

Out-of-Distribution Generalization via Risk Extrapolation

David Krueger^{1,2} Ethan Caballero^{1,2} Joern-Henrik Jacobsen^{3,4} Amy Zhang^{1,5,6} Jonathan Binas^{1,2}
Dinghuai Zhang^{1,2} Remi Le Priol^{1,2} Aaron Courville^{1,2}

Abstract

Distributional shift is one of the major obstacles when transferring machine learning prediction systems from the lab to the real world. To tackle this problem, we assume that variation across training domains is representative of the variation we might encounter at test time, but also that shifts at test time may be more extreme in magnitude. In particular, we show that reducing differences in risk across training domains can reduce a model’s sensitivity to a wide range of extreme distributional shifts, including the challenging setting where the input contains both causal and anti-causal elements. We motivate this approach, **Risk Extrapolation (REx)**, as a form of robust optimization over a perturbation set of extrapolated domains (MM-REx), and propose a penalty on the variance of training risks (V-REx) as a simpler variant. We prove that variants of REx can recover the causal mechanisms of the targets, while also providing some robustness to changes in the input distribution (“covariate shift”). By trading-off robustness to causally induced distributional shifts and covariate shift, REx is able to outperform alternative methods such as Invariant Risk Minimization in situations where these types of shift co-occur.

1. Introduction

While neural networks often exhibit super-human generalization on the training distribution, they can be extremely sensitive to distributional shift, presenting a major roadblock for their practical application (Su et al., 2019; Engstrom et al., 2017; Recht et al., 2019; Hendrycks & Dietterich, 2019). This sensitivity is often caused by relying on “spurious” features unrelated to the core concept we are trying to learn (Geirhos et al., 2018). For instance, Beery et al. (2018) give the example of an image recognition model failing to correctly classify cows on the beach, since it has learned to

make predictions based on the features of the background (e.g. a grassy field) instead of just the animal.

In this work, we consider **out-of-distribution (OOD) generalization**, also known as **domain generalization**, where a model must generalize appropriately to a new test domain for which it has neither labeled nor unlabeled training data. Following common practice (Ben-Tal et al., 2009), we formulate this as optimizing the worst-case performance over a perturbation set of possible test domains, \mathcal{F} :

$$\mathcal{R}_{\mathcal{F}}^{\text{OOD}}(\theta) = \max_{e \in \mathcal{F}} \mathcal{R}_e(\theta) \quad (1)$$

Since generalizing to arbitrary test domains is impossible, the choice of perturbation set encodes our assumptions about which test domains might be encountered. Instead of making such assumptions *a priori*, we assume access to data from multiple training domains, which can inform our choice of perturbation set. A classic approach for this setting is **group distributionally robust optimization (DRO)** (Sagawa et al., 2019), where \mathcal{F} contains all mixtures of the training distributions. This is mathematically equivalent to considering convex combinations of the training risks.

However, we aim for a more ambitious form of OOD generalization, over a larger perturbation set. Our method **minimax Risk Extrapolation (MM-REx)** is an extension of DRO where \mathcal{F} instead contains affine combinations of training risks, see Figure 1. Under specific circumstances, MM-REx can be thought of as DRO over a set of extrapolated domains.¹ But MM-REx also unlocks fundamental new generalization capabilities unavailable to DRO.

In particular, focusing on supervised learning, we show that Risk Extrapolation can uncover **invariant relationships** between inputs X and targets Y . Intuitively, an **invariant relationship** is a statistical relationship which is maintained across all domains in \mathcal{F} . Returning to the cow-on-the-beach example, the relationship between the animal and the label is expected to be invariant, while the relationship between the background and the label is not. A model which bases its predictions on such an invariant relationship is said to perform **invariant prediction**.²

¹Mila ²University of Montreal ³Vector ⁴University of Toronto
⁵McGill University ⁶Facebook AI Research. Correspondence to: <david.scott.krueger@gmail.com>.

¹We define “extrapolation” to mean “outside the convex hull”, see Appendix B for more.

²Note this is different from learning an invariant representation

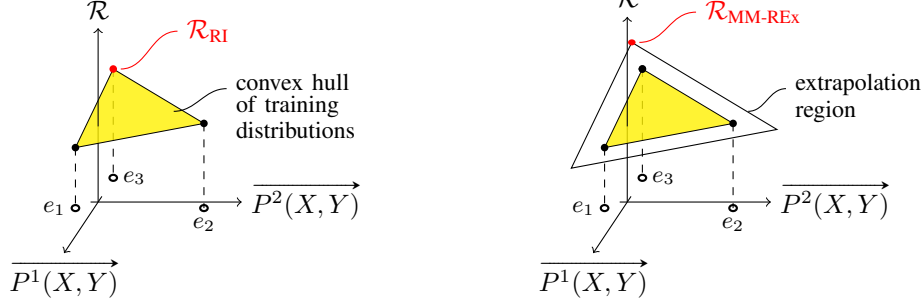


Figure 1. **Left:** Robust optimization optimizes worst-case performance over the convex hull of training distributions. **Right:** By extrapolating risks, REx encourages robustness to larger shifts. Here e_1 , e_2 , and e_3 represent training distributions, and $P^1(X, Y)$, $P^2(X, Y)$ represent some particular directions of variation in the affine space of quasiprobability distributions over (X, Y) .

Many domain generalization methods assume $P(Y|X)$ is an invariant relationship, limiting distributional shift to changes in $P(X)$, which are known as **covariate shift** (Ben-David et al., 2010b). This assumption can easily be violated, however. For instance, when Y causes X , a more sensible assumption is that $P(X|Y)$ is fixed, with $P(Y)$ varying across domains (Schölkopf et al., 2012; Lipton et al., 2018). In general, invariant prediction may involve an aspect of causal discovery. Depending on the perturbation set, however, other, more predictive, invariant relationships may also exist (Koyama & Yamaguchi, 2020).

The first method for invariant prediction to be compatible with modern deep learning problems and techniques is **Invariant Risk Minimization (IRM)** (Arjovsky et al., 2019), making it a natural point of comparison. Our work focuses on explaining how REx addresses OOD generalization, and highlighting differences (especially advantages) of REx compared with IRM and other domain generalization methods, see Table 1. Broadly speaking, **REx optimizes for robustness to the forms of distributional shift that have been observed to have the largest impact on performance in training domains**. This can be a significant advantage over the more focused (but also limited) robustness that IRM targets. **For instance, unlike IRM, REx can also encourage robustness to covariate shift** (see Section 3 and Figure 3.2).

Our experiments show that REx significantly outperforms IRM in settings that involve covariate shift and require invariant prediction, including modified versions of CMNIST and simulated robotics tasks from the Deepmind control suite. On the other hand, because **REx does not distinguish between underfitting and inherent noise**, IRM has an advantage in settings **where some domains are intrinsically harder than others**. Our contributions include:

1. MM-REx, a novel domain generalization problem for (Ganin et al., 2016); see Section 2.3.

mulation suitable for invariant prediction.

2. Demonstrating that REx solves invariant prediction tasks where **IRM fails due to covariate shift**.
3. Proving that **equality of risks can be a sufficient criteria for discovering causal structure**.

2. Background & Related work

We consider multi-source domain generalization, where our goal is to find parameters θ that perform well on unseen domains, given a set of m training domains, $\mathcal{E} = \{e_1, \dots, e_m\}$, sometimes also called **environments**. We assume the loss function, ℓ is fixed, and domains only differ in terms of their data distribution $P_e(X, Y)$ and dataset D_e . The **risk function** for a given domain/distribution e is:

$$\mathcal{R}_e(\theta) \doteq \mathbb{E}_{(x,y) \sim P_e(X,Y)} \ell(f_\theta(x), y) \quad (2)$$

We refer to members of the set $\{\mathcal{R}_e | e \in \mathcal{E}\}$ as the **training risks** or simply **risks**. Changes in $P_e(X, Y)$ can be categorized as either changes in $P(X)$ (**covariate shift**), changes in $P(Y|X)$ (**concept shift**), or a combination. The standard approach to learning problems is **Empirical Risk Minimization (ERM)**, which minimizes the average loss across all the training examples from all the domains:

$$\mathcal{R}_{\text{ERM}}(\theta) \doteq \mathbb{E}_{(x,y) \sim \bigcup_{e \in \mathcal{E}} D_e} \ell(f_\theta(x), y) \quad (3)$$

$$= \sum_e |D_e| \mathbb{E}_{(x,y) \sim D_e} \ell(f_\theta(x), y) \quad (4)$$

2.1. Robust Optimization

An approach more tailored to OOD generalization is **robust optimization** (Ben-Tal et al., 2009), which aims to optimize a model’s worst-case performance over some **perturbation set** of possible data distributions, \mathcal{F} (see Eqn. 1).

Method	Invariant Prediction	Cov. Shift Robustness	Suitable for Deep Learning
DRO	✗	✓	✓
(C-)ADA	✗	✓	✓
ICP	✓	✗	✗
IRM	✓	✗	✓
REx	✓	✓	✓

Table 1. A comparison of approaches for OOD generalization.

When only a single training domain is available (**single-source domain generalization**), it is common to assume that $P(Y|X)$ is fixed, and let \mathcal{F} be all distributions within some f -divergence ball of the training $P(X)$ (Hu et al., 2016; Bagnell, 2005). As another example, adversarial robustness can be seen as instead using a Wasserstein ball as a perturbation set (Sinha et al., 2017). The assumption that $P(Y|X)$ is fixed is commonly called the “covariate shift assumption” (Ben-David et al., 2010b); however, we assume that covariate shift and concept shift can co-occur, and refer to this assumption as **the fixed relationship assumption (FRA)**.

In **multi-source domain generalization**, test distributions are often **assumed to be mixtures** (i.e. convex combinations) **of the training distributions**; this is equivalent to setting $\mathcal{F} \doteq \mathcal{E}$:

$$\mathcal{R}_{\text{RI}}(\theta) \doteq \max_{\substack{\sum_e \lambda_e = 1 \\ \lambda_e \geq 0}} \sum_{e=1}^m \lambda_e \mathcal{R}_e(\theta) = \max_{e \in \mathcal{E}} \mathcal{R}_e(\theta). \quad (5)$$

We call this objective **Risk Interpolation (RI)**, or, following Sagawa et al. (2019), **(group) Distributionally Robust Optimization (DRO)**. While **single-source** methods classically assume that the probability of each data-point can **vary independently** (Hu et al., 2016), DRO yields a much lower dimensional perturbation set, with at most one direction of variation per domain, regardless of the dimensionality of X and Y . It also does not rely on FRA, and can provide robustness to any form of shift in $P(X, Y)$ which occurs across training domains. **Minimax-REx is an extension of this approach to affine combinations of training risks.**

2.2. Invariant representations vs. invariant predictors

An equipredictive representation, Φ , is a function of X with the property that $P_e(Y|\Phi)$ is equal, $\forall e \in \mathcal{F}$. In other words, the relationship between such a Φ and Y is fixed across domains. **Invariant relationships** between X and Y are then exactly those that can be written as $P(Y|\Phi(x))$ with Φ an equipredictive representation. A model $\hat{P}(Y|X=x)$ that learns such an invariant relationship is called an **invariant predictor**. Intuitively, an invariant predictor works equally well **across all domains in \mathcal{F}** . The principle of risk

extrapolation aims to achieve invariant prediction by **enforcing such equality across training domains \mathcal{E} , and does not rely on explicitly learning an equipredictive representation.**

Koyama & Yamaguchi (2020) prove that a *maximal* equipredictive representation – that is, one that maximizes mutual information with the targets, $\Phi^* \doteq \arg\max_{\Phi} I(\Phi, Y)$ – solves the robust optimization problem (Eqn. 1) under fairly general assumptions.³ When Φ^* is unique, we call the features it ignores **spurious**. The result of Koyama & Yamaguchi (2020) provides a theoretical reason for **favoring invariant prediction over the common approach of learning invariant representations** (Pan et al., 2010), which make $P_e(\Phi)$ or $P_e(\Phi|Y)$ equal $\forall e \in \mathcal{E}$. Popular methods here include **adversarial domain adaptation (ADA)** (Ganin et al., 2016) and **conditional ADA (C-ADA)** (Long et al., 2018). Unlike invariant predictors, invariant representations can easily fail to generalize OOD: ADA forces the predictor to have the **same marginal predictions $\hat{P}(Y)$, which is a mistake when $P(Y)$ in fact changes across domains** (Zhao et al., 2019); C-ADA suffers from more subtle issues (Arjovsky et al., 2019).

2.3. Invariance and causality

The relationship between cause and effect is a paradigmatic example of an invariant relationship. Here, we summarize definitions from causal modeling, and discuss causal approaches to domain generalization. We will refer to these definitions for the statements of our theorems in Section 3.2.

Definitions. A **causal graph** is a directed acyclic graph (DAG), where nodes represent variables and edges point from causes to effects. In this work, we use **Structural Causal Models (SCMs)**, which also specify how the value of a variable is computed given its parents. An SCM, \mathcal{C} , is defined by specifying the **mechanism**, $f_Z : Pa(Z) \rightarrow$

³The first formal definition of an equipredictive representation we found was by Koyama & Yamaguchi (2020), who use the term “(maximal) invariant predictor”. We prefer our terminology since: 1) it is more consistent with Arjovsky et al. (2019), and 2) Φ is a representation, not a predictor.

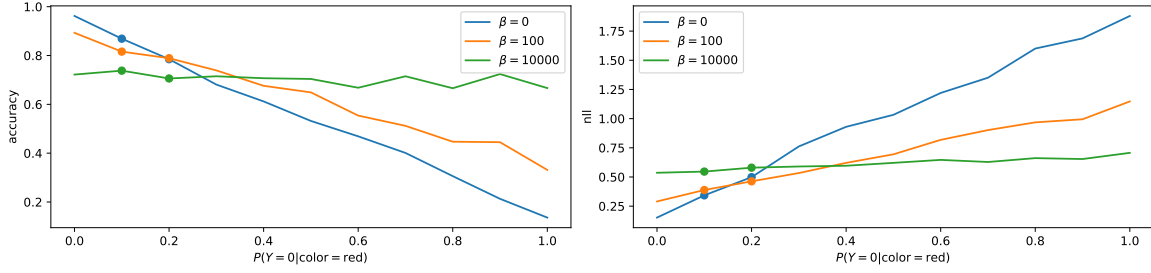


Figure 2. Training accuracies (left) and risks (right) on colored MNIST domains with varying $P(Y=0|\text{color}=\text{red})$ after 500 epochs. Dots represent training risks, lines represent test risks on different domains. Increasing the V-REx penalty (β) leads to a flatter “risk plane” and more consistent performance across domains, as the model learns to ignore color in favor of shape-based invariant prediction. Note that $\beta=100$ gives the best worst-case risk across the 2 training domains, and so would be the solution preferred by DRO (Sagawa et al., 2019). This demonstrates that REx’s counter-intuitive propensity to *increase* training risks can be necessary for good OOD performance.

$\text{dom}(Z)$ for each variable Z .⁴ Mechanisms are *deterministic*; noise in Z is represented explicitly via a special noise variable N_Z , and these noise variables are jointly independent. **An intervention, ι is any modification to the mechanisms of one or more variables**; an intervention can introduce new edges, so long as it does not introduce a cycle. $\text{do}(X_i = x)$ denotes an intervention which sets X_i to the constant value x (removing all incoming edges). Data can be generated from an SCM, \mathcal{C} , by sampling all of the noise variables, and then using the mechanisms to compute the value of every node whose parents’ values are known. This sampling process defines an **entailed distribution**, $P^{\mathcal{C}}(\mathbf{Z})$ over the nodes \mathbf{Z} of \mathcal{C} . We overload f_Z , letting $f_Z(\mathbf{Z})$ refer to the conditional distribution $P^{\mathcal{C}}(Z|\mathbf{Z} \setminus \{Z\})$.

2.3.1. CAUSAL APPROACHES TO DOMAIN GENERALIZATION

Instead of assuming $P(Y|X)$ is fixed (FRA), works that take a causal approach to domain generalization **often assume that the mechanism for Y is fixed; we call this the fixed mechanism assumption (FMA)**. Meanwhile, they assume X may be subject to different (e.g. arbitrary) interventions in different domains (Bühlmann, 2018). We call changes in $P(X, Y)$ resulting from **interventions on X interventional shift**. **Interventional shift can involve both covariate shift and/or concept shift**. In their seminal work on **Invariant Causal Prediction (ICP)**, Peters et al. (2016) leverage this invariance to learn which elements of X cause Y . ICP and its nonlinear extension (Heinze-Deml et al., 2018) use statistical tests to detect whether the residuals of a linear model are equal across domains. Our work differs from ICP in that:

1. Our method is model agnostic and scales to deep networks.

⁴Our definitions follow *Elements of Causal Inference* (Peters et al., 2017); our notation mostly does as well.

2. Our goal is OOD generalization, not causal inference. These are not identical: **invariant prediction can sometimes make use of non-causal relationships, but when deciding which interventions to perform, a truly causal model is called for**.
3. Our learning principle **only requires invariance of risks, not residuals**. Nonetheless, we prove that this can ensure invariant causal prediction.

A more similar method to REx is **Invariant Risk Minimization (IRM)** (Arjovsky et al., 2019), which shares properties (1) and (2) of the list above. Like REx, IRM also uses a weaker form of invariance than ICP; namely, **they insist that the optimal linear classifier must match across domains**.⁵ Still, REx differs significantly from IRM. While IRM specifically aims for invariant prediction, REx seeks robustness to *whichever* forms of distributional shift are present. Thus, REx is more directly focused on the problem of OOD generalization, and can provide robustness to a wider variety of distributional shifts, including covariate shift. Also, unlike REx, IRM seeks to match $\mathbb{E}(Y|\Phi(X))$ across domains, not the full $P(Y|\Phi(X))$. This, combined with IRM’s indifference to covariate shift, **make it more effective in cases where different domains or examples are inherently more noisy**.

2.4. Fairness

Equalizing risk across different groups (e.g. male vs. female) **has been proposed as a definition of fairness** (Donini et al., 2018), generalizing the equal opportunity definition of fairness (Hardt et al., 2016). Williamson & Menon (2019) propose using the **absolute difference of risks to measure deviation from this notion of fairness**; this corresponds to our MM-REx, in the case of only two domains, and is similar to V-REx, which uses the variance of risks. However, in the context of fairness, equalizing the risk of training groups is

⁵In practice, IRMv1 replaces this bilevel optimization problem with a gradient penalty on classifier weights.

the goal. Our work goes beyond this by showing that it can serve as a method for OOD generalization.

3. Risk Extrapolation

Before discussing algorithms for REx and theoretical results, we first expand on our high-level explanations of what REx does, what kind of OOD generalization it promotes, and how. The principle of Risk Extrapolation (REx) has two aims:

1. Reducing training risks
2. Increasing similarity of training risks

In general, these goals can be at odds with each other; decreasing the risk in the domain with the lowest risk also decreases the overall similarity of training risks. Thus methods for REx may seek to increase risk on the best performing domains. While this is counter-intuitive, it can be necessary to achieve good OOD generalization, as Figure 2 demonstrates. From a geometric point of view, encouraging equality of risks flattens the “risk plane” (the affine span of the training risks, considered as a function of the data distribution, see Figures 1 and 2). While this can result in higher training risks, it also means that the risk changes less if the distributional shifts between training domains are magnified at test time.

Figure 2 illustrates how flattening the risk plane can promote OOD generalization on real data, using the Colored MNIST (CMNIST) task as an example (Arjovsky et al., 2019). In the CMNIST training domains, the color of a digit is more predictive of the label than the shape is. But because the correlation between color and label is not invariant, predictors that use the color feature achieve different risk on different domains. By enforcing equality of risks, REx prevents the model from using the color feature enabling successful generalization to the test domain where the correlation between color and label is reversed.

Probabilities vs. Risks. Figure 3 depicts how the extrapolated risks considered in MM-REx can be translated into a corresponding change in $P(X, Y)$, using an example of pure covariate shift. Training distributions can be thought of as points in an affine space with a dimension for every possible value of (X, Y) ; see Appendix C.1 for an example. Because the risk is linear w.r.t. $P(x, y)$, a convex combination of risks from different domains is equivalent to the risk on a domain given by the mixture of their distributions. The same holds for the affine combinations used in MM-REx, with the caveat that the negative coefficients may lead to negative probabilities, making the resulting $P(X, Y)$ a *quasiprobability distribution*, i.e. a signed measure with integral 1. We explore the theoretical implications of this in Appendix E.

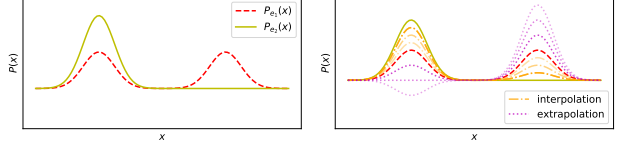


Figure 3. Extrapolation can yield a distribution with *negative* $P(x)$ for some x . **Left:** $P(x)$ for domains e_1 and e_2 . **Right:** Point-wise interpolation/extrapolation of $P^{e_1}(x)$ and $P^{e_2}(x)$. Since MM-REx target worst-case robustness across extrapolated domains, it can provide robustness to such shifts in $P(X)$ (covariate shift).

Covariate Shift. When only $P(X)$ differs across domains (i.e. FRA holds), as in Figure 3, then $\Phi(x) = x$ is already an equipredictive representation, and so *any* predictor is an invariant predictor. Thus methods which only promote invariant prediction – such as IRM – are not expected to improve OOD generalization (compared with ERM). Indeed, Arjovsky et al. (2019) recognize this limitation of IRM in what they call the “realizable” case. Instead, what is needed is robustness to covariate shift, which REx, but not IRM, can provide. Robustness to covariate shift can improve OOD generalization by ensuring that low-capacity models spend sufficient capacity on low-density regions of the input space; we show how REx can provide such benefits in Appendix C.2. But even for high capacity models, $P(X)$ can have a significant influence on what is learned; for instance Sagawa et al. (2019) show that DRO can significantly improves the performance on rare groups in their Waterbirds dataset. Pursuing robustness to covariate shift also comes with drawbacks for REx, however: REx does not distinguish between underfitting and inherent noise in the data, and so can force the model to make equally bad predictions everywhere, even if some examples are less noisy than others.

3.1. Methods of Risk Extrapolation

We now formally describe the **Minimax REx (MM-REx)** and **Variance-REx (V-REx)** techniques for risk extrapolation. Minimax-REx performs robust learning over a perturbation set of *affine* combinations of training risks with **bounded coefficients**:

$$\begin{aligned} \mathcal{R}_{\text{MM-REx}}(\theta) &\doteq \max_{\substack{\sum_e \lambda_e = 1 \\ \lambda_e \geq \lambda_{\min}}} \sum_{e=1}^m \lambda_e \mathcal{R}_e(\theta) \\ &= (1 - m\lambda_{\min}) \max_e \mathcal{R}_e(\theta) + \lambda_{\min} \sum_{e=1}^m \mathcal{R}_e(\theta), \end{aligned} \quad (6)$$

$$(7)$$

where m is the number of domains, and the hyperparameter λ_{\min} controls how much we extrapolate. For negative values of λ_{\min} , MM-REx places negative weights on the risk of all but the worst-case domain, and as $\lambda_{\min} \rightarrow -\infty$,

this criterion enforces strict equality between training risks; $\lambda_{\min} = 0$ recovers risk interpolation (RI). Thus, like RI, MM-REx aims to be robust in the direction of variations in $P(X, Y)$ between test domains. However, negative coefficients allow us to extrapolate to more extreme variations. Geometrically, larger values of λ_{\min} expand the perturbation set farther away from the convex hull of the training risks, encouraging a flatter “risk-plane” (see Figure 2).

While MM-REx makes the relationship to RI/RO clear, we found using the variance of risks as a regularizer (V-REx) simpler, stabler, and more effective:

$$\mathcal{R}_{\text{V-REx}}(\theta) \doteq \beta \text{Var}(\{\mathcal{R}_1(\theta), \dots, \mathcal{R}_m(\theta)\}) + \sum_{e=1}^m \mathcal{R}_e(\theta) \quad (8)$$

Here $\beta \in [0, \infty)$ controls the balance between reducing average risk and enforcing equality of risks, with $\beta = 0$ recovering ERM, and $\beta \rightarrow \infty$ leading V-REx to focus entirely on making the risks equal. See Appendix for the relationship between V-REx and MM-REx and their gradient vector fields.

3.2. Theoretical Conditions for REx to Perform Causal Discovery

We now prove that **exactly equalizing training risks** (as incentivized by REx) leads a model to learn the causal mechanism of Y under assumptions similar to those of Peters et al. (2016), namely:

1. The causes of Y are observed, i.e. $Pa(Y) \subseteq X$.
2. Domains correspond to interventions on X .
3. Homoskedasticity (a slight generalization of the additive noise setting assumed by Peters et al. (2016)). We say an SEM \mathfrak{C} is **homoskedastic** (with respect to a loss function ℓ), if the **Bayes error rate** of $\ell(f_Y(x), f_Y(x))$ is the same for all $x \in \mathcal{X}$.⁶

The contribution of our theory (vs. **ICP**) is to prove that equalizing risks is sufficient to learn the causes of Y . In contrast, **they insist that the entire distribution of error residuals** (in predicting Y) **be the same across domains**. We provide proof sketches here and complete proofs in the appendix.

Theorem 1 demonstrates a practical result: we can identify a linear SEM model using REx with a number of domains linear in the dimensionality of X .

⁶ Note that our definitions of **homoskedastic/heteroskedastic** do *not* correspond to the types of domains constructed in Arjovsky et al. (2019), Section 5.1, but rather are a generalization of the definitions of these terms as commonly used in statistics. Specifically, for us, *heteroskedasticity* means that the “predicatability” (e.g. variance) of Y differs across inputs x , whereas for Arjovsky et al. (2019), it means the predicatability of Y at a given input varies across *domains*; we refer to this second type as *domain-homo/heteroskedasticity* for clarity.

Theorem 1. *Given a Linear SEM, $X_i \leftarrow \sum_{j \neq i} \beta_{(i,j)} X_j + \varepsilon_i$, with $Y \doteq X_0$, and a predictor $f_\beta(X) \doteq \sum_{j:j>0} \beta_j X_j + \varepsilon_j$ that satisfies REx (with mean-squared error) over a perturbation set of domains that contains 3 distinct $do()$ interventions for each $X_i : i > 0$. Then $\beta_j = \beta_{0,j}, \forall j$.*

Proof Sketch. We adapt the proof of Theorem 4i from Peters et al. (2016). They show that matching the residual errors across observational and interventional domains forces the model to learn f_Y . We use the weaker condition of matching risks to derive a quadratic equation that the $do()$ interventions must satisfy for any model other than f_Y . Since there are at most 2 solutions to a quadratic equation, insisting on equality of risks across 3 distinct $do()$ interventions forces the model to learn f_Y .

Given the assumption that a predictor satisfies REx over *all* interventions that do not change the mechanism of Y , we can prove a much more general result. We now consider an arbitrary SCM, \mathfrak{C} , generating Y and X , and let \mathcal{E}^I be the set of domains corresponding to arbitrary interventions on X , similarly to Peters et al. (2016).

Theorem 2. *Suppose ℓ is a (strictly) proper scoring rule. Then a predictor that **satisfies REx** for a over \mathcal{E}^I uses $f_Y(x)$ as its predictive distribution on input x for all $x \in \mathcal{X}$.*

Proof Sketch. Since the distribution of Y given its parents doesn’t depend on the domain, f_Y can make reliable point-wise predictions across domains. This translates into equality of risk across domains when the overall difficulty of the examples is held constant across domains, e.g. by assuming homoskedasticity.⁷ While a different predictor might do a better job on *some* domains, **we can always find a domain where it does worse than f_Y , and so f_Y is both unique and optimal**.

Remark. Theorem 2 is only meant to provide insight into how the REx principle relates to causal invariance; the perturbation set in this theorem is uncountably infinite. Note, however, that even in this setting, the ERM principle does *not*, in general, recover the causal mechanism for Y . Rather, the ERM solution depends on the distribution over domains. For instance, if all but an $\epsilon \rightarrow 0$ fraction of the data comes from the CMNIST training domains, then ERM will learn to use the color feature, just as in original the CMNIST task.

4. Experiments

We evaluate REx and compare with IRM on a range of tasks requiring OOD generalization. REx provides generalization benefits and outperforms IRM on a wide range of tasks, including: i) variants of the Colored MNIST (CMNIST) dataset (Arjovsky et al., 2019) with covariate shift, ii) continuous control tasks with partial observability and

⁷Note we could also assume no covariate shift in order to fix the difficulty, but this seems hard to motivate in the context of interventions on X , which can change $P(X)$.

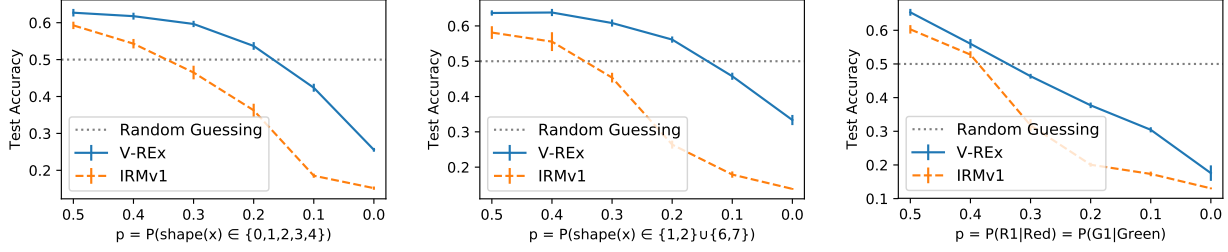


Figure 4. REx outperforms IRM on Colored MNIST variants that include covariate shift. The x-axis indexes increasing amount of shift between training distributions, with $p = 0$ corresponding to disjoint supports. **Left:** class imbalance, **Center:** shape imbalance, **Right:** color imbalance.

Method	train acc	test acc
V-REx (ours)	71.5 ± 1.0	68.7 ± 0.9
IRM	70.8 ± 0.9	66.9 ± 2.5
MM-REx (ours)	72.4 ± 1.8	66.1 ± 1.5
RI	88.9 ± 0.3	22.3 ± 4.6
ERM	87.4 ± 0.2	17.1 ± 0.6
Grayscale oracle	73.5 ± 0.2	73.0 ± 0.4
Optimum	75	75
Chance	50	50

Table 2. Accuracy (percent) on Colored MNIST. REx and IRM learn to ignore the spurious color feature. Strikethrough results achieved via tuning on the test set.

spurious features, iii) domain generalization tasks from the DomainBed suite (Gulrajani & Lopez-Paz, 2020). On the other hand, **when the inherent noise in Y varies across environments, IRM succeeds and REx performs poorly.**

4.1. Colored MNIST

Arjovsky et al. (2019) construct a binary classification problem (with 0-4 and 5-9 each collapsed into a single class) based on the MNIST dataset, using color as a spurious feature. Specifically, digits are either colored red or green, and there is a strong correlation between color and label, which is reversed at test time. The goal is to learn the causal “digit shape” feature and ignore the anti-causal “digit color” feature. The learner has access to three domains:

1. A training domain where green digits have a 80% chance of belonging to class 1 (digits 5-9).
2. A training domain where green digits have a 90% chance of belonging to class 1.
3. A test domain where green digits have a 10% chance of belonging to class 1.

We use the exact same hyperparameters as Arjovsky et al. (2019), only replacing the IRMv1 penalty with MM-REx or V-REx penalty.⁸ These methods all achieve similar perfor-

mance, see Table 2.

CMNIST with covariate shift. To test our hypothesis that REx should outperform IRM under covariate shift, we construct 3 variants of the CMNIST dataset. Each variant represents **a different way of inducing covariate shift to ensure differences across methods are consistent**. These experiments combine covariate shift with interventional shift, since $P(\text{Green}|Y = 1)$ still differs across training domains as in the original CMNIST.

1. **Class imbalance:** varying $p = P(\text{shape}(x) \in \{0, 1, 2, 3, 4\})$; as in Wu et al. (2020).
2. **Digit imbalance:** varying $p = P(\text{shape}(x) \in \{1, 2\} \cup \{6, 7\})$; digits 0 and 5 are removed.
3. **Color imbalance:** We use 2 versions of each color, for 4 total channels: R_1, R_2, G_1, G_2 . We vary $p = P(R_1|\text{Red}) = P(G_1|\text{Green})$.

While (1) also induces change in $P(Y)$, (2) and (3) induce **only covariate shift in the causal shape and anti-causal color features (respectively)**. We compare across several levels of imbalance, $p \in [0, 0.5]$, using the same hyperparameters from Arjovsky et al. (2019), and plot the mean and standard error over 3 trials.

V-REx significantly outperforms IRM in every case, see Figure 3.2. In order to verify that these results are not due to bad hyperparameters for IRM, we perform a random search that samples 340 unique hyperparameter combinations for each value of p , and **compare the the number of times each method achieves better than chance-level (50% accuracy)**. Again, V-REx outperforms IRM; in particular, for small values of p , IRM never achieves better than random in 4.4%/23.7%/2.0% of trials, respectively, in the class/digit/color imbalance scenarios for $p = 0.1/0.1/0.2$. This indicates that REx can achieve good OOD generalization in settings involving both covariate and interventional shift, whereas IRM struggles to do so.

penalty on the Mean Absolute Error (MAE), see Appendix F.2.2.

⁸When there are only 2 domains, MM-REx is equivalent to a

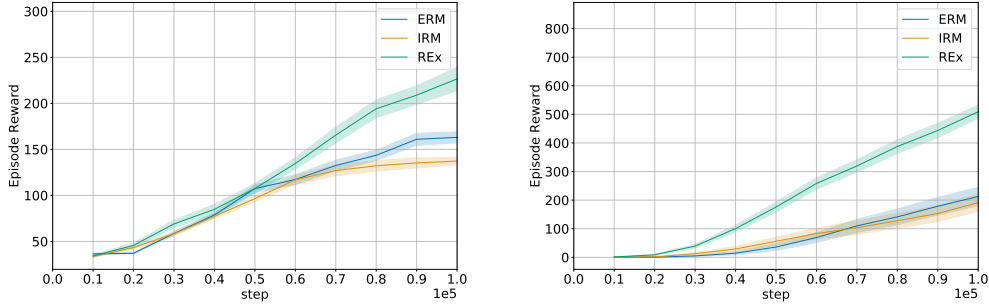


Figure 5. Performance and standard error on walker_walk (top), finger_spin (bottom).

Algorithm	ColoredMNIST	VLCS	PACS	OfficeHome
ERM	52.0 ± 0.1	77.4 ± 0.3	85.7 ± 0.5	67.5 ± 0.5
IRM	51.8 ± 0.1	78.1 ± 0.0	84.4 ± 1.1	66.6 ± 1.0
V-REx	52.1 ± 0.1	77.9 ± 0.5	85.8 ± 0.6	66.7 ± 0.5

Table 3. REx, IRM, and ERM all perform comparably on a set of domain generalization benchmarks.

4.2. Toy Structural Equation Models (SEMs)

REx’s sensitivity to covariate shift can also be a weakness when reallocating capacity towards domains with higher risk does not help the model reduce their risk, e.g. due to irreducible noise. We illustrate this using the linear-Gaussian structural equation model (SEM) tasks introduced by Arjovsky et al. (2019). Like CMNIST, these SEMs include spurious features by construction. They also introduce 1) heteroskedasticity, 2) hidden confounders, and/or 3) elements of X that contain a mixture of causes and effects of Y . These three properties highlight advantages of IRM over ICP (Peters et al., 2016), as demonstrated empirically by Arjovsky et al. (2019). REx is also able to handle (2) and (3), but it performs poorly in the heteroskedastic tasks. See Appendix G.2 for details and Table 5 for results.

4.3. Domain Generalization in the DomainBed Suite

Methodologically, it is inappropriate to assume access to the test environment in domain generalization settings, as the goal is to find methods which generalize to unknown test distributions. Gulrajani & Lopez-Paz (2020) introduced the DomainBed evaluation suite to rigorously compare existing approaches to domain generalization, and found that no method reliably outperformed ERM. We evaluate V-REx on DomainBed using the most commonly used training-domain validation set method for model selection. Due to limited computational resources, we limited ourselves to the 4 cheapest datasets. Results of baseline are taken from Gulrajani & Lopez-Paz (2020), who compare with more methods. Results in Table 3 give the average over 3 different train/valid splits.

4.4. Reinforcement Learning with partial observability and spurious features

Finally, we turn to reinforcement learning, where covariate shift (potentially favoring REx) and heteroskedasticity (favoring IRM) both occur naturally as a result of randomness in the environment and policy. In order to show the benefits of invariant prediction, we modify tasks from the Deepmind Control Suite (Tassa et al., 2018) to include spurious features in the observation, and train a Soft Actor-Critic (Haarnoja et al., 2018) agent. REx outperforms both IRM and ERM, suggesting that REx’s robustness to covariate shift outweighs the challenges it faces with heteroskedasticity in this setting, see Figure 5. We average over 10 runs on finger_spin and walker_walk, using hyperparameters tuned on cartpole_swingup (to avoid overfitting). See Appendix for details and further results.

5. Conclusion

We have demonstrated that REx, a method for robust optimization, can provide robustness and hence out-of-distribution generalization in the challenging case where X contains both causes and effects of Y . In particular, like IRM, REx can perform causal identification, but REx can also perform more robustly in the presence of covariate shift. Covariate shift is known to be problematic when models are misspecified, when training data is limited, or does not cover areas of the test distribution. As such situations are inevitable in practice, REx’s ability to outperform IRM in scenarios involving a combination of covariate shift and interventional shift makes it a powerful approach.

References

- Albuquerque, I., Naik, N., Li, J., Keskar, N., and Socher, R. Improving out-of-distribution generalization via multi-task self-supervised pretraining, 2020.
- Arjovsky, M., Bottou, L., Gulrajani, I., and Lopez-Paz, D. Invariant risk minimization. *arXiv preprint arXiv:1907.02893*, 2019.
- Bachman, P., Hjelm, R. D., and Buchwalter, W. Learning representations by maximizing mutual information across views, 2019.
- Bagnell, J. A. Robust supervised learning. In *Proceedings of the 20th National Conference on Artificial Intelligence - Volume 2, AAAI’05*, pp. 714–719. AAAI Press, 2005. ISBN 157735236x.
- Beery, S., Van Horn, G., and Perona, P. Recognition in terra incognita. *Lecture Notes in Computer Science*, pp. 472–489, 2018. ISSN 1611-3349.
- Ben-David, S., Blitzer, J., Crammer, K., Kulesza, A., Pereira, F., and Vaughan, J. W. A theory of learning from different domains. *Machine learning*, 79(1-2):151–175, 2010a.
- Ben-David, S., Lu, T., Luu, T., and Pál, D. Impossibility theorems for domain adaptation. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pp. 129–136, 2010b.
- Ben-Tal, A., El Ghaoui, L., and Nemirovski, A. *Robust optimization*, volume 28. Princeton University Press, 2009.
- Bühlmann, P. Invariance, causality and robustness, 2018.
- Carlucci, F. M., D’Innocente, A., Bucci, S., Caputo, B., and Tommasi, T. Domain generalization by solving jigsaw puzzles. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2229–2238, 2019.
- Cubuk, E. D., Zoph, B., Mane, D., Vasudevan, V., and Le, Q. V. Autoaugment: Learning augmentation policies from data, 2018.
- Desjardins, G., Simonyan, K., Pascanu, R., et al. Natural neural networks. In *Advances in Neural Information Processing Systems*, pp. 2071–2079, 2015.
- Donini, M., Oneto, L., Ben-David, S., Shawe-Taylor, J., and Pontil, M. Empirical risk minimization under fairness constraints, 2018.
- Engstrom, L., Tran, B., Tsipras, D., Schmidt, L., and Madry, A. Exploring the landscape of spatial robustness. *arXiv preprint arXiv:1712.02779*, 2017.
- Ganin, Y., Ustinova, E., Ajakan, H., Germain, P., Larochelle, H., Laviolette, F., Marchand, M., and Lempitsky, V. Domain-adversarial training of neural networks. *The Journal of Machine Learning Research*, 17(1):2096–2030, 2016.
- Geirhos, R., Rubisch, P., Michaelis, C., Bethge, M., Wichmann, F. A., and Brendel, W. Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. *arXiv preprint arXiv:1811.12231*, 2018.
- Goodfellow, I. J., Shlens, J., and Szegedy, C. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- Gowal, S., Qin, C., Huang, P.-S., Cemgil, T., Dvijotham, K., Mann, T., and Kohli, P. Achieving robustness in the wild via adversarial mixing with disentangled representations. *arXiv preprint arXiv:1912.03192*, 2019.
- Gulrajani, I. and Lopez-Paz, D. In search of lost domain generalization, 2020.
- Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In Dy, J. and Krause, A. (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 1861–1870, Stockholmsmässan, Stockholm Sweden, 10–15 Jul 2018. PMLR.
- Haffner, P. Escaping the convex hull with extrapolated vector machines. In Dietterich, T. G., Becker, S., and Ghahramani, Z. (eds.), *Advances in Neural Information Processing Systems 14*, pp. 753–760. MIT Press, 2002.
- Hardt, M., Price, E., and Srebro, N. Equality of opportunity in supervised learning, 2016.
- Hastie, T., Tibshirani, R., and Friedman, J. *The elements of statistical learning: data mining, inference, and prediction*. Springer Science & Business Media, 2009.
- He, Y., Shen, Z., and Cui, P. Towards non-i.i.d. image classification: A dataset and baselines, 2019.
- Heinze-Deml, C., Peters, J., and Meinshausen, N. Invariant causal prediction for nonlinear models. *Journal of Causal Inference*, 6(2), Sep 2018. ISSN 2193-3685. doi: 10.1515/jci-2017-0016. URL <http://dx.doi.org/10.1515/jci-2017-0016>.
- Hendrycks, D. and Dietterich, T. Benchmarking neural network robustness to common corruptions and perturbations. *arXiv preprint arXiv:1903.12261*, 2019.

- Hendrycks, D. and Gimpel, K. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *arXiv preprint arXiv:1610.02136*, 2016.
- Hendrycks, D., Mazeika, M., and Dietterich, T. Deep anomaly detection with outlier exposure. *arXiv preprint arXiv:1812.04606*, 2018.
- Hendrycks, D., Mazeika, M., Kadavath, S., and Song, D. Using self-supervised learning can improve model robustness and uncertainty, 2019a.
- Hendrycks, D., Mu, N., Cubuk, E. D., Zoph, B., Gilmer, J., and Lakshminarayanan, B. Augmix: A simple data processing method to improve robustness and uncertainty, 2019b.
- Hjelm, R. D., Fedorov, A., Lavoie-Marchildon, S., Grewal, K., Bachman, P., Trischler, A., and Bengio, Y. Learning deep representations by mutual information estimation and maximization, 2018.
- Hu, W., Niu, G., Sato, I., and Sugiyama, M. Does distributionally robust supervised learning give robust classifiers?, 2016.
- Ilse, M., Tomczak, J. M., and Forré, P. Designing data augmentation for simulating interventions. *arXiv preprint arXiv:2005.01856*, 2020.
- Johansson, F. D., Sontag, D., and Ranganath, R. Support and invertibility in domain-invariant representations, 2019.
- Koyama, M. and Yamaguchi, S. Out-of-distribution generalization with maximal invariant predictor, 2020.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pp. 1097–1105, 2012.
- Li, D., Yang, Y., Song, Y.-Z., and Hospedales, T. M. Deeper, broader and artier domain generalization. In *Proceedings of the IEEE international conference on computer vision*, pp. 5542–5550, 2017.
- Li, Y., Tian, X., Gong, M., Liu, Y., Liu, T., Zhang, K., and Tao, D. Deep domain generalization via conditional invariant adversarial networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 624–639, 2018.
- Lipton, Z. C., Wang, Y.-X., and Smola, A. Detecting and correcting for label shift with black box predictors. *arXiv preprint arXiv:1802.03916*, 2018.
- Long, M., Cao, Z., Wang, J., and Jordan, M. I. Conditional adversarial domain adaptation. In *Advances in Neural Information Processing Systems*, pp. 1640–1650, 2018.
- Meinshausen, N., Bühlmann, P., et al. Maximin effects in inhomogeneous large-scale data. *The Annals of Statistics*, 43(4):1801–1830, 2015.
- Pan, S. J., Tsang, I. W., Kwok, J. T., and Yang, Q. Domain adaptation via transfer component analysis. *IEEE Transactions on Neural Networks*, 22(2):199–210, 2010.
- Peters, J., Bühlmann, P., and Meinshausen, N. Causal inference by using invariant prediction: identification and confidence intervals. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 78(5):947–1012, 2016.
- Peters, J., Janzing, D., and Schölkopf, B. *Elements of causal inference: foundations and learning algorithms*. 2017.
- Recht, B., Roelofs, R., Schmidt, L., and Shankar, V. Do imagenet classifiers generalize to imagenet? *arXiv preprint arXiv:1902.10811*, 2019.
- Sagawa, S., Koh, P. W., Hashimoto, T. B., and Liang, P. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization, 2019.
- Sahoo, S. S., Lampert, C. H., and Martius, G. Learning equations for extrapolation and control, 2018.
- Schölkopf, B., Janzing, D., Peters, J., Sgouritsa, E., Zhang, K., and Mooij, J. On causal and anticausal learning. In *Proceedings of the 29th International Conference on Machine Learning, ICML’12*, pp. 459–466, Madison, WI, USA, 2012. Omnipress. ISBN 9781450312851.
- Shorten, C. and Khoshgoftaar, T. M. A survey on image data augmentation for deep learning. *Journal of Big Data*, 6(1):60, 2019.
- Sinha, A., Namkoong, H., Volpi, R., and Duchi, J. Certifying some distributional robustness with principled adversarial training, 2017.
- Su, J., Vargas, D. V., and Sakurai, K. One pixel attack for fooling deep neural networks. *IEEE Transactions on Evolutionary Computation*, 23(5):828–841, 2019.
- Tassa, Y., Doron, Y., Muldal, A., Erez, T., Li, Y., de Las Casas, D., Budden, D., Abdolmaleki, A., Merel, J., Lefrancq, A., Lillicrap, T., and Riedmiller, M. DeepMind control suite. Technical report, DeepMind, January 2018.
- Tian, Y., Krishnan, D., and Isola, P. Contrastive multiview coding, 2019.
- Torralba, A. and Efros, A. A. Unbiased look at dataset bias. In *CVPR 2011*, pp. 1521–1528. IEEE, 2011.

- Tzeng, E., Hoffman, J., Saenko, K., and Darrell, T. Adversarial discriminative domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7167–7176, 2017.
- van den Oord, A., Li, Y., and Vinyals, O. Representation learning with contrastive predictive coding, 2018.
- Wang, H., He, Z., Lipton, Z. C., and Xing, E. P. Learning robust representations by projecting superficial statistics out. *arXiv preprint arXiv:1903.06256*, 2019.
- Williamson, R. C. and Menon, A. K. Fairness risk measures, 2019.
- Wu, X., Guo, Y., Chen, J., Liang, Y., Jha, S., and Chalasani, P. Representation bayesian risk decompositions and multi-source domain adaptation, 2020.
- Zhang, C., Bengio, S., Hardt, M., Recht, B., and Vinyals, O. Understanding deep learning requires rethinking generalization. *arXiv preprint arXiv:1611.03530*, 2016.
- Zhang, H., Cisse, M., Dauphin, Y. N., and Lopez-Paz, D. mixup: Beyond empirical risk minimization, 2017.
- Zhao, H., des Combes, R. T., Zhang, K., and Gordon, G. J. On learning invariant representation for domain adaptation, 2019.

Appendices

A. Appendix Overview

Our code is available online at: <https://anonymous.4open.science/r/12747e81-8505-43cb-b54e-e75e2344a397/>. The sections of our appendix are as follows:

- A) Appendix Overview
- B) Definition and discussion of extrapolation in machine learning
- C) Illustrative examples of how REx works in toy settings
- D) A summary of different types of causal model
- E) Theory
- F) The relationship between MM-REx vs. V-REx, and the role each plays in our work
- G) Further results and details for experiments mentioned in main text
- H) Experiments not mentioned in main text
- I) Overview of other topics related to OOD generalization

B. Definition and discussion of extrapolation in machine learning

We define interpolation and extrapolation as follows: **interpolation** refers to making decisions or predictions about points *within* the convex hull of the training examples and **extrapolation** refers to making decisions or predictions about points *outside* their convex hull.⁹ This generalizes the familiar sense of these terms for one-dimensional functions. An interesting consequence of this definition is: for data of high intrinsic dimension, generalization *requires* extrapolation (Hastie et al., 2009), even in the i.i.d. setting. This is because the volume of high-dimensional manifolds concentrates near their boundary; see Figure 6.

Extrapolation in the space of risk functions. The same geometric considerations apply to extrapolating to new domains. Domains can be highly diverse, varying according to high dimensional attributes, and thus requiring extrapolation to generalize across. Thus Risk Extrapolation might often do a better job of including possible test domains in its perturbation set than Risk Interpolation does.

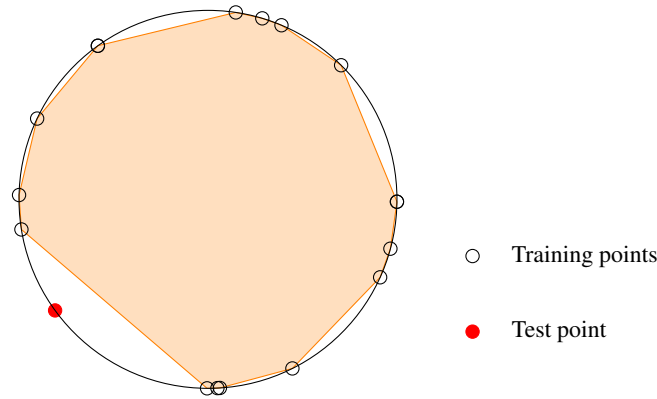


Figure 6. Illustration of the importance of extrapolation for generalizing in high dimensional space. In high dimensional spaces, mass concentrates near the boundary of objects. For instance, the uniform distribution over a ball in $N + 1$ -dimensional space can be approximated by the uniform distribution over the N -dimensional hypersphere. We illustrate this in 2 dimensions, using the 1-sphere (i.e. the unit circle). Dots represent a finite training sample, and the shaded region represents the convex hull of all but one member of the sample. Even in 2 dimensions, we can see why any point from a finite sample from such a distribution remains outside the convex hull of the other samples, with probability 1. The only exception would be if two points in the sample coincide *exactly*.

⁹Surprisingly, we were not able to find any existing definition of these terms in the machine learning literature. They have been used in this sense (Hastie et al., 2009; Haffner, 2002), but also to refer to strong generalization capabilities more generally (Sahoo et al., 2018).

C. Illustrative examples of how REx works in toy settings

Here, we work through two examples to illustrate:

1. How to understand extrapolation in the space of probability density/mass functions (PDF/PMFs)
2. How REx encourages robustness to covariate shift via distributing capacity more evenly across possible input distributions.

C.1. 6D example of REx

Here we provide a simple example illustrating how to understand extrapolations of probability distributions. Suppose $X \in \{0, 1, 2\}$ and $Y \in \{0, 1\}$, so there are a total of 6 possible types of examples, and we can represent their distributions in a particular domain as a point in 6D space: $(P(0, 0), P(0, 1), P(1, 0), P(1, 1), P(2, 0), P(2, 1))$. Now, consider three domains e_1, e_2, e_3 given by

1. (a, b, c, d, e, f)
2. $(a, b, c, d, e - k, f + k)$
3. $(2a, 2b, c(1 - \frac{a+b}{c+d}), d(1 - \frac{a+b}{c+d}), e, f)$

The difference between e_1 and e_2 corresponds to a shift in $P(Y|X = 2)$, and suggests that Y cannot be reliably predicted across different domains when $X = 2$. Meanwhile, the difference between e_1 and e_3 tells us that the relative probability of $X = 0$ vs. $X = 1$ can change, and so we might want our model to be robust to these sorts of covariate shifts. Extrapolating risks across these 3 domains effectively tells the model: “don’t bother trying to predict Y when $X = 2$ (i.e. aim for $\hat{P}(Y = 1|X = 2) = .5$), and split your capacity equally across the $X = 0$ and $X = 1$ cases”. By way of comparison, IRM would also aim for $\hat{P}(Y = 1|X = 2) = .5$, whereas ERM would aim for $\hat{P}(Y = 1|X = 2) = \frac{3f+k}{3e+3f}$ (assuming $|D_1| = |D_2| = |D_3|$). And unlike REx, both ERM and IRM would split capacity between $X = 0/1/2$ cases according to their empirical frequencies.

C.2. Covariate shift example

We now give an example to show how REx provides robustness to covariate shift. Covariate shift is an issue when a model has limited capacity or limited data.

Viewing REx as robust learning over the affine span of the training distributions reveals its potential to improve robustness to distribution shifts. Consider a situation in which a model encounters two types of inputs: COSTLY inputs with probability q and CHEAP inputs with probability $1 - q$. The model tries to predicts the input – it outputs COSTLY with probability p and CHEAP with probability $1 - p$. If the model predicts right its risk is 0, but if it predicts COSTLY instead of CHEAP it gets a risk $u = 2$, and if it predicts CHEAP instead of COSTLY it gets a risk $v = 4$. The risk has expectation $\mathcal{R}_q(p) = (1 - p)(1 - q)u + pqv$. We have access to two domains with different input probabilities $q_1 < q_2$. This is an example of pure covariate shift.

We want to guarantee the minimal risk over the set of all possible domains:

$$\min_{p \in [0, 1]} \max_{q \in [0, 1]} \mathcal{R}_q(p) = (1 - p)(1 - q)u + pqv$$

as illustrated in Figure 7. The saddle point solution of this problem is $p = \omega = u/(u+v)$ and $\mathcal{R}_q(p) = uv/(u+v), \forall q$. From the figure we see that $\mathcal{R}_{q_1}(p) = \mathcal{R}_{q_2}(p)$ can only happen for $p = \omega$, so the risk extrapolation principle will return the minimax optimal solution.

If we use ERM to minimize the risk, we will pool together the domains into a new domain with COSTLY input probability $\bar{q} = (q_1 + q_2)/2$. ERM will return $p = 0$ if $\bar{q} > \omega$ and $p = 1$ otherwise. Risk interpolation (RI) $\min_p \max_{q \in \{q_1, q_2\}} \mathcal{R}_q(p)$ will predict $p = 0$ if $q_1, q_2 > \omega$, $p = 1$ if $q_1, q_2 < \omega$ and $p = \omega$ if $q_1 < \omega < q_2$. We see that only REx finds the minimax optimum for arbitrary values of q_1 and q_2 .

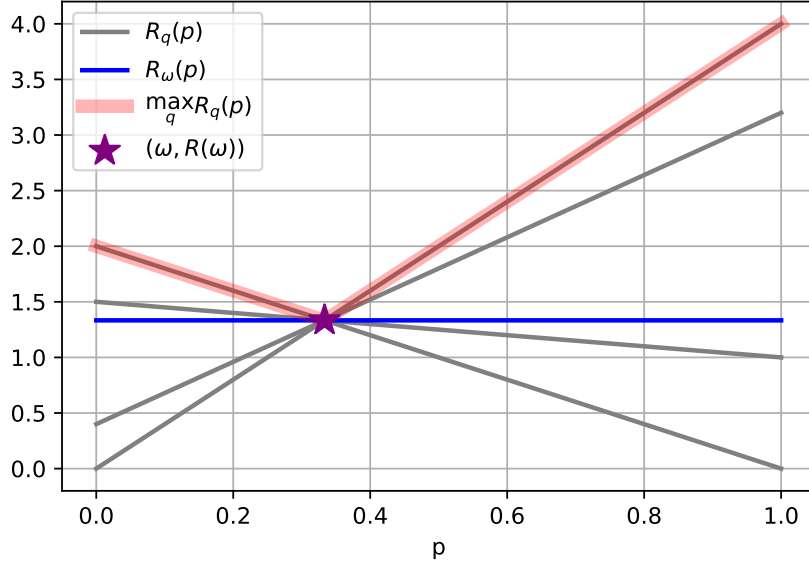


Figure 7. Each grey line is a risk $\mathcal{R}_q(p)$ as functions of p for a specific value of q . The blue line is when $q = \omega$. We highlight in red the curve $\max_q \mathcal{R}_q(p)$ whose minimum is the saddle point marked by a purple star in $p = \omega$.

D. A summary of different types of causal models

Here, we briefly summarize the differences between 3 different types of causal models, see Table 4. Our definitions and notation follow *Elements of Causal Inference: Foundations and Learning Algorithms* (Peters et al., 2017).

A **Causal Graph** is a directed acyclic graph (DAG) over a set of nodes corresponding to random variables \mathbf{Z} , where edges point from causes (including noise variables) to effects. A **Structural Causal Model (SCM)**, \mathfrak{C} , additionally specifies a *deterministic* mapping f_Z for every node Z , which computes the value of that node given the values of its parents, which include a special noise variable N_Z , which is sampled independently from all other nodes. This f_Z is called the **mechanism**, **structural equation**, or **structural assignment** for Z . Given an SCM, \mathfrak{C} , the **entailed distribution** of \mathfrak{C} , $P^\mathfrak{C}(\mathbf{Z})$ is defined via ancestral sampling. Thus for any $Z \in \mathbf{Z}$, we have that the marginal distribution $P^\mathfrak{C}(Z|\mathbf{Z} \setminus Z) = P^\mathfrak{C}(Z|Pa(Z))$. A **Causal Graphical Model (CGM)** can be thought of as specifying these marginal distributions *without* explicitly representing noise variables N_Z . We can draw rough analogies with (non-causal) statistical models. Roughly speaking, Causal Graphs are analogous to Graphical Models, whereas SCMs and CGMs are analogous to joint distributions.

Model	Independences	Distributions	Interventions	Counterfactuals
Graphical Model	✓	✗	✗	✗
Joint Distribution	✓	✓	✗	✗
Causal Graph	✓	✗	✓	✗
Causal Graphical Model	✓	✓	✓	✗
Structural Causal Model	✓	✓	✓	✓

Table 4. A comparison of causal and non-causal models.

E. Theory

E.1. Proofs of theorems 1 and 2

The REX principle (Section 3) has two goals:

1. Reducing training risks
2. Increasing similarity of training risks.

In practice, it may be advantageous to trade-off these two objectives, using a hyperparameter (e.g. β for V-REx or λ_{\min} for MM-REx). However, in this section, we assume the 2nd criteria takes priority; i.e. we define “satisfying” the REX principle as selecting a minimal risk predictor among those that achieve *exact* equality of risks across all the domains in a set \mathcal{E} .

Recall our assumptions from Section 3.2 of the main text:

1. The causes of Y are observed, i.e. $Pa(Y) \subseteq X$.
2. Domains correspond to interventions on X .
3. Homoskedasticity (a slight generalization of the additive noise setting assumed by Peters et al. (2016)). We say an SEM \mathfrak{C} is **homoskedastic** (with respect to a loss function ℓ), if the Bayes error rate of $\ell(f_Y(x), f_Y(x))$ is the same for all $x \in \mathcal{X}$.

And see Section 2.3 for relevant definitions and notation.

We begin with a theorem based on the setting explored by Peters et al. (2016). Here, $\varepsilon_i \doteq N_i$ are assumed to be normally distributed.

Theorem 1. *Given a Linear SEM, $X_i \leftarrow \sum_{j \neq i} \beta_{(i,j)} X_j + \varepsilon_i$, with $Y \doteq X_0$, and a predictor $f_\beta(X) \doteq \sum_{j:j>0} \beta_j X_j + \varepsilon_j$ that satisfies REX (with mean-squared error) over a perturbation set of domains that contains 3 distinct $do(\cdot)$ interventions for each $X_i : i > 0$. Then $\beta_j = \beta_{0,j}, \forall j$.*

Proof. We adapt the proof of Theorem 4i from Peters et al. (2016) to show that REX will learn the correct model under similar assumptions. Let $Y \leftarrow \gamma X + \varepsilon$ be the mechanism for Y , assumed to be fixed across all domains, and let $\hat{Y} = \beta X$ be our predictor. Then the residual is $R(\beta) = (\gamma - \beta)X + \varepsilon$. Define $\alpha_i \doteq \gamma_i - \beta_i$, and consider an intervention $do(X_j = x)$ on the youngest node X_j with $\alpha_j \neq 0$. Then as in eqn 36/37 of Peters et al. (2016), we compare the residuals R of this intervention and of the observational distribution:

$$R^{\text{obs}}(\beta) = \alpha_j X_j + \sum_{i \neq j} \alpha_i X_i + \varepsilon \qquad R^{do(X_j=x)}(\beta) = \alpha_j x + \sum_{i \neq j} \alpha_i X_i + \varepsilon \quad (9)$$

We now compute the MSE risk for both domains, set them equal, and simplify to find a quadratic formula for x :

$$\mathbb{E} \left[(\alpha_j X_j + \sum_{i \neq j} \alpha_i X_i + \varepsilon)^2 \right] = \mathbb{E} \left[(\alpha_j x + \sum_{i \neq j} \alpha_i X_i + \varepsilon)^2 \right] \quad (10)$$

$$0 = \alpha_j^2 x^2 + 2\alpha_j \mathbb{E} \left[\sum_{i \neq j} \alpha_i X_i + \varepsilon \right] x - \mathbb{E} \left[(\alpha_j X_j)^2 - 2\alpha_j X_j \left(\sum_{i \neq j} \alpha_i X_i + \varepsilon \right) \right] \quad (11)$$

Since there are at most two values of x that satisfy this equation, any other value leads to a violation of REX, so that α_j needs to be zero – contradiction. In particular having domains with 3 different do -interventions on every X_i guarantees that the risks are not equal across all domains. \square

Given the assumption that a predictor satisfies REx over *all* interventions that do not change the mechanism of Y , we can prove a much more general result. We now consider an arbitrary SCM, \mathfrak{C} , generating Y and X , and let \mathcal{E}^I be the set of domains corresponding to arbitrary interventions on X , similarly to Peters et al. (2016).

We emphasize that the predictor is not restricted to any particular class of models, and is a generic function $f : \mathcal{X} \rightarrow \mathcal{P}(Y)$, where $\mathcal{P}(Y)$ is the set of distributions over Y . Hence, we drop θ from the below discussion and simply use f to represent the predictor, and $\mathcal{R}(f)$ its risk.

Theorem 2. *Suppose ℓ is a (strictly) proper scoring rule. Then a predictor that satisfies REx for \mathcal{E}^I uses $f_Y(x)$ as its predictive distribution on input x for all $x \in \mathcal{X}$.*

Proof. Let $\mathcal{R}^e(f, x)$ be the loss of predictor f on point x in domain e , and $\mathcal{R}^e(f) = \int_{P^e(x)} \mathcal{R}^e(f, x)$ be the risk of f in e . Define $\iota(x)$ as the domain given by the intervention $do(X = x)$, and note that $\mathcal{R}^{\iota(x)}(f) = \mathcal{R}^{\iota(x)}(f, x)$. We additionally define $X_1 \doteq \text{Par}(Y)$.

The causal mechanism, f_Y , satisfies the REx principle over \mathcal{E}^I . For every $x \in \mathcal{X}$, $f_Y(x) = P(Y|do(X = x)) = P(Y|do(X_1 = x_1)) = P(Y|X_1 = x_1)$ is invariant (meaning ‘independent of domain’) by definition; $P(Y|do(X = x)) = P(Y|do(X_1 = x_1)) = P(Y|X_1 = x_1)$ follows from the semantics of SEM/SCMs, and the fact that we don’t allow f_Y to change across domains. Specifically Y is always generated by the same ancestral sampling process that only depends on X_1 and N^Y . Thus the risk of the predictor $f_Y(x)$ at point x , $\mathcal{R}^e(f_Y, x) = \ell(f_Y(x), f_Y(x))$ is also invariant, so is $\mathcal{R}(f_Y, x)$. Thus $\mathcal{R}^e(f_Y) = \int_{P^e(x)} \mathcal{R}^e(f_Y, x) = \int_{P^e(x)} \mathcal{R}(f_Y, x)$ is invariant whenever $\mathcal{R}(f_Y, x)$ does not depend on x , and the homoskedasticity assumption ensures that this is the case. This establishes that setting $f = f_Y$ will produce equal risk across domains.

No other predictor satisfies the REx principle over \mathcal{E}^I . We show that any other g achieves higher risk than f_Y for at least one domain. This demonstrates both that f_Y achieves minimal risk (thus satisfying REx), and that it is the unique predictor which does so (and thus no other predictors satisfy REx). We suppose such a g exists and construct a domain where it achieves higher risk than f_Y . Specifically, if $g \neq f_Y$ then let $x \in \mathcal{X}$ be a point such that $g(x) \neq f_Y(x)$. And since ℓ is a strictly proper scoring rule, this implies that $\ell(g(x), f_Y(x)) > \ell(f_Y(x), f_Y(x))$. But $\ell(g(x), f_Y(x))$ is exactly the risk of g on the domain $\iota(do(X = x))$, and thus g achieves higher risk than f_Y in $\iota(do(X = x))$, a contradiction. \square

E.2. REx as DRO

We note that MM-REx is also performing robust optimization over a convex hull, see Figure 1. The corners of this convex hull correspond to “extrapolated domains” with coefficients $(\lambda_{\min}, \lambda_{\min}, \dots, (1 - (m - 1)\lambda_{\min}))$ (up to some permutation). However, these domains do not necessarily correspond to valid probability distributions; in general, they are quasidistributions, which can assign negative probabilities to some examples. This means that, even if the original risk functions were convex, the extrapolated risks need not be. However, in the case where they *are* convex, then existing theorems, such as the convergence rate result of (Sagawa et al., 2019). This raises several important questions:

1. When is the affine combination of risks convex?
2. What are the effects of negative probabilities on the optimization problem REx faces, and the solutions ultimately found?

Negative probabilities: Figure 8 illustrates this for a case where $\mathcal{X} = \mathbb{Z}_2^2$, i.e. x is a binary vector of length 2. Suppose x_1, x_2 are independent in our training domains, and represent the distribution for a particular domain by the point $(P(X_1 = 1), P(X_2 = 1))$. And suppose our 4 training distributions have $(P(X_1 = 1), P(X_2 = 1))$ equal to $\{(.4, .1), (.4, .9), (.6, .1), (.6, .9)\}$, with $P(Y|X)$ fixed.

F. The relationship between MM-REx vs. V-REx, and the role each plays in our work

The MM-REx and V-REx methods play different roles in our work:

- We use MM-REx to illustrate that REx can be instantiated as a variant of robust optimization, specifically a generalization of the common Risk Interpolation approach. We also find MM-REx provides a useful geometric intuition, since we can visualize its perturbation set as an expansion of the convex hull of the training risks or distributions.

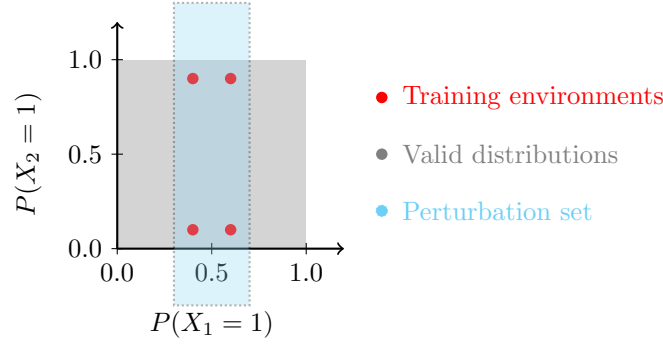


Figure 8. The perturbation set for MM-REx can include “distributions” which assign invalid (e.g. negative) probabilities to some data-points. The range of valid distributions $P(X)$ is shown in grey, and $P(X)$ for 4 different training domains are shown as red points. The interior of the dashed line shows the perturbation set for $\lambda_{\min} = -1/2$.

- We expect V-REx to be the more practical algorithm. It is simple to implement. And it performed better in our CMNIST experiments; we believe this may be due to V-REx providing a smoother gradient vector field, and thus more stable optimization, see Figure F.

Either method recovers the REx principle as a limiting case, as we prove in Section F.1. We also provide a sequence of mathematical derivations that sheds light on the relationship between MM-REx and V-REx in Section F.2 we can view these as a progression of steps for moving from the robust optimization formulation of MM-REx to the penalty term of V-REx:

1. **From minimax to closed form:** We show how to arrive at the closed-form version of MM-REx provided in Eqn. 7.
2. **Closed form as mean absolute error:** The closed form of MM-REx is equivalent to a mean absolute error (MAE) penalty term when there are only two training domains.
3. **V-REx as mean squared error:** V-REx is exactly equivalent to a mean squared error penalty term (always). Thus in the case of only two training domains, the difference between MM-REx and V-REx is just a different choice of norm.

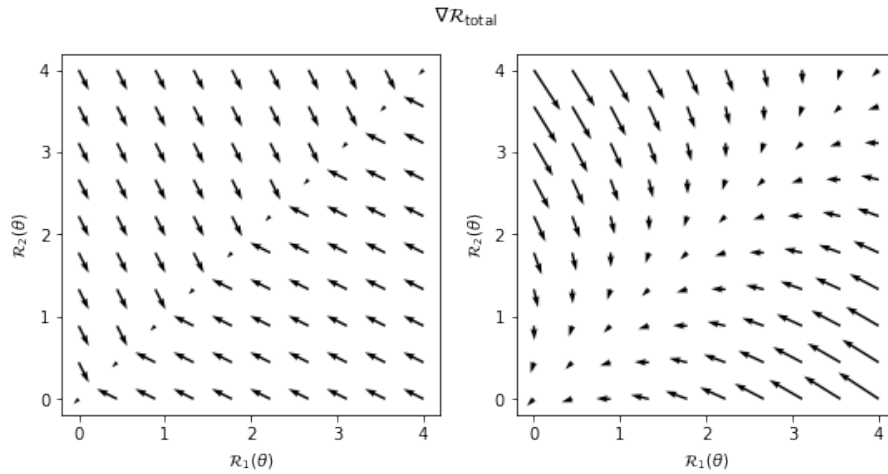


Figure 9. Vector fields of the gradient evaluated at different values of training risks $\mathcal{R}_1(\theta)$, $\mathcal{R}_2(\theta)$. We compare the gradients for $\mathcal{R}_{\text{MM-REx}}$ (left) and $\mathcal{R}_{\text{V-REx}}$ (right). Note that for $\mathcal{R}_{\text{V-REx}}$, the gradient vectors curve smoothly towards the direction of the origin, as they approach the diagonal (where training risks are equal); this leads to a smoother optimization landscape.

F.1. V-REx and MM-REx enforce the REX principle in the limit

We prove that both MM-REx and V-REx recover the constraint of perfect equality between risks in the limit of $\lambda_{\min} \rightarrow -\infty$ or $\beta \rightarrow \infty$, respectively. For both proofs, we assume all training risks are finite.

Proposition 1. *The MM-REx risk of predictor f_θ , $\mathcal{R}_{\text{MM-REx}}(\theta) \rightarrow \infty$ as $\lambda_{\min} \rightarrow -\infty$ unless $\mathcal{R}^d = \mathcal{R}^e$ for all training domains d, e .*

Proof. Suppose the risk is not equal across domains, and let the largest difference between any two training risks be $\epsilon > 0$. Then $\mathcal{R}_{\text{MM-REx}}(\theta) = (1 - m\lambda_{\min}) \max_e \mathcal{R}_e(\theta) + \lambda_{\min} \sum_{i=1}^m \mathcal{R}_i(\theta) = \max_e \mathcal{R}_e(\theta) - m\lambda_{\min} \max_e \mathcal{R}_e(\theta) + \lambda_{\min} \sum_{i=1}^m \mathcal{R}_i(\theta) \geq \max_e \mathcal{R}_e(\theta) - \lambda_{\min} \epsilon$, with the inequality resulting from matching up the m copies of $\lambda_{\min} \max_e \mathcal{R}_e$ with the terms in the sum and noticing that each pair has a non-negative value (since $\mathcal{R}_i - \max_e \mathcal{R}_e$ is non-positive and λ_{\min} is negative), and at least one pair has the value $-\lambda_{\min} \epsilon$. Thus sending $\lambda \rightarrow -\infty$ sends this lower bound on $\mathcal{R}_{\text{MM-REx}}$ to ∞ and hence $\mathcal{R}_{\text{MM-REx}} \rightarrow \infty$ as well. \square

Proposition 2. *The V-REx risk of predictor f_θ , $\mathcal{R}_{\text{V-REx}}(\theta) \rightarrow \infty$ as $\beta \rightarrow \infty$ unless $\mathcal{R}^d = \mathcal{R}^e$ for all training domains d, e .*

Proof. Again, let $\epsilon > 0$ be the largest difference in training risks, and let μ be the mean of the training risks. Then there must exist an e such that $|\mathcal{R}_e - \mu| \geq \epsilon/2$. And thus $\text{Var}_i(\mathcal{R}_i(\theta)) = \sum_i (\mathcal{R}_i - \mu)^2 \geq (\epsilon/2)^2$, since all other terms in the sum are non-negative. Since $\epsilon > 0$ by assumption, the penalty term is positive and thus $\mathcal{R}_{\text{V-REx}}(\theta) \doteq \sum_i \mathcal{R}_i(\theta) + \beta \text{Var}_i(\mathcal{R}_i(\theta))$ goes to infinity as $\beta \rightarrow \infty$. \square

F.2. Connecting MM-REx to V-REx

F.2.1. CLOSED FORM SOLUTIONS TO RISK INTERPOLATION AND MINIMAX-REx

Here, we show that risk interpolation is equivalent to the robust optimization objective of Eqn. 5. Without loss of generality, let \mathcal{R}_1 be the largest risk, so $\mathcal{R}_e \leq \mathcal{R}_1$, for all e . Thus we can express $\mathcal{R}_e = \mathcal{R}_1 - d_e$ for some non-negative d_e , with $d_1 = 0 \geq d_e$ for all e . And thus we can write the weighted sum of Eqn. 7 as:

$$\mathcal{R}_{\text{MM}}(\theta) \doteq \max_{\substack{\sum_e \lambda_e = 1 \\ \lambda_e \geq \lambda_{\min}}} \sum_{e=1}^m \lambda_e \mathcal{R}_e(\theta) \quad (12)$$

$$= \max_{\substack{\sum_e \lambda_e = 1 \\ \lambda_e \geq \lambda_{\min}}} \sum_{e=1}^m \lambda_e (\mathcal{R}_1(\theta) - d_e) \quad (13)$$

$$= \mathcal{R}_1(\theta) + \max_{\substack{\sum_e \lambda_e = 2 \\ \lambda_e \geq \lambda_{\min}}} \sum_{e=1}^m -\lambda_e (d_e) \quad (14)$$

$$(15)$$

Now, since d_e are non-negative, $-d_e$ is non-positive, and the maximal value of this sum is achieved when $\lambda_e = \lambda_{\min}$ for all $e \geq 2$, which also implies that $\lambda_1 = 1 - (m-1)\lambda_{\min}$. This yields the closed form solution provided in Eqn. 7. The special case of Risk Interpolation, where $\lambda_{\min} = 0$, yields Eqn. 5.

F.2.2. MINIMAX-REx AND MEAN ABSOLUTE ERROR REX

In the case of only two training risks, MM-REx is equivalent to using a penalty on the mean absolute error (MAE) between training risks. However, penalizing the pairwise absolute errors is not equivalent when there are $m > 2$ training risks, as we show below. Without loss of generality, assume that $\mathcal{R}_1 < \mathcal{R}_2 < \dots < \mathcal{R}_m$. Then (1/2 of) the \mathcal{R}_{MAE} penalty term is:

$$\sum_i \sum_{j \leq i} (\mathcal{R}_i - \mathcal{R}_j) = m\mathcal{R}_m - \sum_{j \leq m} \mathcal{R}_j + (m-1)\mathcal{R}_{m-1} - \sum_{j \leq m-1} \mathcal{R}_j \dots \quad (16)$$

$$= \sum_j j\mathcal{R}_j - \sum_j \sum_{i \leq j} \mathcal{R}_i \quad (17)$$

$$= \sum_j j\mathcal{R}_j - \sum_j (m-j+1)\mathcal{R}_j \quad (18)$$

$$= \sum_j (2j-m-1)\mathcal{R}_j \quad (19)$$

For $m = 2$, we have $1/2\mathcal{R}_{\text{MAE}} = (2*1-2-1)\mathcal{R}_1 + (2*2-2-1)\mathcal{R}_2 = \mathcal{R}_2 - \mathcal{R}_1$. Now, adding this penalty term with some coefficient β_{MAE} to the ERM term yields:

$$\mathcal{R}_{\text{MAE}} \doteq \mathcal{R}_1 + \mathcal{R}_2 + \beta_{\text{MAE}}(\mathcal{R}_2 - \mathcal{R}_1) = (1 - \beta_{\text{MAE}})\mathcal{R}_1 + (1 + \beta_{\text{MAE}})\mathcal{R}_2 \quad (20)$$

$$(21)$$

We wish to show that this is equal to \mathcal{R}_{MM} for an appropriate choice of learning rate γ_{MAE} and hyperparameter β_{MAE} . Still assuming that $\mathcal{R}_1 < \mathcal{R}_2$, we have that:

$$\mathcal{R}_{\text{MM}} \doteq (1 - \lambda_{\min})\mathcal{R}_2 + \lambda_{\min}\mathcal{R}_1 \quad (22)$$

Choosing $\gamma_{\text{MAE}} = 1/2\gamma_{\text{MM}}$ is equivalent to multiplying \mathcal{R}_{MM} by 2, yielding:

$$2\mathcal{R}_{\text{MM}} \doteq 2(1 - \lambda_{\min})\mathcal{R}_2 + 2\lambda_{\min}\mathcal{R}_1 \quad (23)$$

Now, in order for $\mathcal{R}_{\text{MAE}} = 2\mathcal{R}_{\text{MM}}$, we need that:

$$2 - 2\lambda_{\min} = 1 + \beta_{\text{MAE}} \quad (24)$$

$$2\lambda_{\min} = 1 - \beta_{\text{MAE}} \quad (25)$$

$$(26)$$

And this holds whenever $\beta_{\text{MAE}} = 1 - 2\lambda_{\min}$. When $m > 2$, however, these are not equivalent, since \mathcal{R}_{MM} puts equal weight on all but the highest risk, whereas \mathcal{R}_{MAE} assigns a different weight to each risk.

F.2.3. PENALIZING PAIRWISE MEAN SQUARED ERROR (MSE) YIELDS V-REX

The V-REx penalty (Eqn. 8) is equivalent to the average pairwise mean squared error between all training risks (up to a constant factor of 2). Recall that \mathcal{R}_i denotes the risk on domain i . We have:

$$\frac{1}{2n^2} \sum_i \sum_j (\mathcal{R}_i - \mathcal{R}_j)^2 = \frac{1}{2n^2} \sum_i \sum_j (\mathcal{R}_i^2 + \mathcal{R}_j^2 - 2\mathcal{R}_i\mathcal{R}_j) \quad (27)$$

$$= \frac{1}{2n} \sum_i \mathcal{R}_i^2 + \frac{1}{2n} \sum_j \mathcal{R}_j^2 - \frac{1}{n^2} \sum_i \sum_j \mathcal{R}_i\mathcal{R}_j \quad (28)$$

$$= \frac{1}{n} \sum_i \mathcal{R}_i^2 - \left(\frac{1}{n} \sum_i \mathcal{R}_i \right)^2 \quad (29)$$

$$= \text{Var}(\mathcal{R}). \quad (30)$$

G. Further results and details for experiments mentioned in main text

G.1. CMNIST with covariate shift

Here we present the following additional results:

1. Figure 1 of the main text with additional results using MM-REx, see G.1. These results used the “default” parameters from the code of Arjovsky et al. (2019).
2. A plot with results on these same tasks after performing a random search over hyperparameter values similar to that performed by Arjovsky et al. (2019).
3. A plot with the percentage of the randomly sampled hyperparameter combinations that have satisfactory ($> 50\%$) accuracy, which we count as “success” since this is better than random chance performance.

These results show that REx is able to handle greater covariate shift than IRM, given appropriate hyperparameters. Furthermore, when appropriately tuned, REx can outperform IRM in situations with covariate shift. The lower success rate of REx for high values of p is because it produces degenerate results (where training accuracy is less than test accuracy) more often.

The hyperparameter search consisted of a uniformly random search of 340 samples over the following intervals of the hyperparameters:

1. HiddenDim = $[2^{**}7, 2^{**}12]$
2. L2RegularizerWeight = $[10^{**-2}, 10^{**-4}]$
3. Lr = $[10^{**-2.8}, 10^{**-4.3}]$
4. PenaltyAnnealIters = $[50, 250]$
5. PenaltyWeight = $[10^{**2}, 10^{**6}]$
6. Steps = $[201, 601]$

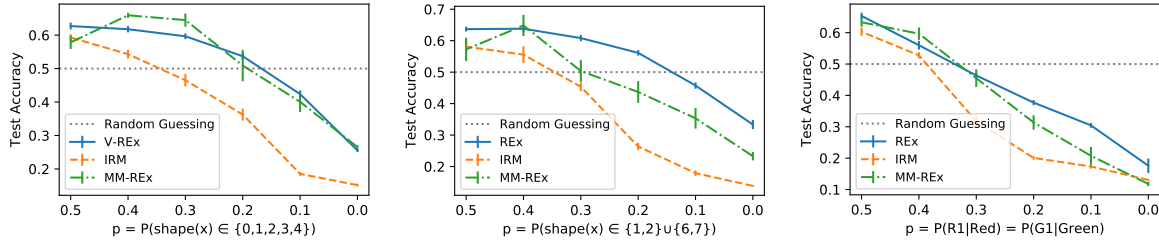


Figure 10. This is Figure 3.2 of main text with additional results using MM-REx. For each covariate shift variant (class imbalance, digit imbalance, and color imbalance from left to right as described in “CMNIST with covariate shift” subsection of Section 4.1 in main text) of CMNIST, the standard error (the vertical bars in plots) is higher for MM-REx than for V-REx.

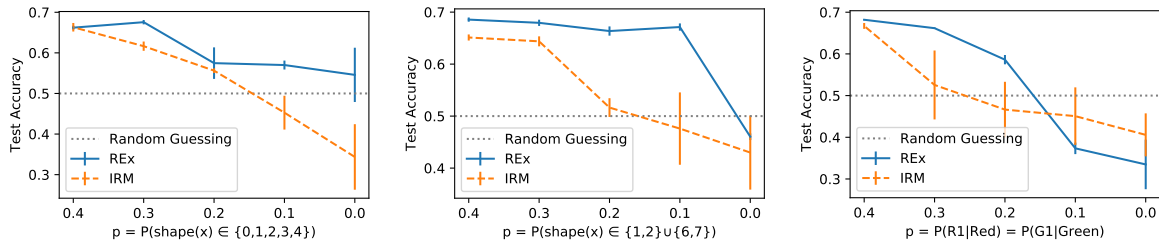


Figure 11. This is Figure 3.2 of main text (class imbalance, digit imbalance, and color imbalance from left to right as described in “CMNIST with covariate shift” subsection of Section 4.1 in main text), but with hyperparameters of REx and IRM each tuned to perform as well as possible for each value of p for each covariate shift type.

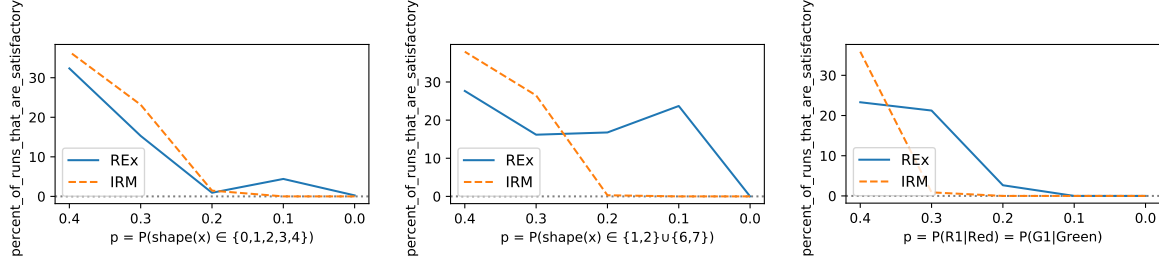


Figure 12. This also corresponds to class imbalance, digit imbalance, and color imbalance from left to right as described in "CMNIST with covariate shift" subsection of Section 4.1 in main text; but now the y-axis refers to what percentage of the randomly sampled hyperparameter combinations we deemed to be satisfactory. We define satisfactory as simultaneously being better than random guessing and having train accuracy greater than test accuracy. For p less than .5, a larger percentage of hyperparameter combinations are often satisfactory for REX than for IRM; for p greater than .5, a larger percentage of hyperparameter combinations are often satisfactory for IRM than for REX because train accuracy is greater than test accuracy for more hyperparameter combinations for IRM. We stipulate that train accuracy must be greater than test accuracy because test accuracy being greater than train accuracy usually means the model has learned a degenerate prediction rule such as "not color".

G.2. SEMs from "Invariant Risk Minimization"

Here we present experiments on the (linear) structural equation model (SEM) tasks introduced by Arjovsky et al. (2019). Arjovsky et al. (2019) construct several varieties of SEM where the task is to predict targets Y from inputs X_1, X_2 , where X_1 are (non-anti-causal) causes of Y , and X_2 are (anti-causal) effects of Y . We refer the reader to Section 5.1 and Figure 3 of Arjovsky et al. (2019) for more details. We use the same experimental settings as Arjovsky et al. (2019) (except we only run 7 trials), and report results in Table 5.

These experiments include several variants of a simple SEM, given by:

$$\begin{aligned} X_1 &= N_1 \\ Y &= W_{1 \rightarrow Y} X_1 + N_Y \\ X_2 &= W_{Y \rightarrow 2} Y + N_2 \end{aligned}$$

Where N_1, N_Y, N_2 are all sampled i.i.d. from normal distributions. The variance of these distributions may vary across domains.

While REX achieves good performance in the **domain-homoskedastic** case, it performs poorly in the **domain-heteroskedastic** case, where the amount of intrinsic noise, σ_y^2 in the target changes across domains.¹⁰ Intuitively, this is because the irreducible error varies across domains in these tasks, meaning that the risk will be larger on some domains than others, even if the model's predictions match the expectation $\mathbb{E}(Y|Pa(Y))$. We tried using a "baseline" (see Eqn. 5) of $r_e = \text{Var}(Y_e)$ (Meinshausen et al., 2015) to account for the different noise levels in Y , but this did not work.

We include a mathematical analysis of the simple SEM given above in order to better understand why REX succeeds in the domain-homoskedastic, but not the domain-heteroskedastic case. Assuming that Y, X_1, X_2 are scalars, this SEM becomes

$$\begin{aligned} X_1 &= N_1 \\ Y &= w_{1 \rightarrow y} N_1 + N_Y \\ X_2 &= w_{y \rightarrow 2} w_{1 \rightarrow y} N_1 + w_{y \rightarrow 2} N_Y + N_2 \end{aligned}$$

We consider learning a model $\hat{Y} = \alpha X_1 + \beta X_2$. Then the residual is:

$$\hat{Y} - Y = (\alpha + w_{1 \rightarrow y}(\beta w_{y \rightarrow 2} - 1))N_1 + (\beta w_{y \rightarrow 2} - 1)N_Y + \beta N_2$$

Since all random variables have zero mean, the MSE loss is the variance of the residual. Using the fact that the noise N_1, N_Y, N_2 are independent, this equals:

$$\mathbb{E}[(\hat{Y} - Y)^2] = (\alpha + w_{1 \rightarrow y}(\beta w_{y \rightarrow 2} - 1))^2 \sigma_1^2 + (\beta w_{y \rightarrow 2} - 1)^2 \sigma_Y^2 + \beta^2 \sigma_2^2$$

¹⁰See Footnote 6.

Out-of-Distribution Generalization via Risk Extrapolation

	FOU(c)	FOU(nc)	FOS(c)	FOS(nc)
IRM	0.001±0.000	0.001±0.000	0.001±0.000	0.000±0.000
REx, $r_e = 0$	0.001±0.000	0.008±0.002	0.007±0.002	0.000±0.000
REx, $r_e = \mathbb{V}(Y_e)$	0.816±0.149	1.417±0.442	0.919±0.091	0.000±0.000
	POU(c)	POU(nc)	POS(c)	POS(nc)
IRM	0.004±0.001	0.006±0.003	0.002±0.000	0.000±0.000
REx, $r_e = 0$	0.004±0.001	0.004±0.001	0.002±0.000	0.000±0.000
REx, $r_e = \mathbb{V}(Y_e)$	0.915±0.055	1.113±0.085	0.937±0.090	0.000±0.000
	FEU(c)	FEU(nc)	FES(c)	FES(nc)
IRM	0.0053±0.0015	0.1025±0.0173	0.0393±0.0054	0.0000±0.0000
REx, $r_e = 0$	0.0390±0.0089	19.1518±3.3012	7.7646±1.1865	0.0000±0.0000
REx, $r_e = \mathbb{V}(Y_e)$	0.7713±0.1402	1.0358±0.1214	0.8603±0.0233	0.0000±0.0000
	PEU(c)	PEU(nc)	PES(c)	PES(nc)
IRM	0.0102±0.0029	0.0991±0.0216	0.0510±0.0049	0.0000±0.0000
REx, $r_e = 0$	0.0784±0.0211	46.7235±11.7409	8.3640±2.6108	0.0000±0.0000
REx, $r_e = \mathbb{V}(Y_e)$	1.0597±0.0829	0.9946±0.0487	1.0252±0.0819	0.0000±0.0000

Table 5. Average mean-squared error between true and estimated weights on causal (X_1) and non-causal (X_2) variables. **Top 2:** When the level of noise in the anti-causal features varies across domains, REx performs well (FOU, FOS, POU, POS). **Bottom 2:** When the level of noise in the targets varies instead, REx performs poorly (FEU, FES, PEU, PES). Using the baselines $r_e = \mathbb{V}(Y)$ does not solve the problem, and indeed, hurts performance on the homoskedastic domains.

Thus when (only) σ_2 changes, the only way to keep the loss unchanged is to set the coefficient in front of σ_2 to 0, meaning $\beta = 0$. By minimizing the loss, we then recover $\alpha = w_{1 \rightarrow y}$; i.e. in the domain-homoskedastic setting, the loss equality constraint of REx yields the causal model. On the other hand, if (only) σ_Y changes, then REx enforces $\beta = 1/w_{y \rightarrow 2}$, which then induces $\alpha = 0$, recovering the *anticausal* model.

While REx (like ICP (Peters et al., 2016)) assumes the mechanism for Y is fixed across domains (meaning $P(Y|Pa(Y))$ is independent of the domain, e), IRM makes the somewhat weaker assumption that $\mathbb{E}(Y|Pa(Y))$ is independent of domain. While it is plausible that an appropriately designed variant of REx could work under this weaker assumption, we believe forbidding interventions on Y is not overly restrictive, and such an extension for future work.

G.3. Reinforcement Learning Experiments

Here we provide details and further results on the experiments in Section 4.1. We take tasks from the Deepmind Control Suite (Tassa et al., 2018) and modify the original state, \mathbf{s} , to produce observation, $\mathbf{o} = (\mathbf{s} + \epsilon, \eta \mathbf{s}')$ including noise ϵ and spurious features $\eta \mathbf{s}'$, where \mathbf{s}' contains 1 or 2 dimensions of \mathbf{s} . The scaling factor takes values $\eta = 1/2/3$ for the two training and test domains, respectively. The agent takes \mathbf{o} as input and learns a representation using Soft Actor-Critic (Haarnoja et al., 2018) and an auxiliary reward predictor, which is trained to predict the next 3 rewards conditioned on the next 3 actions. Since the spurious features are copied from the state before the noise is added, they are more informative for the reward prediction task, but they do not have an invariant relationship with the reward because of the domain-dependent η .

The hyperparameters used for training Soft Actor-Critic can be found in Table 6. We used `cartpole_swingup` as a development task to tune the hyperparameters of penalty weight (chosen from $[0.01, 0.1, 1, 10]$) and number of iterations before the penalty is turned up (chosen from $[5000, 10000, 20000]$), both for REx and IRM. The plots with the hyperparameter sweep are in Figure 13.

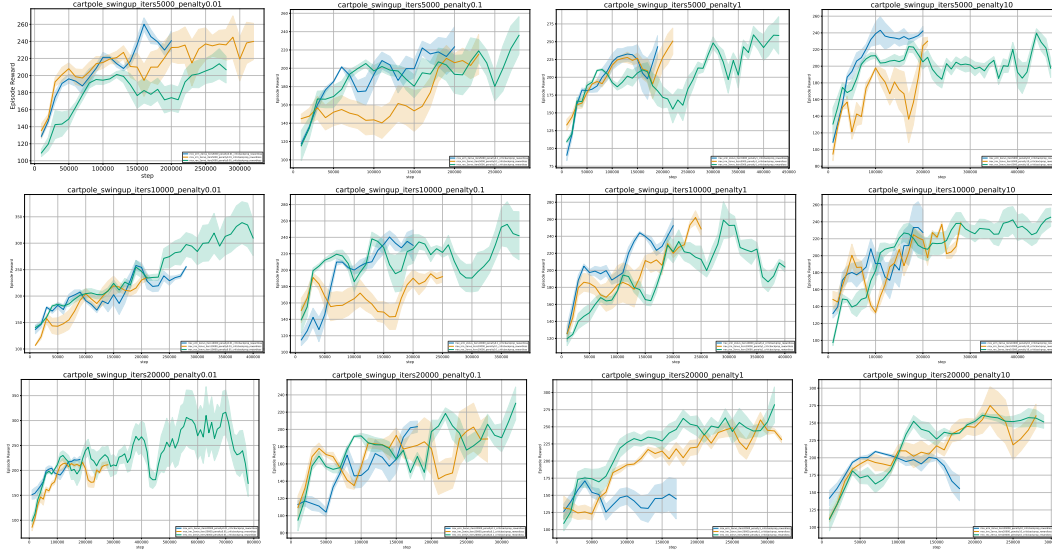


Figure 13. Hyperparameter sweep for IRM and REX on `cartpole_swingup`. Green, blue, and orange curves correspond to REX, ERM, and IRM, respectively. The subfigure titles state the penalty strength (“penalty”) and after how many iterations the penalty strength was increased (“iters”). We chose a penalty factor of 1 and 10k iterations.

Parameter name	Value
Replay buffer capacity	1000000
Batch size	1024
Discount γ	0.99
Optimizer	Adam
Critic learning rate	10^{-5}
Critic target update frequency	2
Critic Q-function soft-update rate τ_Q	0.005
Critic encoder soft-update rate τ_{enc}	0.005
Actor learning rate	10^{-5}
Actor update frequency	2
Actor log stddev bounds	$[-5, 2]$
Encoder learning rate	10^{-5}
Decoder learning rate	10^{-5}
Decoder weight decay	10^{-7}
L1 regularization weight	10^{-5}
Temperature learning rate	10^{-4}
Temperature Adam’s β_1	0.9
Init temperature	0.1

Table 6. A complete overview of hyperparameters used for reinforcement learning experiments.

H. Experiments not mentioned in main text

We include several other experiments which do not contribute directly to the core message of our paper. Here is a summary of the take-aways from these experiments:

1. Our experiments in the CMNIST domain suggest that the IRM/V-REx penalty terms should be amplified exactly when the model starts overfitting training distributions.
2. Our financial indicators experiments suggest that IRM and REx often perform remarkably similarly in practice.

H.1. A possible approach to scheduling IRM/REx penalties

We’ve found that REx and IRM are quite sensitive to the choice of hyperparameters. In particular, hyperparameters controlling the scheduling of the IRM/V-REx penalty terms are of critical importance. For the best performance, the penalty should be increased the relative weight of the penalty term after approximately 100 epochs of training (using a so-called “waterfall” schedule (Desjardins et al., 2015)). See Figure 14(b) for a comparison. We also tried an exponential decay schedule instead of the waterfall and found the results (not reported) were significantly worse, although still above 50% accuracy.

Given the methodological constraints of out-of-distribution generalization mentioned in (Gulrajani & Lopez-Paz, 2020), this could be a significant practical issue for applying these algorithms. We aim to address this limitation by providing a guideline for when to increase the penalty weight, based only on the training domains. We hypothesize that successful learning of causal features using REx or IRM should proceed in two stages:

1. In the first stage, predictive features are learned.
2. In the second stage, causal features are selected and/or predictive features are fine-tuned for stability.

This viewpoint suggests that we could use overfitting on the *training* tasks as an indicator for when to apply (or increase) the IRM or REx penalty.

The experiments presented in this section provide *observational* evidence consistent with this hypothesis. However, since the hypothesis was developed by observing patterns in the CMNIST training runs, it requires further experimental validation on a different task, which we leave for future work.

H.1.1. RESULTS AND INTERPRETATION

In Figure 14, we demonstrate that the optimal point to apply the waterfall in the CMNIST task is after predictive features have been learned, but before the model starts to memorize training examples. Before predictive features are available, the penalty terms push the model to learn a constant predictor, impeding further learning. And after the model starts to memorize, it become difficult to distinguish anti-causal and causal features. This second effect is because neural networks often have the capacity to memorize all training examples given sufficient training time, achieving and near-0 loss (Zhang et al., 2016). In the limits of this memorization regime, the differences between losses become small, and gradients of the loss typically do as well, and so the REx and IRMv1 penalties no longer provide a strong or meaningful training signal, see Figure 15.

H.2. Domain Generalization: VLCS and PACS

Here we provide earlier experiments on the VLCS and PACS dataset. We removed these experiments from the main text of our paper in favor of the more complete DomainBed results.

To test whether REx provides a benefit on more realistic domain generalization tasks, we compared REx, IRM and ERM performance on the VLCS (Torralba & Efros, 2011) and PACS (Li et al., 2017) image datasets. Both datasets are commonly-used for multi-source domain generalization. The task is to train on three domains and generalize to a fourth one at test time.

Since every domain in PACS is used as a test set when training on the other three domains, it is not possible to perform a methodologically sound evaluation on PACS after examining results on *any* of the data. Thus to avoid performing any

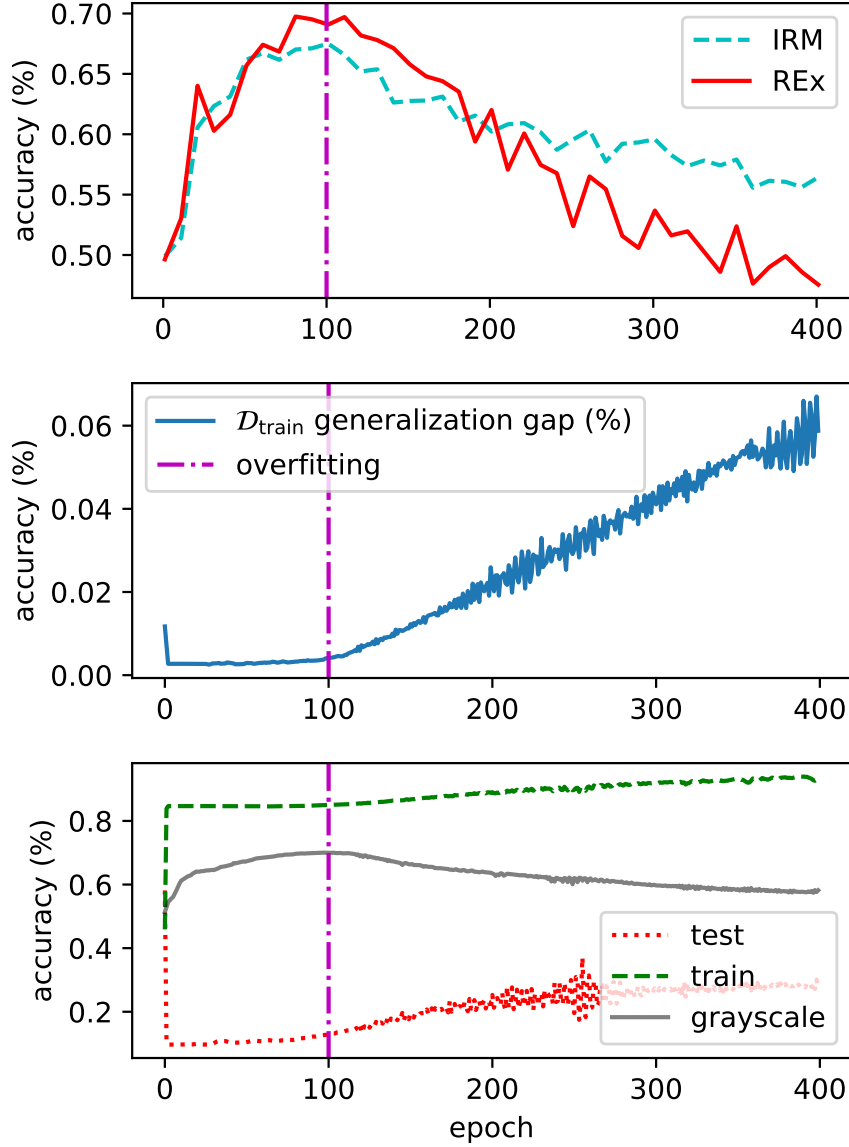


Figure 14. Stability penalties should be applied around when traditional overfitting begins, to ensure that the model has learned predictive features, and that penalties still give meaningful training signals. **Top:** Test accuracy as a function of epoch at which penalty term weight is increased (learning rate is simultaneously decreased proportionally). Choosing this hyperparameter correctly is essential for good performance. **Middle:** Generalization gap on a validation set with 85% correlation between color and label (the same as the average training correlation). The best test accuracy is achieved by increasing the penalty when the generalization gap begins to increase. The increase clearly indicates memorization because color and shape are only 85%/75% correlated with the label, and so cannot be used to make predictions with higher than 85% accuracy. **Bottom:** Accuracy on training/test sets, as well as an auxiliary grayscale set. Training/test performance reach 85%/15% after a few epochs of training, but grayscale performance improves, showing that meaningful features are still being learned.

tuning on test distributions, we use VLCS to tune hyperparameters and then apply these exact same settings to PACS and report the final average over 10 runs on each domain.

We use the same architecture, training procedure and data augmentation strategy as the (formerly) state-of-the-art Jigsaw Puzzle approach (Carlucci et al., 2019) (except with IRM or V-REx instead of Jigsaw as auxiliary loss) for all three methods. As runs are very noisy, we ran each experiment 10 times, and report average test accuracies extracted at the time of the

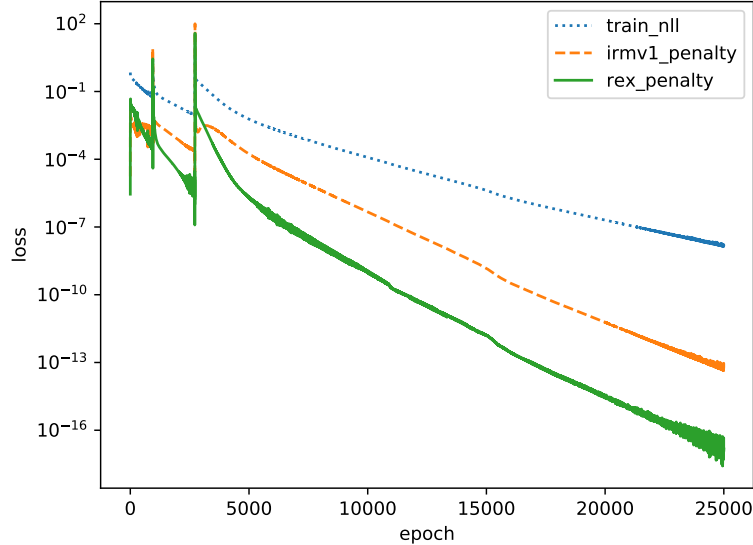


Figure 15. Given sufficient training time, empirical risk minimization (ERM) minimizes both REx and IRMv1 penalty terms on Colored MNIST (*without* including either term in the loss function). This is because the model (a deep network) has sufficient capacity to fit the training sets almost perfectly. This prevents these penalties from having the intended effect, once the model has started to overfit. The y-axis is in log-scale.

highest validation accuracy on each run. Results on PACS are in Table 8. On PACS we found that REx outperforms IRM and IRM outperforms ERM on average, while all are worse than the state-of-the-art Jigsaw method.

We use all hyperparameters from the original Jigsaw codebase.¹¹ We use Imagenet pre-trained AlexNet features and chose batch-size, learning rate, as well as penalty weights based on performance on the VLCS dataset where test performance on the holdout domain was used for the set of parameters producing the highest validation accuracy. The best performing parameters on VLCS were then applied to the PACS dataset without further changes. We searched over batch-sizes in $\{128, 384\}$, over penalty strengths in $\{0.0001, 0.001, 0.01, 0.1, 1, 10\}$, learning rates in $\{0.001, 0.01\}$ and used average performance over all 4 VLCS domains to pick the best performing hyperparameters. Table 7 shows results on VLCS with the best performing hyperparameters.

The final parameters for all methods on PACS were a batch size of 384 with 30 epochs of training with Adam, using a learning rate of 0.001, and multiplying it by 0.1 after 24 epochs (this step schedule was taken from the Jigsaw repo). The penalty weight chosen for Jigsaw was 0.9; for IRM and REx it was 0.1. We used the same data-augmentation pipeline as the original Jigsaw code for ERM, IRM, Jigsaw and REx to allow for a fair comparison.

VLCS	CALTECH	SUN	PASCAL	LABELME	Average
REx (ours)	96.72	63.68	72.41	60.40	73.30
IRM	95.99	62.85	71.71	59.61	72.54
ERM	94.76	61.92	69.03	60.55	71.56
Jigsaw (SOTA)	96.46	63.84	70.49	60.06	72.71

Table 7. Accuracy (percent) of different methods on the VLCS task. Results are test accuracy at the time of the highest validation accuracy, averaged over 10 runs. On VLCS REx outperforms all other methods. Numbers are shown in strike-through because we selected our hyperparameters based on highest test set performance; the goal of this experiment was to find suitable hyperparameters for the PACS experiment.

¹¹<https://github.com/fmcarlucci/JigenDG>

PACS	Art Painting	Cartoon	Sketch	Photo	Average
REx (ours)	66.27 \pm 0.46	68.8 \pm 0.28	59.57 \pm 0.78	89.60 \pm 0.12	71.07
IRM	66.46 \pm 0.31	68.60 \pm 0.40	58.66 \pm 0.73	89.94 \pm 0.13	70.91
ERM	66.01 \pm 0.22	68.62 \pm 0.36	58.38 \pm 0.60	89.40 \pm 0.18	70.60
Jigsaw (SOTA)	66.96 \pm 0.39	66.67 \pm 0.41	61.27 \pm 0.73	89.54 \pm 0.19	71.11

Table 8. Accuracy (percent) of different methods on the PACS task. Results are test accuracy at the time of the highest validation accuracy, averaged over 10 runs. REx outperforms ERM on average, and performs similar to IRM and Jigsaw (the state-of-the-art).



Figure 16. Financial indicators tasks. The left panel indicates the set of training domains; the middle and right panels show the test accuracy on the respective domains relative to ERM (a black dot corresponds to a training domain; a colored patch indicates the test accuracy on the respective domain.)

H.3. Financial indicators

We find that IRM and REx seem to perform similarly across different splits of the data in a prediction task using financial data. The dataset is split into five years, 2014–18, containing 37 publicly reported financial indicators of several thousand publicly listed companies each. The task is to predict if a company’s value will increase or decrease in the following year (see Appendix for dataset details.) We consider each year a different domain, and create 20 different tasks by selecting all possible combinations of domains where three domains represent the training sets, one domain the validation set, and another one the test set. We train an MLP using the validation set to determine an early stopping point, with $\beta = 10^4$. The per-task results summarized in fig. 16 indicate substantial differences between ERM and IRM, and ERM and REx. The predictions produced by IRM and REx, however, only differ insignificantly, highlighting the similarity of IRM and REx. While performance on specific tasks differs significantly between ERM and IRM/REx, performance averaged over tasks is not significantly different.

H.3.1. EXPERIMENT DETAILS

We use $\nabla 1$ of the dataset published on ¹² and prepare the data as described in ¹³. We further remove all the variables that are not shared across all 5 years, leaving us with 37 features, and whiten the data through centering and normalizing by the standard deviation.

On each subtask, we train an MLP with two hidden layers of size 128 with tanh activations and dropout ($p=0.5$) after each layer. We optimize the binary cross-entropy loss using Adam (learning rate 0.001, $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 10^{-8}$), and an L2 penalty (weight 0.001). In the IRM/REx experiments, the respective penalty is added to the loss ($\beta = 1$) and the original loss is scaled by a factor 10^{-4} after 1000 iterations. Experiments are run for a maximum of 9000 training iterations with early stopping based on the validation performance. All results are averaged over 3 trials. The overall performance of the different models, averaged over all tasks, is summarized in Tab. 9. The difference in average performance between ERM, IRM, and REx is not statistically significant, as the error bars are very large.

¹²<https://www.kaggle.com/cnic92/200-financial-indicators-of-us-stocks-20142018>

¹³<https://www.kaggle.com/cnic92/explore-and-clean-financial-indicators-dataset>

	Overall accuracy	Min acc.	Max acc.
ERM	54.6 ± 4.6	47.6	66.2
IRM	55.3 ± 5.9	45.9	67.5
REx	55.5 ± 6.0	47.2	68.0

Table 9. Test accuracy of models trained on the financial domain dataset, averaged over all 20 tasks, as well as min./max. accuracy across the tasks.

I. Overview of other topics related to OOD generalization

Domain adaptation (Ben-David et al., 2010a) shares the goal of generalizing to new distributions at test time, but allows some access to the test distribution. A common approach is to make different domains have a similar distribution of features (Pan et al., 2010). A popular deep learning method for doing so is Adversarial Domain Adaptation (ADA) (Ganin et al., 2016; Tzeng et al., 2017; Long et al., 2018; Li et al., 2018), which seeks a “invariant representation” of the inputs, i.e. one whose distribution is domain-independent. Recent works have identified fundamental shortcomings with this approach, however (Zhao et al., 2019; Johansson et al., 2019; Arjovsky et al., 2019; Wu et al., 2020).

Complementary to the goal of domain generalization is **out-of-distribution detection** (Hendrycks & Gimpel, 2016; Hendrycks et al., 2018), where the goal is to recognize examples as belonging to a new domain. Three common deep learning techniques that can improve OOD generalization are **adversarial training** (Goodfellow et al., 2014; Hendrycks & Dietterich, 2019), **self-supervised learning** (van den Oord et al., 2018; Hjelm et al., 2018; Hendrycks et al., 2019b; Albuquerque et al., 2020) and **data augmentation** (Krizhevsky et al., 2012; Zhang et al., 2017; Cubuk et al., 2018; Shorten & Khoshgoftaar, 2019; Hendrycks et al., 2019a; Carlucci et al., 2019). These methods can also be combined effectively in various ways (Tian et al., 2019; Bachman et al., 2019; Goyal et al., 2019). Data augmentation and self-supervised learning methods typically use prior knowledge such as 2D image structure. Several recent works also use prior knowledge to design augmentation strategies for invariance to superficial features that may be spuriously correlated with labels in object recognition tasks (He et al., 2019; Wang et al., 2019; Goyal et al., 2019; Ilse et al., 2020). In contrast, REx can discover which features have invariant relationships with the label without such prior knowledge.