

# ALIGNMENT, CONVEXITY AND COMPLETENESS: MECHANISMS BEHIND GROUPO

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Models trained with Empirical Risk Minimization (ERM) often fail to generalize under spurious correlations. Group Robustness Methods (GRMs)—notably Group DRO (GDRO)—mitigate this by reweighting losses across groups defined by labels and spurious attributes, yet why they work remains only partially understood. We study the learning dynamics of GRMs and their effects on both the classifier head and the representation. Theoretically, in a fine-tuning setting (fixed features), we analyze the classifier learned by GDRO and show: (i) GDRO aligns less with a spurious classifier and more with an oracle non-spurious classifier than ERM; (ii) when group losses are  $\mu$ -strongly convex, the alignment gap controls performance, yielding an upper bound on the worst-group performance gap between ERM and GDRO; and (iii) for convex losses, adding L2 regularization induces  $\mu$ -strong convexity, so the same guarantees apply—providing an explanation for the empirical gains of GDRO with L2 reported in prior work. Empirically, across standard image and text benchmarks, we confirm the predicted alignment behavior. Beyond the head, under end-to-end training GDRO also reshapes the representation: through a measure called Completeness, we show that task-relevant information is spread across multiple dimensions in GDRO while ERM tends to concentrate it in fewer, making it more susceptible to rely on spurious attributes for prediction. Together, our theory and measurements clarify the mechanisms by which GDRO outperforms ERM.

## 1 INTRODUCTION

Deep learning has achieved remarkable progress in the past decade, yet models trained under the conventional Empirical Risk Minimization (ERM) paradigm (Vapnik, 1991) often fail when facing spurious correlations. For example, a model distinguishing cows from camels may rely on background cues (grass vs. beach) rather than animal shape, leading to failures under distribution shift. This phenomenon—studied under names such as *spurious correlations* (Arjovsky et al., 2020), *simplicity bias* (Shah et al., 2020), or *shortcut learning* (Geirhos et al., 2020)—appears not only in vision but also across text and other modalities (Williams et al., 2018a; Pavlopoulos et al., 2020). To address this, Group Robustness Methods (GRMs) have emerged as state-of-the-art approaches. By partitioning data into groups defined by labels and spurious attributes, they optimize for worst-group performance (Hu et al., 2018). Among them, Group Distributionally Robust Optimization (GDRO) (Sagawa\* et al., 2020) has proven particularly effective, inspiring a broad line of follow-up work (Seo et al., 2022; Sohoni et al., 2020; Liu et al., 2021; Idrissi et al., 2022; Kirichenko et al., 2023). Yet, despite their empirical success, the mechanisms by which GRMs, and GDRO in particular, outperform ERM remain only partially understood. Newer methods emphasize the role of the classifier layer (e.g., Deep Feature Reweighting) (Kirichenko et al., 2023; Idrissi et al., 2022) and disregard any value in the representation learning (Bengio et al., 2014) capabilities of GDRO over ERM. These methods, however, tend to work well only when the classifier is finetuned over *unseen data*, while GDRO works well using only training data. Additionally, why GDRO requires strong L2 regularization to perform (Sagawa\* et al., 2020) is still a mystery. We believe a systematic understanding of the learning dynamics of GDRO is due.

In this work, we investigate the learning dynamics of GRMs both theoretically and empirically. We focus first on the classifier, showing that GDRO induces less alignment with spurious classifiers and more with oracle non-spurious ones. We then show that this difference in spurious alignment upper bounds the difference in worst-group performance for  $\mu$ -strongly convex losses. Using this result,

we show that L2 regularization makes cross-entropy losses  $\mu$ -strongly convex—providing a mechanism by which L2 regularization is needed by GDRO. We then extend the analysis to representations, demonstrating that GDRO does not discard spurious features, but rather distributes task-relevant information across multiple dimensions, reducing reliance on any single spurious attribute. Together, these findings clarify the mechanisms underlying the success of GDRO and highlight broader principles for robust learning under spurious correlations.

**Contributions.** Our contributions can be summarized as follows:

- **Theoretical analysis of the classifier:** We show that GDRO produces classifier heads that align less with spurious classifiers and more with oracle non-spurious classifiers than ERM. For  $\mu$ -strongly convex losses, we prove that this alignment gap yields an upper bound on the worst-group performance difference between ERM and GDRO.
- **Role of regularization:** Using the earlier link between alignment and performance for  $\mu$ -strong losses, we show that GDRO needs L2 regularization on the (non- $\mu$ -strong) Cross-Entropy loss because this regularization makes the loss  $\mu$ -strong.
- **Representation analysis:** Beyond the classifier, we show that GDRO reshapes learned representations. We find that GDRO distributes task-relevant information across multiple dimensions, making classifiers less dependent on individual spurious attributes.
- **Empirical validation:** Across standard vision and text benchmarks, we confirm the predicted alignment dynamics and representation effects, clarifying the mechanisms underlying GDRO’s superior performance over ERM.

## 2 PRELIMINARIES

To begin our analysis we will first state the specific problem we will be studying and provide definitions for key concepts such as *spurious* and *non-spurious classifier*, which will be used during the rest of the article. Then, we proceed to formally define the methods we will be studying along this article which are 4: Empirical Risk Minimization (ERM), Reweighting (RW), Group Distributionally Robust Optimization (GDRO) and Subsampling (SUBG). Finally, we describe datasets and models used for our experiments.

### 2.1 PROBLEM FORMULATION

We’ll work on a binary classification problem with input features  $\vec{x} = \{\vec{x}_{inv}, \vec{x}_{sp}\}$ , where  $\vec{x}_{inv}$  represents invariant features and  $\vec{x}_{sp}$  is a spurious attribute  $A$  correlated with correlation  $c_A$  with the ground truth label  $y$ . We will assume  $A$  may take two values:  $\{\mathcal{B}, -\mathcal{B}\}$ . From  $A$ , we will partition our dataset into  $G = 4$  groups with respect to  $\vec{x}$  and  $y$ . Let  $\mathcal{G}_A = \{s0_A, s1_A, n0_A, n1_A\}$ . These are groups defined by whether the spurious correlation is predictive (“s”) of the ground truth label or not (“n”), and the value of the ground truth label (0,1). Let  $N_g$  be the size of group  $g$  and  $N = \sum_{g \in \mathcal{G}} N_g$  is the total amount of samples.

We will also assume that we will train a model on this problem using SGD, and that this model can be separated into an arbitrary non-linear feature extractor  $\Phi$  and a linear classifier  $C(\vec{x}) = W \cdot \vec{x} + \vec{b}$ . The output of the feature extractor will be a set of features that depend on the invariant and spurious features of the data  $= \Phi(\vec{x}_{inv}, \vec{x}_{sp})$ . Our final model then takes the form  $M(\vec{x}) = C(\Phi(\vec{x}))$

**Definition 2.1.** *Spurious Classifier.* Let  $\alpha = \vec{v} \cdot \vec{x}$ . Define classifier  $\vec{v}$  as spurious w.r.t  $\mathcal{G}_A$  if:

$$\vec{x} \in \{s0_A, n1_A\} \implies \alpha > 0, \vec{x} \in \{s1_A, n0_A\} \implies \alpha < 0$$

**Definition 2.2.** *Non-Spurious Classifier.* Let  $\alpha = \vec{v} \cdot \vec{x}$ . Define classifier  $\vec{v}$  as non-spurious w.r.t  $\mathcal{G}_A$  if:

$$\vec{x} \in \{s1_A, n1_A\} \implies \alpha > 0, \vec{x} \in \{s0_A, n0_A\} \implies \alpha < 0$$

The former definition indicates a classifier that perfectly classifies using  $A$ , while the latter does so according to  $y$ . The theoretical analyses done in Section 3 will be done according to these two definitions.

## 2.2 ROBUSTNESS METHODS

Group Robustness Methods (GRMs) are based around the idea of reweighting examples from the training set. The most successful methods like Group DRO and Reweighting assume the training distribution to be a mixture of  $G$  groups. In practice, what they do is to create  $G = |\mathcal{A}| \times |\mathcal{Y}|$  groups. Where  $|\mathcal{A}|$  is the number of possible values an input spurious attribute might have, while  $|\mathcal{Y}|$  is the number of classes in the classification problem. To define these groups requires previous knowledge of spurious correlations in the data. Plenty of methods attempt to forego this requirement for group labelling in the training phase by finding proxies for it through clustering (Seo et al., 2022) or other means (Liu et al., 2021). The most general form of these methods has the following Loss Function:

$$\mathcal{L}^{GEN} = \sum_{i=1}^N p_i \cdot \mathcal{L}(x_i, y_i)$$

Where  $p_i$  is simply a weighting factor for the  $i$ -th example which can be either a constant or more complex function. Our mathematical definitions for these methods come from the actual implementations of these methods, in particular, GDRO’s definition is taken from the implementation of GDRO used in JTT (Liu et al., 2021).

### 2.2.1 EMPIRICAL RISK MINIMIZATION (ERM)

This is the traditional training loss where the average loss across the training data is minimized. This method suffers heavily from spurious correlations. Its loss function becomes:

$$\mathcal{L}^{ERM} = \sum_{i=1}^N \frac{\mathcal{L}(x_i, y_i)}{N}$$

### 2.2.2 REWEIGHTING

This GRM (Shimodaira, 2000) involves reweighting the loss of each group in the training data, so as to mitigate the impact of each group’s size. In particular, the loss function becomes simply the average loss of each group’s average loss.  $G$  is the number of groups in the dataset.

$$\mathcal{L}^{RW} = \frac{1}{G} \sum_{g \in \mathcal{G}} \sum_{x_i \in \mathbb{X}_g} \frac{\mathcal{L}(x_i, y_i)}{N_g}$$

### 2.2.3 GROUP DRO (GDRO)

GDRO (Sagawa\* et al., 2020) tackles worst-group error by optimizing the following function ( $\epsilon$  is a hyperparameter) during training, which we use for empirical evaluations:

$$\mathcal{L}^{GDRO} = \sum_{g \in \mathcal{G}} p_g \sum_{x_i \in \mathbb{X}_g} \frac{\mathcal{L}(x_i, y_i)}{N_g}, p_i = \frac{e^{\epsilon \cdot \mathcal{L}_i}}{\sum_{i=1}^G e^{\epsilon \cdot \mathcal{L}_i}}$$

### 2.2.4 SUBSAMPLING (SUBG-FT)

An alternate and simple baseline for inducing robustness is by using Subsampling (Idrissi et al., 2022), which consists of finetuning a model using ERM on a subsampled version of the dataset which enforces balance between groups. Usually, each group is subsampled to the size of the smallest group.

Finally, we do not include other methods such as JTT (Liu et al., 2021), DFR (Kirichenko et al., 2023), AFR (Qiu et al., 2023) or LFR (Ghaznavi et al., 2023) in our analysis because these methods in one way or another are a combination of the simpler methods studied here; they estimate weights for each sample either through proxies for group labels and/or they use Subsampling for finetuning a classifier.

**GRM implementations** Our implementation of the methods used in our experiments is mostly based off of Liu et al. (2021). For Subsampling we followed the procedure used in Kirichenko et al. (2023) but without training the 10 classifiers.

## 2.3 DATASETS

**MNIST-CIFAR (Shah et al., 2020)** MNIST-CIFAR consists of images from MNIST (LeCun & Cortes, 2010) and CIFAR-10 (Krizhevsky et al.) concatenated vertically while a spurious correlation is induced by associating CIFAR-10 labels to MNIST labels with a tunable correlation parameter. Classes 0 and 1 from MNIST are correlated with corresponding classes of CIFAR-10. We train our models on versions of this dataset with correlations of 0.0, 0.25, 0.5, 0.75 and 0.9.

**Waterbirds (Sagawa\* et al., 2020)** Waterbirds is a dataset of real images of birds where a spurious correlation is induced with the background. The task consists of identifying if the bird is a land or sea bird, while the background is land or sea. This dataset shows 0.9 correlation between class label and spurious attribute and generate versions for correlations of 0.0, 0.25, 0.5, 0.75 and 0.9.

**CelebA (Liu et al., 2015)** CelebA is a dataset of real images of celebrities carefully annotated with different attributes (gender, hair color, facial hair, attractiveness, etc.) which allows for extensive creation of spuriously correlated datasets. We use the same splits and attributes as in Sagawa\* et al. (2020), where the target attribute is Blonde Hair and the spurious attribute is gender. We only use the original training split for this dataset, which is 0.3. However, for ease of display on tables and figures, we list all results for CelebA under correlation = 0.9.

**MultiNLI (Williams et al., 2018b)** We use the MultiNLI corpus, a large-scale natural language inference benchmark spanning multiple domains. To study spurious correlations, we concentrate on the association between the gold label (*entailment/contradiction/neutral*) and the presence of negation words in the hypothesis (*sentence2\_has\_negation*).

**CivilComments (Borkan et al., 2019; Koh et al., 2020)** We use the CivilComments-WILDS dataset, a subset of the Civil Comments platform annotated for toxicity and identity references. Following prior fairness work, we focus on spurious correlations between toxicity and the presence of demographic identity mentions (*identity\_any*).

## 2.4 MODELS

For MNIST-CIFAR we use a simple convolutional network with 3 convolutional layers of 32, 64, 128 filters and a linear layer at the end, with ReLU activations. No Max-Pooling was used. For Waterbirds and CelebA we use a Resnet-50 (He et al., 2016) pretrained on ImageNet (Deng et al., 2009). For MultiNLI and CivilComments, we use the BERT architecture used by Liu et al. (2021).

## 3 ANALYSIS OF CLASSIFIER ALIGNMENT FOR ERM AND GDRO

Our theoretical analysis is carried out in the *fine-tuning setting*, where the feature extractor is frozen and only the linear classifier is trained. Under this assumption, the only way GDRO can differ from ERM is through how the final classifier aligns with different predictive directions: GDRO should align less with the spurious direction and more with an oracle non-spurious classifier, whereas ERM shows the opposite tendency. We formalize this intuition by tracking the dot product between the learned classifier and the spurious/non-spurious directions throughout training. Furthermore, we prove that when the loss is  $\mu$ -strongly convex, the alignment gap directly bounds the worst-group performance difference between ERM and GDRO. Since cross-entropy in this setting is convex but not  $\mu$ -strongly convex, we also show that adding L2 regularization induces  $\mu$ -strong convexity, extending our guarantees and yielding explicit bounds. We use the Min-Max formulation of GDRO for this analysis.

In this section  $\theta^*$  and  $\theta_{ERM}$  refer to the optimal worst group classifier and the optimal ERM classifier respectively. Our first result talks about how GDRO tends to align less with spurious classifiers and more with non-spurious ones. The full derivation of all results are in the Appendix Section B.

### 3.1 ALIGNMENT TO SPURIOUS AND NON-SPURIOUS DIRECTIONS FOR ERM AND GDRO

Our first set of results show that for even loss functions:  $\alpha_{sp}(\theta^*) \leq \alpha_{sp}(\theta_{ERM})$  and  $\alpha_{ns}(\theta^*) \geq \alpha_{ns}(\theta_{ERM})$ .

**Proposition (B.2 (summary)).** Let  $u = v_{sp}/\|v_{sp}\|$  and write  $\theta = \theta_{\perp} + t u$ ,  $t = \langle u, \theta \rangle$ . Assume  $L_g(\theta) = R(\theta_{\perp}) + \phi(t - a_g)$  with the same  $R$  for all  $g$ , where  $\phi \in C^1(\mathbb{R})$  is even, strictly convex,  $\phi(0) = 0$ ,  $\phi'$  odd, and  $\phi$  strictly increasing on  $[0, \infty)$ . Let  $a_- := \min_g a_g < 0 < \max_g a_g =: a_+$  with  $a_+ = -a_- > 0$ . Define  $\theta_{\text{ERM}} \in \arg \min_{\theta} \sum_g p_g L_g(\theta)$  and  $\theta^* \in \arg \min_{\theta} \max_g L_g(\theta)$  and set  $\alpha_{sp}(\theta) := |\langle u, \theta \rangle|$ . If  $P_{\pm}(a) := \sum_{g: a_g = \pm a} p_g$ :

$$\sum_{a>0} (P_-(a) - P_+(a)) \phi'(a) \neq 0 \implies \alpha_{sp}(\theta^*) < \alpha_{sp}(\theta_{\text{ERM}})$$

**Proposition (B.3 (summary)).** Let  $u$  be a non-spurious direction. Write  $\theta = \theta_{\perp} + t u$  with  $t = \langle u, \theta \rangle$  and  $\theta_{\perp} \in \{u\}^{\perp}$ . Assume  $L_g(\theta) = R(\theta_{\perp}) + \phi(t - a_g)$ , where  $R$  is the same convex function for all groups, and  $\phi \in C^1(\mathbb{R})$  is even, strictly convex, satisfies  $\phi(0) = 0$ , and has odd, strictly increasing derivative  $\phi'$ . Suppose that

$$0 < a_{\min} := \min_g a_g \leq a_g \leq a_{\max} := \max_g a_g,$$

and let  $p_g > 0$  with  $\sum_g p_g = 1$ . Define the ERM and GroupDRO optimizers  $t_{\text{ERM}} \in \arg \min_{t \in \mathbb{R}} \sum_g p_g \phi(t - a_g)$  and  $t^* \in \arg \min_{t \in \mathbb{R}} \max_g \phi(t - a_g)$ , and set  $\alpha_{\text{ns}}(\theta) := |\langle u, \theta \rangle| = |t|$ . Then:

1. Both ERM and GroupDRO choose the same  $\theta_{\perp}^* \in \arg \min_{\theta_{\perp}} R(\theta_{\perp})$ , so the comparison reduces to the one-dimensional problems in  $t$  above.
2. Let  $m := \frac{1}{2} (a_{\min} + a_{\max})$ . Then  $t^* = m$  and  $t_{\text{ERM}}$  is the unique root of

$$f'(t) := \sum_g p_g \phi'(t - a_g) = 0,$$

with  $t_{\text{ERM}} \in [a_{\min}, a_{\max}]$  and  $t_{\text{ERM}} > 0$ .

3. We have

$$\alpha_{\text{ns}}(\theta^*) \begin{cases} > \alpha_{\text{ns}}(\theta_{\text{ERM}}) & \text{if } f'(m) > 0, \\ = \alpha_{\text{ns}}(\theta_{\text{ERM}}) & \text{if } f'(m) = 0, \\ < \alpha_{\text{ns}}(\theta_{\text{ERM}}) & \text{if } f'(m) < 0. \end{cases}$$

In particular, when  $f'(m) > 0$ , GroupDRO pushes further towards  $|t^*| > |t_{\text{ERM}}|$ .

In the Appendix, Propositions B.4 and B.5 show something stronger: that the relation between alignments holds also along the optimization trajectory. Our next result shows how the difference in alignment relates to the difference in worst-group performance.

### 3.2 RELATION BETWEEN ALIGNMENT, PERFORMANCE AND REGULARIZATION FOR GDRO

We assume losses are  $\mu$ -strongly convex, which enables relating classifier alignment to the achieved losses.

**Proposition (B.6 (summary)).** Assume each  $L_g$  is  $\mu$ -strongly convex. Then  $L_g(\theta) \geq L_g(\theta_{\text{ERM}}) + \frac{\mu}{2} \|\theta - \theta_{\text{ERM}}\|^2$ , and consequently  $\alpha_{sp}(\theta_{\text{ERM}}) - \alpha_{sp}(\theta^*) \geq \sqrt{\frac{2(L_{\max}(\theta_{\text{ERM}}) - L_{\max}(\theta^*))}{\mu}}$ .

This implies that for GDRO to attain lower loss than ERM, its alignment with the spurious direction must be smaller. Our previous results confirm this behavior, highlighting alignment as a key mechanism behind GDRO's performance gains. A limitation is that Cross-Entropy is not  $\mu$ -strongly convex. Empirically, GDRO requires L2 regularization (Sagawa\* et al., 2020), but a theoretical justification was missing. We show that adding an L2 penalty of weight  $\frac{\lambda}{2}$  yields  $\lambda$ -strong convexity, allowing us to connect  $\lambda$ , alignments and losses via a lower bound on the alignment gap.

**Proposition (B.7 (summary)).** Let each group risk  $L_g: \Theta \rightarrow \mathbb{R}$  be convex, and fix a spurious unit vector  $u = v_{sp}/\|v_{sp}\|$ . For a regularization parameter  $\lambda > 0$ , define  $R_{\lambda}(\theta) = \frac{\lambda}{2} \|\theta\|^2$ ,  $\tilde{L}_g(\theta) = L_g(\theta) + R_{\lambda}(\theta)$ , and write  $\theta_{\text{ERM}}^{\lambda} = \arg \min_{\theta} \sum_{g=1}^G p_g \tilde{L}_g(\theta)$ ,  $\theta^{\lambda,*} = \arg \min_{\theta} \max_g \tilde{L}_g(\theta)$ , with

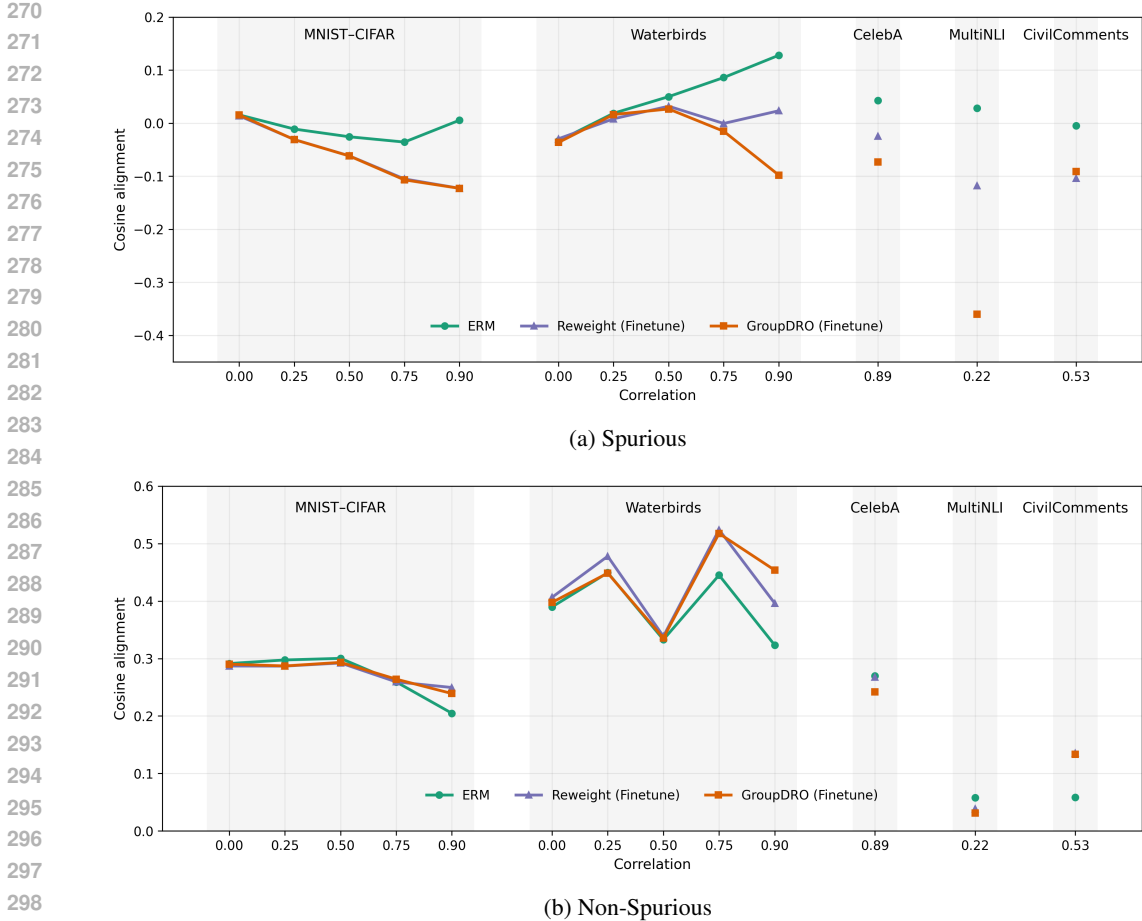


Figure 1: Cosine similarity between classifiers learned by models fine-tuned with different methods on ERM features and two reference classifiers: one predicting the spurious attribute and one oracle non-spurious classifier, across multiple datasets. (a) Consistent with our theory, GDRO exhibits lower similarity to the spurious classifier, in some cases even aligning in the opposite direction. (b) For the non-spurious classifier, GDRO is generally more aligned than ERM, except in CelebA and MultiNLI. Overall, the dominant effect is GDRO’s reduced alignment with the spurious classifier.

$\alpha_{sp}(\theta) = |u^T \theta|$ ,  $\tilde{L}_{\max}(\theta) = \max_g \tilde{L}_g(\theta)$ . Then:

$$\frac{\lambda(\alpha_{sp}(\theta_{ERM}^\lambda) - \alpha_{sp}(\theta^{\lambda,*}))^2}{2} \geq \tilde{L}_{\max}(\theta_{ERM}^\lambda) - \tilde{L}_{\max}(\theta^{\lambda,*})$$

We further show in Proposition B.8 in the Appendix that the result extends to quadratic regularizers of the form  $R_M(\theta) = \frac{1}{2} \theta^T M \theta$  with  $M$  positive definite matrix. In the next section, we show whether this alignment effect actually occurs in practice and how important it is.

#### 4 EMPIRICAL EVALUATION

We first validate our theoretical alignment results empirically across multiple datasets and three methods: ERM, RW, and GDRO. To approximate the spurious classifier, we train a linear model on pretrained features to predict the spurious attribute. For the non-spurious classifier, we approximate an oracle by training on both train and test distributions to predict the class label. Results are shown in Figure 1. For the spurious classifier, both RW and GDRO consistently display lower alignment than ERM, with GDRO almost always showing the least alignment (the only exception being CivilComments, where RW is slightly lower). For the non-spurious classifier, GDRO and RW tend to be more aligned, except in CelebA and MultiNLI. We attribute these exceptions to the fact that the oracle non-spurious classifier is only well-defined in linearly separable settings, which is

Table 1: Worst Group Accuracy (%) across five benchmarks at varying levels of spurious correlation. We compare (i) end-to-end ERM, RW, and GDRO, (ii) finetuned models over ERM features (RW-FT, GDRO-FT, SUBG), and (iii) ablations combining SUBG-FT with RW or GDRO losses (RW+, GDRO+). GDRO-FT generally matches or exceeds end-to-end GDRO, RW-FT is less consistent, and SUBG-FT reliably improves over ERM. GDRO+ provides consistent gains over SUBG across all datasets. The classifier effect thus explains much of GDRO’s success, but in MNIST-CIFAR and Waterbirds it falls short, suggesting that representation-level factors also play a crucial role. MC=MNIST-CIFAR, WB=Waterbirds, CA=CelebA, MNLI=MultiNLI, CC=CivilComments.

DS	CORR	END-TO-END			FINETUNED			SUBG+DRO	
		ERM	RW	GDRO	RW-FT	GDRO-FT	SUBG-FT	RW+	GDRO+
MC	0.25	82.11% (0.4)	86.72% (2.2)	86.59% (2.0)	86.86% (1.3)	86.72% (1.2)	87.95% (0.7)	<b>88.22%</b> (0.5)	<b>88.22%</b> (0.6)
	0.5	78.08% (0.7)	88.50% (0.7)	88.34% (0.6)	<b>88.85%</b> (1.0)	88.72% (0.8)	87.15% (1.2)	87.42% (0.9)	87.01% (2.1)
	0.75	54.86% (5.0)	84.26% (1.1)	83.52% (0.9)	84.38% (0.2)	84.13% (0.6)	83.00% (3.0)	84.34% (0.4)	<b>85.54%</b> (0.0)
	0.9	29.91% (3.2)	<b>84.72%</b> (1.1)	82.96% (0.7)	81.70% (1.0)	77.95% (2.0)	79.83% (1.6)	78.71% (1.8)	80.59% (0.6)
WB	0.25	89.56% (0.7)	89.87% (1.5)	89.50% (3.1)	86.45% (0.4)	86.45% (0.8)	87.33% (0.7)	88.32% (1.8)	<b>90.58%</b> (0.6)
	0.5	87.44% (1.0)	88.84% (1.0)	<b>90.13%</b> (0.5)	82.71% (0.2)	84.48% (0.4)	87.80% (0.5)	85.15% (2.7)	<b>90.13%</b> (0.7)
	0.75	80.74% (0.8)	88.06% (1.6)	87.59% (0.4)	74.35% (1.1)	78.50% (1.2)	84.58% (1.6)	81.25% (2.2)	<b>88.32%</b> (0.6)
	0.9	71.96% (0.8)	<b>86.31%</b> (0.3)	85.87% (0.9)	65.21% (1.9)	73.68% (3.1)	80.54% (1.7)	70.63% (1.2)	82.83% (2.0)
CA	0.89	46.48% (1.2)	86.48% (2.3)	88.52% (0.3)	86.48% (0.9)	<b>90.37%</b> (0.3)	85.00% (1.5)	85.56% (0.6)	89.71% (0.3)
MN	0.22	69.45% (2.2)	67.42% (3.6)	76.13% (1.2)	71.04% (0.3)	<b>76.83%</b> (1.4)	68.73% (1.3)	68.99% (1.0)	70.24% (2.3)
CC	0.53	66.14% (2.4)	80.00% (2.4)	79.89% (3.3)	81.51% (2.2)	<b>81.87%</b> (1.0)	77.46% (6.0)	76.92% (5.8)	77.86% (5.1)

not the case here, making the approximation imperfect. Overall, the results suggest that the primary driver of performance differences is the reduced alignment with the spurious classifier.

#### 4.1 EMPIRICAL EFFECTS OF FINETUNING FINAL CLASSIFIER LAYER ON GRMS

A natural next step is to assess how substantial the classifier effect really is. If robustness gains can be largely attributed to the final linear classifier, as suggested by recent work (Kirichenko et al., 2023), then group-based feature learning may not be essential. To test this, we compare three groups of models: (i) end-to-end models trained with ERM, RW, and GDRO; (ii) models finetuned over ERM features (RW-FT, GDRO-FT, and Subsampling—SUBG); and (iii) ablations that combine SUBG with RW or GDRO losses (RW+, GDRO+). Results are summarized in Table 1.

Overall, finetuning using GDRO consistently yields stronger classifiers: in four out of five datasets, it markedly outperforms ERM and in three cases it actually outperforms end-to-end GDRO. RW exhibits similar tendencies but with much higher variance, sometimes surpassing ERM and sometimes falling well below it. SUBG also reliably improves over ERM, and in MNIST-CIFAR and Waterbirds it even outperforms GDRO-FT. In these same two datasets, however, end-to-end GDRO remains superior to GDRO-FT and SUBG. Remarkably, augmenting SUBG with the GDRO loss (GDRO+) consistently improves performance over SUBG alone across all datasets. We suggest using GDRO+ as a better baseline for robustness methods as it is basically a free improvement over SUBG. These findings show that the classifier effect of GDRO is indeed powerful, but not sufficient to explain all robustness gains. In particular, MNIST-CIFAR and Waterbirds highlight that end-to-end GDRO provides additional benefits, pointing to representation-level effects that we study next.

#### 4.2 PROPERTIES OF GDRO FEATURES

To understand why features learned with GDRO may confer greater robustness, we consider two hypotheses: (1) GDRO does not learn spurious information during training, and (2) GDRO produces more disentangled features, making it easier for the classifier to rely on predictive rather than spurious attributes. The first hypothesis, analyzed in the Appendix in Section C, is false: GDRO learns spurious features just as ERM does. We therefore focus on the second hypothesis through a quantitative study of disentanglement. For this, we adopt the DCI framework (Eastwood & Williams, 2018), which evaluates representations along three dimensions: *Disentanglement*, measuring how specific each latent dimension is to a single attribute (0 = fully mixed, 1 = perfectly specific); *Com-*

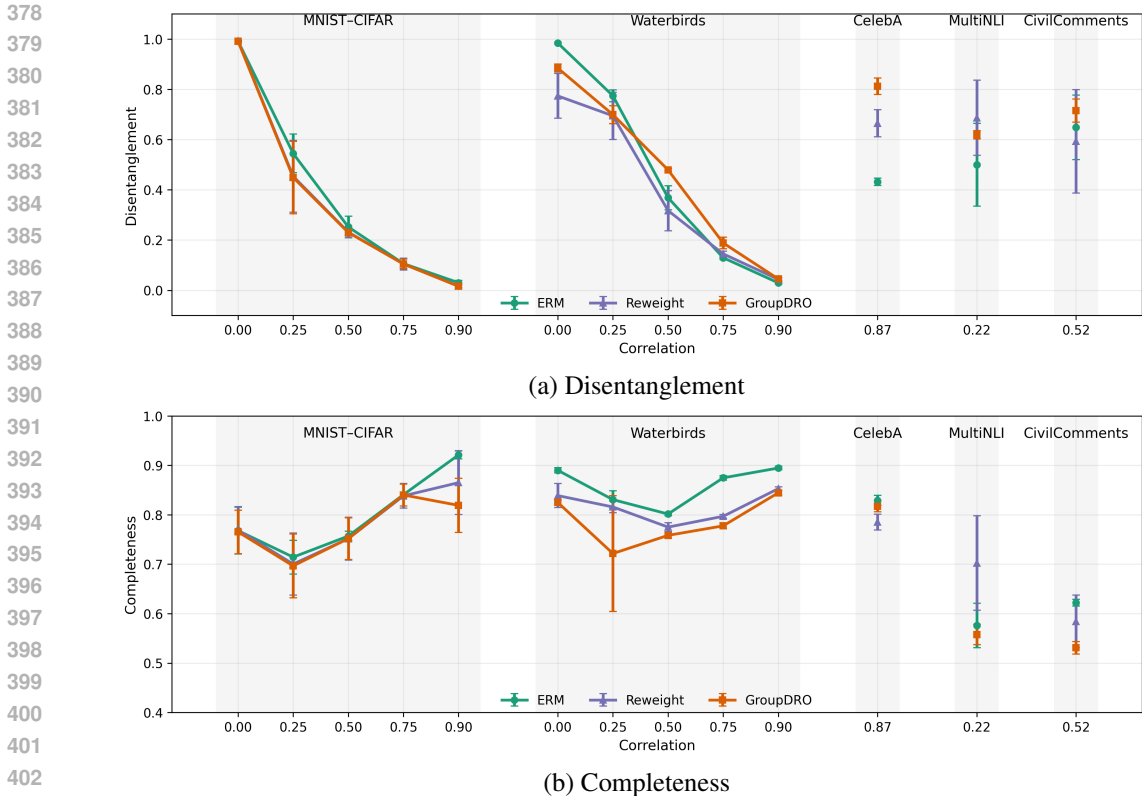


Figure 2: Disentanglement and Completeness of features learned with ERM, Reweighting, and GDRO across several datasets. (a) GDRO frequently increases *disentanglement*, though the effect is not uniform across datasets. (b) In contrast, *completeness* is consistently lower for GDRO, indicating that its robustness partly stems from distributing predictive information across multiple dimensions, thereby reducing reliance on spurious attributes.

*pletteness*, measuring how concentrated each attribute is across dimensions (0 = spread evenly, 1 = concentrated in one dimension); and *Informativeness*, quantifying overall predictive power, typically assessed with linear probes. In theory, a desirable representation should score high on all three. Our results, however, challenge this view in the presence of spurious correlations.

For efficiency, we reduce representations to  $K = 50$  dimensions via PCA, retaining most of the variance. Figure 2 reports results across all datasets. For full numerical results, including informativeness, see Section D.3 in the Appendix. As expected, disentanglement decreases as spurious correlation increases. Unexpectedly, models trained with GDRO show higher disentanglement in more challenging datasets (e.g., CelebA, MultiNLI, CivilComments, high-correlation Waterbirds), but not in simpler cases such as MNIST-CIFAR or low-correlation Waterbirds. Completeness shows a clearer pattern: GRMs consistently achieve lower completeness than ERM across datasets. This indicates that GRM representations distribute predictive information across more dimensions, reducing reliance on any single (potentially spurious) dimension. To the best of our knowledge, this is the first time an empirical link between completeness and robustness has been established. Moreover, new robustness methods may choose to target this property to achieve their goals.

## 5 RELATED WORK

### 5.1 METHODS FOR ROBUST LEARNING

Invariant Risk Minimization (Arjovsky et al., 2020) modifies the loss function to make the model invariant to different environments. Group Reweighting schemes manipulate the importance of samples during training: Reweighting (Shimodaira, 2000) reweights group losses to eliminate the impact of group size on the loss, Group Distributionally Robust Optimization (GDRO) (Sagawa\* et al., 2020) reweights group losses based on their magnitude, while other methods (Seo et al., 2022; So-

honi et al., 2020) try to apply GDRO by finding proxies for the group labels needed. Some have used biased models to reweight losses to train a debiased model (Nam et al., 2020) or used biased models’ representations to train an invariant classifier (Wald et al., 2023). Other methods are based on multiple training passes: some finetune an ERM trained model on balanced data (Kirichenko et al., 2023; Qiu et al., 2023; Ghaznavi et al., 2023) or use that same model to find samples to up-weight to train a new model from scratch (Liu et al., 2021) or finetune a classifier. Others start from a balanced dataset and progressively expand the dataset during training (Deng et al., 2023). Others have worked on creating non-linear classifiers that are orthogonal to a set of attributes (Xu et al., 2022), this requires having a notion of both the train and test distribution of those attributes. Other methods have used contrastive losses (Zhang et al., 2022) to create representations that are invariant to spurious attributes. Finally, some have used Self Supervised pretrained models to estimate a logit adjustment term on the loss (Tsirigotis et al., 2023) and others have proposed adversarial training (Setlur et al., 2023) focusing on features rather than groups.

## 5.2 STUDIES ON ROBUSTNESS METHODS

Several studies highlight why robust methods are needed: networks are highly prone to simplicity bias (Shah et al., 2020), and spurious features are often simpler than core ones (Vasudeva et al., 2024). The most relevant works for us analyze representations from ERM models (Kirichenko et al., 2023; Izmailov et al., 2022), arguing that GRMs mainly affect the final classifier layer and that ERM representations suffice. Our results challenge this view: while the classifier effect is central, GRMs also shape representations in ways ERM does not. Unlike prior empirical-only studies, we provide both theoretical and empirical evidence. Other analyses claim generalized reweighting offers no advantage over ERM (Zhai et al., 2023), though under restrictive assumptions. In contrast, our analysis only assumes a linear classifier over frozen features. Recent theory has studied spurious feature memorization in random/NTK settings (Bombari & Mondelli, 2024), introducing alignment measures across sample pairs. Our alignment notion differs: we quantify alignment between the learned classifier and reference spurious/non-spurious classifiers. Additional works compare methods such as GDRO in systematic generalization settings (Ahmed et al., 2021), but without addressing underlying mechanisms. Others link mutual information between spurious attributes and labels to worst-group error (Zhang et al., 2022), also invoking alignment, though defined in representation space between samples while we focus on classifier-level alignment.

## 6 CONCLUSIONS

We studied why Group Robustness Methods, and GDRO in particular, succeed where ERM fails under spurious correlations. Our theoretical analysis in the fine-tuning setting showed that GDRO produces classifiers less aligned with spurious directions and more aligned with oracle non-spurious ones, and that this alignment gap explains improvements in worst-group performance when the loss is  $\mu$ -strongly convex. Moreover, we proved that L2 regularization induces  $\mu$ -strong convexity in cross-entropy, providing a principled explanation for the necessity of strong regularization in GDRO. Empirically, we confirmed these predictions across vision and text benchmarks: GDRO reduces alignment with spurious classifiers and increases alignment with non-spurious ones compared to ERM. Going beyond the classifier, we showed that GDRO also reshapes the learned representations. Contrary to the intuition that robustness comes from discarding spurious features, we found that GDRO representations still encode them but distribute predictive information across more dimensions (lower completeness), making classifiers less dependent on individual spurious attributes. Taken together, our findings clarify the mechanisms by which GRMs achieve robustness: they act both at the classifier level and at the representation level, with effects that complement each other. We believe these insights can inspire the design of new methods that combine alignment and lower completeness, ultimately leading to more principled and practical approaches for robust learning under spurious correlations.

## 7 REPRODUCIBILITY STATEMENT

We share our code for our experiments, which is modified from Liu et al. (2021), in the supplementary materials. This includes the Jupyter Notebooks used to obtain plots and tables as well as sample scripts to reproduce results. Proof for results on the theoretical section are found in the Appendix, Section B. MNIST-CIFAR, CelebA, Waterbirds, MultiNLI and CivilComments are freely available datasets. Hyperparameters used are detailed in Section F.1.

## REFERENCES

- 486  
487  
488 Faruk Ahmed, Yoshua Bengio, Harm van Seijen, and Aaron Courville. Systematic generalisation  
489 with group invariant predictions. In *International Conference on Learning Representations*, 2021.  
490 URL <https://openreview.net/forum?id=b9P0imzZFJ>.
- 491  
492 Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization.  
493 *stat*, 1050:27, 2020.
- 494  
495 Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new  
496 perspectives, 2014. URL <https://arxiv.org/abs/1206.5538>.
- 497  
498 Simone Bombari and Marco Mondelli. How spurious features are memorized: Precise analysis  
499 for random and NTK features. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian  
500 Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the 41st  
501 International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning  
502 Research*, pp. 4267–4299. PMLR, 21–27 Jul 2024. URL [https://proceedings.mlr.  
503 press/v235/bombari24a.html](https://proceedings.mlr.press/v235/bombari24a.html).
- 504  
505 Daniel Borkan, Lucas Dixon, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. Nuanced  
506 metrics for measuring unintended bias with real data for text classification. In *Compan-  
507 ion Proceedings of The 2019 World Wide Web Conference*, WWW ’19, pp. 491–500, New  
508 York, NY, USA, 2019. Association for Computing Machinery. ISBN 9781450366755. doi:  
509 10.1145/3308560.3317593. URL <https://doi.org/10.1145/3308560.3317593>.
- 510  
511 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hi-  
512 erarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,  
513 pp. 248–255. Ieee, 2009.
- 514  
515 Yihe Deng, Yu Yang, Baharan Mirzasoleiman, and Quanquan Gu. Robust learning with progressive  
516 data expansion against spurious correlation. In *Thirty-seventh Conference on Neural Information  
517 Processing Systems*, 2023. URL <https://openreview.net/forum?id=9QEVJ9qm46>.
- 518  
519 Cian Eastwood and Christopher K. I. Williams. A framework for the quantitative evaluation of  
520 disentangled representations. In *International Conference on Learning Representations*, 2018.  
521 URL <https://openreview.net/forum?id=By-7dz-AZ>.
- 522  
523 R. Geirhos, J.-H. Jacobsen, C. Michaelis, R. Zemel, W. Brendel, M. Bethge, and F. A. Wichmann.  
524 Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673, 2020.  
525 doi: 10.1038/s42256-020-00257-z.
- 526  
527 Mahdi Ghaznavi, Hesam Asadollahzadeh, HamidReza Yaghoubi Araghi, Fahimeh Hosseini  
528 Noohdani, Mohammad Hossein Rohban, and Mahdieh Soleymani Baghshah. Annotation-free  
529 group robustness via loss-based resampling, 2023. URL [https://arxiv.org/abs/2312.  
530 04893](https://arxiv.org/abs/2312.04893).
- 531  
532 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Im-  
533 age Recognition. In *Proceedings of 2016 IEEE Conference on Computer Vision and Pattern  
534 Recognition*, CVPR ’16, pp. 770–778. IEEE, June 2016. doi: 10.1109/CVPR.2016.90. URL  
535 <http://ieeexplore.ieee.org/document/7780459>.
- 536  
537 Weihua Hu, Gang Niu, Issei Sato, and Masashi Sugiyama. Does distributionally robust supervised  
538 learning give robust classifiers? In Jennifer G. Dy and Andreas Krause 0001 (eds.), *Proceed-  
539 ings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan,  
540 Stockholm, Sweden, July 10-15, 2018*, volume 80 of *JMLR Workshop and Conference Proceed-  
541 ings*, pp. 2034–2042. JMLR.org, 2018. URL [http://proceedings.mlr.press/v80/  
542 hul8a.html](http://proceedings.mlr.press/v80/hul8a.html).
- 543  
544 Badr Youbi Idrissi, Martin Arjovsky, Mohammad Pezeshki, and David Lopez-Paz. Simple data  
545 balancing achieves competitive worst-group-accuracy, 2022.

- 540 Pavel Izmailov, Polina Kirichenko, Nate Gruver, and Andrew Gordon Wilson. On feature learning  
541 in the presence of spurious correlations. In *Proceedings of the 36th International Conference on*  
542 *Neural Information Processing Systems, NIPS '22*, Red Hook, NY, USA, 2022. Curran Associates  
543 Inc. ISBN 9781713871088.
- 544 Polina Kirichenko, Pavel Izmailov, and Andrew Gordon Wilson. Last layer re-training is sufficient  
545 for robustness to spurious correlations. In *The Eleventh International Conference on Learning*  
546 *Representations*, 2023. URL <https://openreview.net/forum?id=Zb6c8A-Fghk>.
- 548 Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Bal-  
549 subramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Sara Beery, Jure Leskovec,  
550 Anshul Kundaje, Emma Pierson, Sergey Levine, Chelsea Finn, and Percy Liang. WILDS: A  
551 benchmark of in-the-wild distribution shifts. *CoRR*, abs/2012.07421, 2020. URL <https://arxiv.org/abs/2012.07421>.
- 553 Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-10 (canadian institute for advanced re-  
554 search). URL <http://www.cs.toronto.edu/~kriz/cifar.html>.
- 556 Alexandre Lacoste, Alexandra Luccioni, Victor Schmidt, and Thomas Dandres. Quantifying the  
557 carbon emissions of machine learning. *arXiv preprint arXiv:1910.09700*, 2019.
- 559 Yann LeCun and Corinna Cortes. MNIST handwritten digit database. 2010. URL <http://yann.>  
560 [lecun.com/exdb/mnist/](http://yann.lecun.com/exdb/mnist/).
- 561 Evan Z Liu, Behzad Haghgoo, Annie S Chen, Aditi Raghunathan, Pang Wei Koh, Shiori Sagawa,  
562 Percy Liang, and Chelsea Finn. Just train twice: Improving group robustness without training  
563 group information. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International*  
564 *Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp.  
565 6781–6792. PMLR, 18–24 Jul 2021. URL [https://proceedings.mlr.press/v139/](https://proceedings.mlr.press/v139/liu21f.html)  
566 [liu21f.html](https://proceedings.mlr.press/v139/liu21f.html).
- 568 Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild.  
569 In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.
- 570 Junhyun Nam, Hyuntak Cha, Sungsoo Ahn, Jaeho Lee, and Jinwoo Shin. Learn-  
571 ing from failure: De-biasing classifier from biased classifier. In H. Larochelle,  
572 M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural In-*  
573 *formation Processing Systems*, volume 33, pp. 20673–20684. Curran Associates, Inc.,  
574 2020. URL [https://proceedings.neurips.cc/paper\\_files/paper/2020/](https://proceedings.neurips.cc/paper_files/paper/2020/file/eddc3427c5d77843c2253f1e799fe933-Paper.pdf)  
575 [file/eddc3427c5d77843c2253f1e799fe933-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/eddc3427c5d77843c2253f1e799fe933-Paper.pdf).
- 577 John Pavlopoulos, Jeffrey Sorensen, Lucas Dixon, Nithum Thain, and Ion Androutsopoulos. Toxic-  
578 ity detection: Does context really matter?, 2020.
- 579 Shikai Qiu, Andres Potapczynski, Pavel Izmailov, and Andrew Gordon Wilson. Simple and Fast  
580 Group Robustness by Automatic Feature Reweighting. *International Conference on Machine*  
581 *Learning (ICML)*, 2023.
- 583 Shiori Sagawa\*, Pang Wei Koh\*, Tatsunori B. Hashimoto, and Percy Liang. Distributionally ro-  
584 bust neural networks for group shifts: On the importance of regularization for worst-case gen-  
585 eralization. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=ryxGuJrFvS>.
- 587 Seonguk Seo, Joon-Young Lee, and Bohyung Han. Unsupervised learning of debiased representa-  
588 tions with pseudo-attributes. In *Proceedings of the IEEE/CVF Conference on Computer Vision*  
589 *and Pattern Recognition (CVPR)*, pp. 16742–16751, June 2022.
- 591 Amrith Setlur, Don Dennis, Benjamin Eysenbach, Aditi Raghunathan, Chelsea Finn, Virginia Smith,  
592 and Sergey Levine. Bitrate-constrained DRO: Beyond worst case robustness to unknown group  
593 shifts. In *The Eleventh International Conference on Learning Representations*, 2023. URL  
<https://openreview.net/forum?id=2QzNuaRHn4Z>.

- 594 Harshay Shah, Kaustav Tamuly, Aditi Raghunathan, Prateek Jain, and Praneeth Netrapalli. The  
595 Pitfalls of Simplicity Bias in Neural Networks. (NeurIPS):1–32, 2020. URL <http://arxiv.org/abs/2006.07710>.  
596  
597
- 598 Hidetoshi Shimodaira. Improving predictive inference under covariate shift by weighting  
599 the log-likelihood function. *Journal of Statistical Planning and Inference*, 90(2):227–  
600 244, October 2000. URL [http://www.sciencedirect.com/science/article/  
601 B6V0M-4136355-5/1/6432c256e0be03b1503bbf79e4e91d1a](http://www.sciencedirect.com/science/article/B6V0M-4136355-5/1/6432c256e0be03b1503bbf79e4e91d1a).
- 602 Nimit Sohoni, Jared A. Dunnmon, Geoffrey Angus, Albert Gu, and Christopher Ré. No subclass left  
603 behind: Fine-grained robustness in coarse-grained classification problems. In *Proceedings of the  
604 34th International Conference on Neural Information Processing Systems, NIPS’20*, Red Hook,  
605 NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546.
- 606 Christos Tsirigotis, Joao Monteiro, Pau Rodriguez, David Vazquez, and Aaron Courville. Group  
607 robust classification without any group information. In *Thirty-seventh Conference on Neural  
608 Information Processing Systems, 2023*. URL [https://openreview.net/forum?id=  
609 20cNWFHFpk](https://openreview.net/forum?id=20cNWFHFpk).
- 610  
611 V. Vapnik. Principles of risk minimization for learning theory. In J. Moody, S. Hanson, and  
612 R.P. Lippmann (eds.), *Advances in Neural Information Processing Systems*, volume 4. Morgan-  
613 Kaufmann, 1991. URL [https://proceedings.neurips.cc/paper\\_files/paper/  
614 1991/file/ff4d5fbbafdf976cfdc032e3bde78de5-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/1991/file/ff4d5fbbafdf976cfdc032e3bde78de5-Paper.pdf).
- 615 Bhavya Vasudeva, Kameron Shahabi, and Vatsal Sharan. Mitigating simplicity bias in deep learning  
616 for improved OOD generalization and robustness. *Transactions on Machine Learning Research*,  
617 2024. ISSN 2835-8856. URL <https://openreview.net/forum?id=XccFHGakyU>.
- 618  
619 Yoav Wald, Gal Yona, Uri Shalit, and Yair Carmon. Malign overfitting: Interpolation can provably  
620 preclude invariance. In *International Conference on Learning Representations (ICLR), 2023*.  
621 URL <https://arxiv.org/abs/2211.15724>.
- 622 Adina Williams, Nikita Nangia, and Samuel Bowman. A broad-coverage challenge corpus for sen-  
623 tence understanding through inference. In *Proceedings of the 2018 Conference of the North Amer-  
624 ican Chapter of the Association for Computational Linguistics: Human Language Technologies,  
625 Volume 1 (Long Papers)*, pp. 1112–1122. Association for Computational Linguistics, 2018a. URL  
626 <http://aclweb.org/anthology/N18-1101>.
- 627  
628 Adina Williams, Nikita Nangia, and Samuel Bowman. A broad-coverage challenge corpus for sen-  
629 tence understanding through inference. In Marilyn Walker, Heng Ji, and Amanda Stent (eds.),  
630 *Proceedings of the 2018 Conference of the North American Chapter of the Association for Com-  
631 putational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pp. 1112–  
632 1122, New Orleans, Louisiana, June 2018b. Association for Computational Linguistics. doi:  
10.18653/v1/N18-1101. URL <https://aclanthology.org/N18-1101/>.
- 633  
634 Yilun Xu, Hao He, Tianxiao Shen, and Tommi S. Jaakkola. Controlling directions orthogonal to  
635 a classifier. In *International Conference on Learning Representations, 2022*. URL <https://openreview.net/forum?id=DIjCrIsu6Z>.  
636
- 637 Runtian Zhai, Chen Dan, J Zico Kolter, and Pradeep Kumar Ravikumar. Understanding why  
638 generalized reweighting does not improve over ERM. In *The Eleventh International Confer-  
639 ence on Learning Representations, 2023*. URL [https://openreview.net/forum?id=  
640 ashPce\\_W8F-](https://openreview.net/forum?id=ashPce_W8F-).
- 641 Michael Zhang, Nimit S Sohoni, Hongyang R Zhang, Chelsea Finn, and Christopher Re. Correct-  
642 n-contrast: a contrastive approach for improving robustness to spurious correlations. In Kamalika  
643 Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato  
644 (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of  
645 *Proceedings of Machine Learning Research*, pp. 26484–26516. PMLR, 17–23 Jul 2022. URL  
646 <https://proceedings.mlr.press/v162/zhang22z.html>.  
647

## A DATASET EXAMPLES



Figure 3: Sample images from datasets used in this work. MNIST-CIFAR correlates MNIST digits with CIFAR-10 classes; Waterbirds correlates land/water birds with a land/water background; CelebA consists of heavily annotated images from celebrities. In this work, we use the correlation that appears between gender and hair color.

## B DERIVATION OF THEORETICAL RESULTS

### B.1 PRELIMINARIES

**Theorem B.1.** Let  $\mathcal{G} = \{1, \dots, G\}$  be a finite set of groups. Consider  $\mathcal{D}_g$  the set of the training points in  $g$  with  $|\mathcal{D}_g| = n_g > 0$ . Also, denote  $p_g = n_g/n$  and  $n = \sum_g n_g$ . Let  $\ell: \Theta \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  be a proper, convex and lower bounded loss function in  $\theta \in \Theta \subset \mathbb{R}^d$ . Where  $\Theta$  is a closed convex set. Define the group risk as

$$L_g(\theta) = \frac{1}{n_g} \sum_{(x,y) \in \mathcal{D}_g} \ell(\theta; x, y)$$

and the simplex  $\Delta_G = \left\{ w \in \mathbb{R}_{\geq 0}^G : \sum_{g=1}^G w_g = 1 \right\}$ . Consider the following:

1.  $F(w) = \inf_{\theta \in \Theta} \sum_{g=1}^G w_g L_g(\theta)$
2.  $\theta_{ERM} = \operatorname{argmin}_{\theta} \sum_g p_g L_g(\theta)$
3.  $\theta^* = \operatorname{argmin}_{\theta} \max_g L_g(\theta)$ .

Then

1. the function  $F: \Delta_G \rightarrow \mathbb{R} \cup \{-\infty\}$  is concave,
2.  $\min_{\theta \in \Theta} \max_g L_g(\theta) = \max_{w \in \Delta_G} F(w)$
3. if  $L_{\max}(\theta) = \max_g L_g(\theta)$ , then  $L_{\max}(\theta^*) \leq L_{\max}(\theta_{ERM})$ .

*Proof.* Let's start with a proposition.

**Proposition B.1.** The simplex  $\Delta_G$  with the standard Euclidean topology is a convex, compact subset of  $\mathbb{R}^G$ .

*Proof.* Consider  $u, v \in \Delta_G$  and  $\lambda \in [0, 1]$ . Define

$$w = \lambda u + (1 - \lambda)v = (\lambda u_1 + (1 - \lambda)v_1, \dots, \lambda u_G + (1 - \lambda)v_G).$$

702 Then for each  $g$ , we have  
703

$$\begin{aligned} 704 w_g &= \lambda u_g + (1 - \lambda)v_g \\ 705 &\geq \lambda \cdot 0 + (1 - \lambda) \cdot 0 \\ 706 &= 0, \end{aligned}$$

$$\begin{aligned} 707 \\ 708 \sum_{g=1}^G w_g &= \sum_g [\lambda u_g + (1 - \lambda)v_g] \\ 709 &= \lambda \sum_g u_g + (1 - \lambda) \sum_g v_g \\ 710 &= \lambda \cdot 1 + (1 - \lambda) \cdot 1 \\ 711 &= 1. \end{aligned}$$

712 Hence  $w \in \Delta_G$ , and thus  $\Delta_G$  is convex.  
713

714 Now let

$$715 \Delta_G = \bigcap_{g=1}^G \{w \in \mathbb{R}^G : w_g \geq 0\} \cap \left\{w : \sum_g w_g = 1\right\}.$$

716 Each half-space  $\{w : w_g \geq 0\}$  is closed in  $\mathbb{R}^G$ , and the hyperplane  $\{\sum_g w_g = 1\}$  is the preimage  
717 of the closed set  $\{1\}$  under the continuous map  $w \mapsto \sum_g w_g$ . A finite intersection of closed sets is  
718 closed, so  $\Delta_G$  is closed.  
719

720 Finally, for every  $w \in \Delta_G$  and each coordinate  $g$ ,

$$721 0 \leq w_g \leq \sum_{h=1}^G w_h = 1$$

722 so  $\Delta_G \subset [0, 1]^G$ , which is bounded in  $\mathbb{R}^G$ . Thus, since any closed and bounded subset of  $\mathbb{R}^G$  is  
723 compact, we have that  $\Delta_G$  is compact.  $\square$   
724

725 Let's go back to the proof of Theorem B.1. By the convexity, continuity and coercivity assumptions  
726 each of the objectives

$$727 \theta \mapsto \sum_g p_g L_g(\theta), \quad \theta \mapsto \max_g L_g(\theta),$$

728 admits a minimizer on the closed convex set  $\Theta$ .

729 Fix  $w^{(1)}, w^{(2)} \in \Delta_G$  and  $\lambda \in [0, 1]$ . Since for each  $\theta$  the map  $w \mapsto \sum_g w_g L_g(\theta)$  is affine, we have

$$\begin{aligned} 730 F(\lambda w^{(1)} + (1 - \lambda)w^{(2)}) &= \inf_{\theta} [\lambda f(\theta, w^{(1)}) + (1 - \lambda)f(\theta, w^{(2)})] \\ 731 &\geq \lambda F(w^{(1)}) + (1 - \lambda)F(w^{(2)}), \end{aligned}$$

732 where  $f(\theta, w) = \sum_g w_g L_g(\theta)$  and we used  $\inf(f_1 + f_2) \geq \inf f_1 + \inf f_2$ . Hence  $F$  is concave  
733 and, being an infimum of continuous functions, upper-semi-continuous on the compact simplex  $\Delta_G$ .  
734 Therefore  $\max_{w \in \Delta_G} F(w)$  is attained at some  $w^*$ .  
735

736 Now define the saddle function

$$737 \Phi(\theta, w) = \sum_{g=1}^G w_g L_g(\theta).$$

738 For fixed  $w$ ,  $\Phi(\cdot, w)$  is convex and lower-bounded; for fixed  $\theta$ ,  $\Phi(\theta, \cdot)$  is affine (hence concave).  
739 Then we have

$$740 \min_{\theta} \max_{w \in \Delta_G} \Phi(\theta, w) = \max_{w \in \Delta_G} \min_{\theta} \Phi(\theta, w).$$

741 Since  $\max_{w \in \Delta_G} \Phi(\theta, w) = \max_g L_g(\theta)$  and  $\min_{\theta} \Phi(\theta, w) = F(w)$ , we obtain

$$742 \min_{\theta} \max_g L_g(\theta) = \max_{w \in \Delta_G} F(w).$$

If  $w^*$  attains the maximum, then any minimizer  $\theta^*$  in  $\min_{\theta} \Phi(\theta, w^*)$  also minimizes the worst-group risk, establishing the claimed equivalence.

Observe that for any  $w \in \Delta_G$  and any  $\theta$ ,  $\max_g L_g(\theta) \geq \sum_{g=1}^G w_g L_g(\theta)$ . Taking  $w = w^*$  and applying first at  $\theta_{\text{ERM}}$  and then at  $\theta^*$ , and using optimality of  $\theta^*$  for weights  $w^*$ , yields

$$\begin{aligned} \max_g L_g(\theta_{\text{ERM}}) &\geq \sum_g w_g^* L_g(\theta_{\text{ERM}}) \\ &\geq \sum_g w_g^* L_g(\theta^*) \\ &= \max_g L_g(\theta^*). \end{aligned}$$

If the group risks at  $\theta_{\text{ERM}}$  are not all equal, the first inequality is strict, completing the proof.  $\square$

## B.2 ALIGNMENT TO SPURIOUS AND NON-SPURIOUS DIRECTIONS FOR ERM AND GDRO

**Proposition B.2.** *Let  $u = v_{sp}/\|v_{sp}\|$  and write  $\theta = \theta_{\perp} + t u$ ,  $t = \langle u, \theta \rangle$ . Assume  $L_g(\theta) = R(\theta_{\perp}) + \phi(t - a_g)$  with the same  $R$  for all  $g$ , where  $\phi \in C^1(\mathbb{R})$  is even, strictly convex,  $\phi(0) = 0$ ,  $\phi'$  odd, and  $\phi$  strictly increasing on  $[0, \infty)$ . Let*

$$a_- := \min_g a_g < 0 < \max_g a_g =: a_+$$

with  $a_+ = -a_- > 0$ . Define

$$\begin{aligned} \theta_{\text{ERM}} &\in \arg \min_{\theta} \sum_g p_g L_g(\theta) \\ \theta^* &\in \arg \min_{\theta} \max_g L_g(\theta), \end{aligned}$$

and set  $\alpha_{sp}(\theta) := |\langle u, \theta \rangle|$ . If

$$\sum_{a>0} (P_-(a) - P_+(a)) \phi'(a) \neq 0,$$

where  $P_{\pm}(a) := \sum_{g: a_g = \pm a} p_g$ , then

$$\alpha_{sp}(\theta^*) < \alpha_{sp}(\theta_{\text{ERM}}).$$

*Proof.* With  $R$  identical, the objectives decouple in  $(\theta_{\perp}, t)$ . Both ERM and DRO share the same  $\theta_{\perp}^* \in \arg \min_{\theta_{\perp}} R(\theta_{\perp})$ . The ERM  $t$  minimizes  $f(t) = \sum_g p_g \phi(t - a_g)$ ;  $f$  is strictly convex and  $f'(0) = \sum_{a>0} (P_+(a) - P_-(a)) \phi'(a) \neq 0$ , hence  $t_{\text{ERM}} \neq 0$  and  $|t_{\text{ERM}}| > 0$ . For DRO, we minimize  $\max_g \phi(t - a_g)$ . Since  $\phi$  is even and strictly increasing on  $[0, \infty)$ , this is minimized at the Chebyshev center of  $\{a_g\}$ , which under the symmetric extremes assumption is  $t^* = 0$ . Therefore  $\alpha_{sp}(\theta^*) = 0 < |t_{\text{ERM}}| = \alpha_{sp}(\theta_{\text{ERM}})$ .  $\square$

**Proposition B.3.** *Let  $u$  be a non-spurious direction. Write  $\theta = \theta_{\perp} + t u$  with  $t = \langle u, \theta \rangle$  and  $\theta_{\perp} \in \{u\}^{\perp}$ . Assume  $L_g(\theta) = R(\theta_{\perp}) + \phi(t - a_g)$ , where  $R$  is the same convex function for all groups, and  $\phi \in C^1(\mathbb{R})$  is even, strictly convex, satisfies  $\phi(0) = 0$ , and has odd, strictly increasing derivative  $\phi'$ . Suppose that*

$$0 < a_{\min} := \min_g a_g \leq a_g \leq a_{\max} := \max_g a_g,$$

and let  $p_g > 0$  with  $\sum_g p_g = 1$ . Define the ERM and GroupDRO optimizers

$$\begin{aligned} t_{\text{ERM}} &\in \arg \min_{t \in \mathbb{R}} \sum_g p_g \phi(t - a_g) \\ t^* &\in \arg \min_{t \in \mathbb{R}} \max_g \phi(t - a_g), \end{aligned}$$

and set  $\alpha_{ns}(\theta) := |\langle u, \theta \rangle| = |t|$ .

Then

- 810 1. Both ERM and DRO choose the same  $\theta_{\perp}^* \in \arg \min_{\theta_{\perp}} R(\theta_{\perp})$ , so the comparison reduces  
 811 to the one-dimensional problems in  $t$  above.  
 812  
 813 2. Let  $m := \frac{1}{2} (a_{\min} + a_{\max})$ . Then  $t^* = m$  and  $t_{\text{ERM}}$  is the unique root of

$$814 f'(t) := \sum_g p_g \phi'(t - a_g) = 0,$$

815  
 816 with  $t_{\text{ERM}} \in [a_{\min}, a_{\max}]$  and  $t_{\text{ERM}} > 0$ .  
 817

- 818 3. We have

$$819 \alpha_{\text{ns}}(\theta^*) \begin{cases} > \alpha_{\text{ns}}(\theta_{\text{ERM}}) & \text{if } f'(m) > 0, \\ = \alpha_{\text{ns}}(\theta_{\text{ERM}}) & \text{if } f'(m) = 0, \\ < \alpha_{\text{ns}}(\theta_{\text{ERM}}) & \text{if } f'(m) < 0. \end{cases}$$

820  
 821 In particular, when  $f'(m) > 0$ , GroupDRO pushes further towards  $|t^*| > |t_{\text{ERM}}|$ .  
 822  
 823

824 *Proof.* 1. Since  $R$  is identical across groups, both ERM and DRO minimize  $R(\theta_{\perp})$  indepen-  
 825 dently of  $t$ , hence share  $\theta_{\perp}^* \in \arg \min R$ . The problems decouple to the one-dimensional  
 826  $t$ -objectives.  
 827

- 828 2. Define  $h(t) := \max_g \phi(t - a_g)$ . Because  $\phi$  is even and strictly increasing on  $[0, \infty)$ ,  
 829  $\arg \min_t h(t) = \arg \min_t \max_g |t - a_g|$ . This is attained at the midrange  $m = (a_{\min} +$   
 830  $a_{\max})/2$ , where the two farthest groups (at the extremes) are equidistant; moving  $t$  left or  
 831 right increases the maximal deviation. Thus  $t^* = m$ .  
 832

833 For ERM,  $f(t) := \sum_g p_g \phi(t - a_g)$  is strictly convex, so  $f'$  is strictly increasing and has a  
 834 unique zero  $t_{\text{ERM}}$ . Moreover,  
 835

$$836 f'(a_{\min}) = \sum_g p_g \phi'(a_{\min} - a_g) \leq 0, \quad f'(a_{\max}) = \sum_g p_g \phi'(a_{\max} - a_g) \geq 0,$$

837  
 838 so  $t_{\text{ERM}} \in [a_{\min}, a_{\max}]$ . Since  $a_g > 0$  for all  $g$  and  $\phi'$  is odd and increasing,  
 839

$$840 f'(0) = \sum_g p_g \phi'(-a_g) = - \sum_g p_g \phi'(a_g) < 0,$$

841  
 842 hence, by monotonicity of  $f'$ , the unique root satisfies  $t_{\text{ERM}} > 0$ .  
 843

- 844 3. Because  $f'$  is strictly increasing, the sign of  $f'(m)$  determines the order of  $t_{\text{ERM}}$  relative  
 845 to  $m$ : if  $f'(m) > 0$  then  $t_{\text{ERM}} < m$ ; if  $f'(m) < 0$  then  $t_{\text{ERM}} > m$ ; if  $f'(m) = 0$  then  
 846  $t_{\text{ERM}} = m$ . Since  $a_{\min} > 0$ , both  $m$  and  $t_{\text{ERM}}$  are positive, so  $\alpha_{\text{ns}}(\theta^*) = |t^*| = m$  and  
 847  $\alpha_{\text{ns}}(\theta_{\text{ERM}}) = |t_{\text{ERM}}| = t_{\text{ERM}}$ , yielding the stated trichotomy.  
 848

849  $\square$

850 **Proposition B.4.** Let  $p_g > 0$ ,  $\sum_g p_g = 1$ , and consider the continuous-time gradient flows  
 851

$$852 \dot{\theta}_{\text{ERM}}(t) = - \sum_{g=1}^G p_g \nabla L_g(\theta_{\text{ERM}}(t))$$

$$853 \dot{\theta}_{\text{DRO}}(t) = - \nabla L_{k(t)}(\theta_{\text{DRO}}(t)),$$

854 where  $k(t) \in \arg \max_g L_g(\theta_{\text{DRO}}(t))$  is any measurable selection. Fix a spurious direction  $v_{\text{sp}} \in$   
 855  $\mathbb{R}^d$  and define  $\alpha_{\text{sp}}(t) := v_{\text{sp}}^{\top} \theta(t)$  and  $\Delta \alpha_{\text{sp}}(t) := \frac{d}{dt} \alpha_{\text{sp}}(t) = v_{\text{sp}}^{\top} \dot{\theta}(t)$ .  
 856

857 Assume the following spurious-gradient monotonicity at the point  $\theta$ :

$$858 L_i(\theta) \geq L_j(\theta) \implies v_{\text{sp}}^{\top} \nabla L_i(\theta) \geq v_{\text{sp}}^{\top} \nabla L_j(\theta) \quad \text{for all } i, j \in \{1, \dots, G\}. \quad (1)$$

859 Then, at  $\theta$ ,

$$860 \Delta \alpha_{\text{sp}}|_{\text{ERM}} \geq \Delta \alpha_{\text{sp}}|_{\text{DRO}}.$$

864 *Proof.* Fix  $\theta$  and abbreviate  $x_g := v_{sp}^\top \nabla L_g(\theta)$ . Then

$$865 \Delta\alpha_{sp}|_{\text{ERM}} = v_{sp}^\top \dot{\theta}_{\text{ERM}} = - \sum_g p_g x_g,$$

$$866 \Delta\alpha_{sp}|_{\text{DRO}} = v_{sp}^\top \dot{\theta}_{\text{DRO}} = -x_k,$$

867 with  $k \in \arg \max_g L_g(\theta)$ . Under inequality 1, any maximizer of  $L_g(\theta)$  is also a maximizer of  $x_g$ ;  
868 hence  $x_k = \max_g x_g$ . Since  $\sum_g p_g x_g \leq \max_g x_g$ ,

$$869 \Delta\alpha_{sp}|_{\text{ERM}} = - \sum_g p_g x_g \geq - \max_g x_g = -x_k = \Delta\alpha_{sp}|_{\text{DRO}}.$$

870 If not all  $x_g$  on the support of  $\{p_g\}$  equal  $\max_g x_g$ , then  $\sum_g p_g x_g < \max_g x_g$ , which yields strict  
871 inequality.  $\square$

872 **Proposition B.5.** Let  $u \in \mathbb{R}^d$  be a fixed non-spurious direction (assume  $\|u\| = 1$  w.l.o.g.). Consider  
873 the continuous-time gradient flows started at the same  $\theta(0)$

$$874 \dot{\theta}_{\text{ERM}}(t) = - \sum_{g=1}^G p_g \nabla L_g(\theta_{\text{ERM}}(t)),$$

$$875 \dot{\theta}_{\text{DRO}}(t) = - \nabla L_{k(t)}(\theta_{\text{DRO}}(t)),$$

876 where  $k(t) \in \arg \max_g L_g(\theta_{\text{DRO}}(t))$  is any measurable selection and  $p_g > 0$ ,  $\sum_g p_g = 1$ . Define  
877 the non-spurious projection  $\alpha_{ns}(t) := \langle u, \theta(t) \rangle$  and its instantaneous rate  $\Delta\alpha_{ns}(t) := \frac{d}{dt} \alpha_{ns}(t) =$   
878  $\langle u, \dot{\theta}(t) \rangle$ .

879 Fix a point  $\theta$  and set  $x_g := \langle u, \nabla L_g(\theta) \rangle$ . Assume the following useful-gradient anti-monotonicity  
880 at  $\theta$

$$881 L_i(\theta) \geq L_j(\theta) \implies \langle u, \nabla L_i(\theta) \rangle \leq \langle u, \nabla L_j(\theta) \rangle \quad \text{for all } i, j \in \{1, \dots, G\}. \quad (2)$$

882 Then, at  $\theta$ ,

$$883 \Delta\alpha_{ns}|_{\text{ERM}} \leq \Delta\alpha_{ns}|_{\text{DRO}}.$$

884 *Proof.* At the common point  $\theta$ ,

$$885 \Delta\alpha_{ns}|_{\text{ERM}} = \langle u, - \sum_g p_g \nabla L_g(\theta) \rangle = - \sum_g p_g x_g,$$

$$886 \Delta\alpha_{ns}|_{\text{DRO}} = \langle u, - \nabla L_k(\theta) \rangle = -x_k,$$

887 where  $k \in \arg \max_g L_g(\theta)$ . By inequality 2, any maximizer of  $L_g$  is a minimizer of  $x_g$ , hence  
888  $x_k = \min_g x_g$ . Since a weighted average is lower-bounded by the minimum,  $\sum_g p_g x_g \geq \min_g x_g$ ,  
889 so

$$890 \Delta\alpha_{ns}|_{\text{ERM}} = - \sum_g p_g x_g \leq - \min_g x_g = \Delta\alpha_{ns}|_{\text{DRO}}.$$

891 Strict inequality holds if some  $x_g$  is strictly larger than  $\min_g x_g$  with positive weight.  $\square$

### 892 B.3 RELATION BETWEEN ALIGNMENT, PERFORMANCE AND REGULARIZATION FOR GDRO

893 **Proposition B.6.** Assume each  $L_g$  is  $\mu$ -strongly convex. Then

$$894 L_g(\theta) \geq L_g(\theta_{\text{ERM}}) + \frac{\mu}{2} \|\theta - \theta_{\text{ERM}}\|^2,$$

895 and consequently

$$896 \alpha_{sp}(\theta_{\text{ERM}}) - \alpha_{sp}(\theta^*) \geq \sqrt{\frac{2(L_{\max}(\theta_{\text{ERM}}) - L_{\max}(\theta^*))}{\mu}}.$$

918 *Proof.* By Part (2) of Theorem B.1 we have,

$$919 \quad L_{\max}(\theta^*) = \min_{\theta} \max_g L_g(\theta) = \max_{w \in \Delta_G} \inf_{\theta} \sum_{g=1}^G w_g L_g(\theta).$$

923 In particular, plugging in the ERM weights  $p = (p_1, \dots, p_G)$  gives

$$924 \quad L_{\max}(\theta^*) \geq \inf_{\theta} \sum_{g=1}^G p_g L_g(\theta) = \sum_{g=1}^G p_g L_g(\theta_{ERM}) = L_{ERM}(\theta_{ERM}),$$

926 where  $L_{ERM}(\theta) = \sum_g p_g L_g(\theta)$ . But by definition  $L_{ERM}(\theta_{ERM}) \leq L_{\max}(\theta_{ERM})$ .

929 Consequently

$$930 \quad L_{\max}(\theta_{ERM}) - L_{\max}(\theta^*) \geq L_{\max}(\theta_{ERM}) - L_{ERM}(\theta_{ERM}) \geq 0.$$

932 Since each  $L_g$  is  $\mu$ -strongly convex,

$$933 \quad L_g(\theta_{ERM}) \geq L_g(\theta^*) + \nabla L_g(\theta^*)^\top (\theta_{ERM} - \theta^*) + \frac{\mu}{2} \|\theta_{ERM} - \theta^*\|^2.$$

936 Taking the maximum over  $g$  and noting  $\max_g \nabla L_g(\theta^*)^\top (\theta_{ERM} - \theta^*) \geq 0$  yields

$$937 \quad L_{\max}(\theta_{ERM}) \geq L_{\max}(\theta^*) + \frac{\mu}{2} \|\theta_{ERM} - \theta^*\|^2.$$

939 Rearrange to get

$$940 \quad \|\theta_{ERM} - \theta^*\| \leq \sqrt{\frac{2(L_{\max}(\theta_{ERM}) - L_{\max}(\theta^*))}{\mu}}.$$

943 Finally, since  $\alpha_{sp}(\theta) = |v_{sp}^\top \theta|$ , the one-dimensional distance along  $v_{sp}$  is at most the full Euclidean distance

$$944 \quad \alpha_{sp}(\theta_{ERM}) - \alpha_{sp}(\theta^*) \geq \|(\theta_{ERM} - \theta^*)\| \geq \sqrt{\frac{2(L_{\max}(\theta_{ERM}) - L_{\max}(\theta^*))}{\mu}}.$$

949  $\square$

950 **Proposition B.7.** Let each group risk  $L_g: \Theta \rightarrow \mathbb{R}$  be convex, and fix a spurious unit vector  $u = v_{sp}/\|v_{sp}\|$ . For a regularization parameter  $\lambda > 0$ , define

$$951 \quad R_\lambda(\theta) = \frac{\lambda}{2} \|\theta\|^2, \quad \tilde{L}_g(\theta) = L_g(\theta) + R_\lambda(\theta),$$

952 and write

$$953 \quad \theta_{ERM}^\lambda = \arg \min_{\theta} \sum_{g=1}^G p_g \tilde{L}_g(\theta), \quad \theta^{\lambda,*} = \arg \min_{\theta} \max_g \tilde{L}_g(\theta),$$

954 with  $\alpha_{sp}(\theta) = |u^\top \theta|$ ,  $\tilde{L}_{\max}(\theta) = \max_g \tilde{L}_g(\theta)$ . Then

$$955 \quad \frac{\lambda(\alpha_{sp}(\theta_{ERM}^\lambda) - \alpha_{sp}(\theta^{\lambda,*}))^2}{2} \geq \tilde{L}_{\max}(\theta_{ERM}^\lambda) - \tilde{L}_{\max}(\theta^{\lambda,*})$$

963 *Proof.* Since  $R_\lambda(\theta)$  is  $\lambda$ -strongly convex and each  $L_g$  is convex, each  $\tilde{L}_g = L_g + R_\lambda$  is  $\lambda$ -strongly convex. We then proceed as follows.

965 By  $\lambda$ -strong convexity, for any  $\theta$  and any group  $g$ ,

$$966 \quad \tilde{L}_g(\theta) \geq \tilde{L}_g(\theta_{ERM}^\lambda) + \nabla \tilde{L}_g(\theta_{ERM}^\lambda)^\top (\theta - \theta_{ERM}^\lambda) + \frac{\lambda}{2} \|\theta - \theta_{ERM}^\lambda\|^2.$$

969 Taking maximum over  $g$  and using that  $\max_g \nabla \tilde{L}_g(\theta_{ERM}^\lambda)^\top (\theta - \theta_{ERM}^\lambda) \geq 0$  (since  $\sum_g p_g \nabla \tilde{L}_g(\theta_{ERM}^\lambda) = 0$ ) yields

$$970 \quad \tilde{L}_{\max}(\theta) \geq \tilde{L}_{\max}(\theta_{ERM}^\lambda) + \frac{\lambda}{2} \|\theta - \theta_{ERM}^\lambda\|^2.$$

By the minimax equality,

$$\tilde{L}_{\max}(\theta^{\lambda,*}) = \min_{\theta} \max_g \tilde{L}_g(\theta) = \max_{w \in \Delta_G} \inf_{\theta} \sum_g w_g \tilde{L}_g(\theta) \geq \inf_{\theta} \sum_g p_g \tilde{L}_g(\theta) = \sum_g p_g \tilde{L}_g(\theta_{ERM}^{\lambda}),$$

so in particular  $\tilde{L}_{\max}(\theta_{ERM}^{\lambda}) - \tilde{L}_{\max}(\theta^{\lambda,*}) \geq 0$ .

Combining with the quadratic growth at  $\theta = \theta^{\lambda,*}$  gives

$$\frac{\lambda}{2} \|\theta^{\lambda,*} - \theta_{ERM}^{\lambda}\|^2 \leq \tilde{L}_{\max}(\theta_{ERM}^{\lambda}) - \tilde{L}_{\max}(\theta^{\lambda,*}),$$

hence

$$\|\theta_{ERM}^{\lambda} - \theta^{\lambda,*}\| \leq \sqrt{\frac{2(\tilde{L}_{\max}(\theta_{ERM}^{\lambda}) - \tilde{L}_{\max}(\theta^{\lambda,*}))}{\lambda}}.$$

Finally, since  $\alpha_{sp}(\theta) = |u^T \theta| \leq \|\theta\|$ , we have

$$\alpha_{sp}(\theta_{ERM}^{\lambda}) - \alpha_{sp}(\theta^{\lambda,*}) = |u^T(\theta_{ERM}^{\lambda} - \theta^{\lambda,*})| \geq \|\theta_{ERM}^{\lambda} - \theta^{\lambda,*}\| \geq \sqrt{\frac{2(\tilde{L}_{\max}(\theta_{ERM}^{\lambda}) - \tilde{L}_{\max}(\theta^{\lambda,*}))}{\lambda}}.$$

This completes the proof under  $\ell_2$  regularization.  $\square$

**Proposition B.8.** Assume each group risk  $L_g: \Theta \rightarrow \mathbb{R}$  is convex and continuously differentiable. Let  $M \in \mathbb{R}^{d \times d}$  be symmetric positive definite, and define the regularizer

$$R_M(\theta) = \frac{1}{2} \theta^T M \theta.$$

Fix any unit vector  $w \in \mathbb{R}^d$ , and write

$$\alpha_w(\theta) = |w^T \theta|.$$

For a parameter set of probabilities  $\{p_g\}$ , define

$$\tilde{L}_g(\theta) = L_g(\theta) + R_M(\theta), \quad \theta_{ERM} = \arg \min_{\theta} \sum_{g=1}^G p_g \tilde{L}_g(\theta), \quad \theta^* = \arg \min_{\theta} \max_g \tilde{L}_g(\theta),$$

and set

$$\Delta L = \tilde{L}_{\max}(\theta_{ERM}) - \tilde{L}_{\max}(\theta^*), \quad \tilde{L}_{\max}(\theta) = \max_g \tilde{L}_g(\theta).$$

Let  $\mu_w = w^T M w > 0$ . Then the following directional quadratic bound holds:

$$\frac{\mu_w}{2} (\alpha_w(\theta_{ERM}) - \alpha_w(\theta^*))^2 \leq \Delta L.$$

*Proof.* Since  $M \succ 0$ , the function  $\theta \mapsto R_M(\theta)$  is  $\mu_w$ -strongly convex along  $w$ . In particular, for each  $g$  and any  $\theta$  the strong convexity inequality in direction  $w$  gives

$$R_M(\theta) \geq R_M(\theta_{ERM}) + \nabla R_M(\theta_{ERM})^T (\theta - \theta_{ERM}) + \frac{\mu_w}{2} (w^T \theta - w^T \theta_{ERM})^2.$$

Adding the convex term  $L_g$  preserves this inequality:

$$\tilde{L}_g(\theta) \geq \tilde{L}_g(\theta_{ERM}) + \nabla \tilde{L}_g(\theta_{ERM})^T (\theta - \theta_{ERM}) + \frac{\mu_w}{2} (w^T \theta - w^T \theta_{ERM})^2.$$

Set  $\theta = \theta^*$  and take the maximum over  $g$ . Using the minimax identity

$$\tilde{L}_{\max}(\theta^*) = \min_{\theta} \max_g \tilde{L}_g(\theta) = \max_{v \in \Delta_G} \inf_{\theta} \sum_g v_g \tilde{L}_g(\theta) \geq \sum_g p_g \tilde{L}_g(\theta_{ERM}),$$

we obtain

$$\tilde{L}_{\max}(\theta_{ERM}) - \tilde{L}_{\max}(\theta^*) \geq \sum_g p_g \left[ \nabla \tilde{L}_g(\theta_{ERM})^T (\theta^* - \theta_{ERM}) + \frac{\mu_w}{2} (w^T \theta^* - w^T \theta_{ERM})^2 \right].$$

Stationarity of  $\theta_{ERM}$  implies  $\sum_g p_g \nabla \tilde{L}_g(\theta_{ERM}) = 0$ , so the linear term vanishes. Hence

$$\Delta L \geq \frac{\mu_w}{2} (w^T \theta^* - w^T \theta_{ERM})^2 = \frac{\mu_w}{2} (\alpha_w(\theta_{ERM}) - \alpha_w(\theta^*))^2.$$

This completes the proof.  $\square$

## C DO GRMS ELIMINATE INFORMATION ABOUT SPURIOUS FEATURES?

We will analyze first if the representations learned by GRMs contain information about spurious features. One hypothesis about their success might relate to them discarding this information during training. To do this, we obtain the singular vectors of a PCA decomposition of representations of the training set. For each singular vector we train a logistic regression to predict the spurious label on the training set and then evaluate on the test split. Figure 4 shows these results for all datasets. For models trained on MNIST-CIFAR, the first direction is surprisingly predictive of the MNIST label both for the training set and the test set, independently of which method was used to train. What is remarkable is that this happens even if the spurious label is not actually useful for the task: models trained on datasets with no correlation between the spurious and target label still retain information about the spurious label. This behaviour happens also in more complex datasets like Waterbirds and CelebA, with the first and second directions, respectively being highly predictive of the spurious label. This behavior is consistent across all methods, which suggests that all of them store spurious information in their representations and there is not much of a mechanism to curb this.

1080  
 1081  
 1082  
 1083  
 1084  
 1085  
 1086  
 1087  
 1088  
 1089  
 1090  
 1091  
 1092  
 1093  
 1094  
 1095  
 1096  
 1097  
 1098  
 1099  
 1100  
 1101  
 1102  
 1103  
 1104  
 1105  
 1106  
 1107  
 1108  
 1109  
 1110  
 1111  
 1112  
 1113  
 1114  
 1115  
 1116  
 1117  
 1118  
 1119  
 1120  
 1121  
 1122  
 1123  
 1124  
 1125  
 1126  
 1127  
 1128  
 1129  
 1130  
 1131  
 1132  
 1133

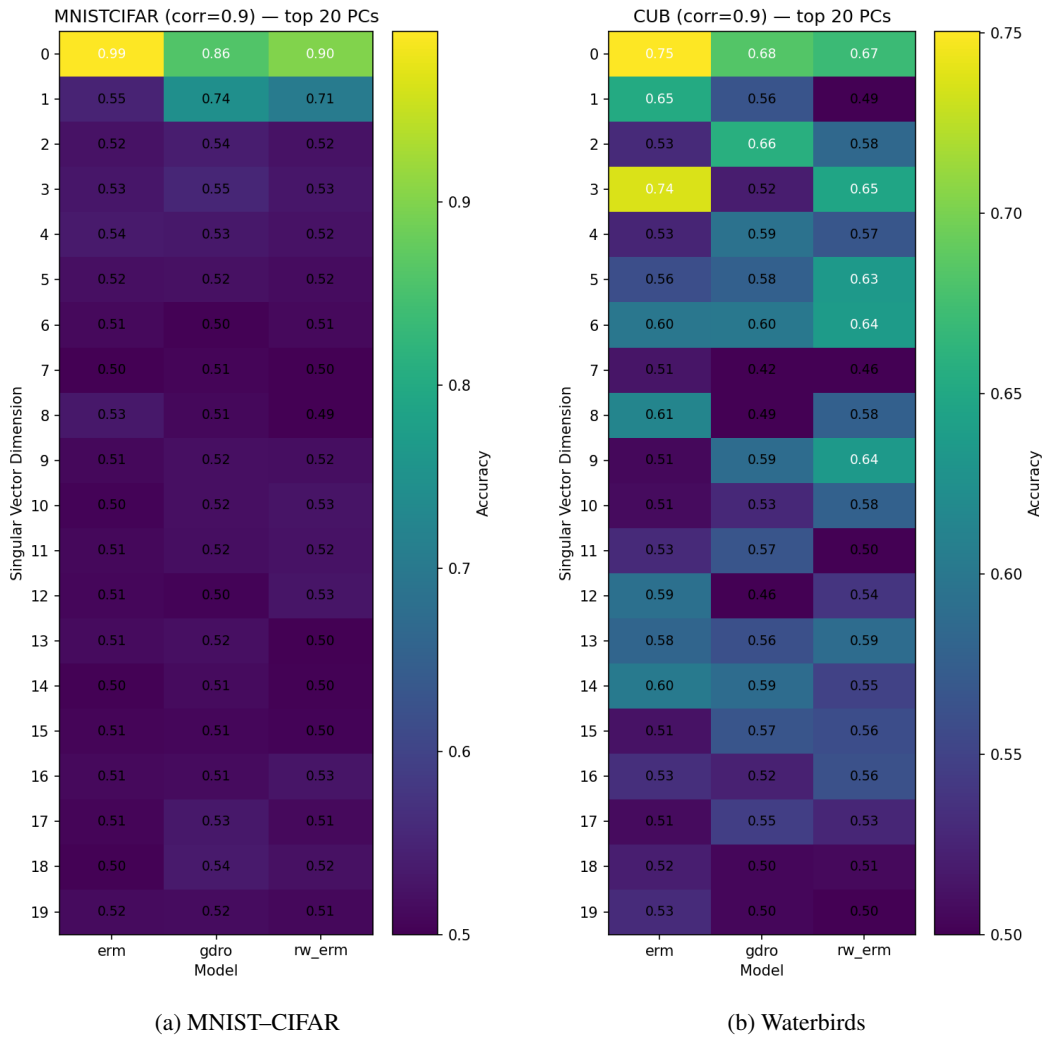


Figure 4: Test accuracy for predicting the spurious label (corr=0.9) across the top-20 singular vectors for models trained on five datasets. For all methods we find highly predictive spurious directions, often within the top-3 singular vectors, suggesting their relevance.

1134  
 1135  
 1136  
 1137  
 1138  
 1139  
 1140  
 1141  
 1142  
 1143  
 1144  
 1145  
 1146  
 1147  
 1148  
 1149  
 1150  
 1151  
 1152  
 1153  
 1154  
 1155  
 1156  
 1157  
 1158  
 1159  
 1160  
 1161  
 1162  
 1163  
 1164  
 1165  
 1166  
 1167  
 1168  
 1169  
 1170  
 1171  
 1172  
 1173  
 1174  
 1175  
 1176  
 1177  
 1178  
 1179  
 1180  
 1181  
 1182  
 1183  
 1184  
 1185  
 1186  
 1187

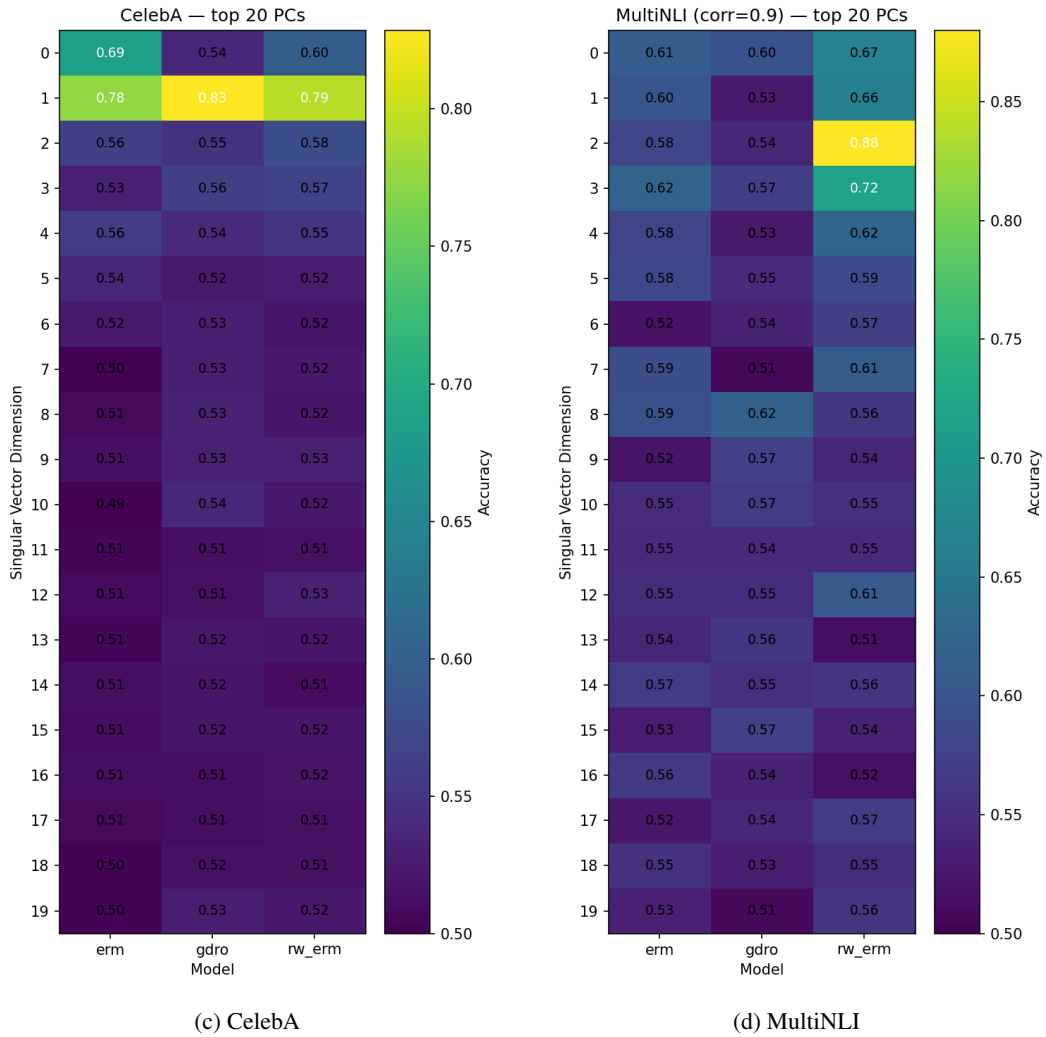
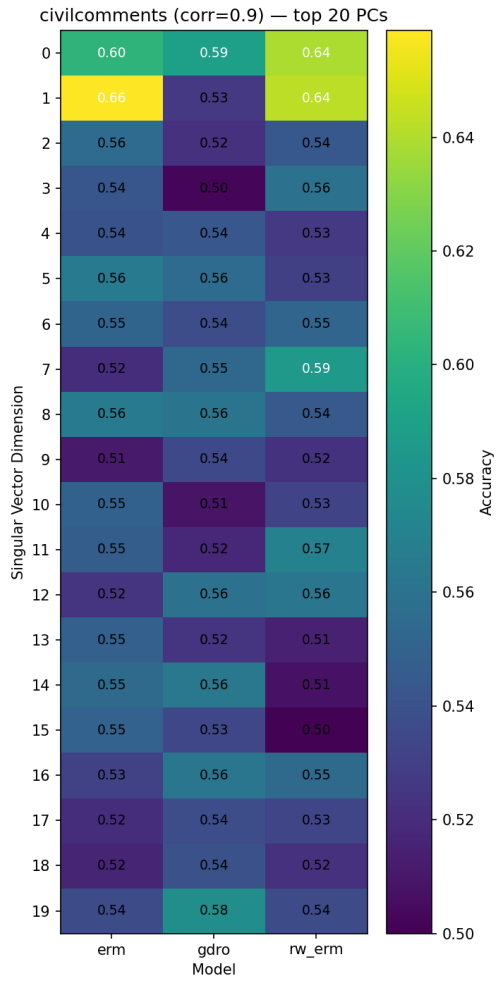


Figure 4: (continued)

1188  
 1189  
 1190  
 1191  
 1192  
 1193  
 1194  
 1195  
 1196  
 1197  
 1198  
 1199  
 1200  
 1201  
 1202  
 1203  
 1204  
 1205  
 1206  
 1207  
 1208  
 1209  
 1210  
 1211  
 1212  
 1213  
 1214  
 1215  
 1216  
 1217  
 1218  
 1219  
 1220  
 1221  
 1222  
 1223  
 1224  
 1225  
 1226  
 1227  
 1228  
 1229  
 1230  
 1231  
 1232  
 1233  
 1234  
 1235  
 1236  
 1237  
 1238  
 1239  
 1240  
 1241



(e) CivilComments

Figure 4: (continued)

## D ADDITIONAL EXPERIMENTS

### D.1 ABLATION ON DATASET SIZE FOR RW+ AND GDRO+

DATASET	METHOD	0.1	0.25	0.5	0.75	0.9
MNIST-CIFAR	RW-FT	<b>80.63%</b> (1.19)	79.33% (1.20)	78.88% (2.79)	79.28% (2.07)	<b>80.40%</b> (0.81)
	GDRO-FT	78.77% (2.54)	78.53% (1.46)	78.01% (2.11)	78.96% (0.75)	78.22% (1.51)
	RW+	75.85% (5.13)	79.79% (1.81)	<b>80.59%</b> (1.01)	<b>80.72%</b> (0.40)	79.12% (3.14)
	GDRO+	76.73% (5.59)	<b>81.39%</b> (0.61)	<b>80.59%</b> (0.23)	80.19% (0.46)	79.25% (2.91)
WATERBIRDS	RW-FT	49.43% (1.04)	56.28% (2.16)	62.88% (3.74)	66.20% (2.36)	66.25% (0.80)
	GDRO-FT	53.32% (1.72)	62.36% (3.28)	67.68% (0.70)	72.95% (3.74)	71.86% (1.26)
	RW+	79.22% (7.19)	83.84% (2.16)	86.64% (0.65)	88.00% (0.82)	<b>88.99%</b> (0.55)
	GDRO+	<b>79.85%</b> (7.42)	<b>84.33%</b> (2.66)	<b>86.75%</b> (0.82)	<b>88.21%</b> (0.78)	88.92% (0.35)
CELEBA	RW-FT	83.15% (0.64)	84.63% (1.95)	85.74% (0.85)	86.48% (1.70)	86.85% (1.40)
	GDRO-FT	<b>87.41%</b> (2.57)	<b>89.60%</b> (1.10)	<b>89.09%</b> (1.18)	<b>89.26%</b> (1.79)	<b>90.44%</b> (0.53)
	RW+	83.70% (0.85)	84.63% (1.95)	85.56% (0.56)	85.00% (0.96)	85.19% (0.85)
	GDRO+	80.74% (4.10)	86.30% (2.10)	88.15% (0.64)	88.70% (0.85)	88.71% (0.85)

Table 2: Worst Group Accuracy results for various methods trained on all datasets (correlation=0.9) on different percentages of training data. All methods finetuned a classifier from frozen features derived from ERM training. Note that RW+/GDRO+ already uses around 5-10% of the original training data. Even with 50% of their training data, very competitive results are obtained with GDRO+.

1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

## D.2 EXPLAINED VARIANCE OF PCA FOR ALL METHODS

DATASET	CORRELATION	METHOD	EXPLAINED RATIO	MIN	MAX	PCA DIMS
WATERBIRDS	0.00	ERM	91.34% (0.22)	91.14	91.58	200
WATERBIRDS	0.00	GDRO	92.05% (0.95)	91.37	93.13	200
WATERBIRDS	0.00	RW	89.45% (0.16)	89.31	89.63	200
WATERBIRDS	0.25	ERM	90.69% (0.57)	90.14	91.28	200
WATERBIRDS	0.25	GDRO	95.06% (4.26)	92.14	99.95	200
WATERBIRDS	0.25	RW	90.10% (1.08)	89.40	91.35	200
WATERBIRDS	0.50	ERM	91.25% (0.16)	91.06	91.38	200
WATERBIRDS	0.50	GDRO	93.23% (2.62)	91.40	96.23	200
WATERBIRDS	0.50	RW	90.07% (1.02)	89.32	91.23	200
WATERBIRDS	0.75	ERM	91.31% (0.02)	91.30	91.33	200
WATERBIRDS	0.75	GDRO	95.11% (0.28)	94.89	95.42	200
WATERBIRDS	0.75	RW	89.46% (0.05)	89.42	89.52	200
WATERBIRDS	0.90	ERM	90.99% (0.12)	90.85	91.06	200
WATERBIRDS	0.90	GDRO	89.45% (0.43)	89.08	89.92	200
WATERBIRDS	0.90	RW	89.39% (0.06)	89.32	89.44	200
CELEBA	0.25	ERM	89.21% (0.08)	89.15	89.30	200
CELEBA	0.25	GDRO	90.24% (0.21)	90.02	90.43	200
CELEBA	0.25	RW	89.92% (0.26)	89.69	90.20	200
MNISTCIFAR	0.00	ERM	87.56% (0.56)	87.06	88.16	50
MNISTCIFAR	0.00	GDRO	87.75% (0.72)	87.22	88.56	50
MNISTCIFAR	0.00	RW	87.58% (0.64)	87.00	88.26	50
MNISTCIFAR	0.25	ERM	88.53% (0.29)	88.20	88.72	50
MNISTCIFAR	0.25	GDRO	87.78% (0.63)	87.11	88.37	50
MNISTCIFAR	0.25	RW	87.74% (0.65)	87.07	88.37	50
MNISTCIFAR	0.50	ERM	89.37% (0.40)	88.98	89.77	50
MNISTCIFAR	0.50	GDRO	87.27% (0.54)	86.75	87.82	50
MNISTCIFAR	0.50	RW	87.30% (0.62)	86.77	87.99	50
MNISTCIFAR	0.75	ERM	90.89% (0.96)	90.07	91.95	50
MNISTCIFAR	0.75	GDRO	87.33% (1.07)	86.36	88.48	50
MNISTCIFAR	0.75	RW	87.24% (1.10)	86.20	88.39	50
MNISTCIFAR	0.90	ERM	92.52% (0.97)	91.70	93.59	50
MNISTCIFAR	0.90	GDRO	89.62% (1.40)	88.13	90.90	50
MNISTCIFAR	0.90	RW	89.39% (1.27)	87.96	90.39	50
MULTINLI	0.90	ERM	99.09% (0.35)	98.87	99.49	50
MULTINLI	0.90	GDRO	99.42% (0.13)	99.28	99.52	50
MULTINLI	0.90	RW	99.07% (0.21)	98.92	99.32	50
CIVILCOMMENTS	0.90	ERM	99.23% (0.18)	99.03	99.33	50
CIVILCOMMENTS	0.90	GDRO	99.32% (0.06)	99.26	99.37	50
CIVILCOMMENTS	0.90	RW	99.36% (0.10)	99.25	99.44	50

Table 3: Percentage of Explained Variance for all PCA decompositions used for calculating DCI metrics.

## D.3 DCI METRICS IN FULL

DATASET	CORR	DISENTANGLEMENT			COMPLETENESS			INFORMATIVENESS		
		ERM	GDRO	RW	ERM	GDRO	RW	ERM	GDRO	RW
MNIST-CIFAR	0.00	0.992	0.991	0.993	0.768	0.765	0.769	0.957	0.957	0.957
	0.25	0.545	0.449	0.454	0.714	0.697	0.700	0.955	0.957	0.957
	0.50	0.252	0.230	0.231	0.757	0.752	0.752	0.958	0.963	0.963
	0.75	0.106	0.105	0.105	0.841	0.840	0.838	0.968	0.973	0.974
	0.90	0.030	0.016	0.017	0.921	0.819	0.865	0.984	0.983	0.984
WATERBIRDS	0.00	0.984	0.884	0.774	0.890	0.825	0.839	0.974	0.971	0.967
	0.25	0.774	0.699	0.695	0.831	0.722	0.816	0.974	0.959	0.971
	0.50	0.369	0.479	0.317	0.802	0.758	0.775	0.979	0.974	0.976
	0.75	0.130	0.188	0.145	0.875	0.778	0.797	0.990	0.984	0.980
	0.90	0.030	0.045	0.044	0.895	0.845	0.853	0.992	0.983	0.985
CELEBA	0.87	0.432	0.812	0.665	0.829	0.817	0.785	0.940	0.927	0.927
MULTINLI	0.22	0.500	0.619	0.687	0.576	0.558	0.702	0.954	0.933	0.920
CIVILCOMMENTS	0.52	0.649	0.716	0.593	0.622	0.531	0.584	0.881	0.858	0.882

Table 4: Disentanglement, Completeness and Informativeness results for the train set for multiple datasets and methods.

DATASET	CORR	DISENTANGLEMENT			COMPLETENESS			INFORMATIVENESS		
		ERM	GDRO	RW	ERM	GDRO	RW	ERM	GDRO	RW
MNIST-CIFAR	0.00	0.992	0.991	0.993	0.768	0.765	0.769	0.939	0.942	0.942
	0.25	0.545	0.449	0.454	0.714	0.697	0.700	0.926	0.930	0.928
	0.50	0.252	0.230	0.231	0.757	0.752	0.752	0.873	0.888	0.886
	0.75	0.106	0.105	0.105	0.841	0.840	0.838	0.755	0.798	0.801
	0.90	0.030	0.016	0.017	0.921	0.819	0.865	0.622	0.647	0.651
WATERBIRDS	0.00	0.984	0.884	0.774	0.890	0.825	0.839	0.937	0.938	0.940
	0.25	0.774	0.699	0.695	0.831	0.722	0.816	0.928	0.901	0.930
	0.50	0.369	0.479	0.317	0.802	0.758	0.775	0.921	0.913	0.915
	0.75	0.130	0.188	0.145	0.875	0.778	0.797	0.897	0.877	0.882
	0.90	0.030	0.045	0.044	0.895	0.845	0.853	0.866	0.815	0.819
CELEBA	0.87	0.432	0.812	0.665	0.829	0.817	0.785	0.935	0.925	0.925
MULTINLI	0.22	0.500	0.619	0.687	0.576	0.558	0.702	0.886	0.879	0.892
CIVILCOMMENTS	0.52	0.649	0.716	0.593	0.622	0.531	0.584	0.868	0.849	0.865

Table 5: Disentanglement, Completeness and Informativeness results for the test set for multiple datasets and methods.

1404 E LIMITATIONS AND SOCIETAL IMPACT

1405  
1406 Our analysis hinges on the following assumptions: groups within the dataset are usually unbalanced  
1407 with the most represented groups benefiting from the spurious correlation. However, this is a stan-  
1408 dard setting in the robustness literature. Our theoretical analysis depends on a specific definition of  
1409 a spurious vector, with which some may disagree.

1410 Our proposed method requires extra data with annotations of both the class and spurious label for  
1411 the finetuning stage. It does not require much of it, 5% for the full method, but Table 3 shows decent  
1412 results with even less. All competing baselines require at least this much extra data.

1413 We believe this work may have a positive societal impact as it is trying to mitigate a problem of spu-  
1414 rious correlations that usually happens on underrepresented subsets of the data, which may represent  
1415 in practice underrepresented parts of society. Moreover, our work seeks to understand the mecha-  
1416 nisms by which GRMs work which could aid in developing better methods that work to achieve  
1417 robustness and fairness.

1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457

## 1458 F EXPERIMENT DETAILS

1459

### 1460 F.1 HYPERPARAMETERS

1461

#### 1462 F.1.1 GENERAL

1463

1464 All methods use SGD with momentum 0.9 as their optimizer for vision datasets. For text datasets  
1465 we used AdamW. It is the same setup as Liu et al. (2021). We use L2 regularization and no data  
1466 augmentation. No learning rate scheduler is used. All methods are run for 3 seeds. Seeds used were:  
1467 {111, 222, 333} Learning rates and weight decay used were taken from (Liu et al., 2021), which  
1468 were based off of (Sagawa\* et al., 2020).

#### 1469 F.1.2 MNIST-CIFAR

1470

1471 We train all models for 5000 epochs/iterations. L2 regularization of  $10^{-4}$ , learning rate of 0.001 for  
1472 all methods. Batch size is 10000.

1473

#### 1474 F.1.3 WATERBIRDS

1475

1476 We train all models for 300 epochs. Batch size is 64. For ERM, we use L2 regularization of  
1477  $10^{-4}$ , learning rate of  $10^{-4}$ ; for GDRO L2 regularization of 1, learning rate of  $10^{-5}$ ; for RW, L2  
1478 regularization of  $10^{-3}$ , learning rate of  $10^{-4}$ .

#### 1479 F.1.4 CELEBA

1480

1481 We train all models for 50 epochs. Batch size is 64. For ERM, we use L2 regularization of  $10^{-4}$ ,  
1482 learning rate of  $10^{-4}$ ; for GDRO L2 regularization of 0.1, learning rate of  $10^{-5}$ ; for RW, L2 regu-  
1483 larization of 0.1, learning rate of  $10^{-5}$ .

#### 1484 F.1.5 MULTINLI

1485

1486 We train all models for 50 epochs. Batch size is 64. For ERM, we use L2 regularization of  $10^{-4}$ ,  
1487 learning rate of  $10^{-4}$ ; for GDRO L2 regularization of 0.1, learning rate of  $10^{-5}$ ; for RW, L2 regu-  
1488 larization of 0.1, learning rate of  $10^{-5}$ .

1489

#### 1490 F.1.6 CIVILCOMMENTS

1491

1492 We train all models for 50 epochs. Batch size is 64. For ERM, we use L2 regularization of  $10^{-4}$ ,  
1493 learning rate of  $10^{-4}$ ; for GDRO L2 regularization of 0.1, learning rate of  $10^{-5}$ ; for RW, L2 regu-  
1494 larization of 0.1, learning rate of  $10^{-5}$ .

1495

## 1496 G CO2 EMISSION RELATED TO EXPERIMENTS

1497

1498 Experiments were conducted using a private infrastructure, which has a carbon efficiency of 0.432  
1499 kgCO<sub>2</sub>eq/kWh. A cumulative of 2275 hours of computation was performed on hardware of type  
1500 GTX 1080 Ti (TDP of 250W).

1501 Total emissions are estimated to be 245.7 kgCO<sub>2</sub>eq of which 0 percents were directly offset.

1502

1503 Estimations were conducted using the MachineLearning Impact calculator presented in Lacoste et al.  
1504 (2019).

1505

## 1506 H LLM USAGE

1507

1508 LLMs were only used in polishing writing and in aiding in creating tables and plots.

1509

1510

1511