
On Invariance Penalties for Risk Minimization

Kia Khezeli
Mayo Clinic
khezeli.kia@mayo.edu

Arno Blaas
Oxford University
arno@robots.ox.ac.uk

Frank Soboczinski
King's College London
frank.soboczinski@kcl.ac.uk

Nicholas Chia
Mayo Clinic
chia.nicholas@mayo.edu

John Kalantari
Mayo Clinic
kalantari.john@mayo.edu

Abstract

The Invariant Risk Minimization (IRM) principle was first proposed by Arjovsky et al. [2019] to address the domain generalization problem by leveraging data heterogeneity from differing experimental conditions. Specifically, IRM seeks to find a data representation under which an optimal classifier remains invariant across all domains. Despite the conceptual appeal of IRM, the effectiveness of the originally proposed invariance penalty has recently been brought into question. In particular, there exists counterexamples for which that invariance penalty can be arbitrarily small for non-invariant data representations. We propose an alternative invariance penalty by revisiting the Gramian matrix of the data representation. We discuss the role of its eigenvalues in the relationship between the risk and the invariance penalty, and demonstrate that it is ill-conditioned for said counterexamples. The proposed approach is guaranteed to recover an invariant representation for linear settings under mild non-degeneracy conditions. Its effectiveness is substantiated by experiments on DomainBed and InvarianceUnitTest, two extensive test beds for domain generalization.

1 Introduction

Under the learning paradigm of Empirical Risk Minimization (ERM) [Vapnik, 1992], data is assumed to consist of independent and identically distributed (iid) samples from an underlying generating distribution. As the data generating distribution is often unknown in practice, ERM seeks predictors with minimal average training error (i.e., empirical risk) over the training set. Despite becoming a ubiquitous paradigm in machine learning, a growing body of literature [Arjovsky et al., 2019, Teney et al., 2020] has revealed that ERM and the the common practice of shuffling data inadvertently results in capturing all correlations found in the training data, whether spurious or causal, and produces models that fail to *generalize* to test data. The potential variation of experimental conditions from training to the utilization in real-world applications, manifests in discrepancy between training and testing distributions. This, in turn, highlights the need for machine learning algorithms to *generalize out-of-distribution (OoD)*.

Shuffling and treating data as iid risks possibly losing important information about the underlying conditions of the data generating process. Instead, partitioning training data into *environments*, e.g., based on the conditions under which data is generated, can exploit these differences to enhance generalization. Based on this observation, Arjovsky et al. [2019] introduce the principle of Invariant Risk Minimization (IRM) with the objective of finding a predictor that is invariant across all training environments (see Definition 1 and Equation 2). Because of the conceptually appealing nature of IRM and its potential to address the OoD-generalization problem, there is a stream of literature scrutinizing various facets of the original framework, e.g., extensions to other settings including online learning

[Javed et al., 2020] and treatment effect estimation [Shi et al., 2020], introducing game-theoretic interpretations [Ahuja et al., 2020], and raising concerns on the drawbacks and limitations of current IRM implementations [Rosenfeld et al., 2021, Kamath et al., 2021].

For an in-depth overview of the broader generalization literature, we refer the interested reader to [Arjovsky, 2020] and the references therein, and for an empirical evaluation of the performance of a number of the state-of-the-art methods on various test cases, we refer the reader to [Gulrajani and Lopez-Paz, 2020].

1.1 Contributions

In this paper, we propose a novel invariance penalty for practical implementation of IRM. Our proposed invariance penalty is directly related to risk. More precisely, we show that the risk in each environment under an arbitrary classifier equals to the risk under the invariant classifier for that environment plus the proposed invariance penalty between the said classifier and the optimal one. Moreover, we show that the proposed framework finds an invariant predictor for the setting in which the data is generated according to a linear Structural Equation Model (SEM) when provided a sufficient number of training environments under a mild non-degeneracy condition, which is similar in nature to the ones considered in [Arjovsky et al., 2019, Rosenfeld et al., 2021].

In addition, this work serves to illustrate the importance of the eigenstructure of the Gram matrix of the data representation for IRM. In particular, we show that the Gram matrix is ill-conditioned in the counterexample of Rosenfeld et al. [2021] where the invariance penalty of Arjovsky et al. [2019] is made arbitrarily small. Moreover, we characterize the difference between our proposed invariance penalty and the one proposed by Arjovsky et al. [2019] in terms of the eigenvalues of the Gram matrix of the data representation. This eigenstructure plays a significant role in the failure of invariance penalties including the one proposed by Arjovsky et al. [2019].

Finally, we evaluate our method on various test cases including DomainBed [Gulrajani and Lopez-Paz, 2020], a test bed including various benchmark data sets for domain generalization, and InvarianceUnitTest [Aubin et al., 2020], a test bed with three synthetically generated data sets capturing different structures of spurious correlations. We demonstrate the competitiveness of our proposed framework with other state-of-the-art methods for addressing OoD generalization problem.

1.2 Organization

The remainder of the paper is organized as follows. In Section 2, we formally define the notion of invariant prediction, the invariant risk minimization principle, and its relaxation proposed by Arjovsky et al. [2019]. In Sections 3 and 4, we introduce our more practical implementation and the rationale for its design. In Section 5, we evaluate the efficacy of our proposed model and compare it with other variations of IRM over a series of experiments. We conclude the paper in Section 6. All mathematical proofs are presented in the Appendix.

2 Background: Invariant Prediction

In this paper, we consider data (X^e, Y^e) collected from multiple training environments \mathcal{E}_{tr} where the distribution of (X^i, Y^i) and (X^j, Y^j) may be different for $i \neq j$ with $i, j \in \mathcal{E}_{\text{tr}}$. We denote by R_e the risk under environment e . That is, for predictor $f : \mathcal{X} \rightarrow \mathcal{Y}$, and loss function $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$, the risk under environment e is defined as

$$R_e(f) = \mathbf{E}_{X^e, Y^e} [\ell(f(X^e), Y^e)]. \tag{1}$$

2.1 Invariant Risk Minimization

Arjovsky et al. [2019] define the notion of invariant predictors under a multi-environment setting as follows.

Definition 1 (Invariant Predictor). A data representation $\varphi : \mathcal{X} \rightarrow \mathcal{H}$ is said to elicit an invariant predictor $w \circ \varphi$ across environments \mathcal{E} if there exists a classifier $w : \mathcal{H} \rightarrow \mathcal{Y}$, which is optimal for all environments, i.e., $w \in \operatorname{argmin}_{\tilde{w} : \mathcal{H} \rightarrow \mathcal{Y}} R_e(\tilde{w} \circ \varphi)$ for all $e \in \mathcal{E}$.

To find such invariant predictors, Arjovsky et al. [2019] introduce the notion of the Invariant Risk Minimization (IRM) principle:

$$\begin{aligned} & \min_{\substack{\varphi: \mathcal{X} \rightarrow \mathcal{H} \\ w: \mathcal{H} \rightarrow \mathcal{Y}}} \sum_{e \in \mathcal{E}_{\text{tr}}} R_e(w \circ \varphi) \\ & \text{subject to } w \in \underset{\tilde{w}: \mathcal{H} \rightarrow \mathcal{Y}}{\text{argmin}} R_e(\tilde{w} \circ \varphi), \forall e \in \mathcal{E}_{\text{tr}}. \end{aligned} \quad (2)$$

As this bi-leveled optimization problem is rather intractable, Arjovsky et al. [2019] propose a practical implementation of IRM by relaxing the invariance constraint (which itself requires solving an optimization problem) to an invariance penalty. We review it in what follows.

2.2 IRMv1: A Relaxation of IRM

In order to provide an implementation of IRM, Arjovsky et al. [2019] restrict the classifier w to linear functions, i.e.,

$$\begin{aligned} & \min_{\substack{\varphi: \mathcal{X} \rightarrow \mathcal{H} \\ w \in \mathbb{R}^{d_\varphi}}} \sum_{e \in \mathcal{E}_{\text{tr}}} R_e(w^\top \varphi) \\ & \text{subject to } w \in \underset{\tilde{w} \in \mathbb{R}^{d_\varphi}}{\text{argmin}} R_e(\tilde{w}^\top \varphi), \forall e \in \mathcal{E}_{\text{tr}}. \end{aligned} \quad (3)$$

To motivate their proposed penalty, Arjovsky et al. [2019] first consider the squared loss, i.e., $\ell(f(x), y) = \|f(x) - y\|^2$ where $\|\cdot\|$ denotes the Euclidean norm. Let the matrix $\mathcal{I}_e(\varphi)$ be defined as

$$\mathcal{I}_e(\varphi) := \mathbf{E}_{X^e} [\varphi(X^e)\varphi(X^e)^\top]. \quad (4)$$

Assuming that $\mathcal{I}_e(\varphi)$ is full rank for a fixed φ , its respective optimal classifier is unique, i.e.,

$$\underset{\tilde{w} \in \mathbb{R}^{d_\varphi}}{\text{argmin}} R_e(\tilde{w}^\top \varphi) = w_e^*(\varphi), \quad (5)$$

where

$$w_e^*(\varphi) = \mathcal{I}_e(\varphi)^{-1} \mathbf{E}_{X^e, Y^e} [\varphi(X^e)Y^e]. \quad (6)$$

To relax the constraint $w - w_e^*(\varphi) = 0$ to a penalty, Arjovsky et al. [2019] first consider the natural choice of $\|w - w_e^*(\varphi)\|^2$. However, they show that this penalty does not capture invariance by constructing an example for which $\|w - w_e^*(\varphi)\|^2$ is not well-behaved (see Section 2.3 for more details). Using the insight from this example, they propose $\|\mathcal{I}_e(\varphi)(w - w_e^*(\varphi))\|^2$ as an invariant penalty. For the squared loss, one can show that

$$\|\mathcal{I}_e(\varphi)(w - w_e^*(\varphi))\|^2 = (1/4) \|\nabla_w R_e(w^\top \varphi)\|^2. \quad (7)$$

Hence, their proposed penalty is given by

$$\rho_e^{\text{IRMv1}}(\varphi, w) := \|\nabla_w R_e(w^\top \varphi)\|^2. \quad (8)$$

Using the penalty (8), the relaxation of IRM is given by

$$\min_{\varphi, w} \sum_{e \in \mathcal{E}_{\text{tr}}} R_e(w^\top \varphi) + \lambda \rho_e^{\text{IRMv1}}(\varphi, w), \quad (9)$$

where $\lambda \geq 0$ is the penalty coefficient. Notice that for a given w and φ the predictor $w \circ \varphi$ can be expressed using different classifiers and data representations, i.e., $w \circ \varphi = \tilde{w} \circ \tilde{\varphi}$ where $\tilde{w} = w \circ \psi^{-1}$ and $\tilde{\varphi} = \psi \circ \varphi$ for some invertible mapping $\psi: \mathcal{H} \rightarrow \mathcal{H}$. Hence, in principle, it is possible to fix w without loss of generality. By relying on this observation, Arjovsky et al. [2019] fix the classifier as a scalar $w = 1$, and, thus, search for invariant data representation of the form $\varphi \in \mathbb{R}^{1 \times d_x}$. Their relaxation of IRM, which they refer to by IRMv1 is given by

$$\min_{\varphi} \sum_{e \in \mathcal{E}_{\text{tr}}} R_e(\varphi) + \lambda \rho_e^{\text{IRMv1}}(\varphi, 1.0). \quad (\text{IRMv1})$$

Although equation (7) only holds for squared loss, Arjovsky et al. [2019, Theorem 4] show that for all differentiable loss functions $(w^\top \Phi)^\top \nabla_w R(w^\top \Phi) = 0$ if and only if w is optimal for all environments. Here, the matrix Φ parameterizes the data representation. Hence, they justify the choice of $\|\nabla_w R_e(w^\top \varphi)\|^2$ as an invariance penalty for other loss functions, e.g., cross-entropy loss. However, recently Rosenfeld et al. [2021] constructed a counterexample by finding a non-invariant data representation for which the penalty $\|\nabla_w R_e(w^\top \varphi)\|^2$ with logistic loss is arbitrarily small.

Recall the assumption of invertability of $\mathcal{I}_e(\varphi)$, which was useful in the derivation of the invariance penalty $\rho_e^{\text{IRMv1}}(\varphi, w)$ for squared loss. In what follows, we investigate the role of the eigenstructure of $\mathcal{I}_e(\varphi)$ in relation to invariance penalization, and in particular, the existing counterexamples for the two penalties considered in this section.

2.3 The Role of $\mathcal{I}_e(\varphi)$

In proposing their invariance penalty, Arjovsky et al. [2019] consider an example in which $\varphi_c(x)$ is parameterized by a variable $c \in \mathbb{R}$ where $c = 0$ for the invariant data representation (see Appendix B.1 for further details). Figure 1 depicts various candidates for invariance penalty at the invariant classifier $w = w_{\text{inv}}$. As Arjovsky et al. [2019] point out, $\|w_{\text{inv}} - w_e^*(\varphi_c)\|^2$ is a poor choice for the invariance penalty as it is discontinuous at the invariant representation with $c = 0$, and vanishes as $c \rightarrow \infty$. Interestingly, $\mathcal{I}_e(\varphi_c)$ is ill-conditioned for both small and large c 's. More precisely, it holds that $\lim_{c \rightarrow 0} \kappa(\mathcal{I}_e(\varphi_c)) = \lim_{c \rightarrow +\infty} \kappa(\mathcal{I}_e(\varphi_c)) = +\infty$ where $\kappa(\cdot)$ denotes the condition number. That is, for a normal matrix A , its condition number is $\kappa(A) := |\lambda_{\max}(A)|/|\lambda_{\min}(A)|$ where λ_{\max} and λ_{\min} denote its maximum and minimum eigenvalues, respectively. Although multiplying $(w_{\text{inv}} - w_e^*(\varphi_c))$ by $\mathcal{I}_e(\varphi_c)$ circumvents the poor behavior of the invariance penalty for this example, it may not appropriately capture invariance in general as argued by Rosenfeld et al. [2021].

We now examine the counterexample introduced by Rosenfeld et al. [2021]. They consider a setting in which the data is generated according to a Structural Equation Model (SEM) (see Section 4.3). They show that for this setting, there exists a non-invariant data representation under which $\|\nabla_w R_e(w^\top \varphi)\|^2$ with logistic loss is arbitrarily small and hence it is poor discrepancy as an invariance penalty. For their counterexample, the matrix $\mathcal{I}_e(\varphi_c)$ is also ill-conditioned.

We provide a detailed derivation of the condition number of $\mathcal{I}_e(\varphi_c)$ for both Arjovsky et al. [2019] and Rosenfeld et al. [2021] in Appendices B.1 and B.2, respectively.

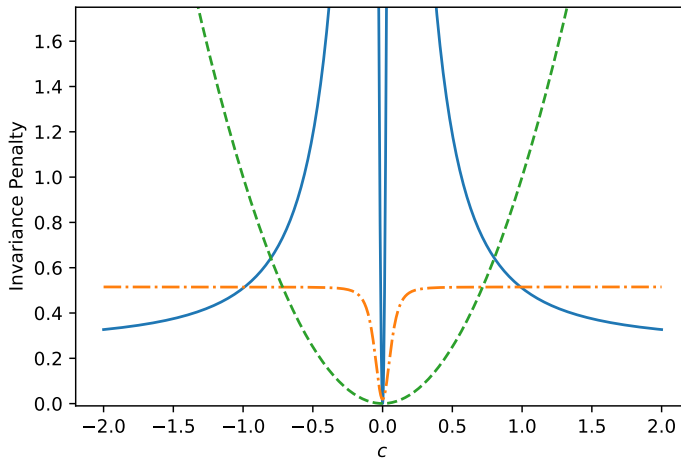


Figure 1: Invariance penalties $\|w_{\text{inv}} - w_e^*(\varphi_c)\|^2$, $\|\mathcal{I}_e(\varphi)^{1/2}(w_{\text{inv}} - w_e^*(\varphi_c))\|^2$, and $\|\mathcal{I}_e(\varphi)(w_{\text{inv}} - w_e^*(\varphi_c))\|^2$ are depicted in solid blue, dashed-dot orange, and dashed green, respectively.

3 IRMv2: An Alternative Penalty

As discussed in Section 2.3, both $\|w - w_e^*(\varphi_c)\|^2$ and $\|\mathcal{I}_e(\varphi)(w - w_e^*(\varphi_c))\|^2$ may be inappropriate choices for the invariance penalty due to their instability in terms of the eigenstructure of $\mathcal{I}_e(\varphi)$. We revisit the structure of the risk in order to propose an alternative penalty. In particular, in the following Lemma, we provide the sub-optimality gap of risk under an arbitrary classifier in comparison to the optimal classifier.

Lemma 1. *Consider squared loss function. Let $w \in \mathbb{R}^{d_\varphi}$ and $w_e^*(\varphi)$ as defined in Equation (6). Then,*

$$R_e(w^\top \varphi) = R_e(w_e^*(\varphi)^\top \varphi) + \left\| \mathcal{I}_e(\varphi)^{1/2} (w - w_e^*(\varphi)) \right\|^2. \quad (10)$$

Based on Lemma 1, we propose an invariance penalty that is directly comparable to risk.

$$\rho_e^{\text{IRMv2}}(\varphi, w) := \left\| \mathcal{I}_e(\varphi)^{1/2} (w - w_e^*(\varphi)) \right\|^2. \quad (11)$$

The relaxation of IRM using the penalty (11) is then given by

$$\min_{\varphi, w} \sum_{e \in \mathcal{E}_{\text{tr}}} R_e(w^\top \varphi) + \lambda \rho_e^{\text{IRMv2}}(\varphi, w). \quad (12)$$

We further simplify the relaxation (12) by finding its optimal classifier for a fixed data representation defined as

$$w^*(\varphi) := \operatorname{argmin}_w \sum_{e \in \mathcal{E}_{\text{tr}}} R_e(w^\top \varphi) + \lambda \rho_e^{\text{IRMv2}}(\varphi, w). \quad (13)$$

In the following Lemma, we leverage on the structure of the squared loss to find $w^*(\varphi)$.

Lemma 2. *Consider squared loss function and fixed φ . Let $w_e^*(\varphi)$ and $w^*(\varphi)$ as defined in Equations (6) and (13), respectively. Then,*

$$w^*(\varphi) = \left(\sum_{e \in \mathcal{E}_{\text{tr}}} \mathcal{I}_e(\varphi) \right)^{-1} \left(\sum_{e \in \mathcal{E}_{\text{tr}}} \mathcal{I}_e(\varphi) w_e^*(\varphi) \right). \quad (14)$$

Moreover,

$$\operatorname{argmin}_w \sum_{e \in \mathcal{E}_{\text{tr}}} R_e(w^\top \varphi) = w^*(\varphi). \quad (15)$$

Based on Lemmas 1 and 2, we propose the following relaxation of IRM, which we refer to by IRMv2.

$$\min_{\varphi} \sum_{e \in \mathcal{E}_{\text{tr}}} R_e(w^*(\varphi)^\top \varphi) + \lambda \rho_e^{\text{IRMv2}}(\varphi, w^*(\varphi)). \quad (\text{IRMv2})$$

We provide the pseudo-code for IRMv2 in Algorithm 1.

There are several factors distinguishing IRMv2 from IRMv1. First, IRMv2 relies on the optimal classifier $w^*(\varphi)$ while $w = 1.0$ in IRMv1. Second, the loss function in IRMv2 is squared loss while IRMv1 allows for utilization of other loss functions. Although this additional flexibility of IRMv1 may seem appealing, the counterexample of Rosenfeld et al. [2021] shows the failure of the penalty of IRMv1 to capture invariance for logistic loss. Finally, $\mathcal{I}_e(\varphi)$ is incorporated differently in the invariance penalty of IRMv1 and IRMv2. We formalize this latter observation in the following Section.

3.1 IRMv1A: An Adaptive Penalty Coefficient

We first bound the invariance penalty of IRMv1 in terms of the penalty of IRMv2 and the eigenvalues of $\mathcal{I}_e(\varphi)$. Then, based on this comparison, we propose an adaptive approach in choosing the penalty coefficient for IRMv1, which we refer to as IRMv1-Adaptive (IRMv1A).

Algorithm 1 IRMv2

- 1: **Input:** Data set: D_e for $e \in \mathcal{E}_{\text{tr}}$. Loss function: Squared loss, Parameters: penalty coefficient $\lambda \geq 0$, data representation parameters $\theta \in \mathbb{R}^{d_\theta}$, learning rate η_t , training horizon T .
 - 2: **Initialize** θ_1 randomly
 - 3: **for** $t = 1, 2, \dots, T$ **do**
 - 4: **for** $e \in \mathcal{E}_{\text{tr}}$ **do**
 - 5: compute the LSE $w_e^*(\varphi_{\theta_t})$ according to Eq. (6)
 - 6: compute the optimal classifier $w^*(\varphi_{\theta_t})$ according to Eq. (13)
 - 7: $\mathcal{L}_t(\varphi_{\theta_t}) \leftarrow \sum_{e \in \mathcal{E}_{\text{tr}}} \mathcal{R}_e(w^*(\varphi_{\theta_t})^\top \varphi_{\theta_t}) + \lambda \rho_e^{\text{IRMv2}}(\varphi_{\theta_t}, w^*(\varphi_{\theta_t}))$
 - 8: $\theta_{t+1} \leftarrow \theta_t - \eta_t \nabla_{\theta_t} \mathcal{L}_t(\varphi_{\theta_t})$
 - 9: **Output** predictor $w^*(\varphi_{\theta_T})^\top \varphi_{\theta_T}$.
-

Lemma 3. Let $\rho_e^{\text{IRMv1}}(\varphi, w)$ and $\rho_e^{\text{IRMv2}}(\varphi, w)$ be the invariance penalties of the IRMv1 and IRMv2 defined in Equations (8) and (11), respectively. Then,

$$\lambda_{\min}(\mathcal{I}_e(\varphi)) \rho_e^{\text{IRMv2}}(\varphi, w) \leq \rho_e^{\text{IRMv1}}(\varphi, w) \leq \lambda_{\max}(\mathcal{I}_e(\varphi)) \rho_e^{\text{IRMv2}}(\varphi, w). \quad (16)$$

The proof of Lemma 3 directly follows from the definition of the invariance penalties $\rho_e^{\text{IRMv1}}(\varphi, w)$ and $\rho_e^{\text{IRMv2}}(\varphi, w)$, and the fact that for a symmetric matrix $A \in \mathbb{R}^{d \times d}$ and a vector $u \in \mathbb{R}^d$, it holds that $\lambda_{\min}(A) \|u\|^2 \leq u^\top A u \leq \lambda_{\max}(A) \|u\|^2$.

Using Lemma 3, we suggest the following rule for the penalty coefficient of IRMv1.

$$\lambda_e := \frac{1}{\lambda_0 + \lambda_{\min}(\mathcal{I}_e(\varphi))} \quad (17)$$

for a user-specified $\lambda_0 \geq 0$. Note that this is an adaptive rule, as φ may change throughout training.

4 Theoretical Results

In this section, we consider the setting introduced by Rosenfeld et al. [2021] and provide theoretical guarantees that IRM with linear classifier and squared loss, and subsequently IRMv2 recover an invariant predictor.

4.1 Problem Setup

We consider a setting in which the data is generated according to a Structural Equation Model [Pearl, 2009]. More precisely, for each environment e , (X^e, Y^e) is generated as

$$X^e = S \begin{bmatrix} Z_c \\ Z_e \end{bmatrix}, \quad Y^e = \begin{cases} 1, & \text{with prob. } \eta, \\ -1, & \text{with prob. } 1 - \eta, \end{cases} \quad (18)$$

where $\eta \in [0, 1]$, and $S \in \mathbb{R}^{d \times (d_c + d_e)}$ is a left invertible matrix, i.e., there exists S^\dagger such that $S^\dagger S = I$. In this model, Z_c captures the causal variables that are invariant across environments, and Z_e captures the spurious environment dependent variables.

The variables Z_c and Z_e are generated as follows

$$Z_c = \mu_c Y + W_c \text{ where } W_c \sim \mathcal{N}(0, \sigma_c^2 I), \quad (19)$$

$$Z_e = \mu_e Y + W_e \text{ where } W_e \sim \mathcal{N}(0, \sigma_e^2 I). \quad (20)$$

Here, $\mu_c \in \mathbb{R}^{d_c}$, $\mu_e \in \mathbb{R}^{d_e}$, and $\mathcal{N}(\mu, \Sigma)$ denotes multi-variate Gaussian distribution with mean equal to μ and covariance matrix equal to Σ . We further assume that W_c , W_e , and Y^e are independent for all environments.

4.2 Invariant Representation under IRM

For the setting introduced in 4.1, the invariant data representation is linear. In particular, for any $d \geq d_c$, $\varphi(X^e) = \Phi_d X^e = Z_c$ is an invariant data representation, where

$$\Phi_d := \begin{bmatrix} I_{d_c \times d_c} & \mathbf{0}_{d_e \times d_c} \\ \mathbf{0}_{d_c \times (d-d_c)} & \mathbf{0}_{d_e \times (d-d_c)} \end{bmatrix} S^\dagger. \quad (21)$$

Naturally, the possibility of finding an invariant predictor depends on the number and the diversity of training environments. We now introduce non-degeneracy conditions on the training environment under which IRM is guaranteed to find an invariant predictor, provided sufficient number of training environments.

Let $|\mathcal{E}_{\text{tr}}| > d_e$. As $\text{span}(\{\mu_e\}_{e \in \mathcal{E}_{\text{tr}}}) \leq d_e$, for each $e \in \mathcal{E}_{\text{tr}}$ there exists a set of coefficients α_i^e for $i \in \mathcal{E}_{\text{tr}} \setminus e$ such that

$$\mu_e = \sum_{i \in \mathcal{E}_{\text{tr}} \setminus e} \alpha_i^e \mu_i. \quad (22)$$

We say that \mathcal{E}_{tr} is a *non-degenerate set of environments* if for all $e \in \mathcal{E}_{\text{tr}}$ it holds that

$$\sum_{i \in \mathcal{E}_{\text{tr}} \setminus e} \alpha_i^e \neq 1, \quad (23)$$

$$\text{rank}(\Gamma_e) = d_e, \quad (24)$$

where Γ_e is defined as

$$\Gamma_e := \frac{1}{1 - \sum_{i \in \mathcal{E}_{\text{tr}} \setminus e} \alpha_i^e} \left(\sigma_e^2 I + \mu_e \mu_e^\top - \sum_{i \in \mathcal{E}_{\text{tr}} \setminus e} (\sigma_i^2 I + \mu_i \mu_i^\top) \alpha_i^e \right).$$

The conditions (23) and (24) specify that the span of covariance matrices of Z_e is R^{d_e} . This is a natural requirement to eliminate the degrees of freedom on the dependency of the data representation on the environment dependent features. We note that the non-degeneracy conditions considered in Rosenfeld et al. [2021] are similar to (23) and (24) with the difference that instead of depending on covariance matrices of Z_e as in (24), their assumption relies on the variances σ_e^2 . This difference in the non-degeneracy requirements is due to the fact that they consider logistic loss and we consider squared loss.

Theorem 1. *Assume that $|\mathcal{E}_{\text{tr}}| > d_e$ where (X^e, Y^e) generated according to (18). Consider a linear data representation $\Phi X = AZ_c + BZ_e$ and a classifier $w(\Phi)$ on top of Φ that is invariant, i.e., $w(\Phi) = w_e^*(\Phi)$ for all $e \in \mathcal{E}_{\text{tr}}$. If non-degeneracy conditions Eqs. (22-24) holds, then either $w(\Phi) = 0$ or $B = 0$.*

4.3 Eigenstructure of $\mathcal{I}_e(\varphi)$

In this section, we compare the penalties of IRMv1 and IRMv2 for the counterexample of Rosenfeld et al. [2021]. They consider a data representation φ_ϵ where $\epsilon > 1$ determines the extent to which $\varphi_\epsilon(X^e)$ depends on Z_e . More specifically, φ_ϵ is defined as

$$\varphi_\epsilon(X^e) := \begin{bmatrix} Z_c \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ Z_e \end{bmatrix} \mathbf{1}_{\{Z_e \notin \mathcal{Z}_e\}}, \quad (25)$$

where $\{Z_e \notin \mathcal{Z}_e\}$ is an event with $\mathbf{P}(Z_e \in \mathcal{Z}_e) \leq p_{e,\epsilon}$ where $p_{e,\epsilon} := \exp(-d_e \min\{\epsilon - 1, (\epsilon - 1)^2\}/8)$. They show that the invariance penalty of IRMv1 decays at a rate faster than $p_{e,\epsilon}^2$ as ϵ grows. Thus, the penalty may be arbitrarily small for a large enough ϵ .

An invariant data representation for this setting is $\varphi_\epsilon(X^e)$ with $\epsilon = 1$. Moreover, in Appendix B.2, we show that $\kappa(\mathcal{I}_e(\varphi)) \geq c/p_{e,\epsilon}$ for some constant c that is independent of ϵ . Thus, $\mathcal{I}_e(\varphi)$ is ill-conditioned when the penalty of IRMv1 is small.

5 Experiments

In this section, we empirically evaluate the efficacy of our proposed implementations of IRM, namely, IRMv2 and IRMv1A with IRMv1. We demonstrate the competitiveness of our approach on *InvarianceUnitTests* [Aubin et al., 2020] and *DomainBed* [Gulrajani and Lopez-Paz, 2020], two recently proposed test beds for evaluation of domain generalization methods. In particular, we show that our approach generalizes in one of the InvarianceUnitTests where all other methods failed (i.e., exhibited tests accuracies that are comparable to random guessing).

5.1 InvarianceUnitTests

In this section, we evaluate the efficacy of our proposed approaches for invariance discovery on the *InvarianceUnitTests* recently proposed by Aubin et al. [2020]. These unit-tests entail three classes of low-dimensional linear problems, each capturing a different structure for inducing spurious correlations. For completeness, we provide a brief overview of these InvarianceUnitTests before providing a performance comparison between IRMv2, IRMv1A, IRMv1, ERM, Inter-environmental Gradient Alignment (IGA) [Koyama and Yamaguchi, 2020], and AND-Mask [Parascandolo et al., 2020]. The IGA method seeks to elicit invariant predictors by an invariance penalty in terms of the variance of the risk under different environments. The AND-Mask method, at each step of the training process, updates the model using the direction where gradient (of the loss) signs agree across environments.

The data set for each problem falls within the multi-environment setting described in Section 2 with $n_e = 10^4$. For all problems, the input $x^e \in \mathbb{R}^d$ is constructed as $x^e = (x_{\text{inv}}^e, x_{\text{spu}}^e)$ where $x_{\text{inv}}^e \in \mathbb{R}^{d_{\text{inv}}}$ and $x_{\text{spu}}^e \in \mathbb{R}^{d_{\text{spu}}}$ denote the invariant and the spurious features, respectively. To make the problems more realistic, Aubin et al. [2020] repeat each experiment and *scramble* the inputs by multiplying x^e by a rotation matrix. In each problem, the spurious correlations that exist in the training environments are discarded in the test environment by random shuffling. As a basis for comparison, similar to Aubin et al. [2020], we implement an Oracle ERM where the spurious correlations are shuffled in the training data sets as well, and hence, ERM can readily identify them.

Example 1 considers a *regression problem* based on Structural Equation Models Pearl [2009] where the target variable is a linear function of the invariant variables and the spurious variables are linear functions of the target variable. Example 2 considers a *classification problem* inspired by the infamous cow vs. camel example Beery et al. [2018] where spurious correlations are interpreted as background color. Example 3 is based on a classification experiment in Parascandolo et al. [2020] where the spurious correlations provide a *shortcut* in minimizing the training error while the invariant classifier takes a more complex form.

We summarize the test errors of all methods on the three examples and their scrambled variations in Table 1. We observe that on these structured unit-tests, most non-ERM methods are only successful in eliciting an invariant predictor in the linear regression case (Example 1). In particular, other than IRMv2 on Example 2 and IRMv1 on Example 3, all methods fail on these cases, i.e., exhibit test errors comparable to random guessing. As the structure of the spurious correlation is different in each of these examples, these mixed results highlight the challenge of constructing methods that generalize well with minimal reliance on the underlying causal structure.

5.2 DomainBed

DomainBed is an extensive framework released by Gulrajani and Lopez-Paz [2020] to test domain generalization algorithms for image classification tasks on various benchmark data sets. In a series of experiments, Gulrajani and Lopez-Paz [2020] show that enabled by data augmentation various state-of-the-art generalization methods perform similar to each other and ERM on several benchmark data sets.

Although the integration of additional data sets and algorithms to DomainBed is straightforward, we note that performing an extensive set of experiments requires significant computational resources as also pointed out by Krueger et al. [2020] (see the supplementary material for further details on the computational details). For this reason, we limit the scope of our experiments to the comparison of ERM, IRMv1, IRMv1A, and IRMv2.

| | ANDMask | ERM | IGA | IRMv1A | IRMv2 | IRMv1 | Oracle |
|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Example1.E0 | 0.07 ± 0.01 | 1.52 ± 0.50 | 6.61 ± 2.77 | 0.07 ± 0.01 | 0.06 ± 0.01 | 0.11 ± 0.02 | 0.05 ± 0.00 |
| Example1.E1 | 11.54 ± 0.31 | 14.10 ± 1.33 | 21.61 ± 3.85 | 13.53 ± 1.29 | 13.15 ± 1.29 | 12.46 ± 1.25 | 11.26 ± 0.17 |
| Example1.E2 | 20.56 ± 0.66 | 23.98 ± 2.21 | 33.43 ± 5.09 | 24.11 ± 2.35 | 23.53 ± 2.26 | 22.18 ± 2.27 | 19.99 ± 0.29 |
| Example1s.E0 | 0.07 ± 0.01 | 1.66 ± 0.49 | 6.49 ± 3.08 | 0.07 ± 0.01 | 0.06 ± 0.01 | 0.10 ± 0.02 | 0.05 ± 0.00 |
| Example1s.E1 | 12.71 ± 0.93 | 14.49 ± 1.45 | 21.83 ± 4.78 | 13.31 ± 1.34 | 13.36 ± 1.28 | 12.61 ± 1.20 | 11.25 ± 0.19 |
| Example1s.E2 | 22.59 ± 1.77 | 24.59 ± 2.40 | 34.01 ± 6.55 | 23.68 ± 2.35 | 23.73 ± 2.21 | 22.40 ± 2.16 | 20.01 ± 0.34 |
| Example2.E0 | 0.43 ± 0.01 | 0.43 ± 0.00 | 0.43 ± 0.00 | 0.43 ± 0.01 | 0.45 ± 0.01 | 0.47 ± 0.01 | 0.42 ± 0.01 |
| Example2.E1 | 0.50 ± 0.00 | 0.50 ± 0.00 | 0.50 ± 0.00 | 0.50 ± 0.01 | 0.35 ± 0.01 | 0.52 ± 0.01 | 0.50 ± 0.01 |
| Example2.E2 | 0.42 ± 0.01 | 0.42 ± 0.01 | 0.42 ± 0.01 | 0.42 ± 0.01 | 0.23 ± 0.01 | 0.47 ± 0.01 | 0.42 ± 0.01 |
| Example2s.E0 | 0.43 ± 0.01 | 0.43 ± 0.01 | 0.43 ± 0.01 | 0.43 ± 0.01 | 0.51 ± 0.07 | 0.44 ± 0.02 | 0.41 ± 0.02 |
| Example2s.E1 | 0.50 ± 0.00 | 0.50 ± 0.00 | 0.50 ± 0.00 | 0.49 ± 0.01 | 0.41 ± 0.06 | 0.50 ± 0.01 | 0.48 ± 0.02 |
| Example2s.E2 | 0.42 ± 0.00 | 0.42 ± 0.01 | 0.42 ± 0.01 | 0.42 ± 0.01 | 0.27 ± 0.04 | 0.44 ± 0.02 | 0.41 ± 0.02 |
| Example3.E0 | 0.41 ± 0.13 | 0.48 ± 0.07 | 0.45 ± 0.12 | 0.35 ± 0.11 | 0.48 ± 0.09 | 0.22 ± 0.13 | 0.01 ± 0.00 |
| Example3.E1 | 0.44 ± 0.12 | 0.49 ± 0.05 | 0.48 ± 0.07 | 0.35 ± 0.11 | 0.47 ± 0.10 | 0.22 ± 0.13 | 0.01 ± 0.00 |
| Example3.E2 | 0.42 ± 0.14 | 0.47 ± 0.10 | 0.47 ± 0.09 | 0.35 ± 0.11 | 0.47 ± 0.09 | 0.22 ± 0.13 | 0.01 ± 0.00 |
| Example3s.E0 | 0.49 ± 0.05 | 0.49 ± 0.06 | 0.50 ± 0.02 | 0.57 ± 0.18 | 0.49 ± 0.06 | 0.44 ± 0.16 | 0.03 ± 0.01 |
| Example3s.E1 | 0.50 ± 0.04 | 0.50 ± 0.00 | 0.50 ± 0.02 | 0.57 ± 0.18 | 0.49 ± 0.07 | 0.44 ± 0.16 | 0.03 ± 0.01 |
| Example3s.E2 | 0.49 ± 0.05 | 0.48 ± 0.09 | 0.49 ± 0.07 | 0.57 ± 0.18 | 0.48 ± 0.07 | 0.44 ± 0.17 | 0.03 ± 0.01 |

Table 1: Test errors for all algorithms and examples with $(d_{\text{inv}}, d_{\text{spu}}, d_{\text{env}}) = (5, 5, 3)$. The errors for Examples 1.E0 through 1s.E2 are in MSE and all others are classification error. The empirical mean and the standard deviation are computed using 10 independent experiments. An ‘s’ indicates the scrambled variation of its corresponding problem setting, e.g. Example 1s is the scrambled variation of the Example 1 regression setting.

| Algorithm | ColoredMNIST | RotatedMNIST | PACS | VLCS | Avg |
|-----------|--------------|--------------|------------|------------|------|
| ERM | 51.7 ± 0.1 | 96.7 ± 0.0 | 81.1 ± 0.1 | 78.8 ± 0.4 | 77.0 |
| IRMv1 | 51.8 ± 0.2 | 95.2 ± 0.4 | 78.6 ± 1.0 | 76.0 ± 0.5 | 75.4 |
| IRMv1A | 50.9 ± 0.1 | 64.7 ± 20.1 | 80.9 ± 0.0 | 77.3 ± 0.2 | 68.4 |
| IRMv2 | 50.8 ± 0.4 | 97.1 ± 0.0 | 82.6 ± 0.9 | 76.5 ± 0.4 | 76.8 |

Table 2: The test accuracy of ERM and different implementations of IRM on benchmark data-sets. Model selection of the DomainBed is chosen as training-domain validation set.

Similar to Gulrajani and Lopez-Paz [2020], we observe that no method significantly outperforms others on any of the benchmark data sets (see Table 2). For a complete set of results on DomainBed with various model selection methods, we refer the reader to Appendix C. As these data sets are image based and equipped with data augmentation, they may not provide comprehensive insight on the strengths and weaknesses of domain generalization algorithms on other modes of data, e.g., gathered in real-world applications.

6 Conclusion

In this paper, we have presented IRMv2, an alternative implementation of the IRM principle that aims to enable out-of-distribution generalization by finding environment invariant predictors. We establish theoretical results on the effectiveness of our approach in the linear setting. In doing so, we bring forward the importance of the eigenstructure of the Gramian matrix of the data representation. In particular, we show that for the existing counterexample on the potential failure of IRMv1, the aforementioned matrix is ill-conditioned for the invariant representation. This highlights the significance of the span of the data representations in relation to the span of the underlying true invariant features of the data. That is, if the data representation allows for more degrees of freedom than needed to capture invariance, the Gramian matrix of the invariant representation would be ill-conditioned. This observation provides intuition on the underlying reasons why current implementations of IRM may fail. While this work attempts to address some of the limitations of IRM that impede its widespread adoption, it leaves for future work subsequent investigations on data gathered from real-world applications beyond curated benchmark data sets.

References

- Kartik Ahuja, Karthikeyan Shanmugam, Kush Varshney, and Amit Dhurandhar. Invariant risk minimization games. In *Proceedings of the 37th International Conference on Machine Learning*, pages 145–155, 2020.
- Martin Arjovsky. *Out of Distribution Generalization in Machine Learning*. PhD thesis, New York University, 2020.
- Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization. *arXiv preprint arXiv:1907.02893*, 2019.
- Benjamin Aubin, Martin Arjovsky, Leon Bottou, and David Lopez-Paz. Linear unit tests for invariance discovery. In *Causal Discovery and Causality-Inspired Machine Learning Workshop at NeurIPS*, 2020.
- Sara Beery, Grant Van Horn, and Pietro Perona. Recognition in terra incognita. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 456–473, 2018.
- Ishaan Gulrajani and David Lopez-Paz. In search of lost domain generalization. *arXiv preprint arXiv:2007.01434*, 2020.
- Khurram Javed, Martha White, and Yoshua Bengio. Learning causal models online. *arXiv preprint arXiv:2006.07461*, 2020.
- Pritish Kamath, Akilesh Tangella, Danica J Sutherland, and Nathan Srebro. Does invariant risk minimization capture invariance? *arXiv preprint arXiv:2101.01134*, 2021.
- Masanori Koyama and Shoichiro Yamaguchi. Out-of-distribution generalization with maximal invariant predictor. *arXiv preprint arXiv:2008.01883*, 2020.
- David Krueger, Ethan Caballero, Joern-Henrik Jacobsen, Amy Zhang, Jonathan Binas, Remi Le Priol, and Aaron Courville. Out-of-distribution generalization via risk extrapolation (rex). *arXiv preprint arXiv:2003.00688*, 2020.
- Giambattista Parascandolo, Alexander Neitz, Antonio Orvieto, Luigi Gresele, and Bernhard Schölkopf. Learning explanations that are hard to vary. *arXiv preprint arXiv:2009.00329*, 2020.
- Judea Pearl. *Causality*. Cambridge university press, 2009.
- Elan Rosenfeld, Pradeep Kumar Ravikumar, and Andrej Risteski. The risks of invariant risk minimization. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=BbNIbVPJ-42>.
- Claudia Shi, Victor Veitch, and David Blei. Invariant representation learning for treatment effect estimation. *arXiv preprint arXiv:2011.12379*, 2020.
- Damien Teney, Ehsan Abbasnejad, and Anton van den Hengel. Unshuffling data for improved generalization. *arXiv preprint arXiv:2002.11894*, 2020.
- Vladimir Vapnik. Principles of risk minimization for learning theory. In *Advances in neural information processing systems*, pages 831–838, 1992.

A Mathematical Proofs

A.1 Proof of Results in Section 2

Lemma 4. Consider squared loss function. Let $w \in \mathbb{R}^{d_\varphi}$ and $w_e^*(\varphi)$ as defined in Equation (6). Then,

$$R_e(w^\top \varphi) = R_e(w_e^*(\varphi)^\top \varphi) + \left\| \mathcal{I}_e(\varphi)^{1/2} (w - w_e^*(\varphi)) \right\|^2. \quad (26)$$

Proof of Lemma 1: First, the risk under environment e with $w = w_e^*(\varphi)$ is given by

$$R_e(w_e^*(\varphi)^\top \varphi) = \mathbf{E}_{Y^e} [\|Y^e\|^2] - \mathbf{E}_{X^e, Y^e} [\varphi(X^e)Y^e]^\top \mathcal{I}_e(\varphi)^{-1} \mathbf{E}_{X^e, Y^e} [\varphi(X^e)Y^e].$$

Then,

$$\begin{aligned} & R_e(w^\top \varphi) - R_e(w_e^*(\varphi)^\top \varphi) \\ &= w^\top \mathcal{I}_e(\varphi)w - 2w^\top \mathbf{E}_{X^e, Y^e} [\varphi(X^e)Y^e] - \mathbf{E}_{X^e, Y^e} [\varphi(X^e)Y^e]^\top \mathcal{I}_e(\varphi)^{-1} \mathbf{E}_{X^e, Y^e} [\varphi(X^e)Y^e] \\ &= \left\| \mathcal{I}_e(\varphi)^{1/2} (w - w_e^*(\varphi)) \right\|^2. \end{aligned}$$

■

Lemma 5. Consider squared loss function and fixed φ . Let $w_e^*(\varphi)$ and $w^*(\varphi)$ as defined in Equations (6) and (13), respectively. Then,

$$w^*(\varphi) = \left(\sum_{e \in \mathcal{E}_{\text{tr}}} \mathcal{I}_e(\varphi) \right)^{-1} \left(\sum_{e \in \mathcal{E}_{\text{tr}}} \mathcal{I}_e(\varphi) w_e^*(\varphi) \right). \quad (27)$$

Moreover,

$$\operatorname{argmin}_w \sum_{e \in \mathcal{E}_{\text{tr}}} R_e(w^\top \varphi) = w^*(\varphi). \quad (28)$$

Proof of Lemma 1: Recall the definition of $w^*(\varphi)$

$$w^*(\varphi) = \operatorname{argmin}_w \sum_{e \in \mathcal{E}_{\text{tr}}} R_e(w^\top \varphi) + \lambda \rho_e^{\text{IRMv2}}(\varphi, w).$$

That is,

$$\begin{aligned} w^*(\varphi) &= \operatorname{argmin}_w \sum_{e \in \mathcal{E}_{\text{tr}}} \mathbf{E} \left[\|w^\top \varphi(X^e) - Y^e\|^2 \right] + \lambda \left\| \mathcal{I}_e(\varphi)^{1/2} (w - w_e^*(\varphi)) \right\|^2 \\ &= \operatorname{argmin}_w \sum_{e \in \mathcal{E}_{\text{tr}}} (1 + \lambda) w^\top \mathcal{I}_e(\varphi) w - 2w^\top (\mathbf{E} [\varphi(X^e)Y^e] + \lambda \mathcal{I}_e(\varphi) w_e^*(\varphi)) \\ &\quad + w^*(\varphi)^\top \mathcal{I}_e(\varphi) w^*(\varphi) + \mathbf{E} [\|Y^e\|^2]. \end{aligned}$$

Note that the objective function is a convex quadratic function of w . Hence, using the first-order optimality condition we have that

$$(1 + \lambda) \left(\sum_{e \in \mathcal{E}_{\text{tr}}} \mathcal{I}_e(\varphi) \right) w^*(\varphi) - \sum_{e \in \mathcal{E}_{\text{tr}}} (\mathbf{E} [\varphi(X^e)Y^e] + \lambda \mathcal{I}_e(\varphi) w_e^*(\varphi)) = 0.$$

Then,

$$w^*(\varphi) = \frac{1}{1 + \lambda} \left(\sum_{e \in \mathcal{E}_{\text{tr}}} \mathcal{I}_e(\varphi) \right)^{-1} \left(\sum_{e \in \mathcal{E}_{\text{tr}}} (\mathbf{E} [\varphi(X^e)Y^e] + \lambda \mathcal{I}_e(\varphi) w_e^*(\varphi)) \right).$$

Recall that $w_e^*(\varphi) = \mathcal{I}_e(\varphi)^{-1} [\varphi(X^e)Y^e]$. Then, $[\varphi(X^e)Y^e] = \mathcal{I}_e(\varphi) w_e^*(\varphi)$. Thus,

$$w^*(\varphi) = \left(\sum_{e \in \mathcal{E}_{\text{tr}}} \mathcal{I}_e(\varphi) \right)^{-1} \left(\sum_{e \in \mathcal{E}_{\text{tr}}} \mathcal{I}_e(\varphi) w_e^*(\varphi) \right).$$

Finally, using a similar argument, we get

$$\begin{aligned} \operatorname{argmin}_w \sum_{e \in \mathcal{E}_{\text{tr}}} R_e(w^\top \varphi) &= \operatorname{argmin}_w \sum_{e \in \mathcal{E}_{\text{tr}}} \mathbf{E} \left[\|w^\top \varphi(X^e) - Y^e\|^2 \right] \\ &= \operatorname{argmin}_w \sum_{e \in \mathcal{E}_{\text{tr}}} w^\top \mathcal{I}_e(\varphi) w - 2w^\top \mathbf{E} [\varphi(X^e)Y^e] + \mathbf{E} [\|Y^e\|^2] \\ &= \left(\sum_{e \in \mathcal{E}_{\text{tr}}} \mathcal{I}_e(\varphi) \right)^{-1} \left(\sum_{e \in \mathcal{E}_{\text{tr}}} \mathcal{I}_e(\varphi) w_e^*(\varphi) \right). \end{aligned}$$

■

Lemma 6. Let $\rho_e^{\text{IRMv1}}(\varphi, w)$ and $\rho_e^{\text{IRMv2}}(\varphi, w)$ be the invariance penalties of the IRMv1 and IRMv2. Then,

$$\lambda_{\min}(\mathcal{I}_e(\varphi)) \rho_e^{\text{IRMv2}}(\varphi, w) \leq \rho_e^{\text{IRMv1}}(\varphi, w) \leq \lambda_{\max}(\mathcal{I}_e(\varphi)) \rho_e^{\text{IRMv2}}(\varphi, w). \quad (29)$$

Proof of Lemma 1: For a symmetric matrix $A \in \mathbb{R}^{d \times d}$ and a vector $u \in \mathbb{R}^d$, it holds that $\lambda_{\min}(A)\|u\|^2 \leq u^\top A u \leq \lambda_{\max}(A)\|u\|^2$. Let $u = \mathcal{I}_e(\varphi)^{1/2}(w - w_e^*(\varphi))$ and $A = \mathcal{I}_e(\varphi)$. Then,

$$\begin{aligned} \|\mathcal{I}_e(\varphi)(w - w_e^*(\varphi))\|^2 &\leq \lambda_{\max}(\mathcal{I}_e(\varphi)) \left\| \mathcal{I}_e(\varphi)^{1/2}(w - w_e^*(\varphi)) \right\|^2 \\ \|\mathcal{I}_e(\varphi)(w - w_e^*(\varphi))\|^2 &\geq \lambda_{\min}(\mathcal{I}_e(\varphi)) \left\| \mathcal{I}_e(\varphi)^{1/2}(w - w_e^*(\varphi)) \right\|^2. \end{aligned}$$

■

A.2 Proof of Results in Section 4

In order to prove Theorem 2, we first find the optimal classifier for a given data representation for the linear setting in the following Lemma.

Lemma 7. Let $w_e^*(\Phi) = \operatorname{argmin}_w R_e(w^\top \Phi)$ where $R_e(w^\top \Phi)$ is defined as

$$R_e(w^\top \Phi) = \mathbf{E} \left[\operatorname{tr} \left(\left(w^\top \Phi X^e - \begin{bmatrix} \mathbb{1}_{\{Y=+1\}} \\ \mathbb{1}_{\{Y=-1\}} \end{bmatrix} \right) \left(w^\top \Phi X^e - \begin{bmatrix} \mathbb{1}_{\{Y=+1\}} \\ \mathbb{1}_{\{Y=-1\}} \end{bmatrix} \right)^\top \right) \right].$$

Then,

$$w_e^*(\Phi) = [\eta \beta_e^*(\Phi) \quad (1 - \eta) \beta_e^*(\Phi)],$$

where

$$\beta_e^*(\Phi) = \frac{1}{1 + \bar{\mu}_e^\top \bar{\Sigma}_e^{-1} \bar{\mu}_e} \bar{\Sigma}_e^{-1} \bar{\mu}_e. \quad (30)$$

Here, $\bar{\Sigma}_e$ and $\bar{\mu}_e$ are defined as

$$\begin{aligned} \bar{\Sigma}_e &:= \Phi S \begin{bmatrix} \sigma_c^2 I & 0 \\ 0 & \sigma_e^2 I \end{bmatrix} S^\top \Phi^\top, \\ \bar{\mu}_e &:= \Phi S \begin{bmatrix} \mu_c \\ \mu_e \end{bmatrix}. \end{aligned}$$

Proof of Lemma 7: We have that

$$w_e^*(\Phi) = \mathcal{I}_e(\Phi)^{-1} \mathbf{E}_{X^e, Y^e} \left[\Phi X^e \tilde{Y}^{e\top} \right].$$

First, we have that

$$\begin{aligned} \mathbf{E}_{X^e, Y^e} \left[\Phi X^e \tilde{Y}^{e\top} \right] &= \Phi S \mathbf{E}_{X^e, Y^e} \left[\begin{bmatrix} Z_c \\ Z_e \end{bmatrix} \begin{bmatrix} \mathbb{1}_{\{Y=1\}} \\ \mathbb{1}_{\{Y=-1\}} \end{bmatrix}^\top \right] \\ &= \Phi S \begin{bmatrix} \eta \mu_c & (1 - \eta) \mu_c \\ \eta \mu_e & (1 - \eta) \mu_e \end{bmatrix} \\ &= [\eta \bar{\mu}_e \quad (1 - \eta) \bar{\mu}_e]. \end{aligned}$$

Thus,

$$w_e^*(\Phi) = [\eta \beta_e^*(\Phi) \quad (1 - \eta) \beta_e^*(\Phi)], \quad (31)$$

where

$$\beta_e^*(\Phi) = \mathcal{I}_e(\Phi)^{-1} \bar{\mu}_e. \quad (32)$$

We now compute $\mathcal{I}_e(\Phi)^{-1}$ in terms of $\bar{\Sigma}_e^{-1}$ and $\bar{\mu}_e$. We have that

$$\mathcal{I}_e(\Phi) = \mathbf{E}_{X^e} \left[\Phi X^e X^{e\top} \Phi^\top \right] = \Phi S \mathbf{E} \left[\begin{bmatrix} Z_c \\ Z_e \end{bmatrix} \begin{bmatrix} Z_c \\ Z_e \end{bmatrix}^\top \right] S^\top \Phi^\top.$$

From the definition of Z_c and Z_e , it follows that

$$\mathbf{E} \begin{bmatrix} \begin{bmatrix} Z_c \\ Z_e \end{bmatrix} \begin{bmatrix} Z_c \\ Z_e \end{bmatrix}^\top \end{bmatrix} = \begin{bmatrix} \sigma_c^2 I & 0 \\ 0 & \sigma_e^2 I \end{bmatrix} + \begin{bmatrix} \mu_c \\ \mu_e \end{bmatrix} \begin{bmatrix} \mu_c \\ \mu_e \end{bmatrix}^\top.$$

Thus,

$$\mathcal{I}_e(\Phi) = \bar{\Sigma}_e + \bar{\mu}_e \bar{\mu}_e^\top$$

Using Sherman-Morrison formula, we have that

$$\mathcal{I}_e(\Phi)^{-1} = \bar{\Sigma}_e^{-1} - \frac{\bar{\Sigma}_e^{-1} \bar{\mu}_e \bar{\mu}_e^\top \bar{\Sigma}_e^{-1}}{1 + \bar{\mu}_e^\top \bar{\Sigma}_e^{-1} \bar{\mu}_e}. \quad (33)$$

Finally, using Equation (32) it follows that

$$\beta_e^*(\Phi) = \left(\bar{\Sigma}_e^{-1} - \frac{\bar{\Sigma}_e^{-1} \bar{\mu}_e \bar{\mu}_e^\top \bar{\Sigma}_e^{-1}}{1 + \bar{\mu}_e^\top \bar{\Sigma}_e^{-1} \bar{\mu}_e} \right) \bar{\mu}_e = \frac{1}{1 + \bar{\mu}_e^\top \bar{\Sigma}_e^{-1} \bar{\mu}_e} \bar{\Sigma}_e^{-1} \bar{\mu}_e. \quad \blacksquare$$

Theorem 2. Assume that $|\mathcal{E}_{\text{tr}}| > d_e$. Consider a linear data representation $\Phi X = AZ_c + BZ_e$ and a classifier $w(\Phi)$ on top of Φ that is invariant, i.e., $w(\Phi) = w_e^*(\Phi)$ for all $e \in \mathcal{E}_{\text{tr}}$. If non-degeneracy conditions Eqs. (22-24) hold, then either $w(\Phi) = 0$ or $B = 0$.

The proof of Lemma 2 closely follows from the proof of Rosenfeld et al., 2021, Lemma C.3.. In what follows, we include a proof to keep the manuscript self-contained.

Proof of Theorem 2: First, notice that the decomposition $\varphi(X^e) = AZ_c + BZ_e$ (or $\Phi S = [A \ B]$) is without loss of generality under the assumption of left-invertibility of S . Then,

$$\bar{\Sigma}_e = \sigma_c^2 AA^\top + \sigma_e^2 BB^\top, \quad (34)$$

$$\bar{\mu}_e = A\mu_c + B\mu_e. \quad (35)$$

Recall from Lemma 7 that $\beta_e^*(\Phi) = \bar{\Sigma}_e^{-1} \bar{\mu}_e / (1 + \bar{\mu}_e^\top \bar{\Sigma}_e^{-1} \bar{\mu}_e)$. If $\beta_e^*(\Phi)$ is invariant, then $\beta^* = \beta_e^*(\Phi)$ for all $e \in \mathcal{E}_{\text{tr}}$. Then, by reorganizing terms, we get

$$\bar{\Sigma}_e \beta^* = (1 - \bar{\mu}_e^\top \beta^*) \bar{\mu}_e.$$

Thus, using Equation (34) and (35), we have that

$$(\sigma_c^2 AA^\top + \sigma_e^2 BB^\top) \beta^* = (1 - (A\mu_c + B\mu_e)^\top \beta^*) (A\mu_c + B\mu_e).$$

Let $v := -\sigma^2 AA^\top \beta^* + (1 - \mu_c^\top A \beta^*) A \mu_c$. Then,

$$B (\sigma_e^2 I + \mu_e \mu_e^\top) B^\top \beta^* - v = (1 - \mu_c^\top A^\top \beta^*) B \mu_e + \mu_e^\top B^\top \beta^* A \mu_c. \quad (36)$$

Similar to proof of Lemma C.3. in [Rosenfeld et al., 2021], we show that for all fixed β^* and A Eq. (36) for all environments only holds (with probability 1) if $B = 0$. If $|\mathcal{E}_{\text{tr}}| > d_e$. Then, by the degeneracy assumption of the training sets, there exists at least one environment for which Eq. (22) holds. Let $\tilde{\mu}$ and $\tilde{\sigma}^2$ be the mean of Z_e and variance of W_e for this environment. Then, we have that $\tilde{\mu} = \sum_{i=1}^{d_e} \alpha_i \mu_i$. By applying this linear combination to Eq. (36) for this environment, we get

$$\begin{aligned} B (\tilde{\sigma}^2 I + \tilde{\mu} \tilde{\mu}^\top) B^\top \beta^* - v &= (1 - \mu_c^\top A^\top \beta^*) B \sum_{i=1}^{d_e} \alpha_i \mu_i + \left(\sum_{i=1}^{d_e} \alpha_i \mu_i \right)^\top B^\top \beta^* A \mu_c \\ &= \sum_{i=1}^{d_e} \alpha_i ((1 - \mu_c^\top A^\top \beta^*) B \mu_i + \mu_i^\top B^\top \beta^* A \mu_c) \\ &= \sum_{i=1}^{d_e} \alpha_i (B (\sigma_i^2 I + \mu_i \mu_i^\top) B^\top \beta^* - v), \end{aligned} \quad (37)$$

where in the last identity, we applied Eq. (36) for all $i = 1, \dots, d_e$. By rearranging the terms in Eq. (37), we get

$$B \left(\tilde{\sigma}^2 I + \tilde{\mu} \tilde{\mu}^\top - \sum_{i=1}^{d_e} \alpha_i (\sigma_i^2 I + \mu_i \mu_i^\top) \right) B^\top \beta^* = \left(1 - \sum_{i=1}^{d_e} \alpha_i \right) v. \quad (38)$$

From the non-degeneracy condition (23), Eq. (38) is equivalent to

$$B \Gamma_\alpha B^\top \beta^* = v, \quad (39)$$

where Γ_α is defined as

$$\Gamma_\alpha := \frac{\tilde{\sigma}^2 I + \tilde{\mu} \tilde{\mu}^\top - \sum_{i=1}^{d_e} \alpha_i (\sigma_i^2 I + \mu_i \mu_i^\top)}{1 - \sum_{i=1}^{d_e} \alpha_i}.$$

Note that B , β^* , and v are environment independent and Γ_α is an environment dependent matrix for which it holds that $\text{rank}(\Gamma_\alpha) = d_e$ from the nondegeneracy condition (24). Thus, Eq. (39) holds if and only if $v = B \Gamma_\alpha B^\top \beta^* = 0$. Then, Eq. (36) reduces to

$$(1 - \mu_c^\top A^\top \beta^*) B \mu_e + \beta^{*\top} B \mu_e A \mu_c = 0$$

for all $e \in \mathcal{E}_{\text{tr}}$. Thus, $B \mu_e = 0$ for all $e \in \mathcal{E}_{\text{tr}}$, which holds if and only if $B = 0$ given that the span of μ_e s is \mathbb{R}^{d_e} .

B The Role of the Eigenstructure of $\mathcal{I}_e(\varphi)$

In this section, we elaborate on the discussions on the eigenstructure of $\mathcal{I}_e(\varphi)$, and in particular, its condition number in the examples of [Arjovsky et al., 2019] and [Rosenfeld et al., 2021].

B.1 Example 1 of [Arjovsky et al., 2019]

Arjovsky et al. [2019] consider data that is generated according to the following SEM

$$\begin{aligned} X_1 &\sim \mathcal{N}(0, \sigma^2), \\ Y &= X_1 + Z_1, \\ X_2 &= Y + Z_2, \end{aligned}$$

where $Z_1 \sim \mathcal{N}(0, \sigma^2)$, $Z_2 \sim (0, 1)$, and X_1 are independent, and $X = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix}$. Consider the following data representation

$$\varphi_c(X) = \begin{bmatrix} X_1 \\ cX_2 \end{bmatrix}. \quad (40)$$

Then,

$$\mathcal{I}(\varphi_c) = \mathbf{E} \left[\begin{bmatrix} X_1 \\ cX_2 \end{bmatrix} \begin{bmatrix} X_1 \\ cX_2 \end{bmatrix}^\top \right] = \begin{bmatrix} \sigma^2 & c\sigma^2 \\ c\sigma^2 & c^2(2\sigma^2 + 1) \end{bmatrix}.$$

We now find the eigenvalues of $\mathcal{I}(\varphi_c)$. That is, the solutions to $\det(\mathcal{I}(\varphi_c) - \lambda I) = 0$.

$$\lambda^2 - \lambda(\sigma^2 + c^2(2\sigma^2 + 1)) + c^2\sigma^2(2\sigma^2 + 1) - c^2\sigma^4 = 0.$$

Then,

$$\begin{aligned} \lambda &= \frac{1}{2} \left(\sigma^2 + c^2(2\sigma^2 + 1) \pm \sqrt{(\sigma^2 + c^2(2\sigma^2 + 1))^2 - 4c^2\sigma^2(\sigma^2 + 1)} \right) \\ &= \frac{1}{2} \left(\sigma^2 + c^2(2\sigma^2 + 1) \pm \sqrt{\sigma^4 + c^4(2\sigma^2 + 1)^2 - 2c^2\sigma^2} \right). \end{aligned}$$

Hence,

$$\begin{aligned}
\kappa(\mathcal{I}(\varphi_c)) &= \frac{\lambda_{\max}(\mathcal{I}(\varphi_c))}{\lambda_{\min}(\mathcal{I}(\varphi_c))} \\
&= \frac{\sigma^2 + c^2(2\sigma^2 + 1) + \sqrt{\sigma^4 + c^4(2\sigma^2 + 1)^2 - 2c^2\sigma^2}}{\sigma^2 + c^2(2\sigma^2 + 1) - \sqrt{\sigma^4 + c^4(2\sigma^2 + 1)^2 - 2c^2\sigma^2}} \\
&= \frac{\left(\sigma^2 + c^2(2\sigma^2 + 1) + \sqrt{\sigma^4 + c^4(2\sigma^2 + 1)^2 - 2c^2\sigma^2}\right)^2}{(\sigma^2 + c^2(2\sigma^2 + 1))^2 - (\sigma^4 + c^4(2\sigma^2 + 1)^2 - 2c^2\sigma^2)} \\
&= \frac{1}{2(\sigma^2 + 1)} \left(\frac{1}{c}\sigma^2 + c(2\sigma^2 + 1) + \sqrt{\frac{1}{c^2}\sigma^4 + c^2(2\sigma^2 + 1)^2 - 2\sigma^2} \right)^2.
\end{aligned}$$

Finally,

$$\lim_{c \rightarrow \infty} \kappa(\mathcal{I}(\varphi_c)) = \lim_{c \rightarrow 0} \kappa(\mathcal{I}(\varphi_c)) = \infty.$$

It is worth noting that Arjovsky et al. [2019] discuss that $\|w - w_c^*(\varphi)\|^2$ is poor discrepancy both for the invariant data representation, i.e., $c = 0$, and for a data representation that heavily rely on the spurious features, i.e., large c .

B.2 Example of [Rosenfeld et al., 2021]

Recall that

$$\varphi_\epsilon(X^e) = \begin{bmatrix} Z_c \\ 0 \end{bmatrix} \mathbf{1}_{\{Z_e \in \mathcal{Z}_\epsilon\}} + \begin{bmatrix} Z_c \\ Z_e \end{bmatrix} \mathbf{1}_{\{Z_e \notin \mathcal{Z}_\epsilon\}}.$$

Here, \mathcal{Z}_ϵ is defined as

$$\mathcal{Z}_\epsilon := \bigcup_{e \in \mathcal{E}} (\mathcal{B}_r(\mu_e) \cup \mathcal{B}_r(-\mu_e)), \quad (41)$$

where $r := \sqrt{\epsilon\sigma_e^2 d_e}$ and $\mathcal{B}_r(\mu)$ denotes the $\ell - 2$ ball of radius r centered at μ . Then,

$$\mathcal{I}_e(\varphi_\epsilon) = \mathbf{E} [\varphi_\epsilon(X^e) \varphi_\epsilon(X^e)^\top] = I_c + I_e.$$

where I_c and I_e are defined as

$$\begin{aligned}
I_c &:= \begin{bmatrix} \mathbf{E} [Z_c Z_c^\top | Z_e \in \mathcal{Z}_\epsilon] & 0 \\ 0 & 0 \end{bmatrix} \mathbf{P}(Z_e \in \mathcal{Z}_\epsilon), \\
I_e &:= \mathbf{E} \left[\begin{bmatrix} Z_c \\ Z_e \end{bmatrix} \begin{bmatrix} Z_c \\ Z_e \end{bmatrix}^\top \middle| Z_e \in \mathcal{Z}_\epsilon \right] \mathbf{P}(Z_e \notin \mathcal{Z}_\epsilon).
\end{aligned}$$

Here, we establish a lower bound on the condition number of $\mathcal{I}_e(\varphi_\epsilon)$ in terms of the probability of event $\mathbf{1}_{\{Z_e \notin \mathcal{Z}_\epsilon\}}$. Using Weyl's inequality, we have that

$$\begin{aligned}
\lambda_{\max}(\mathcal{I}_e(\varphi_\epsilon)) &\geq \lambda_{\max}(I_c) + \lambda_{\min}(I_e), \\
\lambda_{\min}(\mathcal{I}_e(\varphi_\epsilon)) &\leq \lambda_{\min}(I_c) + \lambda_{\max}(I_e).
\end{aligned}$$

As I_e is positive semidefinite, $\lambda_{\min}(I_e) \geq 0$. Moreover, $\lambda_{\min}(I_c) = 0$. Then,

$$\begin{aligned}
\lambda_{\max}(\mathcal{I}_e(\varphi_\epsilon)) &\geq \lambda_{\max}(I_c), \\
\lambda_{\min}(\mathcal{I}_e(\varphi_\epsilon)) &\leq \lambda_{\max}(I_e).
\end{aligned}$$

For the first term, we have

$$\begin{aligned}
\lambda_{\max}(I_c) &\geq \frac{1}{d_e + d_c} \text{tr}(I_c) \\
&= \frac{1}{d_e + d_c} \mathbf{E} [\|Z_c\|^2 | Z_e \in \mathcal{Z}_\epsilon] \mathbf{P}(Z_e \in \mathcal{Z}_\epsilon) \\
&= \frac{1}{d_e + d_c} \mathbf{E} [\|\mu_c\|^2 Y^2 + \|W_c\|^2 + 2\mu_c^\top W_c Y | Z_e \in \mathcal{Z}_\epsilon] \mathbf{P}(Z_e \in \mathcal{Z}_\epsilon) \\
&= \frac{1}{d_e + d_c} (\|\mu_c\|^2 + d_c \sigma_c^2) \mathbf{P}(Z_e \in \mathcal{Z}_\epsilon), \quad (42)
\end{aligned}$$

where the last identity follows from the fact that $Y^2 = 1$ almost surely, and the fact that W_c is independent of Y and W_e , and hence is independent of the event $\mathbb{1}_{\{Z_e \notin \mathcal{Z}_\epsilon\}}$.

For the second term, we have

$$\begin{aligned}\lambda_{\max}(I_e) &\leq \text{tr}(I_c) \\ &= \mathbf{E} [\|Z_c\|^2 + \|Z_e\|^2 | Z_e \notin \mathcal{Z}_\epsilon] \mathbf{P}(Z_e \notin \mathcal{Z}_\epsilon) \\ &= (\|\mu_c\|^2 + d_c \sigma_c^2 + \mathbf{E} [\|Z_e\|^2 | Z_e \notin \mathcal{Z}_\epsilon]) \mathbf{P}(Z_e \notin \mathcal{Z}_\epsilon),\end{aligned}$$

where the last identity follows similarly as Equation (42). Then,

$$\begin{aligned}\mathbf{E} [\|Z_e\|^2 | Z_e \notin \mathcal{Z}_\epsilon] &= \mathbf{E} [\|\mu_e\|^2 Y^2 + \|W_e\|^2 + 2\mu_e^\top W_e Y | Z_e \notin \mathcal{Z}_\epsilon] \\ &= \|\mu_e\|^2 + \mathbf{E} [\|W_e\|^2 + 2\mu_e^\top W_e Y | Z_e \notin \mathcal{Z}_\epsilon].\end{aligned}$$

We have that

$$\mathbf{E} [\|W_e\|^2 | Z_e \notin \mathcal{Z}_\epsilon] \leq d_e \sigma_e^2.$$

Moreover, using Cauchy-Schwarz inequality, we get

$$\mathbf{E} [\mu_e^\top W_e Y | Z_e \notin \mathcal{Z}_\epsilon] \leq \|\mu_e\| \mathbf{E} [\|W_e\| | Y | Z_e \notin \mathcal{Z}_\epsilon] \leq \|\mu_e\| \sqrt{d_e \sigma_e^2}.$$

Hence,

$$\lambda_{\max}(I_e) \leq \left(\|\mu_c\|^2 + d_c \sigma_c^2 + (\|\mu_e\| + \sqrt{d_e \sigma_e^2})^2 \right) \mathbf{P}(Z_e \notin \mathcal{Z}_\epsilon).$$

Thus,

$$\begin{aligned}\kappa(\mathcal{I}_e(\varphi_\epsilon)) &\geq \frac{\lambda_{\max}(I_c)}{\lambda_{\max}(I_e)} \\ &\geq \frac{(\|\mu_c\|^2 + d_c \sigma_c^2) \mathbf{P}(Z_e \in \mathcal{Z}_\epsilon) / (d_e + d_c)}{\left(\|\mu_c\|^2 + d_c \sigma_c^2 + (\|\mu_e\| + \sqrt{d_e \sigma_e^2})^2 \right) \mathbf{P}(Z_e \notin \mathcal{Z}_\epsilon)} \\ &= \frac{\|\mu_c\|^2 + d_c \sigma_c^2}{(d_e + d_c) \left(\|\mu_c\|^2 + d_c \sigma_c^2 + (\|\mu_e\| + \sqrt{d_e \sigma_e^2})^2 \right)} \left(\frac{1}{\mathbf{P}(Z_e \notin \mathcal{Z}_\epsilon)} - 1 \right).\end{aligned}$$

Note that [Rosenfeld et al., 2021, Lemma F.3.] show that

$$\mathbf{P}(Z_e \notin \mathcal{Z}_\epsilon) \leq p_{\epsilon, e},$$

where

$$p_{\epsilon, e} := \exp\left(-\frac{1}{8} \min\{\epsilon - 1, (\epsilon - 1)^2\}\right) \quad (43)$$

Then,

$$\kappa(\mathcal{I}_e(\varphi_\epsilon)) \geq \frac{\|\mu_c\|^2 + d_c \sigma_c^2}{(d_e + d_c) \left(\|\mu_c\|^2 + d_c \sigma_c^2 + (\|\mu_e\| + \sqrt{d_e \sigma_e^2})^2 \right)} \left(\frac{1}{p_{\epsilon, e}} - 1 \right).$$

Rosenfeld et al. [2021] show that the invariance penalty of [Arjovsky et al., 2019] is no greater than $O(p_{\epsilon, e}^2)$, which can be made arbitrarily small by choosing appropriately large ϵ . However, for such choices of ϵ , matrix $\mathcal{I}_e(\varphi_\epsilon)$ is ill-conditioned. In particular,

$$\lim_{\epsilon \rightarrow \infty} \kappa(\mathcal{I}_e(\varphi_\epsilon)) = \infty.$$

C Full DomainBed results

C.1 Model selection: training-domain validation set

C.1.1 ColoredMNIST

| Algorithm | +90% | +80% | -90% | Avg |
|-----------|------------|------------|-----------|------|
| ERM | 72.8 ± 0.1 | 72.6 ± 0.2 | 9.8 ± 0.0 | 51.7 |
| IRMv1 | 72.5 ± 0.3 | 72.9 ± 0.1 | 9.9 ± 0.1 | 51.8 |
| IRMv1A | 70.7 ± 0.3 | 72.3 ± 0.5 | 9.7 ± 0.0 | 50.9 |
| IRMv2 | 69.8 ± 0.8 | 72.9 ± 0.3 | 9.8 ± 0.1 | 50.8 |

C.1.2 RotatedMNIST

| Algorithm | 0 | 15 | 30 | 45 | 60 | 75 | Avg |
|-----------|------------|-------------|-------------|-------------|-------------|-------------|------|
| ERM | 93.1 ± 0.1 | 97.8 ± 0.0 | 98.4 ± 0.0 | 98.3 ± 0.1 | 98.2 ± 0.0 | 94.3 ± 0.1 | 96.7 |
| IRMv1 | 89.6 ± 2.1 | 96.8 ± 0.1 | 97.9 ± 0.1 | 97.8 ± 0.1 | 97.5 ± 0.1 | 91.6 ± 0.0 | 95.2 |
| IRMv1A | 75.9 ± 5.9 | 71.1 ± 17.7 | 60.8 ± 25.1 | 60.4 ± 25.8 | 60.2 ± 25.8 | 59.8 ± 20.5 | 64.7 |
| IRMv2 | 94.1 ± 0.0 | 98.2 ± 0.0 | 98.5 ± 0.1 | 98.4 ± 0.1 | 98.3 ± 0.0 | 95.1 ± 0.2 | 97.1 |

C.1.3 PACS

| Algorithm | A | C | P | S | Avg |
|-----------|------------|------------|------------|------------|------|
| ERM | 84.5 ± 1.6 | 77.1 ± 0.8 | 96.9 ± 0.3 | 65.8 ± 1.9 | 81.1 |
| IRMv1 | 77.0 ± 3.0 | 76.7 ± 1.1 | 96.4 ± 0.4 | 64.4 ± 0.3 | 78.6 |
| IRMv1A | 82.6 ± 0.5 | 77.7 ± 0.7 | 96.6 ± 0.4 | 66.7 ± 0.5 | 80.9 |
| IRMv2 | 86.0 ± 0.9 | 76.6 ± 0.7 | 96.9 ± 0.0 | 70.8 ± 2.0 | 82.6 |

C.1.4 VLCS

| Algorithm | C | L | S | V | Avg |
|-----------|------------|------------|------------|------------|------|
| ERM | 97.4 ± 0.1 | 65.0 ± 0.9 | 74.3 ± 1.1 | 78.7 ± 0.1 | 78.8 |
| IRM | 96.3 ± 0.6 | 61.7 ± 0.3 | 70.1 ± 0.1 | 76.0 ± 1.8 | 76.0 |
| IRMA | 96.9 ± 0.8 | 64.8 ± 0.0 | 70.7 ± 1.4 | 77.0 ± 0.4 | 77.3 |
| IRMv2 | 96.6 ± 1.1 | 65.4 ± 1.5 | 73.5 ± 0.5 | 70.6 ± 2.4 | 76.5 |

C.1.5 Averages

| Algorithm | ColoredMNIST | RotatedMNIST | PACS | VLCS | Avg |
|-----------|--------------|--------------|------------|------------|------|
| ERM | 51.7 ± 0.1 | 96.7 ± 0.0 | 81.1 ± 0.1 | 78.8 ± 0.4 | 77.0 |
| IRMv1 | 51.8 ± 0.2 | 95.2 ± 0.4 | 78.6 ± 1.0 | 76.0 ± 0.5 | 75.4 |
| IRMv1A | 50.9 ± 0.1 | 64.7 ± 20.1 | 80.9 ± 0.0 | 77.3 ± 0.2 | 68.4 |
| IRMv2 | 50.8 ± 0.4 | 97.1 ± 0.0 | 82.6 ± 0.9 | 76.5 ± 0.4 | 76.8 |

C.2 Model selection: leave-one-domain-out cross-validation

C.2.1 ColoredMNIST

| Algorithm | +90% | +80% | -90% | Avg |
|-----------|-------------|------------|-------------|------|
| ERM | 30.4 ± 13.4 | 50.5 ± 0.6 | 9.9 ± 0.0 | 30.2 |
| IRMv1 | 50.1 ± 0.4 | 60.6 ± 7.3 | 30.0 ± 14.1 | 46.9 |
| IRMv1A | 69.5 ± 14.4 | 49.8 ± 0.2 | 10.0 ± 0.1 | 43.1 |
| IRMv2 | 10.0 ± 0.1 | 36.4 ± 3.0 | 9.9 ± 0.0 | 18.8 |

C.2.2 RotatedMNIST

| Algorithm | 0 | 15 | 30 | 45 | 60 | 75 | Avg |
|-----------|------------|-------------|-------------|-------------|-------------|-------------|------|
| ERM | 90.3 ± 1.8 | 97.8 ± 0.0 | 98.2 ± 0.1 | 98.2 ± 0.1 | 97.8 ± 0.2 | 93.5 ± 0.4 | 96.0 |
| IRMv1 | 89.6 ± 2.1 | 96.0 ± 0.5 | 97.9 ± 0.1 | 97.2 ± 0.0 | 97.0 ± 0.2 | 90.9 ± 0.5 | 94.8 |
| IRMv1A | 75.9 ± 5.9 | 64.2 ± 22.4 | 59.4 ± 26.1 | 59.8 ± 26.3 | 59.0 ± 26.7 | 55.9 ± 22.5 | 62.4 |
| IRMv2 | 94.1 ± 0.0 | 98.1 ± 0.3 | 98.5 ± 0.1 | 98.2 ± 0.1 | 98.3 ± 0.0 | 94.4 ± 0.3 | 97.0 |

C.2.3 PACS

| Algorithm | A | C | P | S | Avg |
|-----------|------------|------------|------------|-------------|------|
| ERM | 79.7 ± 0.6 | 73.0 ± 3.8 | 97.1 ± 0.9 | 62.1 ± 1.5 | 78.0 |
| IRMv1 | 67.4 ± 6.7 | 72.3 ± 4.3 | 87.7 ± 5.7 | 64.1 ± 4.5 | 72.9 |
| IRMv1A | 78.8 ± 3.2 | 78.9 ± 0.9 | 96.0 ± 0.2 | 42.3 ± 16.3 | 74.0 |
| IRMv2 | 86.3 ± 0.3 | 76.8 ± 0.5 | 97.0 ± 0.4 | 69.7 ± 2.6 | 82.5 |

C.2.4 VLCS

| Algorithm | C | L | S | V | Avg |
|-----------|-------------|------------|------------|------------|------|
| ERM | 97.5 ± 0.8 | 60.3 ± 3.2 | 70.1 ± 4.2 | 75.5 ± 2.8 | 75.9 |
| IRM | 93.3 ± 2.7 | 61.8 ± 0.4 | 72.9 ± 0.9 | 74.1 ± 3.1 | 75.5 |
| IRMA | 79.2 ± 12.8 | 66.8 ± 1.4 | 68.3 ± 4.1 | 73.5 ± 1.3 | 71.9 |
| IRMv2 | 98.2 ± 0.1 | 63.0 ± 0.1 | 74.4 ± 1.2 | 69.9 ± 0.2 | 76.4 |

C.2.5 Averages

| Algorithm | ColoredMNIST | RotatedMNIST | PACS | VLCS | Avg |
|-----------|--------------|--------------|------------|------------|------|
| ERM | 30.2 ± 4.3 | 96.0 ± 0.4 | 78.0 ± 0.9 | 75.9 ± 0.7 | 70.0 |
| IRMv1 | 46.9 ± 2.1 | 94.8 ± 0.2 | 72.9 ± 3.0 | 75.5 ± 1.6 | 72.5 |
| IRMv1A | 43.1 ± 4.8 | 62.4 ± 21.6 | 74.0 ± 5.0 | 71.9 ± 2.8 | 62.8 |
| IRMv2 | 18.8 ± 1.1 | 97.0 ± 0.1 | 82.5 ± 1.0 | 76.4 ± 0.3 | 68.7 |

C.3 Model selection: test-domain validation set (oracle)

C.3.1 ColoredMNIST

| Algorithm | +90% | +80% | -90% | Avg |
|-----------|------------|------------|------------|------|
| ERM | 72.2 ± 0.3 | 72.6 ± 0.4 | 12.4 ± 1.1 | 52.4 |
| IRMv1 | 72.5 ± 0.3 | 72.9 ± 0.1 | 63.3 ± 6.6 | 69.6 |
| IRMv1A | 71.4 ± 0.2 | 72.8 ± 0.3 | 50.2 ± 0.1 | 64.8 |
| IRMv2 | 70.6 ± 1.2 | 71.9 ± 0.6 | 20.9 ± 0.9 | 54.4 |

C.3.2 RotatedMNIST

| Algorithm | 0 | 15 | 30 | 45 | 60 | 75 | Avg |
|-----------|------------|-------------|-------------|-------------|-------------|-------------|------|
| ERM | 92.5 ± 0.6 | 97.8 ± 0.0 | 97.9 ± 0.1 | 97.9 ± 0.2 | 98.3 ± 0.1 | 94.3 ± 0.1 | 96.4 |
| IRMv1 | 89.6 ± 2.1 | 96.8 ± 0.1 | 98.0 ± 0.2 | 97.5 ± 0.3 | 97.5 ± 0.1 | 91.6 ± 0.0 | 95.2 |
| IRMv1A | 77.9 ± 7.4 | 71.1 ± 17.7 | 59.4 ± 26.1 | 59.8 ± 26.3 | 59.0 ± 26.7 | 55.9 ± 22.5 | 63.9 |
| IRMv2 | 94.7 ± 0.4 | 98.0 ± 0.2 | 98.5 ± 0.1 | 98.3 ± 0.0 | 98.3 ± 0.0 | 95.0 ± 0.2 | 97.2 |

C.3.3 PACS

| Algorithm | A | C | P | S | Avg |
|-----------|------------|------------|------------|------------|------|
| ERM | 83.7 ± 0.5 | 82.1 ± 0.2 | 97.5 ± 0.2 | 69.1 ± 0.6 | 83.1 |
| IRMv1 | 66.7 ± 4.3 | 68.5 ± 1.6 | 87.1 ± 5.3 | 67.7 ± 2.0 | 72.5 |
| IRMv1A | 83.5 ± 0.2 | 75.7 ± 3.2 | 96.4 ± 0.3 | 68.6 ± 1.8 | 81.0 |
| IRMv2 | 84.3 ± 0.3 | 76.5 ± 0.7 | 96.8 ± 0.1 | 70.3 ± 2.4 | 82.0 |

C.3.4 VLCS

| Algorithm | C | L | S | V | Avg |
|-----------|----------------|----------------|----------------|----------------|------|
| ERM | 98.4 ± 0.1 | 65.1 ± 1.4 | 72.9 ± 2.2 | 77.1 ± 1.7 | 78.4 |
| IRM | 97.6 ± 1.2 | 61.9 ± 0.6 | 62.9 ± 1.3 | 73.0 ± 0.3 | 73.9 |
| IRMA | 98.0 ± 0.1 | 64.9 ± 1.1 | 71.8 ± 0.8 | 73.9 ± 1.5 | 77.1 |
| IRMv2 | 96.3 ± 1.0 | 67.1 ± 0.1 | 70.9 ± 1.3 | 71.9 ± 1.5 | 76.5 |

C.3.5 Averages

| Algorithm | ColoredMNIST | RotatedMNIST | PACS | VLCS | Avg |
|-----------|----------------|-----------------|----------------|----------------|------|
| ERM | 52.4 ± 0.1 | 96.4 ± 0.1 | 83.1 ± 0.1 | 78.4 ± 0.6 | 77.6 |
| IRMv1 | 69.6 ± 2.3 | 95.2 ± 0.4 | 72.5 ± 2.3 | 73.9 ± 0.2 | 77.8 |
| IRMv1A | 64.8 ± 0.2 | 63.9 ± 21.1 | 81.0 ± 0.4 | 77.1 ± 0.3 | 71.7 |
| IRMv2 | 54.4 ± 0.9 | 97.2 ± 0.1 | 82.0 ± 0.7 | 76.5 ± 1.0 | 77.5 |