

---

# The Risks of Invariant Risk Minimization

---

Elan Rosenfeld, Pradeep Ravikumar, Andrej Risteski

Machine Learning Department

Carnegie Mellon University

elan@cmu.edu, pradeepr@cs.cmu.edu, aristeski@andrew.cmu.edu

## Abstract

Invariant Causal Prediction [30] is a technique for out-of-distribution generalization which assumes that some aspects of the data distribution vary across the training set but that the underlying causal mechanisms remain constant. Recently, Arjovsky et al. [1] proposed Invariant Risk Minimization (IRM), an objective based on this idea for learning deep, invariant features of data which are a complex function of latent variables; many alternatives have subsequently been suggested. However, formal guarantees for all of these works are severely lacking. In this paper, we present the first analysis of classification under the IRM objective—as well as these recently proposed alternatives—under a fairly natural and general model. In the linear case, we show simple conditions under which the optimal solution succeeds or, more often, fails to recover the optimal invariant predictor. We furthermore present the very first results in the non-linear regime: we demonstrate that IRM can fail catastrophically unless the test data are sufficiently similar to the training distribution—this is precisely the issue that it was intended to solve. Thus, in this setting we find that IRM and its alternatives fundamentally *do not improve* over standard Empirical Risk Minimization.

## 1 Introduction

Prediction algorithms are evaluated by their performance on unseen test data. In classical machine learning, it is common to assume that such data are drawn i.i.d. from the same distribution as the data set on which the learning algorithm was trained—in the real world, however, this is often not the case. When this discrepancy occurs, algorithms with strong in-distribution generalization guarantees, such as Empirical Risk Minimization (ERM), can fail catastrophically. In particular, while deep neural networks achieve superhuman performance on many tasks, there is evidence that they rely on statistically informative but non-causal features in the data [3, 12, 16]. As a result, such models are prone to errors under surprisingly minor distribution shift [37, 31]. To address this, researchers have investigated alternative objectives for training predictors which are robust to possibly egregious shifts in the test distribution.

The task of generalizing under such shifts, known as *Out-of-Distribution (OOD) Generalization*, has led to many separate threads of research. One approach is Bayesian deep learning, accounting for a classifier’s uncertainty at test time [27]. Another technique that has shown promise is data augmentation—this includes both automated data modifications which help prevent overfitting [36] and specific counterfactual augmentations to ensure invariance in the resulting features [39, 19].

A strategy which has recently gained particular traction is Invariant Causal Prediction (ICP; [30]), which views the task of OOD generalization through the lens of causality. This framework assumes that the data are generated according to a Structural Equation Model (SEM; [8]), which consists of a set of so-called mechanisms or structural equations that specify variables given their parents. ICP assumes moreover that the data can be partitioned into *environments*, where each environment

corresponds to interventions on the SEM [28], but where the mechanism by which the target variable is generated via its direct parents is unaffected. The justification for this is that the causal mechanism of the data generating process of the target variable is unchanging but other effects can vary. Therefore, learning relationships that are invariant across environments ensures recovery of the invariant features which generalize under arbitrary interventions. In this work, we consider objectives that attempt to learn what we refer to as the “optimal invariant predictor”—this is the classifier which uses and is optimal with respect to only the invariant features in the SEM. By definition, such a classifier does not overfit to environment-specific properties of the data distribution, so it will generalize even under major distribution shift at test time. In particular, we focus our analysis on one of the more popular objectives, Invariant Risk Minimization (IRM; [1]), but our results can easily be extended to similar recently proposed alternatives.

Various works on invariant prediction [26, 13, 15, 32, 38, 9] consider regression in both the linear and non-linear setting, but they exclusively focus on learning with partially observed covariates or instrumental variables. Under such a condition, results from causal inference [24, 29] allow for formal guarantees of the identification of the invariant features, or at least a strict subset of them. With the rise of deep learning, more recent literature has developed objectives for learning invariant representations when the data are a non-linear function of unobserved latent factors, a common assumption when working with complex, high-dimensional data such as images. Causal discovery and inference with unobserved confounders or latents is a much harder problem [29], so while empirical results seem encouraging, these objectives are presented with few formal guarantees. IRM is one such objective for invariant representation learning. The goal of IRM is to learn a feature embedder such that the optimal linear predictor on top of these features is the same for every environment—the idea being that only the invariant features will have an optimal predictor that is invariant. Thus, IRM hopes to “extrapolate” to new test environments, unlike ERM which excels at “interpolating”. Recent works have pointed to shortcomings of IRM and have suggested modifications which they claim prevent these failures. However, these alternatives are compared in broad strokes, with little in the way of theory.

In this work, we present the first formal analysis of classification under the IRM objective under a fairly natural and general model which is similar to the original work. Our results show that despite being inspired by invariant prediction, this objective can frequently be expected to perform *no better than ERM*. In the linear setting, we present simple conditions under which solving to optimality succeeds or, more often, breaks down in recovering the optimal invariant predictor. We additionally demonstrate a fundamental failure case—under mild conditions, there exists a feasible point that uses only non-invariant features and achieves lower empirical risk than the optimal invariant predictor; thus it will appear as a more attractive solution, yet its reliance on non-invariant features mean it will completely fail to generalize. As corollaries, we present similar cases where Risk Extrapolation (REx; [20]) and similar objectives likewise fail. Furthermore, we present the *first results in the non-linear regime*: we demonstrate the existence of a classifier with exponentially small sub-optimality which nevertheless heavily relies on non-invariant features on most test inputs, resulting in worse-than-chance performance on distributions that are sufficiently dissimilar from the training environments. These findings strongly suggest that existing approaches to ICP for high-dimensional latent variable models do not cleanly achieve their stated objective and that future work would benefit from a more formal treatment.

## 2 Related work

Works on learning deep invariant representations vary considerably: some search for a *domain-invariant* representation [26, 11], i.e. invariance of the distribution  $p(\Phi(x))$ , but this is typically used for domain adaptation [5, 10, 41, 23], with assumed access to labeled or unlabeled data from the target distribution. Further, there is evidence that this may not be sufficient for domain adaptation [42, 18]. Other works instead hope to find representations that are *conditionally* domain-invariant, with invariance of  $p(\Phi(x) | y)$  [14, 22].

More recent works, including the objectives discussed in this paper, suggest invariance of the *feature-conditioned label distribution*. In particular, [1] only assume invariance of  $\mathbb{E}[y | \Phi(x)]$ ; follow-up works rely on a stronger assumption of invariance of  $p(y | \Phi(x))$  [20, 40, 17, 25, 4]. Though this approach has become popular in the last year, it is somewhat similar to the existing concept of *covariate shift* [35, 7], which considers the same setting. The main difference is that these more

recent works assume that the shifts in  $p(\Phi(x))$  occur between discrete, labeled environments, as opposed to more generally from train to test distributions.

### 3 Our Results

We consider an SEM with explicit separation of *invariant* features  $z_c$ , whose joint distribution with the label is fixed for all environments, and *environmental* features  $z_e$  (“non-invariant”), whose distribution can vary. We assume that data are drawn from a set of  $E$  training environments  $\mathcal{E} = \{e_1, e_2, \dots, e_E\}$  and that we know from which environment each sample is drawn. For a given environment  $e$ , the data are defined by the following process: first, a label  $y \in \{\pm 1\}$  is drawn according to a fixed probability:

$$y = \begin{cases} 1, & \text{w.p. } \eta \\ -1, & \text{otherwise.} \end{cases} \quad (1)$$

Next, both invariant features and environmental features are drawn according to a Gaussian<sup>1</sup>:

$$z_c \sim \mathcal{N}(y \cdot \mu_c, \sigma_c^2 I), \quad z_e \sim \mathcal{N}(y \cdot \mu_e, \sigma_e^2 I), \quad (2)$$

with  $\mu_c \in \mathbb{R}^{d_c}$ ,  $\mu_e \in \mathbb{R}^{d_e}$ —typically, for complex, high-dimensional data we would expect  $E < d_c \ll d_e$ . Finally, the observation  $x$  is generated as a function of the latent features:

$$x = f(z_c, z_e). \quad (3)$$

We assume  $f$  is injective, so that it is in principle possible to recover the latent features from the observations, i.e. there exists a function  $\Phi$  such that  $\Phi(f(z_c, z_e)) = [z_c, z_e]^T$ . We write the joint and marginal distributions as  $p^e(x, y, z_c, z_e)$ . When clear from context, we omit the specific arguments.

**Remarks on the model.** This model is natural and flexible; it generalizes several existing models used to analyze learning under the existence of adversarial distribution shift or non-invariant correlations [34, 33]. The fundamental facet of this model is the constancy of the invariant parameters  $\eta, \mu_c, \sigma_c^2, f$  across environments—the dependence of  $\mu_e, \sigma_e$  on the environment allows for varying distributions, while the true causal process remains unchanged. Here we make a few clarifying remarks:

- We do not impose any constraints on the model parameters. In particular, we do not assume a prior over the environmental parameters. Observe that  $\mu_c, \sigma_c^2$  are the same for all environments, hence the subscript indicates the invariant relationship. In contrast, with some abuse of notation, the environmental subscript is used to indicate both dependence on the environment and the index of the environment itself (e.g.,  $\mu_i$  represents the mean specific to environment  $i$ ).
- While we have framed the model as  $y$  causing  $z_c$ , the causation can just as easily be viewed in the other direction. The log-odds of  $y$  are a linear function of  $z_c$ —this matches logistic regression with an invariant regression vector  $\beta_c = 2\mu_c/\sigma_c^2$  and bias  $\beta_0 = \log \frac{\eta}{1-\eta}$ . We present the model as above to emphasize that the causal relationships between  $y$  and the  $z_c, z_e$  are a priori indistinguishable, and because we believe this direction is more intuitive.
- Assuming a constant label marginal is necessary—in logistic regression, the optimal bias term is equal to the log-odds of the marginal:  $\log \frac{\eta}{1-\eta}$ . Thus the optimal classifier can only be invariant if this value is the same for all environments.
- Modeling class-conditional means as direct opposites is made for clarity: it matches existing models and it greatly simplifies the calculations and results. None of our proofs require this condition, and it is straightforward to extend our results to the general case of arbitrary means.

<sup>1</sup>Note the deliberate choice to have  $z_e$  depend on  $y$ . Much work on this problem models spurious features which correlate with the label but are not causal. However, the term “spurious” is often applied incongruously; historically, a spurious correlation is one which (a) appears by chance and would disappear given enough samples or (b) is due to an unobserved confounder. In recent work, the term has been co-opted to refer to any feature that correlates with the label but does not cause it. Thus there is a subtle distinction: if we allow for *the label to cause the features* (e.g. as in natural images), the resulting correlation is *not* spurious. We therefore avoid using the term “spurious” in this work, opting instead for “non-invariant” or “environmental”.

We consider the setting where we are given infinite samples from each environment; this allows us to isolate the behavior of the objectives themselves, rather than finite-sample effects. Upon observing samples from this model, our objective is thus to learn a feature embedder  $\Phi$  and classifier<sup>2</sup>  $\hat{\beta}$  to minimize the risk on an unseen environment  $e$ :

$$\mathcal{R}^e(\Phi, \hat{\beta}) := \mathbb{E}_{(x,y) \sim p^e} \left[ \ell(\sigma(\hat{\beta}^T \Phi(x)), y) \right].$$

The function  $\ell$  can be any loss appropriate to classification: in this work we consider the logistic and the 0-1 loss. Note, however, that we are *not* hoping to minimize this risk in expectation over the environments; this is already accomplished via ERM or distributionally robust optimization (DRO; [2, 6]). Rather, we hope to extract and regress on invariant features while ignoring environmental features, such that our predictor generalizes to all unseen environments regardless of their parameters. In other words, the focus is on minimizing risk for the worst-case environment. We refer to the predictor which will minimize worst-case risk under arbitrary distribution shift as the *optimal invariant predictor*. To discuss this formally, we define precisely what we mean by this term.

**Definition 1.** Under the model described by Equations 1-3, the *optimal invariant predictor* is the predictor defined by the composition of a) the featurizer which recovers the invariant features and b) the classifier which is optimal with respect to those features:

$$\Phi^*(x) := \begin{bmatrix} I & 0 \\ 0 & 0 \end{bmatrix} \circ f^{-1}(x) = [z_c], \quad \hat{\beta}^* := \begin{bmatrix} \beta_c \\ \beta_0 \end{bmatrix} := \begin{bmatrix} 2\mu_c/\sigma_c^2 \\ \log \frac{\eta}{1-\eta} \end{bmatrix}.$$

Observe that this definition closely resembles Definition 3 of [1]; the only difference is that here the optimal invariant predictor must recover *all invariant features*. As [1] do not posit a data model, the concept of recovering “all invariant features” is not well-defined for their setting; technically, a featurizer which outputs the empty set would elicit an invariant predictor, but this would not satisfy the above definition. The classifier  $\hat{\beta}^*$  is optimal with respect to the invariant features and so it achieves the minimum possible risk without using environmental features. Observe that *the optimal invariant predictor is distinct from the Bayes classifier*; the Bayes classifier for a given environment uses environmental features which are informative of the label but non-invariant; the optimal invariant predictor *explicitly ignores* these features.

With the model defined, we can informally present our results; we defer the formal statements to first give a background on the IRM objective in the next section. With a slight abuse of notation, we identify a predictor by the tuple  $\Phi, \hat{\beta}$  which parametrizes it. First, we show that the usefulness of IRM exhibits a “thresholding” behavior depending on  $E$  and  $d_e$ :

**Theorem 3.1** (Informal, Linear). *For linear  $f$ , consider solving the IRM objective to learn a linear  $\Phi$  with invariant optimal classifier  $\hat{\beta}$ . If  $E > d_e$ , then  $\Phi, \hat{\beta}$  is precisely the optimal invariant predictor; it uses only invariant features and generalizes to all environments with minimax-optimal risk. If  $E \leq d_e$ , then  $\Phi, \hat{\beta}$  relies upon non-invariant features.*

In fact, when  $E \leq d_e$  it is even possible to learn a classifier *solely* relying on environmental features that achieves lower risk on the training environments than the optimal invariant predictor:

**Theorem 3.2** (Informal, Linear). *For linear  $f$  and  $E \leq d_e$  there exists a linear predictor  $\Phi, \hat{\beta}$  which uses only environmental features, yet achieves lower risk than the optimal invariant predictor.*

Finally, in the non-linear case, we show that IRM fails unless the training environments approximately “cover” the space of possible environments, and therefore it behaves similarly to ERM:

**Theorem 3.3** (Informal, Non-linear). *For non-linear  $f$ , there exists a non-linear predictor  $\Phi, \hat{\beta}$  which is nearly optimal under the penalized objective and furthermore is nearly identical to the optimal invariant predictor on the training distribution. However, for any test environment with a mean sufficiently different from the training means, this predictor will be equivalent to the ERM solution on nearly all test points. For test distributions where the environmental feature correlations with the label are reversed, this predictor has almost trivial performance.*

---

<sup>2</sup>Following the terminology of [1], we refer to the regression vector  $\hat{\beta}$  as a “classifier” and the composition of  $\Phi, \hat{\beta}$  as a “predictor”.

**Extensions to other objectives.** Several follow-up works have suggested alternatives to IRM (see Section 4). Though these objectives perform better on various baselines, there are few formal guarantees and no results beyond the linear case. Due to their similarities, it is simple to extend every theorem in this paper to these objectives, demonstrating that they all suffer from the same shortcomings. Appendix E contains example corollaries for each of the results presented in this work.

## 4 Background on Invariant Risk Minimization and its Alternatives

During training, a classifier will learn to leverage correlations between features and labels in the training data to make its predictions. If a correlation varies with the environment, it may not be present in future test distributions—worse yet, it may be *reversed*—harming the classifier’s predictive ability. IRM [1] is a recently proposed approach to learning environmentally invariant representations to facilitate invariant prediction.

**The IRM objective.** IRM posits the existence of a feature embedder  $\Phi$  such that the optimal classifier on top of these features is the same for every environment. The authors argue that such a function will use only invariant features, since non-invariant features will have different joint distributions with the label and therefore a fixed classifier on top of them won’t be optimal in all environments. To learn this  $\Phi$ , the IRM objective is the following constrained optimization problem:

$$\min_{\Phi, \hat{\beta}} \frac{1}{|\mathcal{E}|} \sum_{e \in \mathcal{E}} \mathcal{R}^e(\Phi, \hat{\beta}) \quad \text{s.t.} \quad \hat{\beta} \in \arg \min_{\beta} \mathcal{R}^e(\Phi, \beta) \quad \forall e \in \mathcal{E}. \quad (4)$$

This bilevel program is highly non-convex and difficult to solve. To find an approximate solution, the authors consider a Lagrangian form, whereby the sub-optimality with respect to the constraint is expressed as the squared norm of the gradients of each of the inner optimization problems:

$$\min_{\Phi, \hat{\beta}} \frac{1}{|\mathcal{E}|} \sum_{e \in \mathcal{E}} \left[ \mathcal{R}^e(\Phi, \hat{\beta}) + \lambda \|\nabla_{\hat{\beta}} \mathcal{R}^e(\Phi, \hat{\beta})\|_2^2 \right]. \quad (5)$$

Assuming the inner optimization problem is convex, achieving feasibility is equivalent to the penalty term being equal to 0. Thus, Equations 4 and 5 are equivalent if we set  $\lambda = \infty$ .

**Alternative objectives.** IRM is motivated by the existence of a featurizer  $\Phi$  such that  $\mathbb{E}[y | \Phi(x)]$  is invariant. Follow-up works have proposed variations on this objective, based instead on the strictly stronger desideratum of the invariance of  $p(y | \Phi(x))$ . [20] suggest penalizing the variance of the risks, while Xie et al. [40] give the same objective but taking the square root of the variance. Many papers have suggested similar alternatives [17, 25, 4]. These objectives are compelling—indeed, it is easy to show that the optimal invariant predictor constitutes a stationary point of each of these objectives:

**Proposition 4.1.** *Suppose the observed data are generated according to Equations 1-3. Then the (parametrized) optimal invariant predictor  $\Phi^*, \hat{\beta}^*$  is a stationary point for Equation 4.*

The stationarity of the optimal invariant predictor for the other objectives is a trivial corollary. However, in the following sections we will demonstrate that such a result is misleading and that a more careful investigation is necessary.

## 5 The Difficulties of IRM in the Linear Regime

In their work proposing IRM, [1] present specific conditions for an upper bound on the number of training environments needed such that a feasible linear featurizer  $\Phi$  will have an invariant optimal regression vector  $\hat{\beta}$ . Our first result is similar in spirit but presents a substantially stronger (and simplified) upper bound in the classification setting, along with a matching lower bound: we demonstrate that observing a large number of environments—linear in the number of environmental features—is *necessary* for generalization in the linear regime.

**Theorem 5.1** (Linear case). *Assume  $f$  is linear. Suppose we observe  $E$  training environments. Then the following hold:*

1. *Suppose  $E > d_e$ . Consider any linear featurizer  $\Phi$  which is feasible under the IRM objective (4), with invariant optimal classifier  $\hat{\beta} \neq 0$ , and write  $\Phi(f(z_c, z_e)) = Az_c + Bz_e$ . Then under mild non-degeneracy conditions, it holds that  $B = 0$ . Consequently,  $\hat{\beta}$  is the optimal classifier for all possible environments.*
2. *If  $E \leq d_e$  and the environmental means  $\mu_e$  are linearly independent, then there exists a linear  $\Phi$ —where  $\Phi(f(z_c, z_e)) = Az_c + Bz_e$  with  $\text{rank}(B) = d_e + 1 - E$ —which is feasible under the IRM objective. Further, both the logistic and 0-1 risks of this  $\Phi$  and its corresponding optimal  $\hat{\beta}$  are strictly lower than those of the optimal invariant predictor.*

Similar to [1], the set of environments which do not satisfy Theorem 5.1 has measure zero under any absolutely continuous density over environmental parameters. Further details, and the full proof, can be found in Appendix C.1. Since the optimal invariant predictor is Bayes with respect to the invariant features, by the data-processing inequality the only way a predictor can achieve lower risk is by relying on environmental features. Thus, Theorem 5.1 directly implies that when  $E \leq d_e$ , the global minimum necessarily uses these non-invariant features and therefore will not universally generalize to unseen environments. On the other hand, in the (perhaps unlikely) case that  $E > d_e$ , any feasible solution will generalize, and the optimal invariant predictor has the minimum (and minimax) risk of all such predictors:

**Corollary 5.2.** *For both logistic and 0-1 loss, the optimal invariant predictor is the global minimum of the IRM objective if and only if  $E > d_e$ .*

Let us compare our theoretical findings to those of [1]. Suppose the observations  $x$  lie in  $\mathbb{R}^d$ . Roughly, their theorem says that for a learned  $\Phi$  of rank  $r$  with invariant optimal coefficient  $\hat{\beta}$ , if the training set contains  $d - r + d/r$  “non-degenerate” environments, then  $\hat{\beta}$  will be optimal for all environments. There are several important issues with this result: first, they present no result tying the rank of  $\Phi$  to their actual objective; their theory thus motivates the objective, but does not provide any performance guarantees for its solution. Next, observe when  $x$  is high-dimensional (i.e.  $d \gg d_e + d_c$ )—in which case  $\Phi$  will be low-rank (i.e.  $r \leq d_e + d_c$ )—their result requires  $\Omega(d)$  environments, which is extreme. For example, think of images on a low-dimensional manifold embedded in very high-dimensional space. Even when  $d = d_c + d_e$ , the “ideal”  $\Phi$  which recovers precisely  $z_c$  would have rank  $d_c$ , and therefore their condition for invariance would require  $E > d_e + d_e/d_c$ , a stronger requirement than ours; this inequality also seems unlikely to hold in most real-world settings. Finally, they give no lower bound on the number of required environments—prior to this work, there were no existing results for the performance of the IRM objective when their conditions are not met. We also run a simple synthetic experiment to verify our theoretical results, drawing samples according to our model and learning a predictor with the IRM objective. Details and results of this experiment can be found in Appendix C.2. We now sketch a constructive proof of part 2 of the theorem for when  $E = d_e$ :

*Proof Sketch.* Since  $f$  has an inverse over its range, we can define  $\Phi$  as a linear function directly over the latents  $[z_c, z_e]$ . Specifically, define  $\Phi(x) = [z_c, p^T z_e]$ . Here,  $p$  is a unit-norm vector such that  $\forall e \in \mathcal{E}$ ,  $p^T \mu_e = \sigma_e^2 \tilde{\mu}$ ;  $\tilde{\mu}$  is a fixed scalar that depends on the geometry of  $\mu_e, \sigma_e^2$ —such a vector exists so long as the means are linearly independent. Observe that this  $\Phi$  also has the desired rank. Since this is a linear function of a multivariate Gaussian, the label-conditional distribution of each environment’s non-invariant latents has a simple closed form:  $p^T z_e \mid y \sim \mathcal{N}(y \cdot p^T \mu_e, \|p\|_2^2 \sigma_e^2) \stackrel{d}{=} \mathcal{N}(y \cdot \sigma_e^2 \tilde{\mu}, \sigma_e^2)$ .

For separating two Gaussians, the optimal linear classifier is  $\Sigma^{-1}(\mu_1 - \mu_0)$ —here, the optimal classifier on  $p^T z_e$  is precisely  $2\tilde{\mu}$ , which does not depend on the environment (and neither do the optimal coefficients for  $z_c$ ). Though the distribution varies across environments, the optimal classifier is the same! Thus,  $\Phi$  directly depends on the environmental features, yet the optimal regression vector  $\hat{\beta}$  for each environment is constant. To see that it has lower risk than the optimal invariant predictor, note that this classifier is Bayes with respect to its features and that the optimal invariant predictor uses a strict subset of these features, and therefore it has less information for its predictions.  $\square$

**A purely environmental predictor.** The precise value of  $\tilde{\mu}$  in the proof sketch above represents how strongly this non-invariant feature is correlated with the label. In theory, a predictor that achieves

a lower objective value could do so by a very small margin—incorporating an arbitrarily small amount of information from a non-invariant feature would suffice. This result would be less surprising, since achieving low empirical risk might still ensure that we are “close” to the optimal invariant predictor. Our next result shows that this is not the case: there exists a feasible solution which uses *only the environmental features* yet performs better than the optimal invariant predictor on all  $e \in \mathcal{E}$  for which  $\tilde{\mu}$  is large enough.

**Theorem 5.3.** *Suppose we observe  $E \leq d_e$  environments, such that all environmental means are linearly independent. Then there exists a feasible  $\Phi, \hat{\beta}$  which uses only environmental features and achieves lower 0-1 risk than the optimal invariant predictor on every environment  $e$  such that  $\sigma_e \tilde{\mu} > \sigma_c^{-1} \|\mu_c\|_2$  and  $2\sigma_e \tilde{\mu} \sigma_c^{-1} \|\mu_c\|_2 \geq |\beta_0|$ .*

The latter of these two conditions is effectively trivial, requiring only a small separation of the means and balance in class labels. From the construction of  $\tilde{\mu}$  in the proof of Lemma C.1, we can see that the former condition is more likely to be met when  $E \ll d_e$  and in environments where some non-invariant features are reasonably correlated with the label—both of which can be expected to hold in the high-dimensional setting. Figure C.2 in the appendix plots the results for a few toy examples for various dimensionalities and variances to see how often this condition holds in practice. For all settings, the number of environments observed before the condition ceases to hold is quite high—on the order of  $d_e - d_c$ .

## 6 The Failure of IRM in the Non-Linear Regime

We’ve demonstrated that OOD generalization is difficult in the linear case, but it is achievable given enough training environments. Our results—and those of [1]—intuitively proceed by observing that each environment reduces a “degree of freedom” of the solution, such that only the invariant features remain feasible if enough environments are seen. In the non-linear case, it’s not clear how to capture this idea of restricting the “degrees of freedom”—and in fact our results imply that this intuition is simply wrong. Instead, we show that the solution generalizes only to test environments that are sufficiently similar to the training environments. Thus, these objectives present no real improvement over ERM or DRO.

Non-linear transformations of the latent variables make it hard to characterize the optimal linear classifier, which makes reasoning about the constrained solution to Equation 4 difficult. Instead we turn our attention to Equation 5, the penalized IRM objective. In this section we demonstrate a foundational flaw of IRM in the non-linear regime: unless we observe enough environments to “cover” the space of non-invariant features, a solution that appears to be invariant can still wildly underperform on a new test distribution. We begin with a definition about the optimality of a coefficient vector  $\hat{\beta}$ :

**Definition 2.** For  $0 < \gamma < 1$ , a coefficient vector  $\hat{\beta}$  is  $\gamma$ -close to optimal for a label-conditional feature distribution  $z \sim \mathcal{N}(y \cdot \mu, \Sigma)$  if

$$\hat{\beta}^T \mu \geq (1 - \gamma) 2\mu^T \Sigma^{-1} \mu.$$

Since the optimal coefficient vector is precisely  $2\Sigma^{-1}\mu$ , being  $\gamma$ -close implies that  $\hat{\beta}$  is reasonably aligned with that optimum. Observe that the definition does not account for magnitude—the set of vectors which is  $\gamma$ -close to optimal is therefore a halfspace which is normal to the optimal vector. One of our results in the non-linear case uses the following assumption, which says that the observed environmental means are sufficiently similar to one another.

**Assumption 1.** There exists a  $0 \leq \gamma < 1$  such that the ERM-optimal classifier for the non-invariant features,

$$\beta_{e;\text{ERM}} := \arg \min_{\hat{\beta}_e} \frac{1}{|\mathcal{E}|} \sum_{e \in \mathcal{E}} \mathbb{E}_{z_c, z_e, y \sim p^e} \left[ \ell(\sigma(\beta_c^T z_c + \hat{\beta}_e^T z_e + \beta_0), y) \right], \quad (6)$$

is  $\gamma$ -close to optimal for every environmental feature distribution in  $\mathcal{E}$ .

This assumption says the environmental distributions are similar enough such that the optimal “average classifier” is reasonably predictive for each environment individually. This is a natural expectation: we are employing IRM precisely *because* we expect the ERM classifier to do well on

the training set but fail to generalize. If the environmental parameters are sufficiently orthogonal, we might instead expect ERM to ignore the features which are not at least moderately predictive across all environments. Finally, we note that if this assumption only holds for a subset of features, our result still applies by marginalizing out the dimensions for which it does not hold.

We are now ready to give our main result in the non-linear regime. We present a simplified version, assuming that  $\sigma_e^2 = 1 \forall e$ . This is purely for clarity of presentation; the full theorem, for which the result still holds, is presented in Appendix D. We make use of two constants in the following proof—the average squared norm of the environmental means,  $\overline{\|\mu\|_2^2} := \frac{1}{E} \sum_{e \in \mathcal{E}} \|\mu_e\|_2^2$ ; and the standard deviation of the response variable of the ERM-optimal classifier,  $\sigma_{\text{ERM}} := \sqrt{\|\beta_c\|_2^2 \sigma_c^2 + \|\beta_{e;\text{ERM}}\|_2^2 \sigma_e^2}$ .

**Theorem 6.1** (Non-linear case, simplified). *Suppose we observe  $E$  environments  $\mathcal{E} = \{e_1, e_2, \dots, e_E\}$ , where  $\sigma_e^2 = 1 \forall e$ . Then, for any  $\epsilon > 1$ , there exists a featurizer  $\Phi_\epsilon$  which, combined with the ERM-optimal classifier  $\hat{\beta} = [\beta_c, \beta_{e;\text{ERM}}, \beta_0]^T$ , satisfies the following properties, where we define  $p_{\epsilon, d_e} := \exp\{-d_e \min((\epsilon - 1), (\epsilon - 1)^2)/8\}$ :*

1. The regularization term of  $\Phi_\epsilon, \hat{\beta}$  as in Equation 5 is bounded as

$$\frac{1}{E} \sum_{e \in \mathcal{E}} \|\nabla_{\hat{\beta}} \mathcal{R}^e(\Phi_\epsilon, \hat{\beta})\|_2^2 \in \mathcal{O}\left(p_{\epsilon, d_e}^2 (c_\epsilon d_e + \overline{\|\mu\|_2^2})\right),$$

for some constant  $c_\epsilon$  that depends only on  $\epsilon$ .

2.  $\Phi_\epsilon, \hat{\beta}$  exactly matches the optimal invariant predictor on at least a  $1 - p_{\epsilon, d_e}$  fraction of the training set. On the remaining inputs, it matches the ERM-optimal solution.

Further, for any test distribution, suppose its environmental mean  $\mu_{E+1}$  is sufficiently far from the training means:

$$\forall e \in \mathcal{E}, \min_{y \in \{\pm 1\}} \|\mu_{E+1} - y \cdot \mu_e\|_2 \geq (\sqrt{\epsilon} + \delta) \sqrt{d_e} \quad (7)$$

for some  $\delta > 0$ , and define  $q := \frac{2E}{\sqrt{\pi}\delta} \exp\{-\delta^2\}$ . Then the following holds:

3.  $\Phi_\epsilon, \hat{\beta}$  is equivalent to the ERM-optimal predictor on at least a  $1 - q$  fraction of the test distribution.

4. Under Assumption 1, suppose it holds that  $\mu_{E+1} = -\sum_{e \in \mathcal{E}} \alpha_e \mu_e$  for some set of coefficients  $\{\alpha_e\}_{e \in \mathcal{E}}$ . Then so long as

$$\sum_{e \in \mathcal{E}} \alpha_e \|\mu_e\|_2^2 \geq \frac{\|\mu_c\|_2^2 / \sigma_c^2 + |\beta_0|/2 + \sigma_{\text{ERM}}}{1 - \gamma}, \quad (8)$$

the 0-1 risk of  $\Phi_\epsilon, \hat{\beta}$  on the new environment is greater than  $.975 - q$ .

Before giving a proof sketch, we give a brief intuition for each of the claims made in this theorem:

1. The first claim says that the predictor we construct will have a gradient squared norm scaling as  $p_{\epsilon, d_e}^2$  which is exponentially small in  $d_e$ . Thus, in high dimensions, it will appear as a perfectly reasonable solution to the objective (5).
2. The second claim says that this predictor is identical to the invariant optimal predictor on all but an exponentially small fraction of the training data; on the remaining fraction, it matches the ERM-optimal solution, which has lower risk. The correspondence between constrained and penalized optimization implies that for large enough  $d_e$ , the “fake” predictor will often be a preferred solution. In the finite-sample setting, we would need exponentially many samples to even distinguish between the two!
3. The third claim is the crux of the theorem; it says that this predictor we’ve constructed will completely fail to use invariant prediction on most environments. Recall, the intent of IRM is to perform well precisely when ERM breaks down: when the test distribution differs greatly from the training distribution. Assuming a Gaussian prior on the training environment means, they will have separation in  $\mathcal{O}(\sqrt{d_e})$  with high probability. Observe that  $q$  will be vanishingly small so

long as  $\delta \geq \text{polylog}(E)$ . Part 3 says that IRM fails to use invariant prediction on any environment that is slightly outside the high probability region of the prior; even a separation of  $\Omega(\sqrt{d_e \log E})$  suffices. If we expect the new environments to be similar, ERM already guarantees reasonable performance at test-time; thus, IRM fundamentally *does not improve* over ERM in this regime.

4. The final statement demonstrates a particularly egregious failure case of this predictor: just like ERM, if the correlation between the non-invariant features and the label reverses at test-time, our predictor will have significantly worse than chance performance.

With this intuition, we now present a sketch of the proof—the full proof can be found in Appendix D.

*Proof Sketch.* The key is a construction which is almost identical to the optimal invariant predictor on the training data, yet behaves like the ERM solution at test time. We partition the environmental feature space into two sets,  $\mathcal{B}, \mathcal{B}^c \subset \mathbb{R}^{d_e}$ .  $\mathcal{B}$  is the union of balls centered at each  $\mu_e$ , each with a large enough radius that it contains most of the samples from that environment; thus  $\mathcal{B}$  represents the vast majority of the training distribution. On this set, define  $\Phi(x) = [z_c]$ , so our construction is equal to the optimal invariant predictor. Now consider  $\mathcal{B}^c = \mathbb{R}^{d_e} \setminus \mathcal{B}$ . We use standard concentration results to upper bound the measure of  $\mathcal{B}^c$  under the training distribution by  $p_{\epsilon, d_e}$ . Next, we show how choosing  $\Phi(x) = f^{-1}(x) = [z_c, z_e]^T$  on this set results in the sub-optimality bound, which is of order  $p_{\epsilon, d_e}^2$ . It is also clear that our constructed predictor is equivalent to the ERM-optimal solution on  $\mathcal{B}^c$ . Thus, our predictor will often have lower empirical risk on  $\mathcal{B}^c$ , countering the regularization penalty.

The second part of the proof shows that while  $\mathcal{B}$  has large measure under the training environments, it will have very small measure under any moderately different test environment. We can see this by considering the separation of means (Equation 7); the measure of each ball in  $\mathcal{B}$  can be bounded by the measure of the halfspace containing it; if each ball is far enough away from  $\mu_{E+1}$ , then the total measure of these halfspaces must be small. At test time, our predictor will therefore match the ERM solution on all but  $q$  of the observations (part 3). Finally, we lower bound the 0-1 risk of the ERM predictor under such a distribution shift by analyzing the distribution of the response variable. The proof is completed by observing that our predictor’s risk can differ from this by at most  $q$ .  $\square$

Theorem 6.1 shows that it’s possible for the IRM solution to perform poorly on environments which differ even moderately from the training data. We can of course guarantee generalization if the training distributions “cover” (or approximately cover) the full space of environments in order to tie down the performance on future distributions. But in such a scenario, there would no longer be a need for ICP; we could expect ERM or DRO to perform just as well. Once more, we find that our result trivially extends to the alternative objectives; we again refer to Appendix E.

## 7 Conclusion

Out-of-distribution generalization is an important direction for future research, and Invariant Causal Prediction remains a promising approach. However, formal results for latent variable models are lacking, particularly in the non-linear setting with fully unobserved covariates. This paper demonstrates that Invariant Risk Minimization and subsequent works have significant under-explored risks and issues with their formulation. This raises the question: what is the correct formulation for invariant prediction when the observations are complex, non-linear functions of unobserved latent factors? It would be interesting to investigate ours or similar models further; some possible directions include (a) characterizing in which settings specifically an “invariance-like” constraint can lead to improved performance over ERM and (b) formulating an objective that encourages invariance in the non-linear setting in a way that can be formally demonstrated and quantified. We hope that this work will inspire further theoretical study on the effectiveness of IRM and similar objectives.

## References

- [1] Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization. *arXiv preprint arXiv:1907.02893*, 2019.
- [2] J Andrew Bagnell. Robust supervised learning. In *Proceedings of the 20th national conference on Artificial intelligence-Volume 2*, pages 714–719, 2005.
- [3] Sara Beery, Grant Van Horn, and Pietro Perona. Recognition in terra incognita. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 456–473, 2018.
- [4] Alexis Bellot and Mihaela van der Schaar. Generalization and invariances in the presence of unobserved confounding. *arXiv preprint arXiv:2007.10653*, 2020.
- [5] Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman Vaughan. A theory of learning from different domains. *Machine learning*, 79(1-2):151–175, 2010.
- [6] Aharon Ben-Tal, Laurent El Ghaoui, and Arkadi Nemirovski. *Robust optimization*, volume 28. Princeton University Press, 2009.
- [7] Steffen Bickel, Michael Brückner, and Tobias Scheffer. Discriminative learning under covariate shift. *Journal of Machine Learning Research*, 10(9), 2009.
- [8] Kenneth A Bollen. Structural equation models. *Encyclopedia of biostatistics*, 7, 2005.
- [9] Rune Christiansen, Niklas Pfister, Martin Emil Jakobsen, Nicola Gnecco, and Jonas Peters. A causal framework for distribution generalization. *arXiv preprint arXiv:2006.07433*, 2020.
- [10] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In *International conference on machine learning*, pages 1180–1189. PMLR, 2015.
- [11] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *The Journal of Machine Learning Research*, 17(1):2096–2030, 2016.
- [12] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. *arXiv preprint arXiv:1811.12231*, 2018.
- [13] AmirEmad Ghassami, Saber Salehkaleybar, Negar Kiyavash, and Kun Zhang. Learning causal structures using regression invariance. In *Advances in Neural Information Processing Systems*, pages 3011–3021, 2017.
- [14] Mingming Gong, Kun Zhang, Tongliang Liu, Dacheng Tao, Clark Glymour, and Bernhard Schölkopf. Domain adaptation with conditional transferable components. In *International conference on machine learning*, pages 2839–2848, 2016.
- [15] Christina Heinze-Deml, Jonas Peters, and Nicolai Meinshausen. Invariant causal prediction for nonlinear models. *Journal of Causal Inference*, 6(2), 2018.
- [16] Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander Madry. Adversarial examples are not bugs, they are features. In H. Wallach, H. Larochelle, A. Beygelzimer, F. dAlché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 125–136. Curran Associates, Inc., 2019.
- [17] Wengong Jin, Regina Barzilay, and Tommi Jaakkola. Domain extrapolation via regret minimization. *arXiv preprint arXiv:2006.03908*, 2020.
- [18] Fredrik D Johansson, David Sontag, and Rajesh Ranganath. Support and invertibility in domain-invariant representations. *arXiv preprint arXiv:1903.03448*, 2019.
- [19] Divyansh Kaushik, Eduard Hovy, and Zachary Lipton. Learning the difference that makes a difference with counterfactually-augmented data. In *International Conference on Learning Representations*, 2020.

- [20] David Krueger, Ethan Caballero, Joern-Henrik Jacobsen, Amy Zhang, Jonathan Binas, Remi Le Priol, and Aaron Courville. Out-of-distribution generalization via risk extrapolation (rex). *arXiv preprint arXiv:2003.00688*, 2020.
- [21] Frank R Kschischang. The complementary error function. 2017. URL <https://www.commtoronto.ca/frank/notes/erfc.pdf>.
- [22] Ya Li, Xinmei Tian, Mingming Gong, Yajing Liu, Tongliang Liu, Kun Zhang, and Dacheng Tao. Deep domain generalization via conditional invariant adversarial networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 624–639, 2018.
- [23] Mingsheng Long, Zhangjie Cao, Jianmin Wang, and Michael I Jordan. Conditional adversarial domain adaptation. In *Advances in Neural Information Processing Systems*, pages 1640–1650, 2018.
- [24] Marloes H Maathuis, Markus Kalisch, Peter Bühlmann, et al. Estimating high-dimensional intervention effects from observational data. *The Annals of Statistics*, 37(6A):3133–3164, 2009.
- [25] Divyat Mahajan, Shruti Tople, and Amit Sharma. Domain generalization using causal matching. *arXiv preprint arXiv:2006.07500*, 2020.
- [26] Krikamol Muandet, David Balduzzi, and Bernhard Schölkopf. Domain generalization via invariant feature representation. volume 28 of *Proceedings of Machine Learning Research*, pages 10–18, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR. URL <http://proceedings.mlr.press/v28/muandet13.html>.
- [27] Radford M Neal. *Bayesian learning for neural networks*, volume 118. Springer Science & Business Media, 2012.
- [28] Judea Pearl. *Causality*. Cambridge university press, 2009.
- [29] J Peters, D Janzing, and B Schölkopf. Elements of causal inference-foundations and learning algorithms. 2017.
- [30] Jonas Peters, Peter Bühlmann, and Nicolai Meinshausen. Causal inference by using invariant prediction: identification and confidence intervals. *Journal of the Royal Statistical Society*, 2016.
- [31] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet classifiers generalize to imagenet? *arXiv preprint arXiv:1902.10811*, 2019.
- [32] Mateo Rojas-Carulla, Bernhard Schölkopf, Richard Turner, and Jonas Peters. Invariant models for causal transfer learning. *The Journal of Machine Learning Research*, 19(1):1309–1342, 2018.
- [33] Shiori Sagawa, Aditi Raghunathan, Pang Wei Koh, and Percy Liang. An investigation of why overparameterization exacerbates spurious correlations. In *Proceedings of the 37th International Conference on Machine Learning*, 2020.
- [34] Ludwig Schmidt, Shibani Santurkar, Dimitris Tsipras, Kunal Talwar, and Aleksander Madry. Adversarially robust generalization requires more data. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages 5014–5026. Curran Associates, Inc., 2018.
- [35] Hidetoshi Shimodaira. Improving predictive inference under covariate shift by weighting the log-likelihood function. *Journal of statistical planning and inference*, 90(2):227–244, 2000.
- [36] Connor Shorten and Taghi M Khoshgoftaar. A survey on image data augmentation for deep learning. *Journal of Big Data*, 6(1):60, 2019.
- [37] Jiawei Su, Danilo Vasconcellos Vargas, and Kouichi Sakurai. One pixel attack for fooling deep neural networks. *IEEE Transactions on Evolutionary Computation*, 23(5):828–841, 2019.

- [38] Adarsh Subbaswamy, Peter Schulam, and Suchi Saria. Preventing failures due to dataset shift: Learning predictive models that transport. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 3118–3127. PMLR, 2019.
- [39] Riccardo Volpi, Hongseok Namkoong, Ozan Sener, John C Duchi, Vittorio Murino, and Silvio Savarese. Generalizing to unseen domains via adversarial data augmentation. In *Advances in neural information processing systems*, pages 5334–5344, 2018.
- [40] Chuanlong Xie, Fei Chen, Yue Liu, and Zhenguo Li. Risk variance penalization: From distributional robustness to causality. *arXiv preprint arXiv:2006.07544*, 2020.
- [41] Kun Zhang, Mingming Gong, and Bernhard Scholkopf. Multi-source domain adaptation: a causal view. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, pages 3150–3157, 2015.
- [42] Han Zhao, Remi Tachet Des Combes, Kun Zhang, and Geoffrey Gordon. On learning invariant representations for domain adaptation. In *International Conference on Machine Learning*, pages 7523–7532, 2019.