{L}}(\Phi_v) \leq (1 + 2\lambda)\epsilon + H(Y|\mathbf{X}_v).$$

Similar to Equation (12), we have

$$\hat{\mathcal{L}}(\Phi_{-v}) \geq -(1 + 2\lambda)\epsilon + H(Y|\mathbf{X}_{-v}).$$

Then according to Assumption 3, we have

$$\begin{aligned} \hat{\mathcal{L}}(\Phi_{-v}) - \hat{\mathcal{L}}(\Phi_v) &\geq -2(1 + 2\lambda)\epsilon + H(y|\mathbf{X}_{-v}) - H(y|\mathbf{X}_v) \\ &\geq -2(1 + 2\lambda)\epsilon + \gamma. \end{aligned}$$

Thus, if  $\epsilon < \frac{\gamma}{2(1+2\lambda)}$ , we have

$$\hat{\mathcal{L}}(\Phi_{-v}) > \hat{\mathcal{L}}(\Phi_v).$$

In conclusion, with  $\lambda \in [\frac{H(Y)+1/2\delta C}{\delta C-4\epsilon} - \frac{1}{2}, \frac{\gamma}{4\epsilon} - \frac{1}{2}]$ , we can get

$$\hat{\mathcal{L}}(\Phi_v) < \hat{\mathcal{L}}(\Phi), \quad \forall \Phi \neq \Phi_v.$$

Notably, there exists a feasible  $\lambda$  if  $\epsilon < \frac{C\gamma\delta}{4\gamma+2C\delta H(Y)}$ . The proof is complete by noticing that  $\epsilon$  can be chosen arbitrarily according to Assumption 1.  $\square$

### B.3 Proof of Proposition 1

*Proof.* Consider the following feature set

$$\mathbf{X}_{\bar{v}} := \max_{|\mathbf{X}'|} \{ \mathbf{X}' \subset \mathbf{X} : H(Y|\mathbf{X}', \rho(\mathbf{Z})) = H(Y|\mathbf{X}') \quad \forall \rho(\cdot) \},$$

and the corresponding feature mask is denoted as  $\Phi_{\bar{v}}$ . It corresponds to the largest subset of  $\mathbf{X}$  that satisfies the invariance constraint.

Note that  $\Phi_{\bar{v}} \neq \Phi_v$ . By Assumption 1, for a given  $\epsilon$  we can get

$$\begin{aligned} \hat{\mathcal{L}}(\Phi_{\bar{v}}) &\leq (1 + 2\lambda)\epsilon + H(Y|\mathbf{X}_{\bar{v}}) + \lambda(H(Y|\mathbf{X}_{\bar{v}}) - H(Y|\mathbf{X}_v, \rho(\mathbf{Z}))) \\ &= (1 + 2\lambda)\epsilon + H(Y|\mathbf{X}_{\bar{v}}) \\ &\leq (1 + 2\lambda)\epsilon + H(Y). \end{aligned}$$

One the other hand, we have

$$\begin{aligned} \mathcal{L}(\Phi_v) &\geq -(1 + 2\lambda)\epsilon + H(Y|\mathbf{X}_v) + \lambda(H(Y|\mathbf{X}_v) - H(Y|\mathbf{X}_v, \rho(\mathbf{Z}))) \\ &\geq -(1 + 2\lambda)\epsilon + \lambda(H(Y|\mathbf{X}_v) - H(Y|\mathbf{X}_v, \rho(\mathbf{Z}))) \\ &\geq -(1 + 2\lambda)\epsilon + \lambda\delta C', \end{aligned}$$

which can be shown similarly to Equation (12). Thus, if we choose  $\epsilon < \delta C'/4$  and  $\lambda > \frac{H(Y)+2\epsilon}{\delta C'-4\epsilon}$ , we would get

$$\mathcal{L}(\Phi_{\bar{v}}) < \mathcal{L}(\Phi_v).$$

□

### B.4 Proof of Corollary 1

*Proof.* (a) Since  $h$  is injective,  $H(Y|\mathbf{X}_v, h(Y)) = H(Y|\mathbf{X}_v, Y) = 0$  for any  $\mathbf{X}_v$ . By Assumption 3, we have  $H(Y|\mathbf{X}_v) \geq H(Y|\mathbf{X}) + \gamma \geq \gamma$ . Then Proposition 1 indicates that we cannot identify all the invariant features. (b) The proof proceeds the same as above by noting  $H(Y|\mathbf{X}_v, h(\mathbf{X}, Y)) = 0$ . (c) This case can be shown similarly to (a), because  $H(Y|\mathbf{X}_v, h(\text{Index}(\mathbf{X}, Y))) = H(Y|\mathbf{X}_v, \mathbf{X}, Y) = 0$ . □

### B.5 Proof of Proposition 2

*Proof.* Denote a feature set  $\mathbf{X}_{v(+k)}$ , which contains the invariant feature set  $\mathbf{X}_v$  as well a spurious feature  $\mathbf{X}_s^k$  in Proposition 2.

Similar to Eq. (11), we can show that

$$\hat{\mathcal{L}}(\Phi_{v(+k)}) \leq (1 + 2\lambda)\epsilon + H(Y|\mathbf{X}_{v(+k)}).$$

And similar to Eq. (12), we also have

$$\hat{\mathcal{L}}(\Phi_v) \geq -(1 + 2\lambda)\epsilon + H(Y|\mathbf{X}_v).$$

Then it follows that

$$\begin{aligned} \hat{\mathcal{L}}(\Phi_v) - \hat{\mathcal{L}}(\Phi_{v(+k)}) &\geq -2(1 + 2\lambda)\epsilon + H(Y|\mathbf{X}_v) - H(y|\mathbf{X}_{v(+k)}) \\ &\geq -2(1 + 2\lambda)\epsilon + \gamma. \end{aligned}$$

If  $\epsilon < \frac{\gamma}{2(1+2\lambda)}$ , we have

$$\hat{\mathcal{L}}(\Phi_v) > \hat{\mathcal{L}}(\Phi_{v(+k)}).$$

Thus, we cannot identify all the invariant features from Problem 6. □

## B.6 Proof of Meeting Condition 1

We show that  $H(Y|\mathbf{X}_v, \rho(\mathbf{Z})) = H(Y|\mathbf{X}_v)$  for all  $\rho(\cdot)$  if  $H(Y|\mathbf{X}_v, \mathbf{Z}) = H(Y|\mathbf{X}_v)$  holds.

*Proof.* On one hand, because  $\rho(\mathbf{Z})$  contains less information than  $\mathbf{Z}$ , we have

$$H(Y|\mathbf{X}_v, \rho(\mathbf{Z})) \geq H(Y|\mathbf{X}_v, \mathbf{Z}) = H(Y|\mathbf{X}_v).$$

On the other hand,  $\mathbf{X}_v$  and  $\rho(\mathbf{Z})$  contain more information than  $\mathbf{X}_v$ , so we can get

$$H(Y|\mathbf{X}_v, \rho(\mathbf{Z})) \leq H(Y|\mathbf{X}_v).$$

Thus, we conclude  $H(Y|\mathbf{X}_v, \rho(\mathbf{Z})) = H(Y|\mathbf{X}_v)$ .  $\square$

## C More Experimental Details

**Implementing Approximated ZIN.** Since Eq. (5) is a challenging minimax formulation, we replace the penalty term given a environment partition with its first order approximation. Specifically, we consider the following surrogate minimax formulation:

$$\min_{\omega, \Phi} \max_{\rho} = \mathcal{R}(\omega, \Phi) + \lambda \sum_{k=1}^K \|\nabla_{\omega} \mathcal{R}_{\rho^{(k)}}(\omega, \Phi)\|^2. \quad (13)$$

Readers can refer to the Appendix B.3 of [57] for more details about the relationship between Eqs. (5) and (13).

We further use a two-stage method to approximate the minimax procedure:

- Minimize  $\mathcal{R}(\omega, \Phi)$  over  $(\omega, \Phi)$  and simultaneously maximize  $\sum_{k=1}^K \|\nabla_{\omega} \mathcal{R}_{\rho^{(k)}}(\omega, \Phi)\|^2$  over  $\rho$ .
- Fix  $\rho$  and minimize Eq. (13) over  $(\omega, \Phi)$ .

**Training Details.** For the synthetic datasets and the house price prediction dataset, we use the full batch gradient in optimization. The ERM method is trained for 4000 epochs. The IRM/ZIN method is also trained for 4000 epochs, with additional annealing in the first 2000 epochs. We train EIIL for 4000 epochs which is divided equally, i.e., 2000 epochs, in each of the two stages. For the CelebA classification task, we use mini-batch training with batch size of 128. All the methods are trained for 50 epochs, with annealing strategy in the first 25 epochs. For the Landcover task, we use mini-batch training with batch size of 1024, and all the methods are trained for 400 epochs with annealing strategy in the first 40 epochs. We use Adam [23] with learning rate 0.001 as our optimizer. Our experiments are run in a Linux workstation with Intel Xeon 3.20GHz CPU, 128GB RAM, and Nvidia GTX 3090 GPU. It takes about 5 GPU hours for the CelebA tasks, while the house price task and Landcover task cost less than 10 minutes for training. Notably, since IRM suffers from overfitting problem when applied to large models like ResNet-18 [27, 57], we fix the feature extraction backbone (that is, use a pre-trained model) in the CelebA experiment to alleviate this issue.

## D More Experiment Results

### D.1 Spatial Heterogeneity

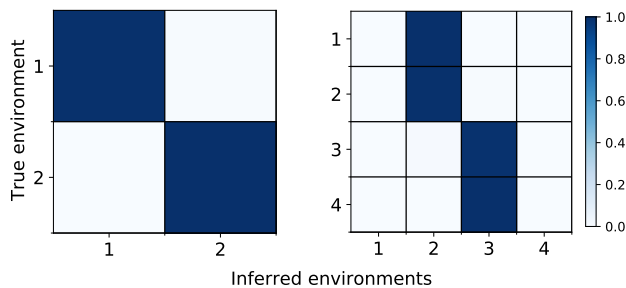
We also consider spatial heterogeneity that is also commonly encountered in practice, e.g., environments may be divided according to locations of latitude and longitude. We simulate spatial heterogeneity in the same way the data generation process for time heterogeneity but use a two-dimensional spatial variable  $\mathbf{r} = [r_1, r_2] \in [0, 1]^2$ . We simulate four environments for training by equally splitting the space into four blocks, i.e.,  $\{[0, 0.5] \times [0, 0.5], [0, 0.5] \times [0.5, 1], [0.5, 1] \times [0, 0.5], [0.5, 1] \times [0.5, 1]\}$ . Similarly, we denote a simulated case by tuple of  $p_s(\mathbf{r})$  in the four elements. We also evaluate the performance on four distinct test environments with  $p_s \in \{0.999, 0.8, 0.2, 0.1\}$  and  $p_v$  being constant. More details regarding the implementations are given in Appendix C.

The experimental results of spatial heterogeneity in Table 6 show similar performances to those of temporal heterogeneity. An interesting observation is that ZIN outperforms IRM with ground-truth environments in this simulation, especially when ‘‘duplicated’’ environments exist, e.g., when

**Table 6:** Test Mean and Worst accuracy (%) on four spatial heterogeneity synthetic datasets.

Env Partition	$p_s(t)$	0.999, 0.999, 0.7, 0.7				0.999, 0.9, 0.8, 0.7				0.999, 0.999, 0.8, 0.8			
	$p_v$	0.9		0.8		0.9		0.8		0.9		0.8	
	Test Acc	Mean	Worst	Mean	Worst	Mean	Worst	Mean	Worst	Mean	Worst	Mean	Worst
No	ERM	76.65	59.48	60.33	27.25	76.59	59.35	60.30	27.25	69.93	44.65	56.23	16.60
	EHL	37.81	16.89	66.46	50.35	37.03	14.01	66.60	50.28	70.18	45.23	71.00	58.72
	HRM	49.98	49.95	49.97	49.92	49.99	49.98	50.00	49.99	49.97	49.95	49.99	49.97
	ZIN	<b>88.66</b>	<b>87.23</b>	<b>79.16</b>	<b>78.04</b>	<b>88.28</b>	<b>86.29</b>	<b>78.92</b>	<b>77.49</b>	<b>88.00</b>	<b>85.75</b>	<b>78.80</b>	<b>77.25</b>
Yes	IRM	83.71	82.24	73.26	71.25	86.73	83.79	75.80	73.33	84.39	81.48	73.15	69.97

$p_s(\mathbf{r}) = (0.999, 0.999, 0.7, 0.7)$ . We conjecture that in this case the “ground-truth” partition may not be the most effective for invariant learning due to the heavily overlapped environments. As shown in the right panel of Figure 2, ZIN automatically merges duplicated environments.



**Figure 2:** Visualization of inferred environments. **Left:** temporal heterogeneity setting (0.999, 0.7). **Right:** spatial heterogeneity setting (0.999, 0.999, 0.7, 0.7).

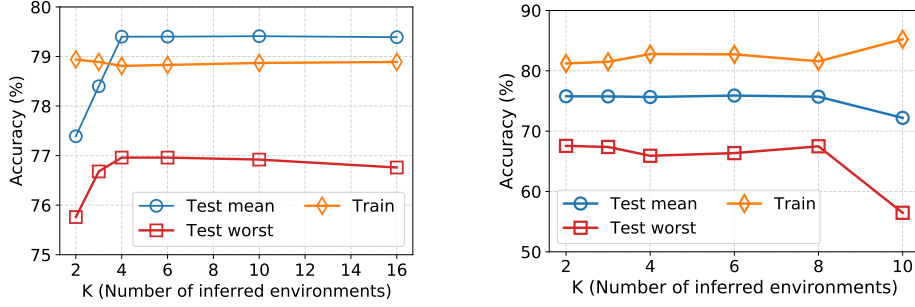
## D.2 Ablation Study on Auxiliary Information and Number of Environments

We now verify the theoretical results in Section 5 by choosing different inputs as  $\mathbf{Z}$ . We adopt a setting of spatial heterogeneity with  $p_s(\mathbf{r}) = (0.999, 0.999, 0.8, 0.8)$  and  $p_v = 0.8$ . Note that the heterogeneity in this setting is only along the second dimension  $r_2$ . The results are shown in Table 7. One can verify that Conditions 1 and 2 are satisfied when  $\mathbf{r}$  or  $r_2$  is chosen as  $\mathbf{Z}$ . The mean and worst test accuracy implies that ZIN based on  $\mathbf{r}$  or  $r_2$  can effectively remove the spurious feature. Since  $p_s(\mathbf{r})$  does not change along  $r_1$ , choosing  $r_1$  to infer environment partition violates Condition 2, which is reflected by the poor performance of ZIN with  $r_1$ . Using  $[\mathbf{X}, Y]$  as input to  $\rho(\cdot)$  is also a violation by Corollary 1, and the corresponding results in Table 7 confirm our analysis. Lastly, notice that  $X_s \in \mathbf{X}$  contains some information of  $Y$  and there may exist a function  $\rho(\cdot)$  so that  $H(Y|X_v) > H(Y|X_v, \rho(\mathbf{X}))$ . This violates Condition 1 and leads to poor results when choosing  $\mathbf{X}$  as input to  $\rho(\cdot)$ .

**Table 7:** Ablation study on choice of  $\mathbf{Z}$ .

$\mathbf{Z}$	Condition 1	Condition2	Test Mean	Test Worst
$\mathbf{r}$	✓	✓	78.80	77.25
$r_1$	✓	✗	56.30	16.84
$r_2$	✓	✓	78.79	77.21
$\mathbf{X}$	✗	✓	59.85	25.99
$\mathbf{X}, Y$	✗	✓	71.09	58.85

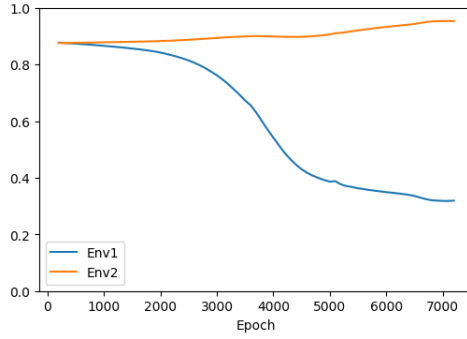
We also empirically verify the choice of hyper-parameter  $K$  using the same simulated setting. Fig. 3 shows that  $K$  has a relatively small impact, especially when  $K \geq 4$ .



**Figure 3:** Ablation study on the choice of  $K$ . **Left:** synthetic dataset. **Right:** CelebA.

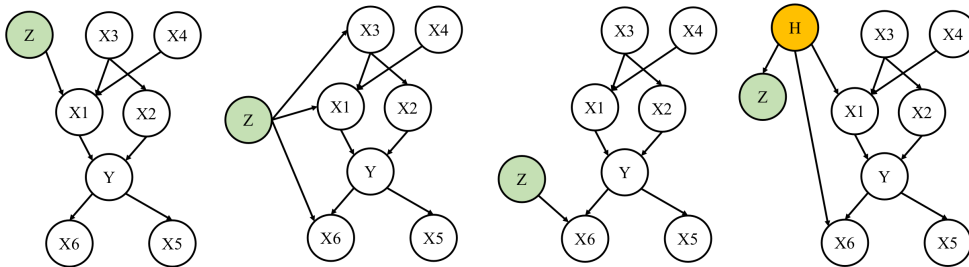
### D.3 About the Inferred Environments on CelebA

We visualize the inferred environments of the CelebA dataset in this section. The target *Smiling* is spuriously correlated with feature *Gender*, i.e., most females are smiling while most males are not. We set  $K$  to be 2 (our framework is insensitive to  $K$  as shown in the right panel of Fig. 3). We visualize the spurious correlations in the two inferred environments during training in Fig. 4. ZIN can generate two environments where the spurious correlations differ. Then we can easily discard the spurious feature using IRM methods on the inferred environments.

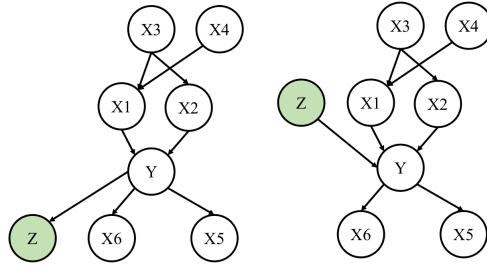


**Figure 4:** Correlation of the spurious feature (Gender) with the target (Smiling). The spurious correlation is calculated as the percentage of samples whose target (Smiling/Not Smiling) aligns with its gender (Female/Male). For example, a smiling female or a non-smiling male is counted as score 1, otherwise as score 0. The average score represents the correlation between Smiling and Gender.

### E Examples of Valid and Invalid Choices of Auxiliary Information



**Figure 5:** Examples of valid choices of  $Z$  satisfying Condition 1. Here “H” in the 4th graph denotes some hidden confounders. The invariant features are  $X_1$  and  $X_2$ , direct causes of  $Y$ .



**Figure 6:** Invalid choices of  $Z$  that violate Condition 1. The invariant features are  $X_1$  and  $X_2$ , the direct causes of  $Y$ .

## F Societal Impact

In this work, we propose to utilize the auxiliary information to aid the invariance learning without environmental indexes. This method is also helpful to fairness issues; see, e.g., the discussion about out-of-distribution generation and algorithmic fairness in [10]. For the additional information, we should also avoid using demographic, private, and sensitive information.