# Does Invariant Risk Minimization Capture Invariance?

Pritish Kamath pritish@ttic.edu

Akilesh Tangella akilesh@ttic.edu

Danica J. Sutherland dsuth@cs.ubc.ca

Nathan Srebro nati@ttic.edu

Toyota Technological Institute at Chicago

# Abstract

We show that the Invariant Risk Minimization (IRM) formulation of Arjovsky et al. (2019) can fail to capture "natural" invariances, at least when used in its practical "linear" form, and even on very simple problems which directly follow the motivating examples for IRM. This can lead to worse generalization on new environments, even when compared to unconstrained ERM. The issue stems from a significant gap between the linear variant (as in their concrete method IRMv1) and the full non-linear IRM formulation. Additionally, even when capturing the "right" invariances, we show that it is possible for IRM to learn a sub-optimal predictor, due to the loss function not being invariant across environments. The issues arise even when measuring invariance on the population distributions, but are exacerbated by the fact that IRM is extremely fragile to sampling.

## 1 INTRODUCTION

Machine learning systems tend to seize on spurious correlations present in the training data, and so when presented with out-of-distribution inputs, they can fail spectacularly. For instance, in the spirit of Beery et al. (2018) and Arjovsky et al. (2019), consider a deep neural network trained to classify images as containing a cow or a camel. Suppose that most pictures of cows in the training set are taken in (green) grassy pastures, and those of camels are mostly in (brown) deserts. Then, the neural network is likely to strongly use background color for its predictions – after all, it is a very easy signal to use, and it barely hurts the loss.

Proceedings of the 24<sup>th</sup> International Conference on Artificial Intelligence and Statistics (AISTATS) 2021, San Diego, California, USA. PMLR: Volume 130. Copyright 2021 by the author(s).

Such a network, however, will perform poorly at recognizing cows on a beach. How, then, can we design a machine learning system to identify key features of interest – face, shape, body color, etc., of animals – and ignore spurious ones, like the background color?

Standard machine learning algorithms assume a training set independently sampled from a *single* distribution, and seek good performance only on new samples from the same distribution. There has been much work on models that can adapt to a new distribution given a small number of labeled samples (see e.g. the survey of Redko et al. 2020), or models that are robust to *nearby* distributions (see e.g. the survey of Rahimian and Mehrotra 2019). Ideally though, we would hope for a model that can handle even *large* changes in distribution, *without* the need for labeled target samples.

In reality, our training data usually does *not* actually come from a single homogeneous source: we may have collected it from different users, on different continents, in different years. We thus may be able to tell which correlations are stable across environments (and hence are more likely to be the "true" correlations we seek), and which behave differently in different environments (and are more likely to be spurious).

One approach, then, is to attempt to learn an *invariant predictor* (e.g. Peters et al. 2015; Heinze-Deml et al. 2018; Rojas-Carulla et al. 2018). We might, for instance, assume that for the *causally relevant* subset S of the input variables X, the conditional distribution  $\{Y|X_S\}$  is invariant across data sampled from different environments. This usually requires assuming a meaningful causal graph relating the observed variables. When classifying cows vs. camels based on image pixels, such assumptions are not likely to hold on the input data, though they could potentially apply to the latent variables underlying these images.

The Invariant Risk Minimization (IRM) framework of Arjovsky et al. (2019) tries to find a data representation  $\varphi$  which discards the spurious correlations, leaving only the "real" signal, by enforcing that the predictor w acting on that representation is simultaneously

optimal in each environment given  $\varphi$ . For instance, in the cows-vs-camels problem,  $\varphi$  might remove the background color. Since this gives a challenging bilevel optimization problem, Arjovsky et al. propose a relaxed version, IRMv1, which assumes w is a linear predictor. (We will overview the framework in Section 2.) For a thorough overview of how this approach fits into the literature on out-of-domain generalization, see the discussion by Arjovsky et al. and in particular Appendix A of Gulrajani and Lopez-Paz (2021). Subsequent work has provided new approaches for training in the IRM paradigm (e.g. Ahuja, Shanmugam, et al. 2020; Teney et al. 2020) and applications in domains such as interpretable language processing models (Chang et al. 2020).

Despite much initial promise, however, many key questions remain about the IRM framework: how well does IRMv1 approximate the exact version of the framework in general settings? Do invariant predictors always generalize well on unseen environments? When does a set of training environments allow us to find representations invariant across a broader set of target environments? How does the framework and/or the algorithm behave on finite samples?

**Our Contributions** We advance the understanding of several core questions about the IRM framework.

In Section 3, we study a simple setting of environments over  $\mathcal{X}=\{0,1\}^2$ , abstracting the Colored-MNIST problem studied by Arjovsky et al. (2019). We show that sometimes IRM with linear w can provably fail to find a "truly" invariant predictor, even when solved with respect to the population loss, and even if we provide infinitely many training environments. In fact, it finds a predictor that is even worse on out-of-distribution environments than unrestricted ERM. This issue persists in the IRMv1 implementation.

In Section 4, we note the population loss of even "truly" invariant predictors need not be invariant. We give a simple setting where IRM, which minimizes loss over training environments, prefers an invariant predictor with worse out-of-distribution generalization.

In Section 5, we study when it is possible to identify invariant predictors for a broad class of environments on the basis of a small range of training environments. Although this is generally impossible, we show conditions on the environments under which it is possible.

Finally, in Section 6, we point out issues that arise when using the IRM paradigm over the distributions of empirical samples rather than the population distributions. Here, even invariant predictors (over the population distributions) might not be invariant when considered over the distribution of empirical samples.

# 2 INVARIANT RISK MINIMIZATION

We now describe the IRM paradigm of Arjovsky et al. (2019). We have a set of environments  $\mathcal{E}$ , where each environment  $e \in \mathcal{E}$  corresponds to a distribution  $\mathcal{D}_e$  over  $\mathcal{X} \times \mathcal{Y}$ , with  $\mathcal{X}$  being the space of inputs and  $\mathcal{Y}$  that of outputs. Our goal is to find a predictor  $f: \mathcal{X} \to \widehat{\mathcal{Y}}$ ; we measure the quality of a prediction with a loss function  $\ell: \widehat{\mathcal{Y}} \times \mathcal{Y} \to \mathbb{R}_{\geq 0}$ , and the quality of a predictor by its population loss on environment  $e \in \mathcal{E}$ , given by  $\mathcal{L}_e(f) := \mathbb{E}_{(x,y) \sim \mathcal{D}_e} \ell(f(x), y)$ . In this paper, we mainly focus on the following special case.

Setting A.  $\mathcal{Y} \subseteq \mathbb{R}$ ,  $\widehat{\mathcal{Y}} = \mathbb{R}$ , and  $\ell$  is either the square loss  $\ell_{sq}(\widehat{y}, y) := \frac{1}{2}(\widehat{y} - y)^2$ , or, when  $\mathcal{Y} = \{-1, 1\}$  (corresponding to binary classification), the logistic loss  $\ell_{log}(\widehat{y}, y) := log(1 + exp(-\widehat{y}y))$ .

Given access to samples from some training environments  $\mathcal{E}_{tr} \subseteq \mathcal{E}$ , our aim to learn a predictor f that minimizes the "out-of-distribution" loss over all environments in  $\mathcal{E}$ , namely

$$\mathcal{L}_{\mathcal{E}}(f) := \sup_{e \in \mathcal{E}} \mathcal{L}_{e}(f). \tag{OOD-Gen}$$

#### 2.1 Notions of Invariance

The IRM paradigm attempts to solve this problem by learning an *invariant* representation  $\varphi: \mathcal{X} \to \mathcal{Z}$ . For instance,  $\varphi$  might "throw away" the spurious background color in the cows-vs.-camels example, if  $e_1 \in \mathcal{E}$  is images from Ireland (where most cow images have grassy backgrounds), and  $e_2 \in \mathcal{E}$  is from India (with many more images of cows on city streets). The formal definition of *invariant* is as follows.

**Definition 1** (Definition 3 of Arjovsky et al. 2019). A representation  $\varphi: \mathcal{X} \to \mathcal{Z}$  is invariant over a set of environments  $\mathcal{E}$  if there exists a  $w: \mathcal{Z} \to \widehat{\mathcal{Y}}$  such that w is simultaneously optimal on  $\varphi$  for all environments  $e \in \mathcal{E}$ , that is,  $w \in \operatorname{argmin}_{\overline{w}: \mathcal{Z} \to \widehat{\mathcal{Y}}} \mathcal{L}_e(\overline{w} \circ \varphi)$ .

This definition is motivated by the following observation of Arjovsky et al. (2019), which corresponds more closely to an intuitive definition of invariance.

**Observation 2.** Under Setting A, a representation  $\varphi : \mathcal{X} \to \mathcal{Z}$  is invariant over  $\mathcal{E}$  if and only if for all  $e_1, e_2 \in \mathcal{E}$ , it holds that

$$\mathbb{E}_{\mathcal{D}_{e_1}}[Y \mid \varphi(X) = z] = \mathbb{E}_{\mathcal{D}_{e_2}}[Y \mid \varphi(X) = z]$$

for all  $z \in \mathcal{Z}_{\varphi}^{e_1} \cap \mathcal{Z}_{\varphi}^{e_2}$ , where  $\mathcal{Z}_{\varphi}^{e}$  are the representations from  $\mathcal{D}_e$ ,  $\mathcal{Z}_{\varphi}^{e} := \{ \varphi(X) \mid (X,Y) \in \operatorname{Supp}(\mathcal{D}_e) \}.$ 

 $<sup>^{1}</sup>$ We always assume  $\varphi$  and w are measurable. For further subtleties with Definitions 1 and 3, see Appendix A.1.

We give a proof in Appendix A.2 for completeness.

Crucially, Definition 1 requires that  $\varphi$  and w are unrestricted in the space of all (measurable) functions. However, we wish to learn  $\varphi$  and w with access to only (finite) training sets  $S_e$  sampled from  $\mathcal{D}_e$ , for only a small subset of training environments  $\mathcal{E}_{\rm tr} \subseteq \mathcal{E}$ . For this to be feasible, it is natural to add a restriction that  $\varphi \in \Phi$  and  $w \in \mathcal{W}$ , for suitable classes  $\Phi$  of functions mapping  $\mathcal{X} \to \mathcal{Z}$  and  $\mathcal{W}$  of functions mapping  $\mathcal{Z} \to \widehat{\mathcal{Y}}$ . Any choice of function classes  $(\Phi, \mathcal{W})$  defines a class of "invariant" predictors for a set of environments  $\mathcal{E}$ .

**Definition 3.** For any  $\Phi$ , W and loss function  $\ell$ , the set of invariant predictors on  $\mathcal{E}$ ,  $\mathcal{I}_{\Phi,\mathcal{W}}^{\ell}(\mathcal{E})$ , is the set of all predictors  $f: \mathcal{X} \to \widehat{\mathcal{Y}}$  such that  $\exists (w, \varphi) \in \mathcal{W} \times \Phi$  satisfying the following:

- $f = w \circ \varphi$ , and
- for all  $e \in \mathcal{E}$ ,  $w \in \operatorname{argmin}_{\overline{w} \in \mathcal{W}} \mathcal{L}_e(\overline{w} \circ \varphi)$ .

For ease of notation, we will keep the loss function  $\ell$  implicit. When  $\Phi$  is the space of all functions  $\mathcal{X} \to \mathcal{Z}$ , we denote  $\mathcal{I}_{\Phi,\mathcal{W}}(\mathcal{E})$  as simply  $\mathcal{I}_{\mathcal{W}}(\mathcal{E})$ . Moreover, when  $\mathcal{W}$  is the space of all functions  $\mathcal{Z} \to \widehat{\mathcal{Y}}$ , we denote  $\mathcal{I}_{\mathcal{W}}(\mathcal{E})$  as  $\mathcal{I}(\mathcal{E})$ , leaving the choice of  $\mathcal{Z}$  implicit.<sup>2</sup>

Because exact optimization over  $\mathcal{W}$  is in general difficult, it is useful to consider some special cases. A natural option is linear invariant predictors, where  $\mathcal{Z} = \mathbb{R}^d$  and  $\mathcal{W} = \mathcal{W}^d_{lin}$  is the space of all linear functions on  $\mathbb{R}^d$ . Arjovsky et al. (2019) argued that linear predictors in fact provide no additional representation advantage over scalar invariant predictors, the linear predictors for d=1,  $\mathcal{W}=\mathcal{S}:=\mathcal{W}^1_{lin}$ . In our notation, this translates to the following lemma, proved in Appendix A.2.

**Lemma 4.** Under Setting A, for all  $\mathcal{E}$  and  $d \geq 1$ ,

$$\mathcal{I}(\mathcal{E}) \subseteq \mathcal{I}_{\mathcal{S}}(\mathcal{E}) = \mathcal{I}_{\mathcal{W}_{\text{lin}}^d}(\mathcal{E}).$$

#### 2.2 Algorithms

Armed with a notion of invariance, we still need a way to pick an invariant predictor based on training environments  $\mathcal{E}_{\rm tr} \subseteq \mathcal{E}$ . Arjovsky et al. (2019) proposed the *Invariant Risk Minimization* objective given by

$$\min_{\substack{\varphi:\mathcal{X}\to\mathcal{Z}\\w:\mathcal{Z}\to\widehat{\mathcal{Y}}}} \sum_{e\in\mathcal{E}_{\mathrm{tr}}} \mathcal{L}_e(w\circ\varphi)$$

s.t.  $\forall e \in \mathcal{E}_{tr}, \ w \in \operatorname{argmin}_{\overline{w}: \mathcal{Z} \to \widehat{\mathcal{Y}}} \mathcal{L}_e(\overline{w} \circ \varphi),$ 

which in our notation is equivalent to

$$\min_{f \in \mathcal{I}(\mathcal{E}_{tr})} \sum_{e \in \mathcal{E}_{tr}} \mathcal{L}_e(f). \tag{IRM}$$

We can analogously define  $\mathsf{IRM}_{\mathcal{W}}$  to choose a predictor  $f \in \mathcal{I}_{\mathcal{W}}(\mathcal{E}_{\mathrm{tr}})$ , and  $\mathsf{IRM}_{\Phi,\mathcal{W}}$  from  $f \in \mathcal{I}_{\Phi,\mathcal{W}}(\mathcal{E}_{\mathrm{tr}})$ .

Characterizing  $\mathcal{I}_{\mathcal{W}}(\mathcal{E}_{\mathrm{tr}})$  is difficult in general; fortunately  $\mathcal{I}_{\mathcal{W}_{\mathrm{lin}}^d}(\mathcal{E}) = \mathcal{I}_{\mathcal{S}}(\mathcal{E})$  affords a simple characterization. Any predictor  $f \in \mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\mathrm{tr}})$  can be written as  $f(x) = w_* \, \varphi_*(x)$  for a scalar  $w_*$ . Without loss of generality, we can simply absorb the scalar  $w_*$  into  $\varphi \coloneqq w_* \, \varphi_*$ , so that  $f = 1 \cdot \varphi$ . In Setting A, where the loss function is convex and differentiable,  $f = 1 \cdot \varphi \in \mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\mathrm{tr}}) = \mathcal{I}_{\mathcal{W}_{\mathrm{lin}}^d}(\mathcal{E}_{\mathrm{tr}})$  if and only if

for all 
$$e \in \mathcal{E}_{tr}$$
,  $\nabla_{w|w=1} \mathcal{L}_e(w \cdot \varphi) = 0$ .  $(\nabla_w)$ 

Yet,  $\mathsf{IRM}_{\mathcal{S}}$  remains a bi-level optimization problem. For practical purposes, Arjovsky et al. (2019) proposed to soften this hard constraint, giving the algorithm  $\mathsf{IRMv1}$  to approximate  $\mathsf{IRM}_{\mathcal{S}}$ :

$$\min_{\varphi:\mathcal{X}\to\mathbb{R}} \sum_{e\in\mathcal{E}_{\mathrm{tr}}} \mathcal{L}_{e}(\varphi) + \lambda \left| \nabla_{w|w=1} \mathcal{L}_{e}(w \cdot \varphi) \right|^{2}.$$
 (IRMv1)

A natural baseline is the ERM algorithm, which simply minimizes the loss over training environments:

$$\min_{f:\mathcal{X}\to\widehat{\mathcal{Y}}} \sum_{e\in\mathcal{E}_{tr}} \mathcal{L}_e(f). \tag{ERM}$$

While we referred to IRM, IRMv1 and ERM as "algorithms" above, there still remain two key details that make these impractical as stated: (i) the loss minimized refers to the *population loss*, to which we do not have direct access, and (ii) we are assuming that  $\varphi$  is unrestricted in the space of all functions. Arjovsky et al. (2019) attempt to remedy these issues in IRMv1 by (i) replacing the population loss by the corresponding empirical loss measured over training sets, and (ii) by optimizing  $\varphi$  over a sufficiently expressive parameterized model, such as a deep neural network, using gradient-based local search methods.

Nevertheless, as we discuss shortly,  $\mathsf{IRM}_{\mathcal{S}}$  does not capture  $\mathsf{IRM}$  even when operating on the population loss with unrestricted  $\varphi$ . Unless otherwise stated, we always consider  $\mathsf{IRM}$ ,  $\mathsf{IRM}_{\mathcal{S}}$ ,  $\mathsf{IRMv1}$  and  $\mathsf{ERM}$  as operating over population losses.

#### 2.3 Related Work

Rosenfeld et al. (2021) demonstrate an example where there exists a near-optimal solution to the IRMv1 objective, that nearly matches performance of IRM on training environments, but does no better than ERM on environments that are "far" away from the training distributions. This example relies on environments which barely overlap, allowing the representation to simply "memorize" the training environments. Indeed, Ahuja, Wang, et al. (2021) argue that IRM can have

<sup>&</sup>lt;sup>2</sup>In defining  $\mathcal{I}(\mathcal{E})$ , the choice of  $\mathcal{Z}$  does not matter, as long as  $\mathcal{Z}$  is large enough compared to  $\mathcal{X}$ ; for instance,  $\mathcal{Z} = \mathcal{X}$  is always a valid choice.

an advantage over ERM only when the support of the different environment distributions have a significant overlap. Gulrajani and Lopez-Paz (2021) find empirically that with current models and data augmentation techniques, ERM achieves state-of-the-art practical performance in domain generalization. Nagarajan et al. (2021), meanwhile, theoretically study the behavior of ERM for domain generalization.

Note that in prior work,  $IRM_S/IRMv1$  and IRM are often referred to interchangeably. As we demonstrate,  $IRM_S$  can behave very differently from IRM, even on simple examples that motivated the IRM approach.

# 3 COLORED-MNIST AND TWO-BIT ENVIRONMENTS

To illustrate the utility of the IRM approach and IRMv1 in particular, Arjovsky et al. (2019) introduced the Colored-MNIST problem, a synthetic task derived from MNIST (LeCun et al. 2010). While MNIST images are grayscale, in Colored-MNIST each image is colored either red or green in a way that correlates strongly (but spuriously) with the class label. Here ERM learns to exploit the color, and fails at test time when the direction of correlation with the color is reversed.

To understand the behavior of IRM<sub>S</sub> and IRMv1 on Colored-MNIST, we study an abstract version based on two bits of input, where Y is the binary label to be predicted,  $X_1$  corresponds to the label of the handwritten digit (0-4 or 5-9), and  $X_2$  corresponds to the color (red or green). We represent each environment e with two parameters  $\alpha_e, \beta_e \in [0, 1]$ . The distribution  $\mathcal{D}_e$  is defined as

$$Y \leftarrow \operatorname{Rad}(0.5),$$
  
 $X_1 \leftarrow Y \cdot \operatorname{Rad}(\alpha_e),$  (Two-Bit-Envs)  
 $X_2 \leftarrow Y \cdot \operatorname{Rad}(\beta_e),$ 

where  $\operatorname{Rad}(\delta)$  is a random variable taking value -1 with probability  $\delta$  and +1 with probability  $1-\delta$ . For convenience, we denote an environment e as  $(\alpha_e, \beta_e)$ .

Following the experiments with Colored-MNIST as done by Arjovsky et al. (2019), we consider a set of environments  $\mathcal{E}_{\alpha} := \{(\alpha, \beta_e) : 0 < \beta_e < 1\}$ . It can be shown that there only two predictors in  $\mathcal{I}(\mathcal{E}_{\alpha})$ , one being the trivial 0-predictor, and another that depends only on  $X_1$  (see proof of Proposition 5 for details).

Motivating example of Arjovsky et al. (2019) Consider  $\mathcal{E} = \mathcal{E}_{0.25}$  and  $\mathcal{E}_{tr} = \{(0.25, 0.1), (0.25, 0.2)\}$ . Focusing on the case of  $\ell_{sq}$ , (ERM) on  $\mathcal{E}_{tr}$  learns the predictor  $f_{\text{ERM}}$  that is (approximately) given by

$f_{ERM}$	$X_2 = 1$	$X_2 = -1$	
$X_1 = 1$	0.8889	-0.3077	
$X_1 = -1$	0.3077	-0.8889	

the prediction clearly depends on  $X_2$  as well as  $X_1$ . On each environment in  $\mathcal{E}_{tr}$ , the signal from  $X_2$  is stronger than that from  $X_1$ , and so the binary predictor here can be summarized as  $sign(f_{\mathsf{ERM}}(X)) = sign(X_2)$ . On the other hand, (IRM) chooses the predictor  $f_{\mathsf{IRM}}$ 

$f_{IRM}$	$X_2 = 1$	$X_2 = -1$	
$X_1 = 1$	0.5	0.5	
$X_1 = -1$	-0.5	-0.5	

whose binary behavior is  $sign(f_{\mathsf{IRM}}(X)) = sign(X_1)$ .

On  $e \in \mathcal{E}_{\mathrm{tr}}$ ,  $f_{\mathsf{ERM}}$  achieves a lower loss than  $f_{\mathsf{IRM}}$ , since it is using the more powerful signal  $X_2$ . But, if we evaluate the ability of these predictors to generalize far out of distribution to a case where the (spurious) correlation of  $X_2$  has flipped entirely, e = (0.25, 0.9),  $f_{\mathsf{ERM}}$  will give the wrong (binary) prediction 90% of the time, and get square loss  $\mathcal{L}_e(f_{\mathsf{ERM}}) = 0.985$ . This is far worse than  $f_{\mathsf{IRM}}$ , which at  $\mathcal{L}_e(f_{\mathsf{IRM}}) = 0.375$  has not suffered at all compared to  $\mathcal{E}_{\mathsf{tr}}$ . It is even worse than the trivial 0-predictor,  $\mathcal{L}_e(f_0) = 0.5$ .

It turns out that  $\mathsf{IRM}_{\mathcal{S}}$  also learns the predictor  $f_{\mathsf{IRM}}$  here, demonstrating the utility of this relaxation of  $\mathsf{IRM}$ . This raises a natural question:

Does  $IRM_S$  always learn the same predictor as IRM?

Arjovsky et al. (2019, Section 4.1) considered a specialized linear family of environments, where they proved that indeed IRM<sub>S</sub> learns an invariant predictor, as learned by IRM, for any  $\mathcal{E}_{\rm tr}$  with a sufficient number of environments in "general position." (See also Rosenfeld et al. 2021, Section 5.) It was left to future work whether IRM<sub>S</sub> learns invariant predictors in the sense of IRM more generally as well.

A failure mode of  $\mathsf{IRM}_{\mathcal{S}}$  and  $\mathsf{IRMv1}$  We show that in fact for a simple set of two-bit environments,  $\mathsf{IRM}_{\mathcal{S}}$  finds a predictor worse than that learned by  $\mathsf{IRM}$ , and even worse than the one learned by  $\mathsf{ERM}$ .

This occurs, e.g., for  $\mathcal{E} = \mathcal{E}_{0.1}$  with training environments  $\mathcal{E}_{tr} = \{e_1 = (0.1, 0.2), e_2 = (0.1, 0.25)\}$ . The learned predictors are (approximately) as follows.

$f_{\sf ERM}$	$X_2 = 1$	$X_2 = -1$
$X_1 = 1$	0.9375	0.4464
$X_1 = -1$	-0.4464	-0.9375

<sup>&</sup>lt;sup>3</sup>The problem (Two-Bit-Envs) does not fit the setting of their Theorem 9, because flipping signs cannot be phrased as independent additive noise.

$f_{IRM}$	$X_2 = 1$	$X_2 = -1$
$X_1 = 1$	0.8	0.8
$X_1 = -1$	-0.8	-0.8

$f_{IRM_{\mathcal{S}}}$	$X_2 = 1$	$X_2 = -1$
$X_1 = 1$	0.9557	0.2943
$X_1 = -1$	-0.2943	-0.9557

 $X_1$  is the stronger signal for Y in this  $\mathcal{E}_{\mathrm{tr}}$ , and all of these predictors make the same binary predictions, but with differing amounts of confidence. Extrapolating to the same kind of test environment where the correlation of  $X_2$  has flipped,  $e_{\mathrm{test}} = (0.1, 0.9)$ , we observe the following (approximate) losses:

	$f_{ERM}$	$f_{IRM}$	$f_{IRM_{\mathcal{S}}}$	$f_0$
$\mathcal{L}_{e_1}(\cdot)$	0.15	0.18	0.15	0.5
$\mathcal{L}_{e_2}(\cdot)$	0.16	0.18	0.17	0.5
$\mathcal{L}_{e_{ ext{test}}}(\cdot)$	0.28	0.18	0.38	0.5

The relation between IRM and ERM is as expected: IRM trades slightly worse loss on the training environments for much better extrapolation to the distant environment  $e_{\text{test}} = (0.1, 0.9) \in \mathcal{E}_{0.1}$ . But while IRM<sub>S</sub> also suffers slightly on the training environments, it is even worse than ERM at extrapolation to  $e_{\text{test}}$ ! The invariant feature  $X_1$  is more correlated with Y than the non-invariant feature  $X_2$  in all of the training environments, and yet IRM<sub>S</sub> depends on  $X_2$  even more seriously than ERM does.

Moreover, this is not a carefully-selected pathological example that would go away with more training environments. In fact, IRM<sub>S</sub> chooses the same predictor even if we include any number of additional training environments  $(0.1, \beta_e)$  for  $\beta_e < 0.28$ . Indeed, we show that for these two-bit environments  $\mathcal{E}_{\alpha}$ , any two training environments are sufficient to recover the set of all invariant predictors (proof in Appendix B).

**Proposition 5.** Under Setting A, for all  $\alpha \in (0,1)$  and  $\mathcal{E}_{tr} = \{e_1, e_2\}$  for any two distinct  $e_1, e_2 \in \mathcal{E}_{\alpha}$ ,

(i) 
$$\mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\mathrm{tr}}) = \mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\alpha})$$
 and (ii)  $\mathcal{I}(\mathcal{E}_{\mathrm{tr}}) = \mathcal{I}(\mathcal{E}_{\alpha})$ .

Thus, the issue is not just that we have don't have enough training environments. Rather, as we will now show, what  $\mathsf{IRM}_{\mathcal{S}}$  determines to be an "invariant predictor" is broader than our intuitive sense – or  $\mathsf{IRM}$ 's notion – of what it means to be invariant.

**Predictors in**  $\mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\alpha})$  Recall a predictor  $f = 1 \cdot \varphi$  is in  $\mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\mathrm{tr}})$  if and only if  $\varphi$  satisfies Equation  $(\nabla_w)$ .

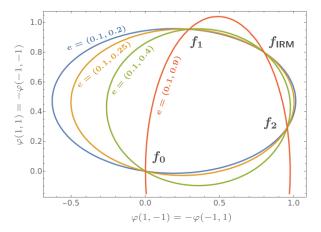


Figure 1: Odd solutions to  $(\nabla_w \text{ for } \ell_{sq})$  for four environments in  $\mathcal{E}_{0,1}$ .

For  $\ell_{sq}$ , this is same as having that for all  $e \in \mathcal{E}_{tr}$ ,

$$\left. \frac{\partial}{\partial w} \left( \mathbb{E}_{(X,Y) \sim \mathcal{D}_e} \left. \frac{(w \cdot \varphi(X) - Y)^2}{2} \right) \right|_{w=1} = 0,$$

or equivalently,

$$\mathbb{E}_{(X,Y)\sim\mathcal{D}_e} (\varphi(X) - Y) \cdot \varphi(X) = 0. \quad (\nabla_w \text{ for } \ell_{\text{sq}})$$

This is a system of quadratic polynomials in four variables  $\{\varphi(x): x \in \{-1,1\}^2\}$ . For ease of visualization, we focus on odd predictors  $f=1\cdot\varphi\in\mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\mathrm{tr}})$ , namely those satisfying f(x)=-f(-x) for all  $x\in\{-1,1\}^2$ . This choice is motivated by the symmetry present in  $\mathcal{D}_e$  and the loss  $\ell_{\mathrm{sq}}$ , along with the observation that the predictors  $f_{\mathrm{ERM}}$ ,  $f_{\mathrm{IRM}}$  and  $f_{\mathrm{IRM}_{\mathcal{S}}}$  are all odd. This allows us to focus on just two variables  $\varphi(1,1)=-\varphi(-1,-1)$  and  $\varphi(1,-1)=-\varphi(-1,1)$ .

Figure 1 shows the solutions of  $(\nabla_w \text{ for } \ell_{\text{sq}})$  among all odd  $\varphi$  for four environments in  $\mathcal{E}_{0.1}$ . There are precisely four odd choices of  $\varphi \in \mathcal{I}_{\mathcal{S}}(\mathcal{E}_{0.1}) = \mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\text{tr}})$ . Two are the expected solutions  $f_0$  and  $f_{\text{IRM}}$  described above; these are the only two predictors in  $\mathcal{I}(\mathcal{E}_{\text{tr}}) = \mathcal{I}(\mathcal{E}_{0.1})$ .  $\mathcal{I}_{\mathcal{S}}(\mathcal{E}_{0.1})$ , however, contains two more odd predictors,  $f_1$  and  $f_2$ , the former being  $f_{\text{IRM}_{\mathcal{S}}}$  from above.  $f_{\text{IRM}_{\mathcal{S}}}$  achieves a smaller loss than the other solutions for the two training environments (0.1, 0.2) and (0.1, 0.25), but higher loss than  $f_{\text{IRM}}$  for environments (0.1, 0.4) or (0.1, 0.9). Figure 2 visualizes the losses of these four odd predictors on environments with varying  $\beta_e$ . Appendix B.1 has more details, including an analysis that explains precisely when these counterexamples arise.

Thus,  $\mathsf{IRM}_{\mathcal{S}}$  can find representations  $\varphi$  which are not invariant in the sense of Definition 1. In particular, for  $\mathcal{E}_{0.1}$  with  $\ell_{\mathrm{sq}}$ ,  $\mathsf{IRM}_{\mathcal{S}}$ 's feasible set of solutions is  $\mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\mathrm{tr}}) \supseteq \mathcal{I}(\mathcal{E}_{\mathrm{tr}})$ , or equivalently  $\mathcal{I}_{\mathcal{S}}(\mathcal{E}_{0.1}) \supseteq \mathcal{I}(\mathcal{E}_{0.1})$ .

<sup>&</sup>lt;sup>4</sup>This analysis was communicated to us by Léon Bottou.

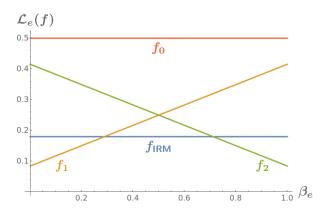


Figure 2: Losses  $\mathcal{L}_e$  (for  $\ell = \ell_{sq}$ ) of odd predictors in  $\mathcal{I}_{\mathcal{S}}(\mathcal{E}_{0.1})$  for various  $e = (0.1, \beta_e)$ .

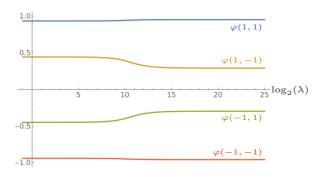


Figure 3: IRMv1 on  $\mathcal{E}_{tr} = \{(0.1, 0.2), (0.1, 0.25)\}$ . The horizontal axis is  $\log_2(\lambda)$ , with -1 representing  $\lambda = 0$ .

As seen from Figure 2,  $f_{\mathsf{IRM}_S} = f_1$  has the lowest loss of those four solutions for  $\beta_e \leq 0.28$ . More training environments will not help  $\mathsf{IRM}_S$  pick  $f_{\mathsf{IRM}}$ , unless the average value of  $\beta_e$  across environments  $e \in \mathcal{E}_{\mathsf{tr}}$  is between 0.29 and 0.71. If the average value of  $\beta_e$  exceeds 0.72,  $\mathsf{IRM}_S$  switches to the other solution  $f_2$ .

We know that IRMv1 becomes exactly ERM when its regularization weight is  $\lambda = 0$ , and IRM<sub>S</sub> for  $\lambda = \infty$ . Figure 3 shows<sup>5</sup> the solution smoothly interpolating between  $f_{\text{ERM}}$  and  $f_{\text{IRM}_S}$ , with the reliance on  $X_2$  increasing as  $\lambda \to \infty$ .

 $\ell_{log}$  loss A similar failure mode occurs for  $\ell_{log}$  on  $\mathcal{E}_{0.05}$  when training on  $\mathcal{E}_{tr} = \{(0.05, 0.1), (0.05, 0.2)\}$ . We give more details in Appendix B.2.

## 3.1 Experiments with Colored-MNIST

We now confirm that the failure mode studied above can also arise in practical training of deep networks based on IRMv1. Colored-MNIST corresponds to the

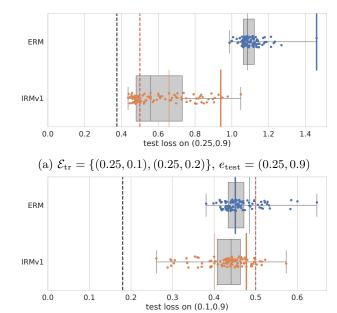


Figure 4: Performance on  $e_{\text{test}}$  when training the given algorithm on  $\mathcal{E}_{\text{tr}}$ , using square loss  $\ell_{\text{sq}}$ , with a fully-connected network. 100 repetitions are shown, using different random hyperparameters and training splits; boxplots show sample quartiles. Black dashed line (left) shows expected loss of the optimal invariant predictor  $f_{\text{IRM}}$ ; red dashed line (right) shows expected loss of the other predictor in  $\mathcal{I}(\mathcal{E}_{\text{tr}})$ , the null predictor  $f_0$ . Shorter, colored vertical lines show the test set performance of the predictor which minimizes the training

objective, (ERM) or (IRMv1) (with  $\lambda = 10^6$ ).

(b)  $\mathcal{E}_{tr} = \{(0.1, 0.2), (0.1, 0.25)\}, e_{test} = (0.1, 0.9)$ 

two-bit environments above, where  $X_1$  is a (grayscale) image from MNIST, and  $X_2$  is a color (red or green) which is assigned to that image.<sup>6</sup> Thus, a learning algorithm which finds global minima of the IRMv1 population-level objective in a model capable of perfectly classifying MNIST digits would behave exactly as described above. In practice, however, we optimize empirical estimates of the risk and gradient penalty, in a model class which may not contain an exactly perfect digit classifier, with an algorithm which may not find the global optimum.

One significant practical issue with IRMv1 is in hyperparameter tuning, since we wish to find models which generalize to environments quite different from  $\mathcal{E}_{\rm tr}$ . Arjovsky et al. (2019) chose hyperparameters arbitrarily for their ERM networks, and for IRMv1 by selecting a network with randomly selected hyperparameters which performed the best on the test set (specifically, the model with the highest minimum accuracy

 $<sup>^5 \</sup>rm{The~IRMv1}$  objective can be non-convex, even for  $\ell_{\rm sq},$  and typical optimization algorithms sometimes find local minima. We instead solved IRMv1 by explicitly enumerating the (odd) stationary points.

<sup>&</sup>lt;sup>6</sup>In practice, we sample the image  $X_1$  first and then flip Y with probability  $\alpha_e$ ; this is equivalent.

on  $\mathcal{E}_{\rm tr} \cup \{e_{\rm test}\}$ ). Since this significantly advantages IRMv1 over ERM, we instead consider the distribution of performances with random hyperparameters from the same proposal distribution as used by Arjovsky et al. We also note which of these models minimized the objective on  $\mathcal{E}_{\rm tr}$  (using a fixed, large  $\lambda$  to compare the objective for IRMv1). Currently, there is no known principled approach for choosing  $\lambda$ ; as noted by Gulrajani and Lopez-Paz (2021), this is often critical to the practical performance of IRM.

Arjovsky et al. (2019) use a fully-connected ReLU network with one hidden layer, operating on the red and green channels of a  $14 \times 14$  image. Running ERM and IRMv1 on this architecture with  $\ell_{\rm sq}$  in the original Colored-MNIST problem shows (Figure 4a) that IRMv1 handily outperforms ERM in test loss, though it does not quite achieve the performance of the best possible  $f_{\rm IRM}$ , and model selection based on  $\mathcal{E}_{\rm tr}$  would choose a predictor notably worse on the test set than the null predictor  $f_0$ . Moving to the example failure mode discussed above, this is no longer the case (Figure 4b): the two algorithms perform about the same in test loss, with model selection on  $\mathcal{E}_{tr}$  selecting a model with performance about the same as  $f_0$  for each algorithm. Although the practical instantiation of IRMv1 clearly suffers here, it is not worse than ERM as we would expect for the population-optimal solutions.

In this representation,  $X_1$  (digit) and  $X_2$  (color) are quite "entangled." In Appendix C, we consider an architecture which processes the grayscale image and total color of the image separately, thus becoming a little closer to the idealized setting (Two-Bit-Envs); here the failure of IRMv1 compared to ERM becomes more apparent. We also explore many variations of the experiment, including experiments with  $\ell_{\log}$ .

Thus,  $IRM_{\mathcal{S}}$ 's surprising failure on the extremely simple problem (Two-Bit-Envs) is essentially reproduced with practical optimization of neural networks on Colored-MNIST.

# 4 CAN IRM FAIL TO CHOOSE THE RIGHT PREDICTOR?

In the previous section, we saw an example where  $\mathsf{IRM}_{\mathcal{S}}$  was able to identify  $\mathcal{I}_{\mathcal{S}}(\mathcal{E})$ , since  $\mathcal{I}_{\mathcal{S}}(\mathcal{E}) = \mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\mathrm{tr}})$  there, but chose a predictor in  $\mathcal{I}_{\mathcal{S}}(\mathcal{E})$  with worse outof-distribution risk for environments "far from"  $\mathcal{E}_{\mathrm{tr}}$ . This happened because the loss  $\mathcal{L}_{e}(f)$  of predictors  $f \in \mathcal{I}_{\mathcal{S}}(\mathcal{E})$  need not be the same (invariant) for all environments  $e \in \mathcal{E}$ , and we pick the "wrong" predictor when optimizing  $\sum_{e \in \mathcal{E}_{\mathrm{tr}}} \mathcal{L}_{e}(f)$  over  $f \in \mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\mathrm{tr}})$ .

Is the same possible for IRM, or does its implicit premise that the optimal invariant predictor on  $\mathcal{E}_{\mathrm{tr}}$ 

will generalize well to  $\mathcal{E}$  hold? IRM can of course fail when  $\mathcal{I}(\mathcal{E}_{\mathrm{tr}}) \supseteq \mathcal{I}(\mathcal{E})$ , when the training environments are not diverse enough to identify the right invariances. But what if we do have  $\mathcal{I}(\mathcal{E}_{\mathrm{tr}}) = \mathcal{I}(\mathcal{E})$ ?

The loss of an invariant predictor  $f \in \mathcal{I}(\mathcal{E})$  need not be invariant for all  $e \in \mathcal{E}$ : consider e.g. varying amounts of inherent additive noise in a regression setting. This would still be acceptable as long as the best invariant predictor with respect to the population loss is the same for all environments  $e \in \mathcal{E}$ . Contrarily, we now give a simple family of environments  $\mathcal{E}$ , training environments  $\mathcal{E}_{tr} \subseteq \mathcal{E}$  satisfying  $\mathcal{I}(\mathcal{E}_{tr}) = \mathcal{I}(\mathcal{E})$ , and two predictors  $f_1, f_2 \in \mathcal{I}(\mathcal{E})$  such that  $\mathcal{L}_e(f_1) > \mathcal{L}_e(f_2)$  for all  $e \in \mathcal{E}_{tr}$ , but  $\mathcal{L}_{\mathcal{E}}(f_1) < \mathcal{L}_{\mathcal{E}}(f_2)$ . Hence IRM prefers  $f_2$ to  $f_1$  based on  $\mathcal{E}_{\mathrm{tr}}$ , but  $f_1$  has better worst-case loss. It is thus generally difficult to handle out-of-distribution prediction in environments with more than one invariant predictor: the invariant predictor which is best on training environments might still perform poorly on unseen test environments, despite being invariant.

Consider environments  $\mathcal{E}$  over  $\mathcal{X} = \{-1, 0, 1\}^3$  and  $\mathcal{Y} = \{-1, 1\}$ , where each environment e is specified by a single parameter  $\theta_e \in (-1/6, 1/3)$  as follows:

$$X_1 \leftarrow \begin{cases} -1 & \text{w.p. } \frac{1}{3} \\ 0 & \text{w.p. } \frac{1}{3} \end{cases}, X_2 \leftarrow \begin{cases} -1 & \text{w.p. } \frac{1}{3} - \theta_e \\ 0 & \text{w.p. } \frac{1}{3} + 2\theta_e \end{cases}, \\ +1 & \text{w.p. } \frac{1}{3} - \theta_e \end{cases}$$

$$\mathbb{E}_{\mathcal{D}_e}[Y|X_1, X_2] = 0.3(X_1 + X_2) + g_{\theta_e}(X_1, X_2) ,$$

where  $g_{\theta_e}(x_1, x_2)$  is given as

$g_{\theta}(x_1, x_2)$	$x_2 = -1$	$x_2 = 0$	$x_2 = +1$
$x_1 = -1$	$\theta(\theta + \frac{2}{3})$	$-\theta(\frac{2}{3}-2\theta)$	$3\theta^2$
$x_1 = 0$	$-\theta(\frac{2}{3}-2\theta)$	0	$\theta(\frac{2}{3}-2\theta)$
$x_1 = +1$	$-3\theta^2$	$\theta(\frac{2}{3}-2\theta)$	$-\theta(\theta+\frac{2}{3})$

While the specific form of  $g_{\theta}$  is a little involved, the main thing to note is that

$$\mathbb{E}_{\mathcal{D}_e}[g_{\theta_e}(X_1, X_2) \mid X_1] = 0 = \mathbb{E}_{\mathcal{D}_e}[g_{\theta_e}(X_1, X_2) \mid X_2]$$

which means that  $\mathbb{E}_{\mathcal{D}_e}[Y|X_1] = 0.3X_1$  as well as  $\mathbb{E}_{\mathcal{D}_e}[Y|X_2] = 0.3X_2$ . Thus for  $\ell_{\text{sq}}$ ,  $\mathcal{I}(\mathcal{E})$  contains the predictors  $f_1(x) = 0.3x_1$  and  $f_2(x) = 0.3x_2$ . In fact, as shown in Appendix D, IRM will indeed pick among these predictors in  $\mathcal{I}(\mathcal{E})$  for almost all  $\mathcal{E}_{\text{tr}}$  containing at least two distinct environments:

**Proposition 6.** In Setting A, for  $\mathcal{E}$  as above, it holds for Lebesgue-almost all  $\mathcal{E}_{tr} \subseteq \mathcal{E}$  with  $|\mathcal{E}_{tr}| \geq 2$  that  $\mathcal{I}(\mathcal{E}) = \mathcal{I}(\mathcal{E}_{tr})$ . Moreover, any  $f \in \mathcal{I}(\mathcal{E})$  depends on at most one of  $x_1$  or  $x_2$ .

Focusing on the case of  $\ell_{sq}$ , the loss of the predictors

can be seen to be<sup>7</sup>, for any  $e \in \mathcal{E}$ ,

$$\mathcal{L}_e(f_1) = 0.47$$
 and  $\mathcal{L}_e(f_2) = 0.47 + 0.09 \cdot \theta_e$ 

Thus, if  $\mathcal{E}_{\mathrm{tr}}$  only contains environments e corresponding to  $\theta_e < 0$ , we will have that  $\mathcal{L}_e(f_2) < \mathcal{L}_e(f_1)$  for all  $e \in \mathcal{E}_{\mathrm{tr}}$ , and yet the invariant predictor that minimizes  $\sup_{e \in \mathcal{E}} \mathcal{L}_e(\cdot)$  is  $f_1$ . See Figure 12 (Appendix D) for an illustration of these loss as a function of  $\theta_e$ .

IRM's notion of invariance ensures  $\mathbb{E}_{\mathcal{D}_e}[Y \mid \varphi(X) =$ z] is invariant across  $\mathcal{E}$ , but allows the loss of the corresponding predictor  $\mathcal{L}_e(f)$  to differ across  $e \in$  $\mathcal{E}$ . Here, in fact the full conditional distribution  $\{Y \mid \varphi(X) = z\}$  is also invariant across  $\mathcal{E}$ , but even so, the loss varies. If we enforced a stronger notion of invariance which requires the entire joint distribution  $\{(Y,\varphi(X))\}_{(X,Y)\sim\mathcal{D}_e}$  to be invariant across all  $e\in\mathcal{E}$ , we would not have faced this issue, since  $\mathcal{L}_e$  would then be invariant, and indeed would pick  $f_1$  in the example above. Yet this joint invariance is clearly too strict for some problems: it is impossible to achieve if the marginal distribution of Y differs across environments, and it is easy to construct other  $\mathcal{E}$  where IRM allows the intuitively-correct predictor but joint invariance allows only a trivial constant predictor.

Thus, IRM is not always guaranteed to achieve optimal out-of-distribution loss, even when all the right invariances are captured by the training environments. The "right" notion of invariance really depends on what we know about the set of all environments  $\mathcal{E}$ .

# 5 WHEN DOES INVARIANCE GENERALIZE?

In the examples of Sections 3 and 4, it held that IRM or IRM<sub>S</sub> were able to identify predictors invariant over all, even unseen, environments: specifically,  $\mathcal{I}_{\mathcal{W}}(\mathcal{E}) = \mathcal{I}_{\mathcal{W}}(\mathcal{E}_{\mathrm{tr}})$ . That this holds is an implicit premise of the IRM framework. Yet it is unclear in general when invariances discovered on training environments will generalize to unseen environments. We now give some partial answers to this question.

For an arbitrary  $\mathcal{E}$ , we of course cannot expect invariances observed across  $\mathcal{E}_{tr}$  to generalize over  $\mathcal{E}$ : simply consider adding a single entirely "irrelevant" e to  $\mathcal{E}$ . To provide some structure, we consider parameterized sets of environments  $\mathcal{E}$ . For simplicity, we focus on finite  $\mathcal{X}$  and  $\mathcal{Y}$ , with  $\mathcal{Y} \subseteq \mathbb{R}$ . Let  $\Delta_{\mathcal{X} \times \mathcal{Y}}$  denote the space of all probability distributions over  $\mathcal{X} \times \mathcal{Y}$ , and let  $\Theta \subseteq \mathbb{R}^d$ . A map  $\Pi : \Theta \to \Delta_{\mathcal{X} \times \mathcal{Y}}$  naturally defines a set of environments  $\mathcal{E}_{\Pi}$  corresponding to the set of distributions  $\{\Pi(\theta) : \theta \in \Theta\}$ . For example, the two-bit

environments  $\mathcal{E}_{\alpha}$  of Section 3 are parameterized by the map  $\Pi: \theta \mapsto e = (\alpha, \theta)$ , for  $\theta \in \Theta = (0, 1)$ .

For 
$$\Theta_{\rm tr} \subseteq \Theta$$
 and  $\mathcal{E}_{\rm tr} = \{\Pi(\theta) \mid \theta \in \Theta_{\rm tr}\},$   
when does it hold that  $\mathcal{I}(\mathcal{E}_{\rm tr}) = \mathcal{I}(\mathcal{E}_{\Pi})$ ?

Note that  $\mathcal{I}(\mathcal{E}_{\Pi}) \subseteq \mathcal{I}(\mathcal{E}_{\mathrm{tr}})$  always holds, but for any hope of  $\mathcal{I}(\mathcal{E}_{\mathrm{tr}}) \subseteq \mathcal{I}(\mathcal{E}_{\Pi})$ , we must assume  $\mathcal{E}_{\mathrm{tr}}$  contains a "representative set" of environments from  $\mathcal{E}_{\Pi}$ .

The most basic assumption to begin with is simply that  $\Pi$  is continuous. This is insufficient to guarantee invariance, even for very large  $\Theta_{\rm tr}$ : the map might simply "change directions" outside of  $\mathcal{E}_{\rm tr}$ . We give a simple example below (proof in Appendix E), where even an uncountable number of environments in  $\mathcal{E}_{\rm tr}$  do not allow us to understand the full behavior of  $\mathcal{E}$ .

**Proposition 7.** There exists a continuous map  $\Pi$ :  $(0,1) \to \Delta_{\mathcal{X} \times \mathcal{Y}}$  such that for  $\Theta_{tr} = (0,\frac{1}{4})$  and  $\mathcal{E}_{tr} = \Pi(\Theta_{tr})$ , it holds that  $\mathcal{I}(\mathcal{E}_{tr}) \neq \mathcal{I}(\mathcal{E}_{\Pi})$ .

On the other hand, if  $\Pi$  is not only continuous but also analytic, we *can* guarantee, under some conditions, that invariances over  $\mathcal{E}_{tr}$  continue to hold over all of  $\mathcal{E}_{\Pi}$ . Let  $\Pi_{(x,y)}(\theta) := \Pr_{(X,Y) \sim \Pi(\theta)}[X = x, Y = y]$  for each  $(x,y) \in \mathcal{X} \times \mathcal{Y}$ . We say the map  $\Pi : \Theta \to \Delta_{\mathcal{X} \times \mathcal{Y}}$  is analytic if, for each  $(x,y) \in \mathcal{X} \times \mathcal{Y}$ ,  $\Pi_{(x,y)} : \Theta \to [0,1]$  is analytic in  $\theta$ .

**Proposition 8.** Let  $\Theta_{tr} \subseteq \Theta \subseteq \mathbb{R}^d$ , where  $\Theta$  is a connected, open set. Suppose  $\Pi : \Theta \to \Delta_{\mathcal{X} \times \mathcal{Y}}$  is analytic,  $\mathcal{X}$  and  $\mathcal{Y}$  are finite and  $\mathcal{E}_{tr} = \Pi(\Theta_{tr})$ . Then, under Setting A,

- (i) For almost all  $\Theta_{\rm tr}$  with  $|\Theta_{\rm tr}| \geq 2$ :  $\mathcal{I}(\mathcal{E}_{\rm tr}) = \mathcal{I}(\mathcal{E}_{\Pi})$ .
- (ii) For all  $\Theta_{\rm tr}$  with non-zero Lebesgue measure:  $\mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\rm tr}) = \mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\Pi}).$

The key step is that when  $\Pi$  is analytic, the conditional expectations  $\mathbb{E}_{\Pi(\theta)}[Y \mid \varphi(X) = z]$  and the gradient  $\nabla_{w \mid w=1} \mathcal{L}_{\Pi(\theta)}(w \cdot \varphi)$  are analytic functions in  $\theta$ ; the result is far stronger, however, for  $\mathcal{I}$  (where the set of representations is finite) than for  $\mathcal{I}_{\mathcal{S}}$ , where our analysis requires uncountably many training environments. A version of Proposition 8 holds even for infinite spaces  $\mathcal{X}$  and  $\mathcal{Y} \subseteq \mathbb{R}$ , under a technical definition of analyticity of  $\Pi$  (details in Appendix E), although in this case our result for  $\mathcal{I}$  also requires  $\mathcal{E}_{\mathrm{tr}}$  to have positive measure.

Recall that the examples studied in Sections 3 and 4 indeed had analytic parameterizations, and hence Proposition 8 implies that  $\mathcal{I}(\mathcal{E}_{\mathrm{tr}}) = \mathcal{I}(\mathcal{E})$  holds for (almost) all  $\mathcal{E}_{\mathrm{tr}}$  with at least two distinct environments.

<sup>&</sup>lt;sup>7</sup>This calculation does not need the specific form of  $g_{\theta}$ .

#### 6 IRM WITH FINITE SAMPLES

Except for Section 3.1, we have so far only discussed algorithms (IRM, IRM $_{\mathcal{S}}$ , and IRMv1) defined in terms of the *population* losses of training environments. In practice, however, we need to work with a finite number of samples from each training environment. If we directly apply IRM or IRM $_{\mathcal{S}}$  as stated in (IRM) to empirical distributions, all correlations will have a small amount of noise, and it is extremely likely that the set of invariant predictors becomes empty.

On the other hand, IRMv1 for a fixed  $\lambda$  could be robust to sampling. We illustrate this in the two-bit environments of Section 3. Consider training environments  $\mathcal{E}_{\rm tr} = \{(0.25, 0.1), (0.25, 0.2), (0.25, 0.3)\}$ : both IRM and IRM<sub>S</sub> are able to learn an invariant predictor. However, when sampling finite datasets, we only have that the empirical distribution of the two environments will be close to – but not exactly the same as – the true distribution; there may not be any exactly-invariant predictors. We illustrate this by evaluating IRM<sub>S</sub> on a set of training environments  $\mathcal{E}'_{\rm tr} = \{(0.245, 0.105), (0.255, 0.195), (0.251, 0.302)\}$ , as a proxy for empirical distributions we see from finite samples. IRM<sub>S</sub> learns the trivial 0 predictor  $f_0$ ; Figure 5 shows the behavior of IRMv1 for increasing  $\lambda$ .

For a fixed empirical distribution, it is likely that as  $\lambda \to \infty$ , IRMv1 approaches IRM<sub>S</sub>, and does not find a good invariant predictor. If we instead take  $n \to \infty$  for a fixed  $\lambda$ , though, we should approach the population version of IRMv1, and hence taking  $\lambda \to \infty$  at an appropriate rate as  $n \to \infty$  may approach the population IRM<sub>S</sub> predictor. Ahuja, Wang, et al. (2021) recently considered a variant of IRM<sub>S</sub> where the constraints  $(\nabla_w)$  defining  $\mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\mathrm{tr}})$  need to hold  $\varepsilon$ -approximately. When training on the objective with finite samples, they bounds the sample complexity to get an out-of-distribution loss close to that of the corresponding population version of this  $\varepsilon$ -IRM<sub>S</sub>.

Given the discrepancy between  $\mathsf{IRM}_{\mathcal{S}}$  and  $\mathsf{IRM}$  as pointed out in Section 3, however, it is important to make  $\mathsf{IRM}$  itself more robust to finite samples. For instance, one possible approach would be to relax the requirement of  $w \in \operatorname{argmin}_{\overline{w}: \mathcal{Z} \to \widehat{\mathcal{V}}} \mathcal{L}_e(\overline{w} \circ \varphi)$  to

$$\mathcal{L}_e(w \circ \varphi) \leq \min_{\overline{w}: \mathcal{Z} \to \widehat{\mathcal{Y}}} \mathcal{L}_e(\overline{w} \circ \varphi) + \varepsilon$$

for a suitable  $\varepsilon > 0$ . How to practically implement a version of this  $\varepsilon$ -IRM remains an open challenge.

## 7 DISCUSSION

The IRM framework of Arjovsky et al. (2019) proposes a promising new paradigm of learning, which attempts

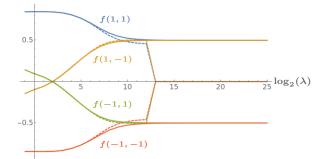


Figure 5: IRMv1 algorithm on exact environments  $\mathcal{E}_{\rm tr}$  (solid lines), and a noisy set  $\mathcal{E}'_{\rm tr}$  (dashed lines; definitions in text). The horizontal axis is  $\log_2(\lambda)$ , with -1 for  $\lambda=0$ . Results are similar for small  $\lambda$ , until the noisy set abruptly gives the 0-predictor.

to exploit information we usually ignore to find models robust to even some quite dramatic changes in the input distribution. We have helped shed light on the applicability of this framework.

We now know that  $\mathsf{IRM}_{\mathcal{S}}$  and  $\mathsf{IRMv1}$  can be surprisingly different from  $\mathsf{IRM}$ , even on very simple environments. This emphasizes the importance of finding practical algorithms to approximate  $\mathsf{IRM}_{\mathcal{W}}$  for some nonlinear class of functions  $\mathcal{W}$ .

We also know that even for IRM, choosing among invariant predictors can also be vital for out-of-domain generalization, and there exist cases where these algorithms choose the wrong one for out-of-distribution robustness. This holds even if we insist on a stronger notion of invariance, namely that of the conditional distribution  $\{Y \mid \varphi(X)\}_{(X,Y) \sim \mathcal{D}_e}$ . To truly handle worst-case out-of-distribution generalization, a stronger notion is needed: for example, it suffices to require invariance of the joint distribution  $\{(Y,\varphi(X))\}_{(X,Y) \sim \mathcal{D}_e}$ , but this seems overly stringent.

We also now know more about the possibility of generalizing invariances learned from  $\mathcal{E}_{\rm tr}$  to a larger set of environments  $\mathcal{E}$ . With significant structure on  $\mathcal{E}$ , it is possible to ensure  $\mathcal{I}(\mathcal{E}) = \mathcal{I}(\mathcal{E}_{\rm tr})$ , but substantial questions remain as to the situation for  $\mathcal{I}_{\mathcal{S}}$  or more realistic assumptions on  $\mathcal{E}$ .

Finally, we demonstrated that IRM and even IRMv1 can be surprisingly brittle when run on samples, rather than populations. Thus more analysis, and perhaps new algorithms, are needed to realize the promise of this framework in practice.

#### Acknowledgments

The authors would like to thank Léon Bottou, Martin Arjovksy, Ishaan Gulrajani, and David Lopez-Paz for useful discussions, particularly the derivation of the form of predictors in Appendix B.1.1.

Work was supported in part by NSF BIGDATA award 1546500 and NSF RI award 1764032. Work done while the authors participated in a special quarter on the Theory of Deep Learning sponsored by NSF TRIPOD award 1934843 (IDEAL) and while the first author participated in the Theory of Reinforcement Learning program at the Simons Institute for the Theory of Computing.

#### References

- Ahuja, Kartik, Karthikeyan Shanmugam, Kush R. Varshney, and Amit Dhurandhar (2020). "Invariant Risk Minimization Games." *International Conference on Machine Learning*. arXiv: 2002.04692.
- Ahuja, Kartik, Jun Wang, Amit Dhurandhar, Karthikeyan Shanmugam, and Kush R. Varshney (2021). "Empirical or Invariant Risk Minimization? A Sample Complexity Perspective." International Conference on Learning Representations. arXiv: 2010.16412.
- Arjovsky, Martin, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz (2019). *Invariant Risk Minimization*. arXiv: 1907.02893.
- Beery, Sara, Grant Van Horn, and Pietro Perona (2018). "Recognition in Terra Incognita." 15th European Conference on Computer Vision. DOI: 10. 1007/978-3-030-01270-0\_28.
- Chang, Shiyu, Yang Zhang, Mo Yu, and Tommi S. Jaakkola (2020). "Invariant Rationalization." *International Conference on Machine Learning*. arXiv: 2003.09772.
- Gulrajani, Ishaan and David Lopez-Paz (2021). "In Search of Lost Domain Generalization." *International Conference on Learning Representations*. arXiv: 2007.01434.
- Heinze-Deml, Christina, Jonas Peters, and Nicolai Meinshausen (2018). "Invariant Causal Prediction for Nonlinear Models." *Journal of Causal Inference* 6.2, p. 20170016. DOI: 10.1515/jci-2017-0016.
- LeCun, Yann, Corinna Cortes, and CJ Burges (2010). "MNIST handwritten digit database." *ATT Labs* [Online] 2. URL: http://yann.lecun.com/exdb/mnist.
- Mityagin, Boris (2015). The Zero Set of a Real Analytic Function. arXiv: 1512.07276.
- Nagarajan, Vaishnavh, Anders Andreassen, and Behnam Neyshabur (2021). "Understanding the failure modes of out-of-distribution generalization."

- International Conference on Learning Representations. arXiv: 2010.15775.
- Peters, Jonas, Peter Bühlmann, and Nicolai Meinshausen (2015). "Causal inference using invariant prediction: identification and confidence intervals." *Journal of the Royal Statistical Society, Series B* 78.5, pp. 947–1012. DOI: 10.1111/rssb.12167. arXiv: 1501.01332.
- Planet Math (Mar. 22, 2013). Differentiation under the Integral Sign. URL: https://planetmath.org/differentiationundertheintegralsign.
- Rahimian, Hamed and Sanjay Mehrotra (2019). Distributionally Robust Optimization: A Review. arXiv: 1908.05659.
- Redko, Ievgen, Emilie Morvant, Amaury Habrard, Marc Sebban, and Younès Bennani (2020). A survey on domain adaptation theory: learning bounds and theoretical quarantees. arXiv: 2004.11829.
- Rojas-Carulla, Mateo, Bernhard Schölkopf, Richard Turner, and Jonas Peters (2018). "Invariant Models for Causal Transfer Learning." *Journal of Machine Learning Research* 19.36, pp. 1–34. arXiv: 1507.05333.
- Rosenfeld, Elan, Pradeep Kumar Ravikumar, and Andrej Risteski (2021). "The Risks of Invariant Risk Minimization." *International Conference on Learning Representations*. arXiv: 2010.05761.
- Teney, Damien, Ehsan Abbasnejad, and Anton van den Hengel (2020). *Unshuffling Data for Improved Generalization*. arXiv: 2002.11894.

## A More details on Section 2

### A.1 Subtleties involving Definitions 1 and 3

Definition 1 implicitly assumes that a minimizer  $w \in \operatorname{argmin}_{\overline{w}: \mathcal{Z} \to \widehat{\mathcal{Y}}} \mathcal{L}_e(\overline{w} \circ \varphi)$  exists. This may not always be the case: for example, if we take logistic loss, there will be no exact maximizer if the problem under  $\varphi$  is separable, i.e.  $\{Y \mid \varphi(X) = z\}$  is constant for each  $z \in \mathcal{Z}_{\varphi}^e$ . To handle such cases, Definition 1 can be modified as follows.

**Definition 9.** A representation  $\varphi : \mathcal{X} \to \mathcal{Z}$  is invariant for a set of environments  $\mathcal{E}$  if for all  $\varepsilon > 0$ , there exists  $a \ w : \mathcal{Z} \to \widehat{\mathcal{Y}}$  such that w is simultaneously  $\varepsilon$ -optimal on  $\varphi$  for all environments  $e \in \mathcal{E}$ : that is, we have that  $\mathcal{L}_e(w \circ \varphi) \leq \inf_{\overline{w} : \mathcal{Z} \to \widehat{\mathcal{Y}}} \mathcal{L}_e(\overline{w} \circ \varphi) + \varepsilon$ .

A related problem arises in Definition 3, where  $\mathcal{L}_e(\overline{w} \circ \varphi)$  may not have a minimizer inside  $\mathcal{W}$ . In addition to the case where the data is separable (and hence we would want w to take values  $\pm \infty$ ), a similar problem can arise even for square loss if  $\mathcal{W}$  contains points arbitrarily close to the conditional expectation function but not the conditional expectation function itself; this can happen, for instance, if  $\mathcal{W}$  is a Gaussian RKHS and the conditional expectation is  $L_2$ -integrable but not in the RKHS. To work around this problem, we can allow w to lie in an appropriate "closure" of  $\mathcal{W}$ .

#### A.2 Proofs of Observation 2 and Lemma 4

The following observation was made by Arjovsky et al. (2019). We include a proof, for completeness and clarity. **Observation 2.** Under Setting A, a representation  $\varphi : \mathcal{X} \to \mathcal{Z}$  is invariant over  $\mathcal{E}$  if and only if for all  $e_1, e_2 \in \mathcal{E}$ , it holds that

$$\mathbb{E}_{\mathcal{D}_{e_1}}[Y \mid \varphi(X) = z] = \mathbb{E}_{\mathcal{D}_{e_2}}[Y \mid \varphi(X) = z]$$

for all  $z \in \mathcal{Z}_{\varphi}^{e_1} \cap \mathcal{Z}_{\varphi}^{e_2}$ , where  $\mathcal{Z}_{\varphi}^{e}$  are the representations from  $\mathcal{D}_e$ ,  $\mathcal{Z}_{\varphi}^{e} := \{ \varphi(X) \mid (X,Y) \in \operatorname{Supp}(\mathcal{D}_e) \}.$ 

*Proof.* Suppose the representation  $\varphi: \mathcal{X} \to \mathcal{Z}$  is invariant for  $\mathcal{E}$ . That is, there exists a predictor  $w: \mathcal{Z} \to \mathbb{R}$  such that  $w \in \operatorname{argmin}_{\overline{w}: \mathcal{Z} \to \mathbb{R}} \mathcal{L}_e(\overline{w} \circ \varphi)$  simultaneously for all environments  $e \in \mathcal{E}$ . In other words, for all  $e \in \mathcal{E}$  and  $z \in \mathcal{Z}_{\varphi}^e$ , it holds that  $w(z) \in \operatorname{argmin}_{\omega \in \mathbb{R}} \mathbb{E}_{\mathcal{D}_e} [\ell(\omega, Y) \mid \varphi(X) = z]$ .

First, consider the case of  $\ell = \ell_{\text{sq}}$ . It follows that  $w(z) = \mathbb{E}_{\mathcal{D}_e}\left[Y \mid \varphi(X) = z\right]$  for all  $e \in \mathcal{E}$  and  $z \in \mathcal{Z}_{\varphi}^e$ . In particular, it holds for all  $e_1, e_2 \in \mathcal{E}$  and  $z \in \mathcal{Z}_{\varphi}^{e_1} \cap \mathcal{Z}_{\varphi}^{e_2}$  that  $\mathbb{E}_{\mathcal{D}_{e_1}}\left[Y \mid \varphi(X) = z\right] = \mathbb{E}_{\mathcal{D}_{e_2}}\left[Y \mid \varphi(X) = z\right] = w(z)$ .

Conversely, suppose that  $\varphi$  is such that  $\mathbb{E}_{\mathcal{D}_{e_1}}[Y \mid \varphi(X) = z] = \mathbb{E}_{\mathcal{D}_{e_2}}[Y \mid \varphi(X) = z]$  for all  $z \in \mathcal{Z}_{\varphi}^{e_1} \cap \mathcal{Z}_{\varphi}^{e_2}$  and all  $e_1, e_2 \in \mathcal{E}$ . Then,  $w(z) := \mathbb{E}_{\mathcal{D}_e}[Y \mid \varphi(X) = z]$  for any e such that  $z \in \mathcal{Z}_{\varphi}^e$  is well-defined and gives a predictor w that is simultaneously optimal for all environments.

The case of  $\ell_{\log}$  is handled similarly by noting that the minimizer of  $\mathbb{E}_{\mathcal{D}_e}[\ell_{\log}(\omega, Y) \mid \varphi(X) = z]$ , given by

$$\omega = \log \left( \frac{\Pr_{\mathcal{D}_e}[Y=1 \mid \varphi(X)=z]}{\Pr_{\mathcal{D}_e}[Y=-1 \mid \varphi(X)=z]} \right) = \log \left( \frac{1+\mathbb{E}_{\mathcal{D}_e}[Y \mid \varphi(X)=z]}{1-\mathbb{E}_{\mathcal{D}_e}[Y \mid \varphi(X)=z]} \right),$$

uniquely corresponds to  $\mathbb{E}_{\mathcal{D}_e}[Y \mid \varphi(X) = z]$ .

The following lemma is implicit in Arjovsky et al. 2019.

**Lemma 4.** Under Setting A, for all  $\mathcal{E}$  and  $d \geq 1$ ,

$$\mathcal{I}(\mathcal{E}) \subseteq \mathcal{I}_{\mathcal{S}}(\mathcal{E}) = \mathcal{I}_{\mathcal{W}_{\mathrm{lin}}^d}(\mathcal{E}).$$

*Proof.* We prove the lemma in the following three parts.

 $\mathcal{I}(\mathcal{E}) \subseteq \mathcal{I}_{\mathcal{S}}(\mathcal{E})$ . Given  $f \in \mathcal{I}(\mathcal{E})$ , let  $\varphi : \mathcal{X} \to \mathcal{Z}$  and  $w : \mathcal{Z} \to \mathbb{R}$  be such that  $f = w \circ \varphi$ , where  $w \in \underset{\overline{w}: \mathcal{Z} \to \mathbb{R}}{\operatorname{argmin}_{\overline{w}: \mathcal{Z} \to \mathbb{R}}} \mathcal{L}_e(\overline{w} \circ \varphi)$  for all  $e \in \mathcal{E}$ . Define  $\varphi' : \mathcal{X} \to \mathbb{R}$  as  $\varphi'(x) := w(\varphi(x))$  and  $w' : \mathbb{R} \to \mathbb{R}$  to be the identity function w'(z) = z. Thus, we have  $w' \circ \varphi' = w \circ \varphi = f$ . Additionally, it holds that for all  $e \in \mathcal{E}$ ,  $w' \in \underset{\overline{w}' \in \mathcal{S}}{\operatorname{argmin}_{\overline{w}' \in \mathcal{S}}} \mathcal{L}_e(\overline{w}' \circ \varphi')$ . (Suppose for contradiction that this is not the case. Then for some environment  $e \in \mathcal{E}$ , there exists  $c \neq 1$  such that  $\mathcal{L}_e(cf) < \mathcal{L}_e(f)$ , corresponding to  $\overline{w}' \in \mathcal{S}$  such that  $\overline{w}'(z) \coloneqq cz$ . Hence  $\mathcal{L}_e(c \cdot w) \circ \varphi > \mathcal{L}_e(w \circ \varphi)$ , which contradicts that  $w \in \underset{\overline{w} \in \mathcal{W}}{\operatorname{argmin}} \mathcal{L}_e(\overline{w} \circ \varphi)$ .) Thus, we get  $f \in \mathcal{I}_{\mathcal{S}}(\mathcal{E})$ .

- $\mathcal{I}_{\mathcal{W}_{\text{lin}}^d}(\mathcal{E}) \subseteq \mathcal{I}_{\mathcal{S}}(\mathcal{E})$ . The proof of the above part shows, more generally, that  $\mathcal{I}_{\mathcal{W}}(\mathcal{E}) \subseteq \mathcal{I}_{\mathcal{S}}(\mathcal{E})$  for any  $\mathcal{W}$  that is closed under scalar multiplications (that is,  $w \in \mathcal{W} \implies c \cdot w \in \mathcal{W}$  for all  $c \in \mathbb{R}$ ), it holds that  $\mathcal{I}_{\mathcal{W}}(\mathcal{E}) \subseteq \mathcal{I}_{\mathcal{S}}(\mathcal{E})$ . Since  $\mathcal{W}_{\text{lin}}^d$  is closed under scalar multiplications, we get  $\mathcal{I}_{\mathcal{W}_{\text{in}}^d}(\mathcal{E}) \subseteq \mathcal{I}_{\mathcal{S}}(\mathcal{E})$ .
- $\mathcal{I}_{\mathcal{S}}(\mathcal{E}) \subseteq \mathcal{I}_{\mathcal{W}_{\mathrm{lin}}^d}(\mathcal{E})$ . Given  $f \in \mathcal{I}_{\mathcal{S}}(\mathcal{E})$ , let  $\varphi : \mathcal{X} \to \mathbb{R}$  and  $w : \mathbb{R} \to \mathbb{R}$  be such that  $f = w \circ \varphi$ , where  $w \in \underset{\overline{w} \in \mathcal{S}}{\operatorname{argmin}}_{\overline{w} \in \mathcal{S}} \mathcal{L}_e(\overline{w} \circ \varphi)$  for all  $e \in \mathcal{E}$ . Define  $\varphi' : \mathcal{X} \to \mathbb{R}^d$  as  $\varphi'(x) \coloneqq \varphi(x) \cdot v$  for any unit vector  $v \in \mathbb{R}^d$  and  $w' : \mathbb{R}^d \to \mathbb{R}$  as  $w'(z) \coloneqq w(\langle z, v \rangle)$ . It is easy to see that  $w' \in \underset{\overline{w}' \in \mathcal{W}_{\mathrm{lin}}^d}{\operatorname{sgmin}}_{\overline{w}' \in \mathcal{W}_{\mathrm{lin}}^d}(\overline{w}' \circ \varphi')$  for all  $e \in \mathcal{E}$  and hence  $w' \circ \varphi' = w \circ \varphi = f$ . Thus,  $f \in \mathcal{I}_{\mathcal{W}_{\mathrm{lin}}^d}(\mathcal{E})$ .

## B More details on Two-Bit Environments (from Section 3)

We show that for all  $\alpha \in (0,1)$ , just two environments in  $\mathcal{E}_{\alpha}$  are sufficient to determine both  $\mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\alpha})$  and  $\mathcal{I}(\mathcal{E}_{\alpha})$ . Thus, the failure of  $\mathsf{IRM}_{\mathcal{S}}$  observed in Section 3 is not due to lack of sufficiently representative training environments, but instead due to the difference between what  $\mathsf{IRM}_{\mathcal{S}}$  deems an "invariant predictor" and the notion of invariance as in Definition 1.

**Proposition 5.** Under Setting A, for all  $\alpha \in (0,1)$  and  $\mathcal{E}_{tr} = \{e_1, e_2\}$  for any two distinct  $e_1, e_2 \in \mathcal{E}_{\alpha}$ ,

(i) 
$$\mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\mathrm{tr}}) = \mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\alpha})$$
 and (ii)  $\mathcal{I}(\mathcal{E}_{\mathrm{tr}}) = \mathcal{I}(\mathcal{E}_{\alpha})$ .

*Proof.* (i). By definition, we have  $\mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\alpha}) \subseteq \mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\mathrm{tr}})$ , since  $\mathcal{E}_{\mathrm{tr}} \subseteq \mathcal{E}_{\alpha}$ . We show the converse. As noted in Section 2, for any convex and differentiable loss  $\ell$  and for any set of environments  $\mathcal{E}$  we have that  $f \in \mathcal{I}_{\mathcal{S}}(\mathcal{E})$  if and only if  $f = 1 \cdot \varphi$  such that  $\nabla_{w|w=1}\mathcal{L}_{e}(w \cdot \varphi) = 0$  for all  $e \in \mathcal{E}$ . The key observation is that for environment  $e = (\alpha, \beta_{e})$ ,

$$\nabla_{w|w=1} \mathcal{L}_{e}(w \cdot \varphi) := \sum_{x_{1}, x_{2}, y} \Pr_{\mathcal{D}_{e}}[X_{1} = x_{1}, X_{2} = x_{2}, Y = y] \cdot \nabla_{w|w=1} \ell(w \cdot \varphi(x_{1}, x_{2}), y)$$

$$= \sum_{x_{1}, x_{2}, y} ((1 - \alpha) \mathbb{1}_{x_{1} = y} + \alpha \mathbb{1}_{x_{1} \neq y}) \cdot ((1 - \beta_{e}) \mathbb{1}_{x_{2} = y} + \beta_{e} \mathbb{1}_{x_{2} \neq y}) \cdot \nabla_{w|w=1} \ell(w \cdot \varphi(x_{1}, x_{2}), y)$$

is affine in  $\beta_e$ . In particular, it can be decomposed as  $\nabla_{w|w=1} \mathcal{L}_e(w \cdot \varphi) = F(\varphi) + \beta_e G(\varphi)$  for some functions F and G. If  $f \in \mathcal{E}_{tr}$ , then we have that  $f = 1 \cdot \varphi$  such that both  $F(\varphi) + \beta_{e_1} G(\varphi) = 0$  and  $F(\varphi) + \beta_{e_2} G(\varphi) = 0$  hold, which happens if and only if  $F(\varphi) = 0 = G(\varphi)$ . This implies  $F(\varphi) + \beta_e G(\varphi) = 0$  for all  $\beta_e \in (0, 1)$ , and hence  $f \in \mathcal{E}_{\alpha}$ .

(ii). By definition, we have that  $\mathcal{I}(\mathcal{E}_{\alpha}) \subseteq \mathcal{I}(\mathcal{E}_{\mathrm{tr}})$ , since  $\mathcal{E}_{\mathrm{tr}} \subseteq \mathcal{E}_{\alpha}$ . We show the converse by establishing that the only invariant predictors in  $\mathcal{I}(\mathcal{E}_{\mathrm{tr}})$  are those that do not depend on  $X_2$ . By Observation 2, we have that  $\varphi$  is invariant over  $\mathcal{E}_{\mathrm{tr}}$  if and only if  $\mathbb{E}_{\mathcal{D}_{e_1}}[Y \mid \varphi(X) = z] = \mathbb{E}_{\mathcal{D}_{e_2}}[Y \mid \varphi(X) = z]$  for all  $z \in \mathcal{Z}_{\varphi}^{e_1} \cap \mathcal{Z}_{\varphi}^{e_2}$ . In other words,  $\varphi$  is invariant over  $\mathcal{E}_{\mathrm{tr}}$  if and only if  $\mathbb{E}_{\mathcal{D}_{e}}[Y \mid (X_1, X_2) \in \varphi^{-1}(z)]$  is identical for  $e_1$  and  $e_2$  as long as  $\Pr_{\mathcal{D}_{e}}[\varphi(X_1, X_2) = z]$  is non-zero in both environments.

In Table 1, we compute  $\mathbb{E}_{\mathcal{D}_e}[Y \mid (X_1, X_2) \in S]$  for all possible non-empty subsets  $S \subseteq \{-1, 1\}^2$ , in terms of the environment parameters  $\alpha$  and  $\beta_e$  and track which of these depend or do not depend on  $\beta_e$ . The ones that depend on  $\beta_e$  can be seen to be distinct for any two distinct values of  $\beta_e$ . Thus, the only invariant representations over  $\mathcal{E}_{tr}$  are those corresponding to the following partitions.

- ▶ {{((1,1),(1,-1),(-1,1),(-1,-1)}}, that is,  $\varphi(x_1,x_2)$  is constant. The predictor  $f \in \mathcal{I}(\mathcal{E}_{tr})$  corresponding to this representation is the identically zero-predictor  $f_0$  (for both  $\ell_{sq}$  and  $\ell_{log}$ ).
- ▶  $\{\{(1,1),(1,-1)\},\{(-1,1),(-1,-1)\}\}$ , that is,  $\varphi(1,1)=\varphi(1,-1)$  and  $\varphi(-1,1)=\varphi(-1,-1)$ , or essentially  $\varphi(x_1,x_2)=x_1$ . The predictor  $f\in\mathcal{I}(\mathcal{E}_{\mathrm{tr}})$  corresponding to this representation is  $f(x)=(1-2\alpha)\cdot x_1$  (for  $\ell=\ell_{\mathrm{sq}}$ ) or  $f(x)=\log\frac{(1-\alpha)}{\alpha}\cdot x_1$  (for  $\ell=\ell_{\mathrm{log}}$ ) see proof of Observation 2 for reference. ▶  $\{\{(1,1),(-1,-1)\},\{(1,-1),(-1,1)\}\}$ , that is,  $\varphi(1,1)=\varphi(-1,-1)$  and  $\varphi(1,-1)=\varphi(-1,1)$ , or essentially
- ▶ {{(1,1), (-1,-1)}, {(1,-1), (-1,1)}}, that is,  $\varphi(1,1) = \varphi(-1,-1)$  and  $\varphi(1,-1) = \varphi(-1,1)$ , or essentially  $\varphi(x_1,x_2) = x_1 \cdot x_2$ . While this representation does depend on  $x_2$ , the predictor  $f \in \mathcal{I}(\mathcal{E}_{tr})$  corresponding to this representation is the identically zero-predictor  $f_0$  (for both  $\ell_{sq}$  and  $\ell_{log}$ ).

In all the above cases, we observe that the invariant representations over  $\mathcal{E}_{tr}$  are also invariant over  $\mathcal{E}_{\alpha}$  and moreover, the corresponding predictors are simultaneously optimal for all  $e \in \mathcal{E}_{\alpha}$  and hence in  $\mathcal{I}(\mathcal{E}_{\alpha})$ . Thus, we have  $\mathcal{I}(\mathcal{E}_{tr}) \subseteq \mathcal{I}(\mathcal{E}_{\alpha})$ .

S	ubset $S$	$\subseteq \{-1, 1\}$	$\lfloor l \rfloor^2$	$\mathbb{E}_{\mathcal{D}_e}[Y (X_1, X_2) \in S]$	Independent of $\beta_e$ ?
(1,1)				$\frac{1-\alpha-\beta_e}{1-\alpha+(2\alpha-1)\beta_e}$	No
	(1,-1)			$\frac{\alpha - \beta_e}{-\alpha + (2\alpha - 1)\beta_e}$	No
		(-1,1)		$\frac{-\alpha + \beta_e}{-\alpha + (2\alpha - 1)\beta_e}$	No
			(-1,-1)	$\frac{\alpha - 1 + \beta_e}{1 - \alpha + (2\alpha - 1)\beta_e}$	No
(1,1)	(1,-1)			$1-2\alpha$	Yes
(1,1)		(-1,1)		$1-2\beta_e$	No
(1,1)			(-1,-1)	0	Yes
	(1,-1)	(-1,1)		0	Yes
	(1,-1)		(-1,-1)	$2\beta_e - 1$	No
		(-1,1)	(-1,-1)	$2\alpha - 1$	Yes
(1,1)	(1,-1)	(-1,1)		$\frac{\alpha - 1 + \beta_e}{-1 - \alpha + (2\alpha - 1)\beta_e}$	No
(1,1)	(1,-1)		(-1,-1)	$\frac{-\alpha + \beta_e}{2 - \alpha + (2\alpha - 1)\beta_e}$	No
(1,1)		(-1,1)	(-1,-1)	$\frac{\alpha - \beta_e}{2 - \alpha + (2\alpha - 1)\beta_e}$	No
	(1,-1)	(-1,1)	(-1,-1)	$\frac{1-\alpha-\beta_e}{-1-\alpha+(2\alpha-1)\beta_e}$	No
(1,1)	(1,-1)	(-1,1)	(-1,-1)	0	Yes

Table 1:  $\mathbb{E}_{\mathcal{D}_e}[Y \mid X \in S]$  for different choices of S in the proof of Part (ii) of Proposition 5.

#### B.1 Case of square loss

We recall the example described in Section 3 that demonstrated the difference between IRM<sub>S</sub> and IRM. We have  $\mathcal{E} = \mathcal{E}_{0.1}$  and  $\mathcal{E}_{tr} = \{(0.1, 0.2), (0.1, 0.25)\}$ . We get from Proposition 5 that  $\mathcal{I}_{\mathcal{S}}(\mathcal{E}) = \mathcal{I}_{\mathcal{S}}(\mathcal{E}_{tr})$ , which can be numerically seen to contain (approximately) the following four predictors, by simultaneously solving  $(\nabla_w \text{ for } \ell_{sq})$  for all  $e \in \mathcal{E}_{tr}$ .

	$f_0$		$f_{H}$	RM	f	, 1	f	2
	$X_2 = +1$	$X_2 = -1$	$X_2 = +1$ $X_2 = -1$		$X_2 = +1$	$X_2 = -1$	$X_2 = +1$	$X_2 = -1$
$X_1 = +1$	0	0	0.8	0.8	0.9557	0.2943	0.2943	0.9557
$X_1 = -1$	0	0	-0.8	-0.8	-0.2943	-0.9557	-0.9557	-0.2943

On the other hand,  $\mathcal{I}(\mathcal{E}_{0.1})$  contains only two of the predictors, namely  $f_0$  and  $f_{\mathsf{IRM}}$ , the latter being the optimal predictor chosen by IRM on  $\mathcal{E}_{\mathsf{tr}}$  — note that this predictor depends only on  $X_1$ .

Figure 2 shows the population square losses  $\mathcal{L}_e$  for each of the predictors in  $\mathcal{I}_{\mathcal{S}}(\mathcal{E}_{0.1})$  for all  $e \in \mathcal{E}_{0.1}$ . It can observed that for  $e = (0.1, \beta_e)$  with  $\beta_e < 0.28$ , it holds that  $\mathcal{L}_e(f_1) < \mathcal{L}_e(f_{\mathsf{IRM}})$ . Thus, no matter how many training environments are present in  $\mathcal{E}_{\mathsf{tr}}$ ,  $\mathsf{IRM}_{\mathcal{S}}$  will choose  $f_1$  as the optimal predictor as long as  $\beta_e < 0.28$  for all  $e \in \mathcal{E}_{\mathsf{tr}}$ . On the other hand,  $\mathsf{IRM}$  with just two environments learns the predictor  $f_{\mathsf{IRM}}$ .

We also note that the value  $\alpha = 0.1$  is not special either. In fact, a similar phenomenon as above is observed in  $\mathcal{E}_{\alpha}$  for any value of  $\alpha < 0.1464$  or  $\alpha > 0.8536$ . The following section explains the meaning of these cutoff values.

## B.1.1 Analytic characterization of odd predictors in $\mathcal{I}_{\mathcal{S}}(\mathcal{E})$

Following the initial version of this paper (which found these constants only by numerically solving certain quadratic systems), Léon Bottou communicated to us the following clean analysis, which provides a closed-form understanding of these  $\mathsf{IRM}_{\mathcal{S}}$  solutions and the range of  $\alpha$  when such examples arise. We are grateful to Léon

for allowing us to include his calculations here.

Firstly, observe that for  $\mathcal{X} = \{-1,1\}^2$ , a representation  $\varphi(x_1,x_2)$  is odd if and only if it is *linear*, namely,  $\varphi(x_1,x_2) := w_1x_1 + w_2x_2$ . Suppose  $\mathcal{E}_{tr}$  consists of two environments  $e_1 = (\alpha,\beta_1)$  and  $e_2 = (\alpha,\beta_2)$ . From  $(\nabla_w \text{ for } \ell_{sq})$ , we for any  $f(x) = 1 \cdot \varphi(x) = w_1x_1 + w_2x_2 \in \mathcal{I}_{\mathcal{S}}(\mathcal{E}_{tr})$  that

$$\mathbb{E}_{\mathcal{D}_{e_1}}(w_1x_1 + w_2x_2 - y)(w_1x_1 + w_2x_2) = 0$$

$$\mathbb{E}_{\mathcal{D}_{e_2}}(w_1x_1 + w_2x_2 - y)(w_1x_1 + w_2x_2) = 0.$$

From the definition (Two-Bit-Envs), we have (i)  $\mathbb{E}_{\mathcal{D}_i}[x_1^2] = \mathbb{E}_{\mathcal{D}_i}[x_2^2] = 1$ , (ii)  $\mathbb{E}_{\mathcal{D}_i}[x_1y] = a$ , (iii)  $\mathbb{E}_{\mathcal{D}_i}[x_2y] = b_i$  and (iv)  $\mathbb{E}_{\mathcal{D}_i}[x_1x_2] = ab_i$ , where  $a := 1 - 2\alpha$  and  $b_i = 1 - 2\beta_i$  for  $i \in \{1, 2\}$ . Thus, we get

$$w_1^2 + w_2^2 + 2w_1w_2ab_1 = w_1a + w_2b_1, (1)$$

$$w_1^2 + w_2^2 + 2w_1w_2ab_2 = w_1a + w_2b_2. (2)$$

By subtracting and using  $b_1 - b_2 \neq 0$ , we get

$$2w_1w_2a = w_2. (3)$$

When  $w_2 = 0$ , we get from (1) (or (2)) that either  $w_1 = 0$  or  $w_1 = a$ . But when  $w_2 \neq 0$ , we have from (3) that  $w_1 = 1/2a$ . Substituting this in (1) (or (2)), we get the two additional solutions given by

$$w_1 = \frac{1}{2a}$$
 and  $w_2 = \pm \sqrt{\frac{1}{2} - \frac{1}{4a^2}}$ .

Note that these additional solutions (with  $w_2 \neq 0$ ) exist only when  $\frac{1}{2} - \frac{1}{4a^2} > 0$ , or  $(1 - 2\alpha)^2 > \frac{1}{2}$ . That is,

$$\alpha < \frac{1}{2} - \frac{1}{2\sqrt{2}} \approx 0.1464$$
 or  $\alpha > \frac{1}{2} + \frac{1}{2\sqrt{2}} \approx 0.8536$ .

In this regime, the four odd (or linear) predictors in  $\mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\alpha})$  are

$$f_0(x) = 0$$

$$f_{\text{IRM}}(x) = (1 - 2\alpha) \cdot x_1$$

$$f_1(x) = \frac{1}{2 - 4\alpha} \cdot x_1 + \sqrt{\frac{1}{2} - \frac{1}{4(1 - 2\alpha)^2}} \cdot x_2$$

$$f_2(x) = \frac{1}{2 - 4\alpha} \cdot x_1 - \sqrt{\frac{1}{2} - \frac{1}{4(1 - 2\alpha)^2}} \cdot x_2.$$

#### B.2 Case of logistic loss

We observe a similar phenomenon with logistic loss as was observed for square loss. We consider  $\mathcal{E} = \mathcal{E}_{0.05}$  and  $\mathcal{E}_{tr} = \{(0.05, 0.1), (0.05, 0.2)\}$ . Again, we get from Proposition 5 that  $\mathcal{I}_{\mathcal{S}}(\mathcal{E}) = \mathcal{I}_{\mathcal{S}}(\mathcal{E}_{tr})$ , which can be numerically seen to contain (approximately) the following predictors, by simultaneously solving Equation  $(\nabla_w)$  for all  $e \in \mathcal{E}_{tr}$ .

	$f_0$		$f_0$ $f_{IRM}$		$f_1$		$f_2$	
$X_2 = +1$		$X_2 = -1$	$X_2 = +1$	$X_2 = -1$	$X_2 = +1$	$X_2 = -1$	$X_2 = +1$	$X_2 = -1$
$X_1 = +1$	0	0	2.9444	2.9444	4.9847	0.9041	0.9041	4.9847
$X_1 = -1$	0	0	-2.9444	-2.9444	-0.9041	-4.9847	-4.9847	-0.90413

On the other hand,  $\mathcal{I}(\mathcal{E}_{0.05})$  contains only two of the predictors, namely  $f_0$  and  $f_{\mathsf{IRM}}$ , the latter being the optimal predictor chosen by IRM on  $\mathcal{E}_{\mathsf{tr}}$  — note that this predictor depends only on  $X_1$ .

Figure 7 shows the population square losses  $\mathcal{L}_e$  for each of the predictors in  $\mathcal{I}_{\mathcal{S}}(\mathcal{E}_{0.05})$  for all  $e \in \mathcal{E}_{0.05}$ . It can observed that for  $e = (0.05, \beta_e)$  with  $\beta_e < 0.25$ , it holds that  $\mathcal{L}_e(f_1) < \mathcal{L}_e(f_{\mathsf{IRM}})$ . Thus, no matter how many training environments are present in  $\mathcal{E}_{\mathsf{tr}}$ ,  $\mathsf{IRM}_{\mathcal{S}}$  will choose  $f_1$  as the optimal predictor as long as  $\beta_e < 0.25$  for all  $e \in \mathcal{E}_{\mathsf{tr}}$ . On the other hand,  $\mathsf{IRM}$  with just two environments learns the predictor  $f_{\mathsf{IRM}}$ .

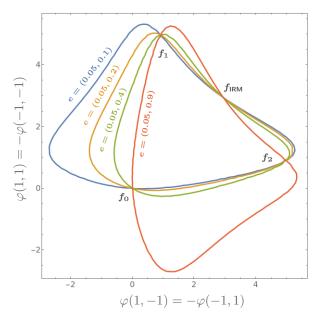


Figure 6: Odd solutions to  $\nabla_{w|w=1}\mathcal{L}_e(w\cdot\varphi)=0$  (with  $\ell=\ell_{\log}$ ) for four environments in  $\mathcal{E}_{0.05}$ . (Compare to Figure 1.)

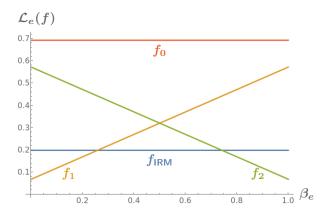


Figure 7: Losses  $\mathcal{L}_e$  (for  $\ell = \ell_{\log}$ ) of odd predictors in  $\mathcal{I}_{\mathcal{S}}(\mathcal{E}_{0.05})$  for various  $e = (0.05, \beta_e)$ . (Compare to Figure 2.)

We also note that the value  $\alpha = 0.05$  is not special; a similar phenomenon as above is observed in  $\mathcal{E}_{\alpha}$  for any value of  $\alpha < 0.077$ .

In the supplementary material, we include the Mathematica code (two-bit/two-bit-irm.nb, and a PDF version two-bit/two-bit-irm.pdf) that was used to compute  $f_{\mathsf{IRM}}$  and  $f_{\mathsf{IRM}_S}$  solutions and plot Figures 1 to 3, 5, 6 and 7.

# C More Colored-MNIST experiments

We now consider more details and variations of the Colored-MNIST experiments of Section 3.1.

The architecture used by Arjovsky et al. (2019) is fully connected, mapping inputs of dimension  $2 \cdot 14 \cdot 14$  to hidden dimension h, from h to h, and then from h to a scalar prediction, with ReLU activations on each layer except the last. The model is optimized with full-batch Adam for 501 steps, with a scaled penalty on the squared (Frobenius) norm of each parameter, and hyperparameters selected as:

- ▶ Hidden dimension  $h: |2^{\text{Uniform}[6,9)}|$ .
- ▶ Weight of  $L_2$  regularization:  $10^{\text{Uniform}[-2,-5)}$ .
- ▶ Learning rate:  $10^{\text{Uniform}[-2.5, -3.5)}$ .

▶ For IRMv1, the gradient penalty weight  $\lambda$  is 1 for Uniform $\{50, 51, \dots, 250\}$  iterations, then  $10^{\text{Uniform}[2,6)}$ .

In Figure 8, we reproduce the results of Figure 4 (left column) but also show results of versions of the architecture forced to depend only on  $X_1$  or  $X_2$  while training via ERM: color-only takes inputs of shape 2, a one-hot indicator for whether the color is red or green, while digit-only receives a flattened grayscale image of dimension 14 · 14. This allows us to see the amount of variation we can expect based purely on changes in the learning process. We also show (in the right column) a flipped version of the problem, where the invariant feature is color rather than the digit identity; this is the same from the point of view of the abstract Two-Bit environment, but allows us to see how much of the behavior depends on the different way that this network processes digit and color information.

As mentioned in Section 3.1, we also consider a "split" variant of the architecture, which is perhaps closer to the abstract two-bit version. Here, the network has two branches: one takes a grayscale  $14 \times 14$  version of its input, which is processed as in the previous architecture down to a scalar. The other branch takes a one-hot (two-dimensional) indicator for the color, and (via a  $2 \times 1$  linear layer) outputs an arbitrary scalar for each color. The top of the network takes in these two scalar values, processes them with an 8-dimensional ReLU layer, then makes a final linear prediction. color-only and digit-only versions simply omit one of those branches. Results for  $\ell_{sq}$  are shown in Figure 9. Here we most clearly see the "average-case" failure of IRMv1 in the color-invariant case.

Similar results for  $\ell_{log}$  are shown in Figures 10 and 11. The expected failure mode is generally less visible here, though it is more evident in the color-invariant settings than the digit-invariant ones.

In the supplementary material, we include the PyTorch code, modified from that of Arjovsky et al. (2019), used to produce these results (colored-mnist directory).

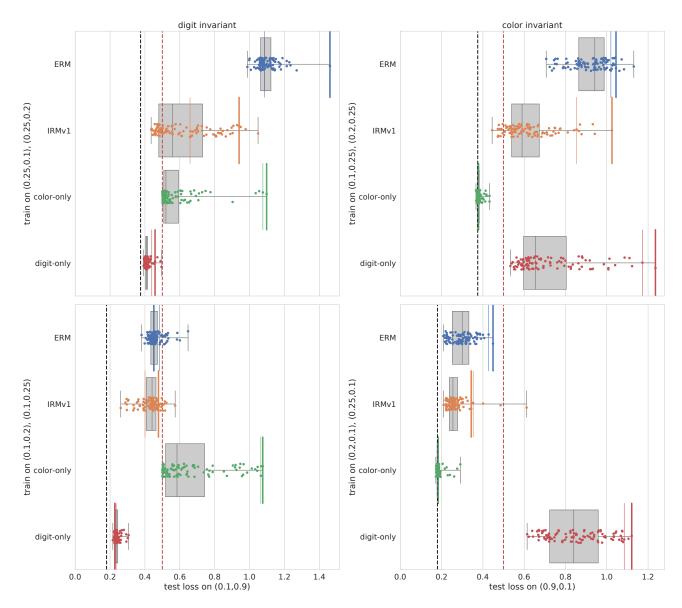


Figure 8: Colored-MNIST using  $\ell_{\rm sq}$  with the architecture of Arjovsky et al. (2019): the same as Figure 4, but additionally showing cases where color is invariant rather than the digit (right column), and performance of networks which receive only grayscale digits as input, or only a one-hot indicator of the color. Thin colored lines show performance of the second-best hyperparameter setting on the training environments.

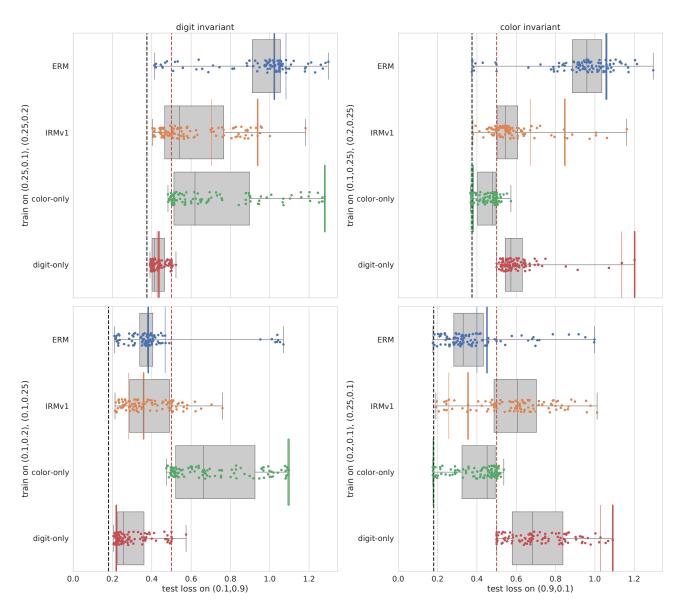


Figure 9: Colored-MNIST using  $\ell_{\rm sq},$  with a "split" architecture.

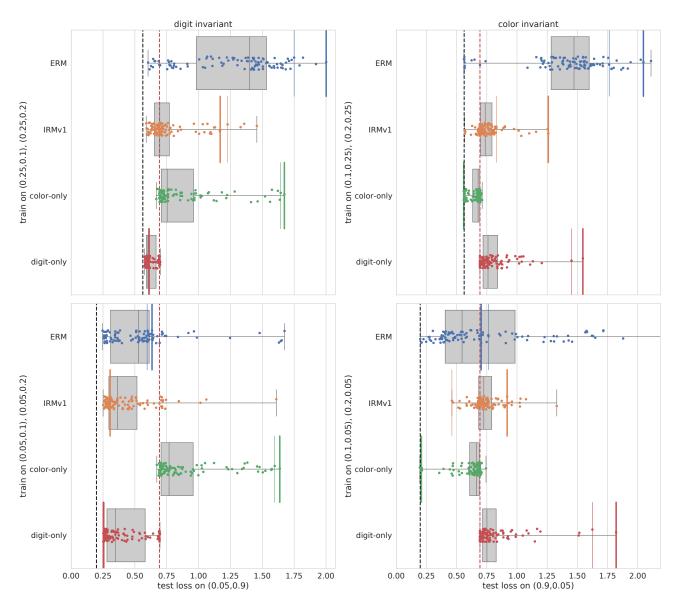


Figure 10: Colored-MNIST using  $\ell_{\log},$  with a "split" architecture.

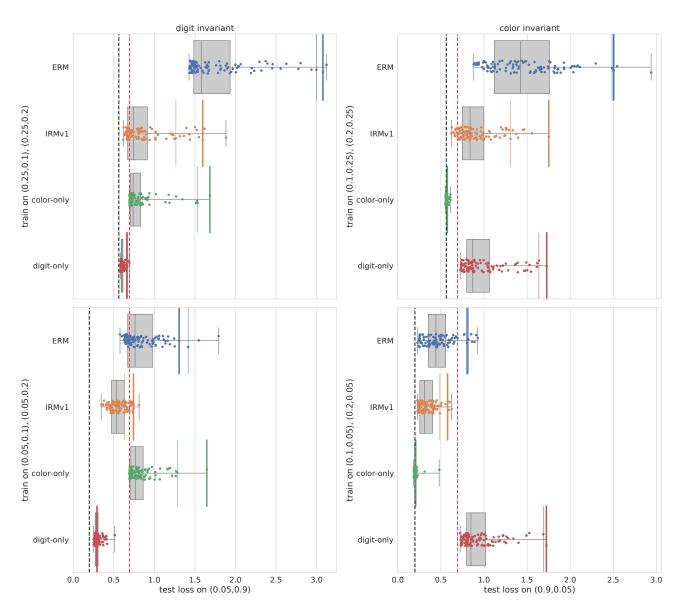


Figure 11: Colored-MNIST using  $\ell_{\log}$ , with the architecture of Arjovsky et al. (2019).

## D More details on failure of IRM (Section 4)

We first prove Proposition 6, restated below for convenience.

**Proposition 6.** In Setting A, for  $\mathcal{E}$  as above, it holds for Lebesgue-almost all  $\mathcal{E}_{tr} \subseteq \mathcal{E}$  with  $|\mathcal{E}_{tr}| \geq 2$  that  $\mathcal{I}(\mathcal{E}) = \mathcal{I}(\mathcal{E}_{tr})$ . Moreover, any  $f \in \mathcal{I}(\mathcal{E})$  depends on at most one of  $x_1$  or  $x_2$ .

*Proof.* Since the parameterization of the environments is analytic and  $\Theta = (-1/6, 1/3)$  is a connected open set, we get from part (i) of Proposition 8 that for almost all  $\mathcal{E}_{tr} \subseteq \mathcal{E}$  with  $|\mathcal{E}_{tr}| \ge 2$ , it holds that  $\mathcal{I}(\mathcal{E}_{tr}) = \mathcal{I}(\mathcal{E})$ . We now establish the second part: any  $f \in \mathcal{I}(\mathcal{E})$  depends on at most one of  $x_1$  or  $x_2$ .

Similar to Table 1, we can compute  $\mathbb{E}_{\mathcal{D}_e}[Y \mid (X_1, X_2) \in S]$  for all possible non-empty subsets  $S \subseteq \{-1, 0, 1\}^2$  and track which of these depend or do not depend on  $\theta_e$ . Since it is cumbersome to enumerate manually over all the 511  $(=2^9-1)$  possible non-empty subsets of  $\{-1,0,1\}^2$ , we enumerate this symbolically, using the SymPy package in Python, to identify all the subsets where  $\mathbb{E}_{\mathcal{D}_e}[Y \mid (X_1, X_2) \in S]$  does not depend on  $\theta_e$ ; note that  $\mathbb{E}_{\mathcal{D}_e}[Y \mid (X_1, X_2) \in S]$  is a rational function in  $\theta_e$  and hence if it is not identically zero, then it is in fact different for almost all pairs of choices for  $\theta_e$ . (Code is in the supplementary material; two-bit/pure-irm-fail-example.py.)

There turn out to be 37 non-empty subsets S for which  $\mathbb{E}_{\mathcal{D}_e}[Y \mid (X_1, X_2) \in S]$  does not depend on  $\theta_e$ ; out of which  $\mathbb{E}_{\mathcal{D}_e}[Y \mid (X_1, X_2) \in S]$  is non-zero for only 6 choices of S as given in Table 2.

	Subset $S \subseteq \{-1, 0, +1\}^2$									Characterization of $S$
		(+1,-1)			(+1,0)			(+1,+1)	0.3	$X_1 = +1$
(-1,-1)			(-1,0)			(-1,+1)			-0.3	$X_1 = -1$
						(-1,+1)	(0,+1)	(+1,+1)	0.3	$X_2 = +1$
(-1,-1)	(0,-1)	(+1,-1)							-0.3	$X_2 = -1$
(-1,-1)	(0,-1)		(-1,0)	(0,0)		(-1,+1)	(0,+1)		-0.15	$X_1 \in \{-1, 0\}$
	(0,-1)	(+1,-1)		(0,0)	(+1,0)		(0,+1)	(+1,+1)	0.15	$X_1 \in \{0, +1\}$

Table 2: Conditional expectations for different choices of  $\varphi$  in the proof of Proposition 6.

For any predictor  $w \circ \varphi \in \mathcal{I}(\mathcal{E}_{tr})$  and any z satisfying  $w(z) \neq 0$ , it must be the case that  $\varphi^{-1}(z)$  is among the ones in Table 2. Thus, it is easy to see that the only predictors in  $\mathcal{I}(\mathcal{E}_{tr})$  are those that depend only on  $x_1$ , or depend only on  $x_2$ , or neither (for the identically zero predictor  $f_0$ ). Clearly, all these predictors are also in  $\mathcal{I}(\mathcal{E})$  and thus, we get  $\mathcal{I}(\mathcal{E}_{tr}) = \mathcal{I}(\mathcal{E})$ .

Moreover, for any environment  $e \in \mathcal{E}$ , it holds in the case of  $\ell = \ell_{\text{sq}}$  that among all the predictors that depend only on  $x_1$ , the one with the lowest loss  $\mathcal{L}_e(\cdot)$  is  $f_1(x) = 0.3x_1$  and similarly, among all the predictors that depend only on  $x_2$ , the one with the lowest loss  $\mathcal{L}_e$  is  $f_2(x) = 0.3x_2$ . (Similar, argument holds for  $\ell = \ell_{\text{log}}$ .) Thus, IRM will always pick one among  $f_1$  and  $f_2$ .

Finally, we visualize the loss of the predictors  $f_1(x) := 0.3x_1$ ,  $f_2(x) = 0.3x_2$  and the zero predictor  $f_0(x) = 0$  over all choices of  $\theta_e \in (-1/6, 1/3)$  in Figure 12. It is easy to see from the figure that if  $\mathcal{E}_{tr}$  only contains environments e corresponding to  $\theta_e < 0$ , we will have that  $\mathcal{L}_e(f_2) < \mathcal{L}_e(f_1)$  for all  $e \in \mathcal{E}_{tr}$ , and yet the invariant predictor that minimizes  $\sup_{e \in \mathcal{E}} \mathcal{L}_e(\cdot)$  is  $f_1$  and in fact  $\sup_{e \in \mathcal{E}} \mathcal{L}_e(f_2) = \sup_{e \in \mathcal{E}} \mathcal{L}_e(f_0) = 0.5$ , that is, worst-case over all environments,  $f_2$  is no better than the identically zero predictor.

## E More details on Generalization of Invariance (from Section 5)

**Proposition 7.** There exists a continuous map  $\Pi:(0,1)\to\Delta_{\mathcal{X}\times\mathcal{Y}}$  such that for  $\Theta_{\mathrm{tr}}=\left(0,\frac{1}{4}\right)$  and  $\mathcal{E}_{\mathrm{tr}}=\Pi(\Theta_{\mathrm{tr}})$ , it holds that  $\mathcal{I}(\mathcal{E}_{\mathrm{tr}})\neq\mathcal{I}(\mathcal{E}_{\Pi})$ .

*Proof.* Consider the two-bit environments of Section 3, denoted  $(\alpha_e, \beta_e)$ . Define  $\Pi$  as the continuous, piecewise-linear map

$$\Pi(\theta) = \begin{cases} \left(\frac{1}{10}, \frac{6\theta}{5}\right) & 0 < \theta \le \frac{1}{4} \\ \left(\frac{6\theta - 1}{5}, \frac{3}{10}\right) & \frac{1}{4} < \theta < 1 \end{cases}.$$

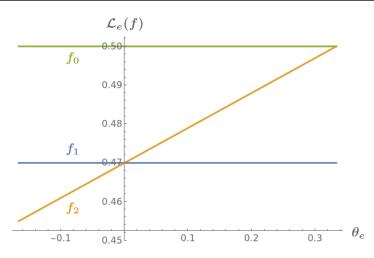


Figure 12: Losses  $\mathcal{L}_e$  (for  $\ell = \ell_{sq}$ ) of predictors  $f_1$ ,  $f_2$  and the zero predictor  $f_0$  in  $\mathcal{I}(\mathcal{E})$  for  $\theta_e \in (-1/6, 1/3)$ .

Consider  $\Theta_{\text{tr}} = \{\theta : 0 < \theta < \frac{1}{4}\}$ . Then the representation  $\varphi_1(X) := X_1$  is invariant across  $\mathcal{E}_{\text{tr}}$ , because  $\mathbb{E}_{\mathcal{D}_e}[Y \mid X_1 = x_1]$  is invariant across  $\mathcal{E}_{\text{tr}}$ . Thus, in the case of  $\ell_{\text{sq}}$ , the predictor  $f_1(X) := 0.8X_1$  is in  $\mathcal{I}(\mathcal{E}_{\text{tr}})$ . However,  $f_1 \notin \mathcal{I}(\mathcal{E})$ , because  $\mathbb{E}_{\mathcal{D}_e}[Y \mid X_1 = x_1]$  changes on environments in  $\mathcal{E} \setminus \mathcal{E}_{\text{tr}}$  when  $\frac{1}{4} < \theta < 1$ .

We now prove Proposition 8, restated below for convenience. First, we recall a basic fact about analytic functions.

**Fact 10** (Mityagin 2015). Let  $\Theta$  be a connected, open subset of  $\mathbb{R}^d$ . The set of zeros  $\{z \in \Theta \mid g(z) = 0\}$  of an analytic function  $g: \Theta \to \mathbb{R}$  has non-zero Lebesgue measure in  $\mathbb{R}^d$  if and only if g is identically 0.

**Proposition 8.** Let  $\Theta_{tr} \subseteq \Theta \subseteq \mathbb{R}^d$ , where  $\Theta$  is a connected, open set. Suppose  $\Pi : \Theta \to \Delta_{\mathcal{X} \times \mathcal{Y}}$  is analytic,  $\mathcal{X}$  and  $\mathcal{Y}$  are finite and  $\mathcal{E}_{tr} = \Pi(\Theta_{tr})$ . Then, under Setting A,

- (i) For almost all  $\Theta_{\rm tr}$  with  $|\Theta_{\rm tr}| \geq 2$ :  $\mathcal{I}(\mathcal{E}_{\rm tr}) = \mathcal{I}(\mathcal{E}_{\Pi})$ .
- (ii) For all  $\Theta_{\rm tr}$  with non-zero Lebesgue measure:  $\mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\rm tr}) = \mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\Pi})$ .

*Proof.* Part (i). We have  $\mathcal{I}(\mathcal{E}_{\Pi}) \subseteq \mathcal{I}(\mathcal{E}_{tr})$  by definition. We establish the converse by showing that  $F := \{(\theta_1, \theta_2) \mid \mathcal{I}(\{\Pi(\theta_1), \Pi(\theta_2)\}) \neq \mathcal{I}(\mathcal{E}_{\Pi})\}$  has measure zero in  $\Theta \times \Theta$ .

For any  $S \subseteq \mathcal{X}$  define the analytic functions  $n_S$  and  $d_S$  as

$$\begin{split} n_S(\theta) := \sum_{x \in S} \sum_{y \in \mathcal{Y}} y \cdot \Pi_{x,y}(\theta), \qquad d_S(\theta) := \sum_{x \in S} \sum_{y \in \mathcal{Y}} \Pi_{x,y}(\theta) \\ \text{so that } & \mathbb{E}_{\Pi(\theta)}[Y \mid X \in S] = \frac{n_S(\theta)}{d_S(\theta)} \text{ whenever } \Pr_{\Pi(\theta)}[X \in S] = d_S(\theta) \neq 0. \end{split}$$

We say that S is "valid" if either (i)  $d_S(\theta) = 0$  for all  $\theta \in \Theta$ , or (ii) there exists  $c_S \in \mathbb{R}$  such that  $n_S(\theta) = c_S \cdot d_S(\theta)$  for all  $\theta \in \Theta$  subject to  $d_S(\theta) \neq 0$ . Note that  $w \circ \varphi \in \mathcal{I}(\mathcal{E}_{\Pi})$  if and only if, (i) for all  $z \in \mathsf{range}(\varphi)$ , it holds that  $\varphi^{-1}(z) \subseteq \mathcal{X}$  is valid, and (ii)  $w(z) = \mathbb{E}_{\Pi(\theta)}[Y|X \in \varphi^{-1}(z)]$  for any  $\theta \in \Theta$  such that  $d_S(\theta) \neq 0$  (in the case of  $\ell_{sq}$ ).

For any invalid set  $S \subseteq \mathcal{X}$ , define  $F_S$  to consist of all pairs  $(\theta_1, \theta_2)$  for which at least one of the following condition holds: (i)  $d_S(\theta_1) = 0$ , or (ii)  $d_S(\theta_2) = 0$  or (iii)  $n_S(\theta_1) \cdot d_S(\theta_2) - n_S(\theta_2) \cdot d_S(\theta_1) = 0$ . Since S is not valid, it follows from Fact 10 that  $F_S$  has zero Lebesgue measure.

Finally, we show that  $F \subseteq \bigcup_{S \subseteq \mathcal{X}: S \text{ is invalid}} F_S$ . For any  $(\theta_1, \theta_2) \in F$  and any  $w \circ \varphi \in I(\{\Pi(\theta_1), \Pi(\theta_2)\}) \setminus I(E)$ ,

there exists  $z \in \mathsf{range}(\varphi)$  such that  $S = \varphi^{-1}(z) \subseteq \mathcal{X}$  is invalid. This implies  $(\theta_1, \theta_2) \in F_S$ . Since there are only finitely many  $S \subseteq \mathcal{X}$ , we get that F also has zero Lebesgue measure, thereby concluding the proof of part (i).

**Part** (ii). We have  $\mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\Pi}) \subseteq \mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\operatorname{tr}})$  by definition. To show the converse, consider any predictor  $f = 1 \cdot \varphi \in \mathcal{E}_{\operatorname{tr}}$ 

 $\mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\mathrm{tr}})$ , and consider

$$g(\theta) \ := \ \nabla_{w|w=1} \mathcal{L}_{\Pi(\theta)}(w \cdot \varphi) \ = \ \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} \Pi_{x,y}(\theta) \cdot \nabla_{w|w=1} \ell(w \cdot \varphi(x), y).$$

 $g(\theta)$  is linear in  $\{\Pi_{(x,y)}(\theta) \mid (x,y) \in \mathcal{X} \times \mathcal{Y}\}$ , each of which is analytic in  $\theta$ ; thus g is analytic in  $\theta$ . Since  $(\nabla_w)$  holds for all  $e \in \mathcal{E}_{tr}$ ,  $g(\theta) = 0$  for all  $\theta \in \Theta_{tr}$ . But since  $\Theta_{tr}$  has non-zero Lebesgue measure in  $\mathbb{R}^d$ , by Fact 10 g is identically 0 on  $\Theta$ , hence  $f \in \mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\Pi})$ .

We show how to extend Proposition 8 to the case of infinite (measurable) spaces  $\mathcal{X}$  and  $\mathcal{Y} \subseteq \mathbb{R}$ , where  $|y| \leq B$  for all  $y \in \mathcal{Y}$  for some known bound B. Similar to before, let  $\Delta_{\mathcal{X} \times \mathcal{Y}}$  be the set of all probability measures over  $\mathcal{X} \times \mathcal{Y}$ . For simplicity, we use  $\Omega$  to denote  $\mathcal{X} \times \mathcal{Y}$ .

**Definition 11.** For  $\Theta \subseteq \mathbb{R}^d$  and a measurable space  $\Omega$ , the parameterization  $\Pi : \Theta \to \Delta_{\Omega}$  is said to be analytic if for every measurable set  $S \subseteq \Omega$  and every measurable function  $g : \Omega \to \mathbb{R}$ , the function

$$\Pi_S^g(\theta) \coloneqq \int_{\omega \in S} g(\omega) \ \mathrm{d}\Pi_{\theta}(\omega)$$

is an analytic function in  $\theta$  (where we use  $\Pi_{\theta}$  to denote the measure  $\Pi(\theta)$  for simplicity).

We now state the extension of Proposition 8 to the case of infinite (measurable) spaces. In the case of  $\mathcal{I}_{\mathcal{S}}(\mathcal{E})$ , we will focus on the representations  $\varphi: \mathcal{X} \to \mathbb{R}$  where  $|\varphi(x)| \leq B$  for all  $x \in \mathcal{X}$ . From the point of view of  $\mathsf{IRM}_{\mathcal{S}}$ , this is without loss of generality because we know that  $|y| \leq B$  for all  $y \in \mathcal{Y}$ .

Proposition 12, however, requires a far stronger condition for  $\mathcal{I}(\mathcal{E}_{tr}) = \mathcal{I}(\mathcal{E}_{\Pi})$ :  $\Theta_{tr}$  needs non-zero Lebesgue measure, rather than simply almost all sets of at least two environments as in Proposition 8. The key step in the proof of Proposition 8 that allowed for this stronger statement was that the number of subsets  $S \subseteq \mathcal{X}$  is finite. We do not know if Proposition 12 can be strengthened to hold for finite  $\Theta_{tr}$ ; if not, it will be interesting to determine other conditions under which we can get generalization of invariance for finite  $\Theta_{tr}$ .

**Proposition 12.** Let  $\Theta_{tr} \subseteq \Theta \subseteq \mathbb{R}^d$ , where  $\Theta$  is a connected, open set and  $\Theta_{tr}$  has non-zero Lebesgue measure, in  $\mathbb{R}^d$ . Suppose  $\Pi : \Theta \to \Delta_{\mathcal{X} \times \mathcal{Y}}$  is analytic (as in Definition 11), and  $\mathcal{E}_{tr} = \Pi(\Theta_{tr})$ . Then for the  $\ell_{sq}$  loss,

(i) 
$$\mathcal{I}(\mathcal{E}_{\mathrm{tr}}) = \mathcal{I}(\mathcal{E}_{\Pi})$$
 and (ii)  $\mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\mathrm{tr}}) = \mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\Pi})$ .

*Proof.* The proof is similar to that of Proposition 8.

**Part (i).** We have  $\mathcal{I}(\mathcal{E}_{\Pi}) \subseteq \mathcal{I}(\mathcal{E}_{\mathrm{tr}})$  by definition. To show the converse, consider any  $f = w \circ \varphi \in \mathcal{I}(\mathcal{E}_{\mathrm{tr}})$ , with  $\varphi$  invariant over  $\mathcal{E}_{\mathrm{tr}}$ . For any z in the range of  $\varphi$ , consider the function

$$g_z(\theta) \ = \ \mathbb{E}_{\Pi(\theta)}[Y \mid \varphi(X) = z] \ = \ \frac{\int_{\varphi^{-1}(z) \times \mathcal{Y}} \ y \, \mathrm{d}\Pi_\theta(x,y)}{\int_{\varphi^{-1}(z) \times \mathcal{Y}} \ \mathrm{d}\Pi_\theta(x,y)}.$$

Let  $n_z(\theta)$  and  $d_z(\theta)$  denote the numerator and denominator of  $g_z(\theta)$ , respectively, both of which are analytic in  $\theta$  by Definition 11 (and boundedness of y). By Observation 2, there exists a constant  $\alpha$  such that for all  $\theta \in \Theta_{\text{tr}}$  that satisfy  $d_z(\theta) \neq 0$ , it holds that

$$g_z(\theta) = \frac{n_z(\theta)}{d_z(\theta)} = \alpha \implies h_z(\theta) := n_z(\theta) - \alpha \cdot d_z(\theta) = 0.$$

Moreover,  $d_z(\theta) = 0$  implies  $n_z(\theta) = 0$ , hence  $h_z(\theta) = 0$  for all  $\theta \in \Theta_{tr}$ . Since  $\Theta_{tr}$  has non-zero Lebesgue measure, it follows from Fact 10 that  $h_z(\theta)$  that is identically zero on  $\Theta$ . This implies for all  $\theta \in \Theta$  such that  $d_z(\theta) \neq 0$ ,  $g_z(\theta) = \alpha$ . Hence, by Observation 2, we get that  $f \in \mathcal{I}(\mathcal{E}_{\Pi})$ .

**Part (ii).** This follows similarly. We have  $\mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\Pi}) \subseteq \mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\mathrm{tr}})$  by definition. To show the converse, consider any predictor  $f = 1 \cdot \varphi \in \mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\mathrm{tr}})$ , and consider the following function of  $\theta$ :

$$g(\theta) := \nabla_{w|w=1} \mathcal{L}_{\Pi(\theta)}(w \cdot \varphi) = \int_{\mathcal{X} \times \mathcal{Y}} \nabla_{w|w=1} \ell(w \cdot \varphi(x), y) \, d\Pi_{\theta}(x, y) \,,$$

## Does Invariant Risk Minimization Capture Invariance?

which by Definition 11 is an analytic function in  $\theta$ . To derive this, we swapped the  $\nabla_{w|w=1}$  with  $\int_{\mathcal{X}\times\mathcal{Y}}$ , possible because |y| and  $|\varphi(x)|$  are uniformly bounded (Planet Math 2013, Theorem 2).

Since  $(\nabla_w)$  holds for all  $e \in \mathcal{E}_{tr}$ ,  $g(\theta) = 0$  for all  $\theta \in \Theta_{tr}$ . But since  $\Theta_{tr}$  has non-zero Lebesgue measure in  $\mathbb{R}^d$ , we have from Fact 10 that g is identically 0 on  $\Theta$ , hence  $f \in \mathcal{I}_{\mathcal{S}}(\mathcal{E}_{\Pi})$ .