# TRAINED MODELS TELL US HOW TO MAKE THEM ROBUST TO SPURIOUS CORRELATION WITHOUT GROUP ANNOTATION

Anonymous authors

Paper under double-blind review

### ABSTRACT

Classifiers trained with Empirical Risk Minimization (ERM) tend to rely on attributes that have high spurious correlation with the target. This can degrade the performance on underrepresented (or *minority*) groups that lack these attributes, posing significant challenges for both out-of-distribution generalization and fairness objectives. Many studies aim to enhance robustness to spurious correlation, but they sometimes depend on group annotations for training. Additionally, a common limitation in previous research is the reliance on group-annotated validation datasets for model selection. This constrains their applicability in situations where the nature of the spurious correlation is not known, or when group labels for certain spurious attributes are not available. To enhance model robustness with minimal group annotation assumptions, we propose Environment-based Validation and Loss-based Sampling (EVaLS). It uses the losses from an ERM-trained model to construct a balanced dataset of high-loss and low-loss samples, mitigating group imbalance in data. This significantly enhances robustness to group shifts when equipped with a simple post-training last layer retraining. By using environment inference methods to create diverse environments with correlation shifts, EVaLS can potentially eliminate the need for group annotation in validation data. In this context, the worst environment accuracy acts as a reliable surrogate throughout the retraining process for tuning hyperparameters and finding a model that performs well across diverse group shifts. EVaLS effectively achieves group robustness, showing that group annotation is not necessary even for validation. It is a fast, straightforward, and effective approach that reaches near-optimal worst group accuracy without needing group annotations, marking a new chapter in the robustness of trained models against spurious correlation.

034 035

037

006

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

027

028

029

031

032

033

# 1 INTRODUCTION

Training deep learning models using Empirical Risk Minimization (ERM) on a dataset, poses the risk of relying on *spurious correlation*. These are correlations between certain patterns in the train-040 ing dataset and the target (e.g., the class label in a classification task) despite lacking any causal 041 relationship. Learning such correlations as shortcuts can negatively impact the models' accuracy 042 on minority groups that do not contain the spurious patterns associated with the target (Kirichenko 043 et al., 2023; LaBonte et al., 2023). This problem leads to concerns regarding fairness (Hashimoto 044 et al., 2018), and can also cause a marked reduction in the performance. This occurs particularly when minority groups, which are underrepresented during training, become overrepresented at the inference time, as a result of shifts within the subpopulations (Yang et al., 2023b). Hence, ensuring 046 robustness to group shifts and developing methods that improve worst group accuracy (WGA) is 047 crucial for achieving both fairness and robustness in the realm of deep learning. 048

Many studies have proposed solutions to address this challenge. A promising line of research focuses on increasing the contribution of minority groups in the model's training (Liu et al., 2021a;
Yang et al., 2023a; Sagawa et al., 2019). A strong assumption that is considered by some previous
works is having access to group annotations for training or fully/partially fine-tuning a pretrained
model (Nam et al., 2021; Sagawa et al., 2019; Kirichenko et al., 2023). The study by Kirichenko et al. (2023) proposes that retraining the last layer of a model on a dataset that is balanced in terms

of group annotation can effectively enhance the model's robustness against shifts in spurious correlation. While these works have shown tremendous robustness performance, their assumption for the availability of the group annotation restricts their usage.

057 In many real-world applications, the process of labeling samples according to their respective groups can be prohibitively expensive, and sometimes impractical, especially when all minority groups may 059 not be identifiable beforehand. A widely adopted strategy in these situations involves the indirect 060 inference of various groups, followed by the training of models using a loss function that is balanced 061 across groups (Liu et al., 2021a; Qiu et al., 2023; Nam et al., 2020; Yang et al., 2023b). The loss 062 value of the model, or its alternatives, are popular signals for recognizing minority groups (Liu et al., 063 2021a; Qiu et al., 2023; Nam et al., 2020; Noohdani et al., 2024). While most of these techniques 064 necessitate full training of a model, Qiu et al. (2023) attempt to adapt the DFR method (Kirichenko et al., 2023) with the aim of preserving computational efficiency while simultaneously improving 065 robustness to the group shift. However, this method still requires group annotations of the validation 066 set for the model selection and hyperparameter tuning. Consequently, this constitutes a restrictive as-067 sumption when adequate annotations for certain groups are not supplied. It also applies to situations 068 where some shortcut attributes are completely unknown. 069

070 In this study, we present a novel strategy that effectively mitigates reliance on spurious correlation, completely eliminating the need for group annotations during both training and retraining. More 071 interestingly, we provide empirical evidence indicating that group annotations are not necessary, 072 even for model selection. We show that assembling a diverse collection of environments for model 073 selection, which reflects group shifts can serve as an effective alternative approach. Our proposed 074 scheme, Environment-based Validation and Loss-based Sampling (EVaLS), strengthens the robust-075 ness of trained models against spurious correlation, all without relying on group annotations. EVaLS 076 is pioneering in its ability to eliminate the need for group annotations at *every phase*, including the 077 model selection step. EVaLS posits that in the absence of group annotations, a set of environments showcasing group shifts is sufficient. Worst Environment Accuracy (WEA) could then be utilized 079 for model selection. We observe that spurious correlations, as a form of subpopulation shifts, cause 080 significant group shifts when using environment inference methods (Creager et al., 2021). Conse-081 quently, the inferred environments-which could be obtained even by simply dividing validation data based on predictions from a random linear layer atop a trained model's feature space—can effectively compare different sets of hyperparameters for tuning. Figure 1 demonstrates the overall 083 procedure of the main parts of EVaLS. 084

Aligned with AFR (Qiu et al., 2023) and DFR (Kirichenko et al., 2023), EVaLS offers a significant advantage by not requiring any modifications to the standard ERM training procedure or the original training data. Moreover, it does not require information from the initial phases of ERM training, such as an early-stopped model. This characteristic is particularly beneficial in enhancing the robustness of ERM-pretrained networks against their potential inherent biases. Specifically, it eliminates the need to retrain the entire model, which may be impractical or infeasible when the original training data is unavailable.

092 Our empirical observations support prior research which suggests that high-loss data points in a 093 trained model may signal the presence of minority groups (Liu et al., 2021a; Qiu et al., 2023; Nam et al., 2020). EVaLS evenly selects from both high-loss and low-loss data to form a balanced dataset 094 that is used for last-layer retraining. We offer theoretical explanations for the effectiveness of this 095 approach in addressing group imbalances, and experimentally show the superiority of our efficient 096 solution to the previous strategies. Comprehensive experiments conducted on spurious correlation benchmarks such as CelebA (Liu et al., 2014), Waterbirds (Sagawa et al., 2019), and UrbanCars (Li 098 et al., 2023), demonstrate that EVaLS achieves optimal accuracy. Moreover, when group annotations are accessible solely for model selection, our approach, EVaLS-GL, exhibits enhanced performance 100 against various distribution shifts, including attribute imbalance, as seen in MultiNLI (Williams 101 et al., 2017), and class imbalance, exemplified by CivilComments (Borkan et al., 2019). We further 102 present a new dataset, Dominoes Colored-MNIST-FashionMNIST, which depicts a situation featur-103 ing multiple independent shortcuts, that group annotations are only available for part of them (see 104 Section 2.2). In this setting, we show that strategies with lower levels of group supervision are paradoxically more effective in mitigating the reliance on both known and unknown shortcuts. 105

106

The main contributions of this paper are summarized as follows:

- We present EVaLS, a simple yet effective post-hoc approach that enhances the robustness of ERM-pretrained models against both known and unknown spurious correlations, without relying on ground-truth group annotations.
  We offer both theoretical and empirical insights on how balanced sampling from high-loss and low-loss samples offers a dataset in which the group imbalance is notably mitigated.
  Using simple environment inference techniques, EVaLS introduces worst environment accuracy as a reliable indicator for model selection.
  EVaLS achieves near-optimal performance in spurious correlation benchmarks with zero group annotations, and delivers state-of-the-art performance when group annotations are available for model selection.
  By utilizing a newly introduced dataset with two spurious attributes, we demonstrate that EVaLS improves robustness to both known and unknown spurious attributes learned by an ERM-trained model better than methods relying on group information.
- 121 122 123

108

110

111

112

113

114

115

116

117

118

119

2 PRELIMINARIES

# 125 2.1 PROBLEM SETTING

126 We assume a general setting of a supervised learning problem with distinct data partitions  $\mathcal{D}^{Tr}$  for 127 training,  $\mathcal{D}^{Val}$  for validation, and  $\mathcal{D}^{Te}$  for final evaluation. Each dataset comprises a set of paired 128 samples (x, y), where  $x \in \mathcal{X}$  represents the data and  $y \in \mathcal{Y}$  denotes the corresponding labels. 129 Conventionally,  $\mathcal{D}^{Tr}$ ,  $\mathcal{D}^{Val}$ , and  $\mathcal{D}^{Te}$  are assumed to be uniformly sampled from the same distribution. 130 However, this idealized assumption does not hold in many real-world problems where distribution 131 shift is inevitable. In this context, we consider the subpopulation shift problem (Yang et al., 2023b). In a general form of this setting, it is assumed that data samples consist of different groups  $\mathcal{G}_i$ , where 132 each group comprises samples that share a property. More specifically, the overall data distribution 133  $p(x,y) = \sum_{i} \alpha_{i} p_{i}(x,y)$  is a composition of individual group distributions  $p_{i}(x,y)$  weighted by 134 their respective proportions  $\alpha_i$ , where  $\sum_i \alpha_i = 1$ . In this work, we assume that  $\mathcal{D}^{\text{Tr}}, \mathcal{D}^{\text{Val}}$ , and  $\mathcal{D}^{\text{Te}}$ 135 are composed of identical groups but with a different set of mixing coefficients  $\{\alpha_i\}$ . It is noteworthy 136 that the validation set may have approximately identical coefficients to those of the training or testing 137 sets, or it may have entirely different coefficients. 138

Several kinds of subpopulation shifts are defined in the literature, including class imbalance, at-139 tribute imbalance, and spurious correlation (Yang et al., 2023b). Class imbalance refers to the cases 140 where there is a difference between the proportion of samples from each class, while attribute imbal-141 ance occurs when instances with a certain attribute are underrepresented in the training data, even 142 though this attribute may not necessarily be a reliable predictor of the label. On the other hand, 143 spurious correlation occurs when various groups are differentiated by spurious attributes that are 144 partially predictive and correlated with class labels but are causally irrelevant. More precisely, we 145 can consider a set of spurious attributes S that partition the data into  $|S| \times |Y|$  groups. When the 146 concurrence of a spurious attribute with a label is significantly higher than its correlation with other 147 labels, that spurious attribute could become predictive of the label, resulting in deep models relying 148 on the spurious attributes as shortcuts instead of the core ones. This is followed by a decrease in the 149 model's performance on groups that do not have this attribute.

150 Given a class, the group containing samples with correlated spurious attributes is referred to as 151 *majority* group of that class, while the other groups are called the *minority* groups. As an example, 152 in the Waterbirds dataset (Sagawa et al., 2019), for which the task is to classify images of birds into 153 landbird and waterbird, there are spurious attributes {water background, land background}. Each 154 background is spuriously correlated with its associated label, decompose the data into two majority 155 groups waterbird on water background, and landbird on land background, and two minority groups waterbird on land background and landbird on water background. Our goal is to make the classifier 156 robust to spurious attributes by increasing performance for all groups. 157

- 158 159
- 2.2 ROBUSTNESS OF A TRAINED MODEL TO UNKNOWN SHORTCUTS
- 161 In scenarios where group annotations are absent, traditional methods that depend on these annotations for training or model selection become infeasible. Moreover, as previously discussed by Li



Figure 1: Overview of the proposed approach. (a) We randomly split the dataset  $\mathcal{D}$  into  $\mathcal{D}^{Tr}$ ,  $\mathcal{D}^{MS}$ ,  $\mathcal{D}^{LL}$  and  $\mathcal{D}^{Te}$ . We train the initial classifier on  $\mathcal{D}^{Tr}$  with empirical risk minimization (ERM). Alternatively, we can assume that an ERM-trained model is given. (b) An environment inference method is utilized to infer diverse environments for each class of  $\mathcal{D}^{MS}$ . (c) We evaluate  $\mathcal{D}^{LL}$  samples on the initial ERM classifier and sort high-loss and low-loss samples of each class for loss-based sampling. (d) Finally, we perform last-layer retraining on the loss-based selected samples  $\mathcal{D}^{Bal}$ . Each retraining setting (e.g. different k for loss-based sampling) is validated based on the worst accuracy of the inferred environments. Note that majority and minority groups are shown with dark and light colors for better visualization, but are not known in our setting.

186 187

et al. (2023), when data contains multiple spurious attributes and annotations are only available for 188 some of them, such methods would make the model robust only to the known spurious attributes. To 189 further explore such complex scenarios, we introduce the Dominoes Colored-MNIST-FashionMNIST 190 (Dominoes CMF) dataset (Figure 4(a)). Drawing inspiration from Pagliardini et al. (2022a) and Ar-191 jovsky et al. (2020), Dominoes CMF merges an image from CIFAR10 (Krizhevsky & Hinton, 2009) 192 at the top with a colored (red or green) MNIST (Deng, 2012) or FashionMNIST (Xiao et al., 2017) 193 image at the bottom. The primary label is derived from the CIFAR10 image, while the bottom part 194 introduces two independent spurious attributes: color (red or green) and style (MNIST or FashionM-195 NIST). Although annotations for shape are provided for training and model selection, color remains 196 an unknown variable until testing. For more details on the dataset refer to the Appendix.

The illustrations in Figure 2(a-c) depict the outlined scenario. A classifier trained using ERM is dependent on both spurious features (Figure 2(b)). Yet, achieving robustness against one spurious correlation (Figure 2(c)), does not ensure robustness against both (Figure 2(a)). In Section 4 we show that our approach, which does not rely on the group annotations of the identified group, achieves enhanced robustness to both spurious correlations, outperforming strategies that depend on the known group's information.

- 203
- 204 205 206

# 3 ENVIRONMENT-BASED VALIDATION AND LOSS-BASED SAMPLING

EVaLS is designed to improve the robustness of ERM-trained deep learning models to group shifts 207 without the need for group annotation. In line with the DFR (Kirichenko et al., 2023) approach, we 208 utilize a classifier defined as  $f = h_{\phi} \circ g_{\theta}$ , where  $g_{\theta}$  represents a deep neural network serving as a 209 feature extractor, and  $h_{\phi}$  denotes a linear classifier. The classifier is initially trained with the ERM 210 objective on the training dataset  $\mathcal{D}^{Tr}$ . Subsequently, we freeze the feature extractor  $g_{\theta}$  and focus 211 solely on retraining the last linear layer  $h_{\phi}$  using the validation dataset  $\mathcal{D}^{Val}$  as a held-out dataset. 212 This scheme helps us make our method available in settings where  $\mathcal{D}^{Tr}$  is not available, or where 213 repeating the training process is infeasible. 214

We randomly divide the validation set  $\mathcal{D}^{Val}$  into two subsets,  $\mathcal{D}^{LL}$  and  $\mathcal{D}^{MS}$  which are used for last layer training and model selection, respectively. In Section 3.1 we explain how to sample a



Figure 2: (a) If all spurious attributes in a dataset are known, they can be utilized to fit a classifier 229 that captures the essential attributes. (b) In the absence of knowledge about all spurious attributes, 230 the model would depend on them for classification, leading to incorrect classification of minority 231 samples. (c) If some spurious attribute is unknown (Spurious 2), the model becomes robust only to the known spurious correlations (Spurious 1), but it still underperforms on minority samples. 232

subset of  $\mathcal{D}^{LL}$  that statistically handles the group shifts inherent in the dataset. In Section 3.2 we 235 describe how  $\mathcal{D}^{MS}$  is divided into different environments that are later used for model selection. The 236 optimal number of selected samples from  $\mathcal{D}^{LL}$  and other hyperparameters is determined based on the worst environment accuracies among environments that are obtained from  $\mathcal{D}^{MS}$ . By combining our 238 sampling and validation strategy, we aim to provide a robust linear classifier  $h_{\phi^*}$  that significantly 239 improves the accuracy of underrepresented groups without requiring group annotations of training 240 or validation sets. Finally in Section 3.3, we provide theoretical support for the loss-based sampling procedure and its effectiveness. Figure 1 illustrates the comprehensive workflow of the EVaLS. 242

#### 3.1 LOSS-BASED INSTANCE SAMPLING

233 234

237

241

243

244

245 Following previous works (Liu et al., 2021a; Nam et al., 2020; Qiu et al., 2023), we use the loss 246 value as an indicator for identifying minority groups. We first evaluate classifier f on samples within  $\mathcal{D}^{LL}$  and choose k samples with the highest and lowest loss values in each class for a given 247 k. By combining these 2k samples from each class, we construct a balanced set  $\mathcal{D}^{\text{Bal}}$ , consisting of 248 high-loss and low-loss samples (see Figure 1(c)).  $\mathcal{D}^{Bal}$  is then used for the training of the last layer 249 of the model. As depicted in Figure 3, the proportion of minority samples among various percentiles 250 of samples with the highest loss values increases as we select a smaller subset of samples with the 251 highest loss. This suggests that high and low-loss samples could serve as effective representatives of minority and majority groups, respectively. In Section 3.3, we offer theoretical insights explaining 253 why this approach could lead to the creation of group-balanced data. 254

255 3.2 PARTITIONING VALIDATION SET INTO ENVIRONMENTS 256

257 Contrary to common assumptions and practices in the field, precise group labels for the validation 258 set are not essential for training models robust to spurious correlations. Our empirical findings, 259 detailed in Section 4, reveal that partitioning the validation set into environments that exhibit sig-260 nificant subpopulation shifts can be used for model selection. Under these conditions, the worst 261 environment accuracy (WEA) emerges as a viable metric for selecting the most effective model and hyperparameters. 262

263 The concept of an *environment*, as frequently discussed in the invariant learning literature, denotes 264 partitions of data that exhibit different distributions. A model that consistently excels across these 265 varied environments, achieving impressive worst environment accuracy (WEA), is likely to perform 266 equally well across different groups in the test set. Several methods for inferring environments with notable distribution shifts have been introduced (Creager et al., 2021; Liu et al., 2021b). Environ-267 ment Inference for Invariant Learning (EIIL) (Creager et al., 2021), leverages the predictions from 268 an earlier trained ERM model to divide the data into two distinct environments that significantly deviate from the invariant learning principle proposed by Arjovsky et al. (2020), thus creating en-



Figure 3: The percentage of samples with the highest (lowest) losses across various thresholds that belong to the minority (majority) group within different classes in  $\mathcal{D}^{LL}$  for (a) the Waterbirds and (b) CelebA datasets. Minority group samples are more prevalent among high-loss samples, while majority group samples dominate the low-loss areas. The error bars are calculated across three ERM models.<sup>1</sup>

vironments with distribution shifts. Initially, EIIL is employed to split  $\mathcal{D}^{MS}$  into two environments. Subsequently, each environment is further divided based on sample labels, resulting in  $2 \times |\mathcal{Y}|$  environments. To measure the difference between the distribution of environments, we define *group shift* of a class as the absolute difference in the proportion of a minority group between two environments. As detailed in the Appendix, environments inferred by EIIL demonstrate an average group shift of 28.7% over datasets with spurious correlation. Further information about EIIL and the group shift quantities for each dataset can be found in the Appendix.

We demonstrate that even more straightforward techniques, such as applying a random linear layer over the feature embedding space and distinguishing environments based on correctly and incorrectly classified samples of each class, can be effective to an extent in several cases (See Appendix F.3). It underscores that the feature space of a trained model is a valuable resource of information for identifying groups affected by spurious correlations. This supports the logic of previous research that employs clustering (Sohoni et al., 2020) or contrastive methods (Zhang et al., 2021) in this space to differentiate between groups.

300 301

302

284

3.3 THEORETICAL ANALYSIS

The environments obtained as described in Section 3.2 are utilized for hyperparameter tuning, specifically for tuning k, which is the number of selected samples from loss tails. It is known that minority samples are more prevalent among high-loss samples, while majority samples dominate the low-loss category. However, the question remains whether loss-based sampling can construct a balanced dataset without introducing spurious correlations. In this section, aligned with our practical approach, we provide theoretical insights into how loss-based sampling within a class can be used to create a group-balanced dataset.

Consider a binary classification problem with a cross-entropy loss function. Let logits be denoted as *L*. We assume a general assumption that in feature space (output of  $g_{\theta}$ ) samples from the minority and majority of a class are derived from Gaussian distributions. As a result, we can consider  $\mathcal{N}(\mu_{\min}, \sigma_{\min}^2)$  and  $\mathcal{N}(\mu_{\max}, \sigma_{\max}^2)$  as the distribution of minority and majority samples in logits space (See Lemma D.1 in Appendix D for details). Because the loss function is a monotonic function of logits, the tails of the distribution of loss across samples are equivalent to that of the logits in each class.

**Proposition 3.1.** [Feasiblity Of Loss-based Group Balancing] Suppose that L is derived from the mixture of two distributions  $\mathcal{N}(\mu_{min}, \sigma_{min}^2)$  and  $\mathcal{N}(\mu_{maj}, \sigma_{maj}^2)$  with proportion of  $\varepsilon$  and  $1 - \varepsilon$ , respectively, where  $\varepsilon \leq \frac{1}{2}$ . If (i)  $\sigma_{min} > \sigma_{maj}$ , or (ii) under sufficient and necessary conditions on  $\mu_{min}, \mu_{maj}, \sigma_{min}$  and  $\sigma_{maj}$  including inequality 1 (see App.D), there exists  $\alpha$  and  $\beta$  such that restricting L to the  $\alpha$ -left and  $\beta$ -right tails of its distribution results in a group-balanced distribution; in which both components are equally represented.

<sup>323</sup> 

<sup>&</sup>lt;sup>1</sup>Note that in the CelebA dataset, only the "blond hair" class includes a minority group.

 $\epsilon$ 

324

327

333

334

335 336

337

$$\geq \operatorname{sigmoid}\left(-\frac{\left(\mu_{\operatorname{maj}} - \mu_{\operatorname{min}}\right)^{2}}{2\left(\sigma_{\operatorname{maj}}^{2} - \sigma_{\operatorname{min}}^{2}\right)} - \log\left(\frac{\sigma_{\operatorname{maj}}}{\sigma_{\operatorname{min}}}\right)\right) \tag{1}$$

We provide an outline for proof of Proposition 3.1 here and leave the complete and formal proof and also exact bounds to Appendix D. We also analyze the conditions and effects of spurious correlation in satisfying these conditions. Practical justifications for Proposition 3.1 can be found in Appendix D.2. To proceed with the outline, we first define a key concept.

**Definition 3.1** (Proportional Density Difference). For any interval I = (a, b] and a mixture distribution  $\varepsilon P_1(x) + (1 - \varepsilon)P_2(x)$ , the proportional density difference is defined as the difference of accumulation of two component distributions in the interval I and is denoted by  $\Delta_{\varepsilon} P_{mixture}(I)$ .

$$\Delta_{\varepsilon} P_{mixture}(I) \stackrel{\Delta}{=} \varepsilon P_1(x \in I) - (1 - \varepsilon) P_2(x \in I)$$
<sup>(2)</sup>

338 Our proof proceeds with three steps. First, we reformulate the theorem as an **Proof outline** 339 equality of left- and right-tail proportional distribution differences. In other words, we show that the more mass the minority distribution has on one tail, the more mass the majority distribution must 340 have on the other tail. Afterward, supposing  $\mu_{min} < \mu_{maj}$  WOLG, we propose a proper range for 341  $\beta$  values on the right tail. We show that when  $\sigma_{mai} \leq \sigma_{min}$ , values for  $\alpha$  trivially exist that can 342 overcome the imbalance between the two distributions. In the last step, for the case in which the 343 variance of the majority is higher than the minority, we discuss a necessary and sufficient condition 344 for the existence of  $\alpha$  and  $\beta$  based on the left-tail proportional density difference using the properties 345 of its derivative with respect to  $\alpha$ . 346

Condition 1 suggests that for a given degree of spurious correlation  $\epsilon$  and variations  $\sigma_{maj}$ ,  $\sigma_{min}$ , an essential prerequisite for the efficacy of loss-based sampling is a sufficiently large disparity between the mean distributions of minority and majority samples, denoted by  $\|\mu_{maj} - \mu_{min}\|^2$ . This indicates that the groups should be distinctly separable in the logits space.

Although the parameters  $\alpha$  and  $\beta$  are theoretically established under certain conditions, their actual values remain undetermined. Therefore, validation data is essential to identify the appropriate tails. For practicality and simplicity, we assume an equal number k of samples for both tails and explore this count (high- and low-loss samples) from a predefined set of values. By leveraging the worst environment accuracy on validation data after last-layer retraining, as detailed in Section 3.2, we identify the optimal candidate that ensures uniform accuracy across all environments.

4 EXPERIMENTS

In this section, we evaluate the effectiveness of the proposed scheme through comprehensive exper iments on multiple datasets and compare it with various methods and baselines. We begin by briefly
 describing evaluation datasets and then introduce baselines and comparative methods. Finally, we
 report and fully explain the results.

364

357 358

359

365 Datasets Our approach, along with other baselines, is evaluated on Waterbirds (Sagawa et al., 2019), CelebA (Liu et al., 2014), UrbanCars (Li et al., 2023), CivilComments (Borkan et al., 2019), and MultiNLI (Williams et al., 2017). As per the study by Yang et al. (2023b), Waterbirds, CelebA, and UrbanCars among these datasets exhibit spurious correlation. Among the rest, CivilComments has class and attribute imbalance, whereas MultiNLI exhibits attribute imbalance. For additional details on the datasets, please refer to the Appendix E.3.

370

Baselines We compare EVaLS with six baselines in addition to standard ERM. Group-DRO (Sagawa et al., 2019) trains a model on the data with the objective of minimizing its average loss on the minority samples. This method requires group labels of both the training and validation sets. DFR (Kirichenko et al., 2023) argues that models trained with ERM are capable of extracting the core features of images. Thus, it first trains a model with ERM, and retrains only the last linear classifier layer on a group-balanced subset of the validation or the held-out training data. While DFR reduces the number of group-annotated samples, it still requires group labels in the training phase. GroupDRO + EIIL (Creager et al., 2021) infers environments of the training set and trains a model

378 with GroupDRO on the inferred environments. JTT (Liu et al., 2021a) first trains a model with ERM 379 on the dataset, and then retrains it on the dataset by upweighting the samples that were misclassified 380 by the initial ERM model. ES Disagreement SELF (LaBonte et al., 2023) selects samples with the 381 highest difference in output when comparing an ERM-trained model to its early-stopped version. 382 Then, they fine-tune the last layer of the ERM-trained model on the selected samples. AFR (Qiu et al., 2023) trains a model with standard ERM, and retrains the classifier on a weighted held-out 383 data. The weights assigned to retraining samples are determined by the probability that the ERM-384 pretrained model assigns to the ground-truth label, leading to an increased weighting of samples 385 from minority groups. 386

387 GroupDRO + EIIL, JTT, ES Disagreement SELF, and AFR eliminate the reliance on group annotations for their (re)training. However, unlike EVaLS, they all require group labels for model selection. 388 JTT, GroupDRO, and GroupDRO + EIIL necessitate training the entire model to apply their meth-389 ods. Additionally, ES Disagreement and SELF require early-stopped versions during training with 390 ERM. In contrast, DFR, AFR, and EVaLS operate in a completely post-training manner without 391 relying on any information from ERM training. This property makes these methods applicable in 392 real-world scenarios when training checkpoints or training data are unavailable, or when it is infea-393 sible to repeat the training due to reasons such as a large training set. 394

395

Setup Similar to all the works mentioned in Section 4, we use ResNet-50 (He et al., 2016) pretrained on ImageNet (Russakovsky et al., 2015) for image classification tasks. We used random crop and random horizontal flip as data augmentation, similar to Kirichenko et al. (2023). For a fair comparison with the baselines, we did not employ any data augmentation techniques in the process of retraining the last layer of the model. For the CivilComments and MultiNLI, we use pretrained BERT (Devlin et al., 2019) and crop sentences to 220 tokens length. In EvaLS, we use the implementation of EIIL by spuco package (Joshi et al., 2023) for environments inference on the model selection set with 20000 steps, SGD optimizer, and learning rate 10<sup>-2</sup> for all datasets.

404 Model selection and hyper-parameter fine-tuning are done according to the worst environment (or 405 group if annotations are assumed to be available) accuracy on the validation set. For each dataset, we 406 assess the performance of our model in two cases: fine-tuning the ERM classifier or retraining it. For 407 all datasets except MultiNLI and Urbancars, retraining yielded better validation results. We report 408 the results of our experiments in two settings: (i) EVaLS, which incorporates loss-based instance 409 sampling for training the last layer, and environment inference for model selection. (ii) EVaLS-GL, 410 similar to EVaLS except in using ground-truth group labels for model selection. For more details on 411 the ERM training and last layer re-training hyperparameters refer to the Appendix.

- 412
- 413 414

# 4.1 Results

415 416

The results of our experiments along with the reported results on GroupDRO (Sagawa et al., 2019), 417 DFR (Kirichenko et al., 2023), JTT (Liu et al., 2021a), ES Disagreement SELF (LaBonte et al., 418 2023), and AFR (Qiu et al., 2023) on five datasets are shown in Table 1. The reported results for 419 GroupDRO, DFR, JTT, and AFR except those for the UrbanCars are taken from Qiu et al. (2023). 420 For EIIL+Group DRO, the results for Waterbirds, CelebA, and CivilComments are reported from 421 Zhang et al. (2021). The results of SELF on CelebA and MultiNLI are reported from the original 422 paper (LaBonte et al., 2023). We report only the worst group accuracy of methods in Table 1. The average group accuracies are documented in the Appendix. The Group Info column shows whether 423 group annotation is required for training or model selection entry for each method. Methods that do 424 not require information regarding ERM training (such as training data or checkpoints) are identified 425 with a star in the table. 426

Overall, our approaches outperform methods that do not require group annotations for (re)training
in 2 out of 3 datasets with spurious correlations. Moreover, EVaLS-GL surpasses other methods
with a similar level of group supervision on MultiNLI (Williams et al., 2017) and achieves state-ofthe-art performance among all methods on UrbanCars (Li et al., 2023). Furthermore, EVaLS and
EVaLS-GL, similar to DFR (Kirichenko et al., 2023) and AFR (Qiu et al., 2023), can be applied to
ERM-trained models without needing further information about their training.

432 Table 1: Comparison of worst group accuracy across various methods, including ours, on five 433 datasets. The Group Info column indicates if each method utilizes group labels of the train-434 ing/validation data, with  $\checkmark$  denoting that group information is employed during both stages. Bold numbers are the highest results overall, while underlined ones are the best among methods that may 435 require group annotation only for model selection. CivilComments is class imbalanced, MultiNLI 436 has imbalanced attributes, and the other three datasets have spurious correlations. The  $\times$  sign indi-437 cates that the dataset is out of the scope of the method. Methods that do not rely on ERM training 438 information are identified with \*. Mean and standard deviation are calculated over three runs. 439

Method	Group Info			Datasets	8	
	Train/Val	Waterbirds	CelebA	UrbanCars	CivilComments	MultiNLI
GDRO (Sagawa et al., 2019) DFR* (Kirichenko et al., 2023)	√  √ X √√	91.4 92.9 <sub>±0.2</sub>	88.9 $88.3_{\pm 1.1}$	73.1 $79.6_{\pm 2.2}$	$\begin{array}{c} 69.9 \\ \textbf{70.1}_{\pm 0.8} \end{array}$	$\begin{array}{c} \textbf{77.7} \\ \textbf{74.7}_{\pm 0.7} \end{array}$
GDRO + EIIL (Creager et al., 2021) JTT (Liu et al., 2021a) SELF (LaBonte et al., 2023) AFR* (Qiu et al., 2023) EVaLS-GL* (Ours)	X √ X √ X √ X √ X √	$\begin{array}{c} 77.2_{\pm 1} \\ 86.7 \\ \underline{91.6_{\pm 1.4}} \\ 90.4_{\pm 1.1} \\ 89.4_{\pm 0.3} \end{array}$	$\begin{array}{c} 81.7_{\pm 0.8}\\ 81.1\\ 83.9_{\pm 0.9}\\ 82.0_{\pm 0.5}\\ 84.6_{\pm 1.6}\end{array}$	$\begin{array}{c} 76.5_{\pm 2.6} \\ 79.5 \\ 83.2_{\pm 0.8} \\ 80.2_{\pm 2.0} \\ \textbf{83.5}_{\pm 1.7} \end{array}$	$ \begin{vmatrix} 67.0_{\pm 2.4} \\ \underline{69.3} \\ 66.0_{\pm 1.7} \\ 68.7_{\pm 0.6} \\ 68.0_{\pm 0.5} \end{vmatrix} $	$\begin{array}{c} 61.2_{\pm 0.5} \\ 72.6 \\ 70.7_{\pm 2.5} \\ 73.4_{\pm 0.6} \\ 75.1_{\pm 1.2} \end{array}$
EVaLS* (Ours) ERM	X/X X/X	$\begin{array}{c} 88.4_{\pm 3.1} \\ 66.4_{\pm 2.3} \end{array}$	$\frac{85.3_{\pm 0.4}}{47.4_{\pm 2.3}}$	$\begin{array}{c} 82.1_{\pm 0.9} \\ 18.67_{\pm 2.0} \end{array}$	$\begin{array}{c} \times \\ 61.2_{\pm 3.6} \end{array}$	$\times 64.8_{\pm 1.9}$

449 450 451

The comparison between EVaLS and GroupDRO + EIIL indicates that when environments are avail able instead of groups, our method, which uses environments solely for model selection and utilizes
 loss-based sampling, is more effective than GroupDRO, a potent invariant learning method.

Regarding the UrbanCars, which contains an un-annotated spurious attribute, Li et al. (2023) has
shown that shortcut mitigation methods often struggle to address multiple shortcuts simultaneously.
Notably, techniques such as DFR (Kirichenko et al., 2023) and GDRO (Sagawa et al., 2019) which
are designed to reduce reliance on a specific shortcut feature, fail to make the model robust to
unknown shortcuts effectively. In contrast, our experiments suggest that annotation-free methods
can mitigate the impact of both labeled and unlabeled shortcut features more effectively.

461 Our evaluation of EVaLS is based on the spurious correlation benchmarks. This is because, in other 462 instances of subpopulation shift, the attributes that differ across groups are not predictive of the label, thereby reducing the visibility of these attributes' effects in the model's final layers (Lee et al., 463 2023). Consequently, EIIL, which depends on output logits for prediction, might not effectively 464 separate the groups. This observation is further supported by our findings related to the degree of 465 group shift between the environments inferred by EIIL for each class in the CivilComments and 466 MultiNLI datasets. The average group shift (defined in the Section 3.2) in the environments of the 467 minority class of CivilComments is only  $0.8_{\pm 0.0}$ %. Also, environments associated with Classes 1 468 and 2 in MultiNLI show only  $1.1_{\pm 0.3}\%$  and  $1.9_{\pm 1.0}\%$  group shift respectively. More results and 469 ablation studies can be found in the Appendix. 470

471 **Mitigating Multiple Spurious Attributes** To evaluate the performance of our method in the case 472 of unknown spurious correlations, we train a ResNet-18 He et al. (2016) model on the Dominoes-473 CMF dataset. We apply DFR Kirichenko et al. (2023), EVaLS-GL, and EVaLS on top of the trained 474 ERMs to assess their ability to mitigate multiple shortcuts. We consider the style (MNIST/Fashion-475 MNIST) feature as the known group label, and the color as the unknown spurious attribute. We set 476 the spurious correlation of the known attribute to 75% and conduct experiments for various amount 477 of unknown spurious correlation. During model selection, we calculate the worst-group accuracy on the validation set considering only the label of the known shortcut, *i.e.*, the lowest accuracy among 478 the four groups based on the combination of the target label and the single known shortcut label. 479 However, the final results on test data are based on the worst group accuracies, taking into account 480 groups defined by the labels of both spurious attributes. The results are shown in Figure 4(b). Note 481 that EVaLS operates without using annotations for either the known or unknown spurious attributes. 482

Our results confirm findings by Li et al. (2023), suggesting that methods using group labels mit igate reliance on the known shortcut but not necessarily on the unknown one. DFR (Kirichenko
 et al., 2023) experiences a significant drop in performance (34.55% under 95% color spurious correlation) when it relies on a single known spurious attribute for grouping, compared to the oracle



Figure 4: (a) The Dominoes-CMF dataset, which contains two spurious attributes. (b) Performance 499 on Dominoes-CMF is measured by worst-group accuracy across varying levels of correlation be-500 tween the target label and the unknown spurious attribute (color). Lower reliance on available group annotations (based on known spurious attributes, i.e., style) results in higher robustness to both attributes. The performance gap between EVaLS and EVaLS-GL with lower group supervi-502 sion compared to DFR (Kirichenko et al., 2023) increases with higher correlations. The oracle uses 503 DFR (Kirichenko et al., 2023) with complete group information regarding both attributes. 504

that uses both attributes for grouping. EVaLS-GL reduces this issue using its loss-based sampling 506 approach, but surprisingly EVaLS even outperforms EVaLS-GL. Combining a loss-based sampling 507 approach for last layer training and environment-based model selection, results in a completely 508 group-annotation-free method in a multi-shortcut setting with unknown spurious correlations, and 509 successfully re-weights features to perform well with respect to multiple spurious attributes. It is 510 also evident that increasing unknown spurious correlation results in a larger gap between the perfor-511 mance of EVaLS and EVaLS-GL compared to DFR (Kirichenko et al., 2023). 512

5 DISCUSSION

514 515

513

501

505

This study presents EVaLS, a novel approach to improve robustness to spurious correlations with 516 zero group annotation. EVaLS uses loss-based sampling to create a balanced training dataset that ef-517 fectively disrupts spurious correlations and employs EIIL to infer environments for model selection. 518 We also explore situations with multiple spurious correlations, some of which are unknown. In this 519 context, we introduce Dominoes-CMF, a dataset in which two factors are spuriously correlated with 520 the label, but only one is identified. Our findings suggest that EVaLS attains near-optimal worst test 521 group accuracy on spurious correlation datasets. We also present EVaLS-GL, which needs group 522 labels only for model selection. Our empirical tests on various datasets demonstrate that EVaLS-GL 523 outperforms state-of-the-art methods requiring group labels during evaluation or training.

524 Note that this paper remains consistent with the findings of Lin et al. (2022). Our approach does 525 not involve identifying spurious attributes without auxiliary information. Instead, the objective is 526 to make a trained model robust against its reliance on shortcuts. Specifically, conditioning on what 527 a trained model learns, we ascertain that both the loss value and the model's feature space are 528 instrumental in mitigating shortcuts.

529 EVaLS and EVaLS-GL may struggle with small datasets due to a low number of selected samples 530 for the last layer training. Also, as environment inference from the last layer features is not effective 531 for all types of subpopulation shifts, EVaLS is limited to datasets with spurious correlation. Similar 532 to other methods in the field, EVaLS prioritizes the worst group accuracy at the cost of less average 533 accuracy. Additionally, a notable variance has been observed in some of our experiments.

534 EVaLS represents a significant advancement in the development of methods for enhancing model fairness and robustness without prior knowledge about group annotations. EVaLS could be simply 536 applied as a plug-and-play solution on various ERM-pretrained models with unknown inherent bi-537 ases to make them robust to possible spurious correlations. Future work could explore developing 538 environment inference methods effective for other types of subpopulation shift, such as attribute and class imbalance.

# 540 REFERENCES

559

561

576

577

578

579 580

581

582

- Faruk Ahmed, Yoshua Bengio, Harm van Seijen, and Aaron Courville. Systematic generalisation
   with group invariant predictions. In *International Conference on Learning Representations*, 2021.
   URL https://openreview.net/forum?id=b9PoimzZFJ.
- Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization, 2020.
- Sara Beery, Grant Van Horn, and Pietro Perona. Recognition in terra incognita. In *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018.
- Daniel Borkan, Lucas Dixon, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. Nuanced metrics for measuring unintended bias with real data for text classification. In *Companion proceedings of the 2019 world wide web conference*, pp. 491–500, 2019.
- Elliot Creager, Joern-Henrik Jacobsen, and Richard Zemel. Environment inference for invariant
   learning. In Marina Meila and Tong Zhang (eds.), Proceedings of the 38th International Con *ference on Machine Learning*, volume 139 of Proceedings of Machine Learning Research, pp.
   2189–2200. PMLR, 18–24 Jul 2021. URL https://proceedings.mlr.press/v139/
   creager21a.html.
  - Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE* Signal Processing Magazine, 29(6):141–142, 2012.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Thamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL https://aclanthology.org/ N19-1423.
- Tatsunori Hashimoto, Megha Srivastava, Hongseok Namkoong, and Percy Liang. Fairness without
   demographics in repeated loss minimization. In *International Conference on Machine Learning*,
   pp. 1929–1938. PMLR, 2018.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 770–778, 2016. doi: 10.1109/CVPR.2016.90.
  - Siddharth Joshi, Yu Yang, Yihao Xue, Wenhan Yang, and Baharan Mirzasoleiman. Towards mitigating spurious correlations in the wild: A benchmark & a more realistic dataset. ArXiv, abs/2306.11957, 2023. URL https://api.semanticscholar.org/ CorpusID:259211935.
    - Polina Kirichenko, Pavel Izmailov, and Andrew Gordon Wilson. Last layer re-training is sufficient for robustness to spurious correlations. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=Zb6c8A-Fghk.
- Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanas Phillips, Irena Gao, Tony Lee, Etienne David, Ian Stavness, Wei Guo, Berton Earnshaw, Imran Haque, Sara M Beery, Jure Leskovec, Anshul Kundaje, Emma Pierson, Sergey Levine, Chelsea Finn, and Percy Liang. Wilds: A benchmark of in-the-wild distribution shifts. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 5637–5664. PMLR, 18–24 Jul 2021. URL https://proceedings.mlr.press/v139/koh21a.html.
- Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images.
   Technical Report 0, University of Toronto, Toronto, Ontario, 2009. URL https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf.

- David Krueger, Ethan Caballero, Joern-Henrik Jacobsen, Amy Zhang, Jonathan Binas, Dinghuai Zhang, Remi Le Priol, and Aaron Courville. Out-of-distribution generalization via risk extrapolation (rex). In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 5815–5826. PMLR, 18–24 Jul 2021. URL https://proceedings.mlr.press/v139/krueger21a.html.
- Tyler LaBonte, Vidya Muthukumar, and Abhishek Kumar. Towards last-layer retraining for group robustness with fewer annotations. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL https://openreview.net/forum?id=kshC3NOP6h.
- Yoonho Lee, Annie S Chen, Fahim Tajwar, Ananya Kumar, Huaxiu Yao, Percy Liang, and Chelsea
   Finn. Surgical fine-tuning improves adaptation to distribution shifts. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?
   id=APuPRxjHvZ.
- <sup>608</sup>
   <sup>609</sup> Zhiheng Li, Ivan Evtimov, Albert Gordo, Caner Hazirbas, Tal Hassner, Cristian Canton Ferrer, Chenliang Xu, and Mark Ibrahim. A whac-a-mole dilemma: Shortcuts come in multiples where mitigating one amplifies others. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 20071–20082, June 2023.
- Yong Lin, Shengyu Zhu, Lu Tan, and Peng Cui. Zin: When and how to learn invariance without environment partition? In S. Koyejo, S. Mohamed, A. Agarwal,
  D. Belgrave, K. Cho, and A. Oh (eds.), Advances in Neural Information Processing Systems, volume 35, pp. 24529–24542. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper\_files/paper/2022/file/ 9b77f07301b1ef1fe810aae96c12cb7b-Paper-Conference.pdf.
- Evan Z Liu, Behzad Haghgoo, Annie S Chen, Aditi Raghunathan, Pang Wei Koh, Shiori Sagawa, Percy Liang, and Chelsea Finn. Just train twice: Improving group robustness without training group information. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp.
  623 6781–6792. PMLR, 18–24 Jul 2021a. URL https://proceedings.mlr.press/v139/ liu21f.html.
- Jiashuo Liu, Zheyuan Hu, Peng Cui, Bo Li, and Zheyan Shen. Heterogeneous risk minimization. In Marina Meila and Tong Zhang (eds.), Proceedings of the 38th International Conference on Machine Learning, volume 139 of Proceedings of Machine Learning Research, pp. 6804–6814.
  PMLR, 18–24 Jul 2021b. URL https://proceedings.mlr.press/v139/liu21h. html.
- Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild.
   2015 IEEE International Conference on Computer Vision (ICCV), pp. 3730–3738, 2014. URL
   https://api.semanticscholar.org/CorpusID:459456.

634

635

636

- Yuzhen Mao, Zhun Deng, Huaxiu Yao, Ting Ye, Kenji Kawaguchi, and James Zou. Last-layer fairness fine-tuning is simple and effective for neural networks. *arXiv preprint arXiv:2304.03935*, 2023.
- Junhyun Nam, Hyuntak Cha, Sungsoo Ahn, Jaeho Lee, and Jinwoo Shin. Learning from failure:
   De-biasing classifier from biased classifier. *Advances in Neural Information Processing Systems*, 33:20673–20684, 2020.
- Junhyun Nam, Jaehyung Kim, Jaeho Lee, and Jinwoo Shin. Spread spurious attribute: Improving worst-group accuracy with spurious attribute estimation. In *International Conference on Learning Representations*, 2021.
- Fahimeh Hosseini Noohdani, Parsa Hosseini, Aryan Yazdan Parast, HamidReza Yaghoubi Araghi, and Mahdieh Soleymani Baghshah. Decompose-and-compose: A compositional approach to mitigating spurious correlation. *CoRR*, abs/2402.18919, 2024. doi: 10.48550/ARXIV.2402.18919. URL https://doi.org/10.48550/arXiv.2402.18919.

- 648 Matteo Pagliardini, Martin Jaggi, François Fleuret, and Sai Praneeth Karimireddy. Agree to dis-649 agree: Diversity through disagreement for better transferability. In The Eleventh International 650 Conference on Learning Representations, 2022a. 651
- Matteo Pagliardini, Martin Jaggi, François Fleuret, and Sai Praneeth Karimireddy. Agree 652 to disagree: Diversity through disagreement for better transferability. arXiv preprint. 653 arXiv:2202.04414, 2022b. 654
- 655 Shikai Qiu, Andres Potapczynski, Pavel Izmailov, and Andrew Gordon Wilson. Simple and fast 656 group robustness by automatic feature reweighting. In International Conference on Machine 657 Learning, pp. 28448–28467. PMLR, 2023.
- 658 Alexandre Rame, Corentin Dancette, and Matthieu Cord. Fishr: Invariant gradient variances for 659 out-of-distribution generalization. In International Conference on Machine Learning, pp. 18347-660 18377. PMLR, 2022. 661
- 662 Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual 663 recognition challenge. International journal of computer vision, 115:211–252, 2015. 664
- 665 Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust 666 neural networks. In International Conference on Learning Representations, 2019. 667
- Seonguk Seo, Joon-Young Lee, and Bohyung Han. Information-theoretic bias reduction via causal 668 view of spurious correlation. In Proceedings of the AAAI Conference on Artificial Intelligence, 669 volume 36, pp. 2180-2188, 2022. 670
- 671 Harshay Shah, Kaustav Tamuly, Aditi Raghunathan, Prateek Jain, and Praneeth Netrapalli. The 672 pitfalls of simplicity bias in neural networks. Advances in Neural Information Processing Systems, 673 33:9573-9585, 2020.
- 674 Zheyan Shen, Jiashuo Liu, Yue He, Xingxuan Zhang, Renzhe Xu, Han Yu, and Peng Cui. Towards 675 out-of-distribution generalization: A survey. ArXiv, abs/2108.13624, 2021. URL https:// 676 api.semanticscholar.org/CorpusID:237364121. 677
- 678 Nimit Sohoni, Jared Dunnmon, Geoffrey Angus, Albert Gu, and Christopher Ré. No subclass left 679 behind: Fine-grained robustness in coarse-grained classification problems. Advances in Neural Information Processing Systems, 33:19339–19352, 2020. 680
- 681 Adina Williams, Nikita Nangia, and Samuel R Bowman. A broad-coverage challenge corpus for 682 sentence understanding through inference. arXiv preprint arXiv:1704.05426, 2017. 683
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset 685 for benchmarking machine learning algorithms, 2017. URL http://arxiv.org/ abs/1708.07747. cite arxiv:1708.07747Comment: Dataset is freely available at 686 https://github.com/zalandoresearch/fashion-mnist Benchmark is available at http://fashionmnist.s3-website.eu-central-1.amazonaws.com/. 688
- 689 Yu Yang, Eric Gan, Gintare Karolina Dziugaite, and Baharan Mirzasoleiman. Identifying spurious 690 biases early in training through the lens of simplicity bias. ArXiv, abs/2305.18761, 2023a. URL 691 https://api.semanticscholar.org/CorpusID:258967752. 692
- Yuzhe Yang, Haoran Zhang, Dina Katabi, and Marzyeh Ghassemi. Change is hard: A closer look at 693 subpopulation shift. In International Conference on Machine Learning, 2023b. 694
- Michael Zhang, Nimit Sharad Sohoni, Hongyang R. Zhang, Chelsea Finn, and Christopher Ré. 696 Correct-n-contrast: A contrastive approach for improving robustness to spurious correlations. In 697 NeurIPS 2021 Workshop on Distribution Shifts: Connecting Methods and Applications, 2021. URL https://openreview.net/forum?id=Q41kl\_DwS3Y.
- 699

687

# A RELATED WORK

703 704

Robustness to spurious correlation is a critical concern across various machine learning subfields. It
 is a form of out-of-distribution generalization (Shen et al., 2021) where the distribution shift arises
 from the disproportionate representation of minority groups—those instances that are devoid of the
 correlated spurious patterns associated with their labels (Yang et al., 2023b). The issue of spurious
 correlation also intersects with the discourse on fairness in machine learning (Seo et al., 2022; Mao
 et al., 2023).

Past studies have proposed a range of strategies to mitigate the models' reliance on spurious correlation. Broadly speaking, these methods can be categorized according to the degree of supervision they require regarding group labels.

Invariant learning (IL) methods (Arjovsky et al., 2020; Krueger et al., 2021; Rame et al., 2022) 714 operate under the assumption of having access to a collection of environments that comprise group 715 shift. By imposing invariant conditions on these environments, IL methods strive to create classifiers 716 robust against group-sensitive features. IRM (Arjovsky et al., 2020) is designed to learn a feature 717 extractor, which, when utilized, guarantees the existence of a classifier that would be optimal in all 718 training environments. VREx (Krueger et al., 2021) aims to decrease the risk variance among dif-719 ferent training environments. PGI (Ahmed et al., 2021) works by minimizing the distance between 720 the expected softmax distribution of labels, conditioned on inputs across both majority and minority 721 environments. Lastly, Fishr (Rame et al., 2022) focuses on bringing the variance of risk gradients 722 closer together across different training environments. For scenarios that the environments are not 723 available, environment inference methods (Creager et al., 2021; Liu et al., 2021b) are used to obtain 724 a set of environments. Creager et al. (2021) introduce environment inference for invariant learning (EIIL), which tries to partition samples into two groups such that the objective of IRM (Arjovsky 725 et al., 2020) is maximized. HRM (Liu et al., 2021b) aims to optimize both an environment inference 726 module and an invariant prediction module jointly, with the goal of achieving an invariant predictor. 727

728 When group annotations are accessible, various methods leverage this information to equalize the 729 impact of different groups on the model's loss. The Group Distributionally Robust Optimization 730 (GDRO) approach (Sagawa et al., 2019), for instance, focuses on optimizing the loss for the worst-731 performing group during training. Kirichenko et al. (2023) has shown that models can still learn and extract core data features even in the presence high spurious correlation. Consequently, They 732 suggest that retraining just the last layer of a model initially trained with Empirical Risk Mini-733 mization (ERM) can effectively reduce reliance on spurious correlation for predicting class labels. 734 This method, termed Deep Feature Re-weighting (DFR), has been validated as not only highly ef-735 fective but also significantly more efficient than earlier techniques that necessitated retraining the 736 full model (Nam et al., 2021; Sagawa et al., 2019). However, availability of group annotations is 737 considered a serious restrictive assumption. 738

Several recent studies have endeavored to enhance model robustness against spurious correlation, 739 even in the absence of group annotations (Liu et al., 2021a; Zhang et al., 2021; Qiu et al., 2023; 740 LaBonte et al., 2023; Yang et al., 2023a). Liu et al. (2021a) introduce a two-stage method that 741 involves training a model using ERM for a number of epochs before retraining it to give more 742 weight to misclassified samples. The study by Zhang et al. (2021) employs the same two-stage 743 training process, but with a twist for the second stage: they utilize contrastive methods. The goal 744 is to bring samples from the same class but with divergent predictions closer in the feature space, 745 while simultaneously increasing the separation between samples from different classes that have 746 similar predictions. Another method, known as automatic feature reweighting (AFR) (Qiu et al., 2023), reweights the last layer of an ERM-pretrained model to favor samples that the original model 747 was less accurate on. LaBonte et al. (2023) refine the last layer of an ERM-trained model through 748 class-balanced finetuning, identifying challenging data points by comparing the classifier's predic-749 tions with those of an early-stopped version. While these methods have significantly reduced the 750 reliance on group annotations, they still required for validation and model selection. This remains a 751 constraint, particularly when the spurious correlation is completely unknown. 752

To make a trained model robust to subpopulation shifts with zero group annotations, LaBonte et al.
(2023) have recently demonstrated that class-balanced retraining of a model pretrained with ERM
can effectively improve the worst-group accuracy (WGA) for certain datasets. While this method effectively reduces the impact of class imbalance, it fails in datasets with spurious correlations.

Table 2: The average and variation percentage (%)(across 3 seeds) of group shift between the inferred environments using EIIL (Creager et al., 2021) for each class, which is the absolute difference
between the proportion of a minority group in the two environments of a class. Higher group shift
indicates better separation of environments. In most cases, a significant group shift is observed between the inferred environments.

Class No.		Dataset						
01000 1101	Waterbirds	CelebA	UrbanCars					
0	$16.6_{\pm 0.7}$	$3.6_{\pm 0.2}$	$17.7_{\pm 1.2}, 23.5_{\pm 0.1}, 62.1_{\pm 1.9}$					
1	$50.5_{\pm 0.3}$	$14.1_{\pm 0.9}$	$40.7_{\pm 7.9}, 13.8_{\pm 0.1}, 19.2_{\pm 3.9}$					

### 

# **B** ENVIRONMENT INFERENCE FOR INVARIANT LEARNING

Consider the training dataset  $\mathcal{D}^{\mathrm{Tr}} = \{(x^{(i)}, y^{(i)}) | x^{(i)} \in \mathcal{X}, y^{(i)} \in \mathcal{Y}\}$ , where  $\mathcal{X}$  and  $\mathcal{Y}$  represent the input and output spaces, respectively. This dataset can be partitioned into different environments  $\mathcal{E}^{tr} = \{e_1, \dots, e_n\}$ , such that for any  $i \neq j$ , the data distribution in  $e_i$  and  $e_j$  differs. The objective of invariant learning is to train a predictor that performs consistently across all environments in  $\mathcal{E}^{tr}$ . Under certain conditions, this predictor is also expected to perform well on  $e^{tst}$ , a test environment with a distribution distinct from the training data. Invariant Risk Minimization (IRM) (Arjovsky et al., 2020) approaches this problem by learning a feature extractor  $\Phi(.)$  such that a classifier  $\omega(.)$ exists, where  $\omega \circ \Phi(.)$  performs consistently across all training environments. The practical imple-mentation of the IRM objective is to minimize 

$$\sum_{e \in \mathcal{E}^{tr}} R^e(\Phi) + \lambda ||\nabla_{\bar{\omega}} R^e(\bar{\omega} \circ \Phi)||^2,$$
(3)

where  $\bar{\omega}$  is a constant scalar with a value of 1.0,  $\lambda$  is a hyperparameter, and  $R^e(f) = \mathbb{E}_{(x,y)\sim p_e}[l(f(x),y)]$  is referred to as the risk on environment e.

In real-world scenarios, training environments might not always be available. To address this, Environment Inference for Invariant Learning (EIIL) (Creager et al., 2021) partitions samples into two environments in a way that maximizes the objective in Eq 3.

During the training phase, the EIIL algorithm replaces the hard assignment of environments to samples with a soft assignment  $\mathbf{q}_i(e) = p(e|(x^{(i)}, y^{(i)}))$ , where  $\mathbf{q}_i$  is learnable. Consequently, the relaxed version of the risk function is defined as  $\tilde{R}^e(\Phi) = \frac{1}{N} \sum_{i}^{N} \mathbf{q}_i(e)[l(\Phi(x^{(i)}), y^{(i)})]$ . Given a model  $\Phi$  that has been trained with ERM on the dataset, EIIL optimizes

$$\mathbf{q}^* = \operatorname*{arg\,max}_{\mathbf{q}} ||\nabla_{\bar{\omega}} \tilde{R}^e(\bar{\omega} \circ \Phi)||. \tag{4}$$

As discussed in Creager et al. (2021), using a biased base model  $\Phi$  could lead to environments exhibiting varying degrees of spurious correlation. During the inference phase, the soft assignment is converted to a hard assignment. The average group shift between the inferred environments using EIIL is illustrated in Table 2.

# <sup>810</sup> C ALGORITHM

## Algorithm 1 EVaLS

1: Input: Held-out dataset  $\mathcal{D}^{Val}$ , ERM-trained model  $f_{ERM}$ , maximum k value  $k_{max}$ 814 2: **Output:** Optimal number of samples  $k^*$ , best model  $f^*$ , best performance wea\* 815 3:  $(\mathcal{D}^{LL}, \mathcal{D}^{MS}) \leftarrow \text{splitDataset}(\mathcal{D}^{Val})$ 816 ▷ Split the held-out dataset 4: Envs[y]  $\leftarrow$  inferEnvs( $\mathcal{D}^{MS}$ )[y]  $\forall y \in \mathcal{Y}$  $\triangleright$  Infer environments from  $\mathcal{D}^{MS}$ 817 5: sortedSamples[y]  $\leftarrow$  sortByLoss( $f_{\text{ERM}}, \mathcal{D}^{\text{LL}}[y]$ )  $\forall y \in \mathcal{Y}$  $\triangleright$  Sort  $\mathcal{D}^{LL}$  samples by their loss 818 6: Initialize wea<sup>\*</sup>  $\leftarrow 0, k^* \leftarrow 0, f^* \leftarrow None$ 819 7: for k = 1 to  $k_{\text{max}}$  do 820 highLossSamples[y]  $\leftarrow$  sortedSamples[y][: k]  $\forall y \in \mathcal{Y} \triangleright$  Select top-k high-loss samples 8: 821 lowLossSamples[y]  $\leftarrow$  sortedSamples[y][-k :]  $\forall y \in \mathcal{Y} \triangleright$  Select top-k low-loss samples 9: 822  $\mathcal{D}^{\text{Bal}} \leftarrow \{\text{highLossSamples}, \text{lowLossSamples}\}$ 10:  $\triangleright$  Combine samples 823  $f \leftarrow \text{retrainLastLayer}(\mathcal{D}^{\text{Bal}})$ ▷ Retrain the last layer with combined samples 11: 824 wea  $\leftarrow$  evaluateWEA(f, Envs) ▷ Evaluate the retrained model 12: 825 13: if wea > wea\* then wea\*  $\leftarrow$  wea,  $f^* \leftarrow f, k^* \leftarrow k$ 826 14: ▷ Record the best configuration 827 15: end if 16: end for 828 17: **Return:**  $k^*$ , wea<sup>\*</sup>,  $f^*$ 829 830

831 832

833 834

835

836

837

839

812

813

# D THEORETICAL ANALYSIS

In this section, we establish a more formal description of loss-based sampling for balanced dataset creation and then prove it. We thoroughly analyze the close relationship between the availability of the balanced dataset and the gap between spurious features of minority and majority groups.

# 838 D.1 FEASIBILITY OF LOSS-BASED GROUP BALANCING

Consider a binary classification problem with a cross-entropy loss function. Let logits be denoted as *L*. Because loss is a monotonic function of logits, the tails of the distribution of loss across samples are equivalent to that of the logits in each class. We assume that in feature space (output of  $g_{\theta}$ ) samples from the minority and majority of a class are derived from Gaussian distributions  $\mathcal{N}(h_{\min}, \Sigma_{\min})$  and  $\mathcal{N}(h_{\max j}, \Sigma_{\max j})$ , respectively. Before diving into the group balance problem we initially show that the distribution of minority and majority samples in the logit space (output of  $h_{\phi}$ ) are Gaussian too.

Lemma D.1. [Gaussain Distribution of Logits] Considering a Gaussian distribution  $Z \sim \mathcal{N}(h, \Sigma)$ in feature space and  $W \in \mathbb{R}^d$ , then the distribution of logits is as follows:  $L = \langle W, Z \rangle \sim \mathcal{N}(Wh, ||W||_{\Sigma}^2)$ .

850 851 Proof. Let  $Z \sim \mathcal{N}(h, \Sigma)$ .

Consider  $L = \langle W, Z \rangle = W^T Z$ , where  $W \in \mathbb{R}^d$ . L is a linear combination of jointly gaussian random variables which makes it an univariate gaussian random variable.

To find the distribution of L, we need to determine its mean and variance.

1. Mean of L

$$\mathbb{E}[L] = \mathbb{E}[\langle W, Z \rangle] = \mathbb{E}[W^T Z] = W^T \mathbb{E}[Z] = W^T h = \langle W, h \rangle.$$

860 Therefore, the mean of L is Wh.

# 8612. Variance of *L*:

863 The variance of L can be computed using the properties of covariance. Recall that if  $Z \sim \mathcal{N}(h, \Sigma)$ , then the covariance matrix of Z is  $\Sigma$ .

The variance of the linear combination  $L = W^T Z$  is given by: 

$$\operatorname{Var}(L) = \operatorname{Var}(W^T Z) = W^T \Sigma W = ||W||_{\Sigma}^2,$$

where  $||W||_{\Sigma}$  denotes the Mahalanobis norm of W.

Thus, we have proved that if  $Z \sim \mathcal{N}(h, \Sigma)$ , then the logits  $L = \langle W, Z \rangle$  follow the distribution  $\mathcal{N}(Wh, ||W||_{\Sigma}^2)$ .

From now on, we consider  $\mathcal{N}(\mu_{\min}, \sigma_{\min}^2)$  and  $\mathcal{N}(\mu_{\max}, \sigma_{\max}^2)$  as the distribution of minority and majority samples in logits space.

Next, we prove the more formal version of the main proposition 3.1, which describes the existence of a balanced dataset, only after we define a key concept, *proportional density difference* (illustrated in figure 5) to outline our proof.
2.1 State Data (Content of the main proposition 2.1) and the state of the

**Definition D.1** (Proportional Density Difference). For any interval I = (a, b] and a mixture distribution  $\varepsilon P_1(x) + (1 - \varepsilon)P_2(x)$ , proportional density difference is defined by the difference of accumulation of two component distributions in the interval I and is denoted by  $\Delta_{\varepsilon}P_{mixture}(I)$ .

$$\Delta_{\varepsilon} P_{mixture}(I) \stackrel{\Delta}{=} \varepsilon P_1(x \in I) - (1 - \varepsilon) P_2(x \in I)$$

**Definition D.2** (Tail Proportional Density Difference). For a mixture distribution  $\varepsilon P_1(x) + (1 - \varepsilon)P_2(x)$ , we define  $tail_L(\alpha)$  as  $\Delta_{\varepsilon}P_{mixture}((-\infty, \alpha])$  and  $tail_R(\beta)$  as  $-\Delta_{\varepsilon}P_{mixture}((\beta, +\infty))$ .

Corollary D.1.

$$tail_L(\alpha) = \varepsilon F^1(\alpha) - (1 - \varepsilon)F^2(\alpha)$$
$$tail_R(\beta) = (1 - \varepsilon)[1 - F^2(\beta)] - \varepsilon[1 - F^1(\beta)]$$

where  $F^1$  and  $F^2$  are CDF of two component distributions.



Figure 5: (a) Illustration of proportion density difference D.1, (b) equation of  $tail_L(\alpha) = tail_R(\beta)$ at D.2.

**Proposition D.1.** [Feasiblity Of Loss-based Group Balancing] Suppose that L is derived from the mixture of two distributions  $\mathcal{N}(\mu_{\min}, \sigma_{\min}^2)$  and  $\mathcal{N}(\mu_{maj}, \sigma_{maj}^2)$  with proportion of  $\varepsilon$  and  $1 - \varepsilon$ , respectively, where  $\varepsilon \leq \frac{1}{2}$ . There exists  $\alpha$  and  $\beta$  such that restricting L to the  $\alpha$ -left and  $\beta$ -right tails of its distribution results in a group-balanced distribution if and only if (i)

$$\sigma_{\min} \ge \sigma_{\max},\tag{5}$$

or (ii)

 $tail_L(\frac{-B+\sqrt{\Delta}}{2A}) > 0 \tag{6}$ 

and

 $\epsilon \geq sigmoid\left(-\frac{\left(\mu_{maj}-\mu_{min}\right)^{2}}{2\left(\sigma_{maj}^{2}-\sigma_{min}^{2}\right)}-\log\left(\frac{\sigma_{maj}}{\sigma_{min}}\right)\right).$ (7)

where 
$$A = \left(\frac{1}{2\sigma_{maj}^2} - \frac{1}{2\sigma_{min}^2}\right)$$
,  $B = \left(\frac{\mu_{min}}{\sigma_{maj}^2} - \frac{\mu_{maj}}{\sigma_{maj}^2}\right)$  and  $\Delta = \frac{(\mu_{min} - \mu_{maj})^2}{\sigma_{min}^2 \sigma_{maj}^2} - 4\left[\log\left(\frac{\sigma_{maj}}{\sigma_{min}}\right) + \log\left(\frac{\epsilon}{1-\epsilon}\right)\right]\left[\frac{1}{2\sigma_{maj}^2} - \frac{1}{2\sigma_{min}^2}\right]$ .

### Proof outline

Our proof proceeds with three steps. First, we reformulate the theorem as an equality of left- and right-tail proportional distribution differences. In other words, we show that the more mass the minority distribution has on one tail, the more mass the majority distribution must have on the other tail. Afterward, supposing  $\mu_{\min} < \mu_{\max}$  WLOG , we propose a proper range for  $\beta$  values on the right tail. We show that when  $\sigma_{maj} \leq \sigma_{min}$ , values for  $\alpha$  trivially exist that can overcome the imbalance between the two distributions. In the last step, for the case in which the variance of the majority is higher than the minority, we discuss a necessary and sufficient condition for the existence of  $\alpha$ and  $\beta$  based on the left-tail proportional density difference using the properties of its derivative with respect to  $\alpha$ .

**Step 1** *Reformulating the problem based on proportional distribution difference.* 

We introduce a utility random variable *Logit Value Tier* as T, which is defined as a function of a random variable L.

$$T_{\alpha,\beta} = \begin{cases} High & \text{if } L \ge \beta \\ Mid & \text{if } \alpha < L < \beta \\ Low & \text{if } L \le \alpha \end{cases}$$
(8)

We can rewrite the problem in formal form as finding an  $\alpha$  and  $\beta$  which satisfies the following equation:

$$P\left(g = \min \left| T_{\alpha,\beta} \neq Mid \right) = P\left(g = \max \left| T_{\alpha,\beta} \neq Mid \right) \right)$$
(9)

972 Equation 7 now can be rewritten to a more suitable form:

/

$$P\left(g = \min \left| T_{\alpha,\beta} \neq Mid \right) = P\left(g = \max \left| T_{\alpha,\beta} \neq Mid \right) \right)$$
(10)

$$\iff \frac{P(T_{\alpha,\beta} \neq Mid|g = \min)P(g = \min)}{P(T_{\alpha,\beta} \neq Mid)} = \frac{P(T_{\alpha,\beta} \neq Mid|g = \max)P(g = \max)}{P(T_{\alpha,\beta} \neq Mid)}$$
(11)

$$\iff P\left(T_{\alpha,\beta} \neq Mid \middle| g = \min\right) P\left(g = \min\right) = P\left(T_{\alpha,\beta} \neq Mid \middle| g = \max\right) P\left(g = \max\right)$$
(12)

$$\iff \qquad \varepsilon P\Big(T_{\alpha,\beta} \neq Mid \Big| g = \min\Big) = (1 - \varepsilon)P\Big(T_{\alpha,\beta} \neq Mid \Big| g = \max j\Big) \qquad (13)$$

$$\iff \varepsilon \left[ P(T_{\alpha,\beta} = Low | g = \min) + P(T_{\alpha,\beta} = High | g = \min) \right] =$$
(14)

$$(1-\varepsilon)\left[P\left(T_{\alpha,\beta} = Low \middle| g = \operatorname{maj}\right) + P\left(T_{\alpha,\beta} = High \middle| g = \operatorname{maj}\right)\right]$$
(15)

$$\iff \varepsilon \left[ P\left( L \le \alpha \middle| g = \min \right) + P\left( L \ge \beta \middle| g = \min \right) \right] =$$
(16)

$$(1-\varepsilon)\left[P\left(L \le \alpha \left| g = \operatorname{maj}\right) + P\left(L \ge \beta \left| g = \operatorname{maj}\right)\right]\right]$$
(17)

$$\iff \qquad \varepsilon \left[ F^{\min}(\alpha) + \left( 1 - F^{\min}(\beta) \right) \right] = (1 - \varepsilon) \left[ F^{\max}(\alpha) + \left( 1 - F^{\max}(\beta) \right) \right] \tag{18}$$

$$\iff \varepsilon F^{\min}(\alpha) - (1-\varepsilon)F^{\max}(\alpha) = (1-\varepsilon)\left[1 - F^{\max}(\beta)\right] - \varepsilon\left[1 - F^{\min}(\beta)\right]$$
(19)

<sup>999</sup> We can see the left side of equation 19 is just a function of alpha. The same goes for the right side of the equation which is a function of  $\beta$ .

Rewriting the left side of the equation as  $tail_L(\alpha)$  and right side as  $tail_R(\beta)$ , the problem is now reduced to finding an  $\alpha$  and  $\beta$  that satisfies

$$tail_L(\alpha) = tail_R(\beta) \tag{20}$$

which is shown in figure 5.

1007 Before reaching out to step two we discuss the properties of  $tail_L$  and  $tail_R$  in Lemma D.2.

Lemma D.2.  $tail_L(\alpha)$  and  $tail_R(\beta)$  are continuous functions and  $\lim_{\alpha \to -\infty} tail_L(\alpha) = 0$ ,  $\lim_{\alpha \to +\infty} tail_L(\alpha) = 2\varepsilon - 1 < 0$ ,  $\lim_{\beta \to +\infty} tail_R(\beta) = 0$  and  $\lim_{\beta \to -\infty} tail_R(\beta) = 1 - 2\varepsilon > 0$ .

*Proof.* Simply proved by the definition of tail functions and properties of CDF.

**Step 2** Solving the equation 20 for simple cases.

**1016** Lemma D.3.  $tail_R(\mu_{maj}) > \frac{1}{2} - \varepsilon \ge 0$ 

1018 Proof.

$$tail_R(\mu_{\rm maj}) = (1 - \varepsilon) \left[ 1 - F^{\rm maj}(\mu_{\rm maj}) \right] - \varepsilon \left[ 1 - F^{\rm min}(\mu_{\rm maj}) \right]$$
(21)

$$= (1 - \varepsilon) \left[ 1 - \phi(0) \right] - \varepsilon \left[ 1 - \phi\left(\frac{\mu_{\text{maj}} - \mu_{\text{min}}}{\sigma_{\text{min}}}\right) \right]$$
(22)

1025  
1024  
1025 
$$> \frac{(1-\varepsilon)}{2} - \varepsilon \left(1 - \frac{1}{2}\right) = \frac{1-2\varepsilon}{2} = \frac{1}{2} - \varepsilon$$
(23)

**Corollary D.2.** Because  $tail_R$  is continuous and  $\lim_{\beta \to +\infty} tail_R(\beta) = 0$ , based on the mean value theorem, any value between zero and  $\frac{(1-2\varepsilon)}{2}$  is obtainable by selecting a  $\beta$  in  $[\mu_2, +\infty)$ .

According to the previous corollary D.2 finding a positive  $tail_L(\alpha)$  will satisfy our need. to find a suitable point, we employ derivatives and properties of relative PDFs to maximize  $tail_L(\alpha)$  and find a positive value.

1033 1034

1035 1036 1037

1038

1039 1040 1041

1043

1045

$$\frac{\mathrm{d}tail_L(\alpha)}{\mathrm{d}\alpha} = \varepsilon f^{\min}(\alpha) - (1-\varepsilon)f^{\max}(\alpha) = \varepsilon f^{\max}(\alpha) \left[\frac{f^{\min}(\alpha)}{f^{\max}(\alpha)} - \frac{1-\varepsilon}{\varepsilon}\right]$$
(24)

The term  $\left[\frac{f^{\min}(\alpha)}{f^{\max}(\alpha)} - \frac{1-\varepsilon}{\varepsilon}\right]$  has the same sign with derivative of  $tail_L(\alpha)$ , also it's roots are critical points of  $tail_L$ , analyzing characteristics of  $\log \frac{f^{\min}(\alpha)}{f^{\max}(\alpha)}$  is the key insight to find a proper  $\alpha$  value.

$$\log f^{\min}(\alpha) - \log f^{\max}(\alpha) = \log \left(\frac{1-\epsilon}{\epsilon}\right)$$

$$\Rightarrow \log\left(\frac{\sigma_{\rm maj}}{\sigma_{\rm min}}\right) - \log\left(\frac{1-\epsilon}{\epsilon}\right) - \frac{(\alpha - \mu_{\rm min})^2}{2\sigma_{\rm min}^2} + \frac{(\alpha - \mu_{\rm maj})^2}{2\sigma_{\rm maj}^2} = 0$$

1046 1047 1048

1051 Because  $\lim_{\alpha \to -\infty} tail_L(\alpha) = 0$  and  $\lim_{\beta \to +\infty} tail_R(\beta) < 0$  to have a positive  $tail_L(\alpha)$ , we need 1052 to have an interval which  $\frac{dtail_L(\alpha)}{d\alpha}$  is positive. For a second degree polynomial like  $ax^2 + bx + c$  to 1054 have positive value, either  $a \ge 0$  or  $\Delta > 0$ , in our case a is  $\left(\frac{1}{\sigma_{mai}^2} - \frac{1}{\sigma_{mai}^2}\right)$ . if  $\sigma_{mai} \ge \sigma_{maj}$  then  $a \ge 0$ 1055 and the minority CDF function will dominate the majority CDF function in the left-side tail and by 1056 choosing a negative number with big enough absolute value for alpha and  $tail_L(\alpha)$  will be positive.



Figure 6: Tail thresholds for three cases: (a) minority group variance is less than majority ( $\sigma_{\min} < \sigma_{\max}$ ), (b) the variance of two groups are equal ( $\sigma_{\min} = \sigma_{\max}$ ) and (c) the variance of the minority group is more than majority ( $\sigma_{\min} > \sigma_{\max}$ ).

1073 Step 3 Solving equation 20 for special case  $\sigma_{min} < \sigma_{maj}$  In case of  $\sigma_{min} \leq \sigma_{maj}$ , having  $\Delta > 0$ 1074 is a necessary condition, also derivative of  $tail_L(\alpha)$  is only positive in  $(\frac{-b-\sqrt{\Delta}}{2a}, \frac{-b+\sqrt{\Delta}}{2a})$  so the 1075 maximum of  $tail_L$  is either in  $-\infty$  or in  $\frac{-b+\sqrt{\Delta}}{2a}$ . Having  $tail_L(\frac{-b+\sqrt{\Delta}}{2a}) > 0$  next to  $\Delta > 0$ 1076 condition, would be the necessary and also sufficient in this case.

0

1078

1072

$$B^2 = \frac{\mu_{\min}^2}{\sigma_{\min}^4} + \frac{\mu_{\max}}{\sigma_{\max}^4} - 2\frac{\mu_{\max}\mu_{\min}}{\sigma_{\max}^2\sigma_{\min}^2}$$

 $4AC = \frac{\mu_{\min}^2}{\sigma_{\min}^4} - \frac{\mu_{\max}^2}{\sigma_{\max}^2 \sigma_{\min}^2} - \frac{\mu_{\max}^2}{\sigma_{\max}^2 \sigma_{\min}^2} + \frac{\mu_{\max}^2}{\sigma_{\max}^4} + 4\left[\log\left(\frac{\sigma_{\max}}{\sigma_{\min}}\right) + \log\left(\frac{\epsilon}{1-\epsilon}\right)\right] \left[\frac{1}{2\sigma_{\max}^2} - \frac{1}{2\sigma_{\min}^2}\right]$   $\Delta = \frac{(\mu_{\min} - \mu_{\max})^2}{\sigma_{\min}^2 \sigma_{\max}^2} - 4\left[\log\left(\frac{\sigma_{\max}}{\sigma_{\min}}\right) + \log\left(\frac{\epsilon}{1-\epsilon}\right)\right] \left[\frac{1}{2\sigma_{\max}^2} - \frac{1}{2\sigma_{\min}^2}\right] \ge 0$   $\Leftrightarrow (\mu_{\min} - \mu_{\max})^2 \ge 2\left[\log\left(\frac{1-\epsilon}{\epsilon}\right) - \log\left(\frac{\sigma_{\max}}{\sigma_{\min}}\right)\right] \left[\sigma_{\max}^2 - \sigma_{\min}^2\right]$   $\Leftrightarrow \epsilon \ge \text{sigmoid}\left(-\frac{(\mu_{\max} - \mu_{\min})^2}{2(\sigma_{\max}^2 - \sigma_{\min}^2)} - \log\left(\frac{\sigma_{\max}}{\sigma_{\min}}\right)\right)$ 

1095

1099

1100

1101

1102

1103

1104

1105

1106

1107 1108

109

Next, we investigate properties of the conditions of the proposition D.1 in case of  $\sigma_{maj} < \sigma_{min}$ . Schematic interpretation of these conditions is presented in figure 7.

• As equation 7 indicates, the minority group is not allowed to be too underrepresented. This especially has a direct relation with the difference of means. The more mean values of groups are different, the more imbalance can be mitigated through loss-based sampling. Mean value difference is especially affected by the spurious correlation, it escalates as the model relies on spurious correlation and also when the spurious features between groups are too different.

• On the other hand condition 6 is more complex and doesn't have a simple closed form, we analytically describe its behaviors by fixating the means and calculating the valid values for  $\varepsilon$ . As the results show in figure 7, most of  $\varepsilon$  are feasible in for  $\sigma_{\min} < \Delta \mu$  as we can see the possible region declines with an increase of  $\sigma_{\min}$  and valid  $\varepsilon$  values cease to exist.

# 1109 D.2 PRACTICAL JUSTIFICATION

1110 As shown in Table 3, the standard deviation ( $\sigma$ ) of the minority group is consistently greater than that of the majority group across all analyzed datasets. Consequently, condition (i) (Eq. 5) of Proposition D.1 is satisfied. Therefore, we theoretically expect the existence of properly balanced left and right tails.

Table 3: Means, standard deviations (STD), and Earth Mover's Distance across WaterBirds andCelebA datasets.

		Waterb		CelebA		
	Cla	ass 1	Class 2		Class 2	
	Min	Maj	Min	Maj	Min	Maj
Mean $(\mu)$	-6.77	-19.17	2.55	11.39	-1.02	6.42
<b>STD</b> $(\sigma)$	6.31	6.23	6.97	4.75	7.64	6.48
Earth Mover's Distance	12.40		8.84		7.43	

1124 1125 1126

1127

# E EXPERIMENTAL DETAILS

1128 E.1 COMPLETE RESULTS

The complete results on Waterbirds, CelebA, and UrbanCars, in addition to complete results on
CivilComments and MultiNLI are reported in Tables 4 and 5 respectively. The results for all methods
except Group DRO + EIIL on all datasets except UrbanCars are reported by Qiu et al. (2023). The
results for Group DRO + EIIL are taken from Zhang et al. (2021). Also, the results of our method and DFR are shown in Table 6



Figure 7: (a) Conditions if  $\sigma_{\min} > \sigma_{\max}$ , (b), (c), (d) minimum, maximum and interval length of feasible  $\varepsilon$  values across ( $\sigma_{\min}, \sigma_{\max}$ ) field for  $\mu_{\min} = 0, \mu_{\max} = 1$ .

Table 4: A comparison of the various methods, ours included, on spurious correlation datasets. The Group Info column indicates if each method utilizes group labels of the training/validation data, with √ denoting that group information is employed during both the training and validation stages. Both the average test accuracy and worst test group accuracy are reported. The mean and standard deviation are calculated over three runs with different seeds. The numbers in bold represent the highest results among all methods, while the underlined numbers represent the best results among methods that may not require group annotation in the training phase.

Method	Group Info	Waterbirds		CelebA		Urbai	ıCa
Mellou	Train/Val	Worst	Average	Worst	Average	Worst	
GDRO (Sagawa et al., 2019) DFR (Kirichenko et al., 2023)	√  √ X √	91.4 92.9 <sub>±0.2</sub>	93.5 $94.2_{\pm 0.4}$	88.9 $88.3_{\pm 1.1}$	92.9 91.3 $_{\pm 0.3}$	73.1 $79.6_{\pm 2.22}$	8
GDRO + EIIL (Creager et al., 2021) JTT (Liu et al., 2021a) SELF (LaBonte et al., 2023) AFR (Qiu et al., 2023) EVaLS-GL (Ours)	X √ X √ X √ X √ X √	$\begin{array}{c} 77.2_{\pm 1} \\ 86.7 \\ \underline{91.6_{\pm 1.4}} \\ 90.4_{\pm 1.1} \\ 89.4_{\pm 0.3} \end{array}$	$\frac{96.5_{\pm 0.2}}{93.3}\\93.6_{\pm 1.1}\\94.2_{1.2}\\95.1_{\pm 0.3}$	$\begin{array}{c} 81.7_{\pm 0.8}\\ 81.1\\ 83.9_{\pm 0.9}\\ 82.0_{\pm 0.5}\\ 84.6_{\pm 1.6}\end{array}$	$\begin{array}{c} 85.7_{\pm 0.1} \\ 88.0 \\ 91.7_{\pm 0.4} \\ 91.3_{\pm 0.3} \\ 91.1_{\pm 0.6} \end{array}$	$\begin{array}{c} 76.5_{\pm 2.6} \\ 79.5 \\ 83.2_{\pm 0.8} \\ 80.2_{\pm 2.0} \\ \textbf{83.5}_{\pm 1.7} \end{array}$	9 
ERM EVaLS (Ours)	x/x x/x	$\begin{array}{c} 66.4_{\pm 2.3} \\ 88.4_{\pm 3.1} \end{array}$	$\begin{array}{c} 90.3_{\pm 0.5} \\ 94.1_{\pm 0.1} \end{array}$	$47.4_{\pm 2.3}$ $85.3_{\pm 0.4}$	$\frac{95.5_{\pm 0.0}}{89.4_{\pm 0.5}}$	$\begin{array}{c} 18.67_{\pm 2.01} \\ 82.1_{\pm 0.9} \end{array}$	1

1188 Table 5: A comparison of the various methods, ours included, on CivilComments and MultiNLI. 1189 The Group Info column indicates if each method utilizes group labels of the training/validation data, 1190 with  $\checkmark$  denoting that group information is employed during both the training and validation stages. Both the average test accuracy and worst test group accuracy are reported. The mean and standard 1191 deviation are calculated over three runs with different seeds. The numbers in bold represent the 1192 highest results among all methods, while the underlined numbers represent the best results among 1193 methods that may not require group annotation in the training phase. 1194

Method	Group Info	CivilCo	omments	MultiNLI	
moniou	Train/Val	Worst	Average	Worst	Average
GDRO (Sagawa et al., 2019)	$\sqrt{1}$	69.9	88.9	77.7	81.4
DFR (Kirichenko et al., 2023)	<b>X</b> /√/	$70.1_{\pm 0.8}$	$87.2_{\pm 0.3}$	$74.7_{\pm 0.7}$	$82.1_{\pm 0.5}$
GDRO + EIIL (Creager et al., 2021)	<b>X</b> /√	$67.0_{\pm 2.4}$	$90.5_{\pm 0.2}$	$61.2_{\pm 0.5}$	$79.4_{\pm 0.}$
JTT (Liu et al., 2021a)	<b>X</b> /√	69.3	91.1	72.6	78.6
SELF (LaBonte et al., 2023)	<b>X</b> /√	$65.9_{\pm 1.7}$	$89.7_{\pm 0.6}$	$70.7_{\pm 2.5}$	$81.2_{\pm 0}$
AFR (Qiu et al., 2023)	<b>X</b> /√	$68.7_{\pm 0.6}$	$89.8_{\pm 0.6}$	$73.4_{\pm 0.6}$	$81.4_{\pm 0}$
EVaLS-GL (Ours)	<b>X</b> /√	$68.0_{\pm 0.5}$	$89.2_{\pm 0.3}$	$75.1_{\pm 1.2}$	$81.6_{\pm 0.0}$
ERM	X/X	$61.2_{\pm 3.6}$	$92.0_{\pm 0.0}$	$64.8_{\pm 1.9}$	$82.6_{\pm 0}$

1207 Table 6: A Comparison of ERM, DFR, EVaLS, and EVaLS-GL on the Dominoes-CMF with dif-1208 ferent spurious correlations for the unknown feature. Both the worst and average of test group 1209 accuracies are presented. The mean and standard deviation are calculated based on runs with three 1210 distinct seeds.

1212		85%	Corr	90%	Corr	95% (	Corr
1213	Method	Worst	Average	Worst	Average	Worst	Average
1214	ERM	$68.3_{\pm 1.5}$	$97.1_{\pm 0.5}$	$50.6_{\pm 1.0}$	$96.1_{\pm 0.0}$	$36.8_{\pm 2.0}$	$95.4_{\pm 1.0}$
1215	DFR	$70.7_{\pm 0.5}$	$86.2_{\pm 0.6}$	$60.2_{\pm 1.2}$	$84.6_{\pm 0.4}$	$42.7_{\pm 2.7}$	$81.5_{\pm 1.2}$
1216	AFR	$65.7_{\pm 0.2}$	$94.2_{\pm 0.8}$	$54.2_{\pm 0.2}^{-}$	$94.9\pm_{2.1}$	$40.3_{\pm 0.5}$	$95.9_{\pm 1.2}^{-}$
1217	AFR + EIIL	$69.1_{\pm 0.1}$	$92_{\pm 1.3}$	$61.5_{\pm 0.2}$	$92.1_{\pm 1.9}$	$40.4_{\pm 0.1}$	$92.9_{\pm 1.5}$
1218	EVaLS-GL	$70.1_{\pm 2.9}$	$82.5_{\pm 1.8}$	$63.6_{\pm 1.3}$	$78.7_{\pm 1.5}$	$48.5_{\pm 0.8}$	$77.0_{\pm 2.0}$
1219	EVaLS	$73.0_{\pm 4.8}$	$81.5_{\pm 1.8}$	$67.1_{\pm 4.2}$	$78.6_{\pm 2.0}$	$51.2_{\pm 1.4}$	$77.5_{\pm 2.5}$

#### 1220 1221

1211

#### 1222 E.2 DOMINOES-COLORED-MNIST-FASHIONMNIST

1223 Dominoes-Colored-MNIST-FashionMNIST (Dominoes-CMF) is a synthetic dataset. We adopt 1224 a similar approach to previous works Pagliardini et al. (2022b); Shah et al. (2020); Kirichenko et al. 1225 (2023) using a modified version of the Dominoes binary classification dataset. This dataset consists 1226 of images with the top half showing CIFAR-10 images Krizhevsky & Hinton (2009), divided into 1227 two meaningful classes: vehicles (airplane, car, ship, truck) and animals (cat, dog, horse, deer). The 1228 bottom half displays either MNIST Deng (2012) images from classes  $\{0 - 3\}$  or Fashion-MNIST 1229 Xiao et al. (2017) images from classes {T-shirt, Dress, Coat, Shirt}. The complex feature (top half) 1230 serves as the core feature and the simple feature (bottom half) is linearly separable and correlated 1231 with the class label at 75%. Furthermore, inspired by the approaches in Zhang et al. (2021); Arjovsky et al. (2020), we intentionally introduce an additional spurious attribute by artificially coloring a 1232 subset of images as follows: for three different datasets, 85%, 90%, and 95% of the images in the 1233 bottom half of class  $c_1$  are randomly assigned a red color in each respective dataset, while 15%, 1234 10%, and 5% of the images are assigned a green color, respectively. The same procedure is applied 1235 inversely for class  $c_2$ . 1236

- 1237 See Table 7 for more details about the dataset statistics.
- 1238
- E.3 DATASETS 1239 1240
- Waterbirds (Sagawa et al., 2019) The dataset comprises images of diverse bird species, classified 1241 into two categories: waterbirds and landbirds. Each image features a bird set against a backdrop of

Top part		Bottom	<b>Part</b> (85% Corr.)	Bottom	Part (90% Corr.)	Bottom Part (95% Corr.)		
CIFAR-10 Class	Color	MNIST	Fashion-MNIST	MNIST	Fashion-MNIST	MNIST	Fashion-MNIST	
c <sub>1</sub> (Vehicle)	Red	12,750	4,250	13,500	4,500	14,250	4,750	
	Green	2,250	750	1,500	500	750	250	
c <sub>2</sub> (Animal)	Red	750	2,250	500	1,500	250	750	
	Green	4,250	12,750	4,500	13,500	4,750	14,250	
Total		40,000		40,000		40,000		

Table 7: Dominoes-CMF Dataset Statistics for 85%, 90%, and 95% Correlation

Table 8: ERM Accuracies on *Dominoes-CMF* Dataset. The mean and standard deviation are reported based on three runs with different seeds.

Top part		Bottom Pa	rt (85% Corr.)	Bottom P	art (90% Corr.)	Bottom Part (95% Corr.)		
CIFAR-10 Class	Color	MNIST	Fashion-MNIST	MNIST	Fashion-MNIST	MNIST	Fashion-MNIST	
$c_1$ (Vehicle)	Red Green	$\begin{array}{c} 98.53_{\pm 0.01}\%\\ 89.33_{\pm 2.4}\%\end{array}$	$\begin{array}{c} 95.61_{\pm 1.1}\% \\ 68.57_{\pm 0.5}\% \end{array}$	${99.2_{\pm 0.01}\%\atop 84.5_{\pm 2.4}\%}$	$95.2_{\pm 1.1}\%$ $54.7_{\pm 0.5}\%$	${}^{99.63_{\pm 0.01}\%}_{63.1_{\pm 1.4}\%}$	$\begin{array}{c} 98.11_{\pm 1.1}\% \\ 36.84_{\pm 0.5}\% \end{array}$	
c <sub>2</sub> (Animal)	Red Green	${}^{68.28 \pm 2.6 \%}_{93.97 \pm 0.5 \%}$	$\frac{86.18_{\pm 2.4}\%}{98.36_{\pm 0.2}\%}$	${}^{56.8_{\pm 5.6}\%}_{96.2_{\pm 0.5}\%}$	$\frac{86.7_{\pm 2.4}\%}{99.3_{\pm 0.2}\%}$	$\begin{array}{c} 39.13_{\pm 1.6}\% \\ 97.92_{\pm 0.5}\% \end{array}$	$68.53_{\pm 2.4}\%$ $99.25_{\pm 0.2}\%$	

either water or land. Interestingly, the background scene acts as a spurious feature in this classification task. Waterbirds are primarily shown against water backgrounds, and landbirds against land backgrounds. Consequently, waterbirds on water and landbirds on land form the minority groups in the training data. It's important to note that the validation dataset for waterbirds is group-balanced, meaning birds from each class are equally represented against both water and land backgrounds. This dataset is mainly categorized as a spurious correlation dataset.

CelebA (Liu et al., 2014) is a widely used dataset in image classification tasks, featuring annotations for 40 binary facial attributes such as hair color, gender, and age. Hair color classification is particularly prominent in literature focusing on spurious correlation robustness. Notably, gender serves as a spurious attribute within this dataset, where a significant majority 94% of individuals with blond hair are women, while men with blond hair represent a minority group. In addition to spurious correlation in the class of blond hair, this dataset also exhibits class imbalance.

MultiNLI (Williams et al., 2017) dataset involves a text classification task focused on determining the relationship between pairs of sentences: contradiction, entailment, or neutral. Sentences containing negation words such as "no" or "never" are under-represented in all three classes, inducing attribute imbalance in the dataset. Figure 8 illustrates the distinct behavior of this dataset compared to other datasets that contain spurious attributes.

CivilComments (Borkan et al., 2019) dataset, as part of the WILDS benchmark, involves a text classification task focused on labeling online comments as either "toxic" or "not toxic". Each comment is associated with 8 attributes, including gender (male, female), sexual orientation (LGBTQ), race (black, white), and religion (Christian, Muslim, or other), based on whether these characteristics are mentioned in the comment. While there is a small attribute imbalance in the dataset, it can categorized into datasets with class imbalance. The detailed proportion of each attribute in each class is described in Table 9. In this paper, we use the implementation of the dataset by the WILDS package (Koh et al., 2021).

 Table 9: Proportion of attributes in each class for CivilComments dataset.

Toxicity (Class)	Male	Female	LGBTQ	Christian	Muslim	Other Religions	Black	White
0	0.11	0.12	0.03	0.10	0.05	0.02	0.03	0.05
1	0.14	0.15	0.08	0.08	0.10	0.03	0.1	0.14



Figure 8: The percentage of samples with the highest (lowest) losses across various thresholds that belong to the minority (majority) group within different classes in  $\mathcal{D}^{LL}$  for (a) MultiNLI and (b) UrbanCars datasets.

UrbanCars (Li et al., 2023) is an image classification dataset with multiple shortcuts. Each image in the dataset consists of a car in the center of the image on a natural scene background, with another object to the right of the image. Images are labeled *Urban* or *City* according to the type of car present in the center. However, each of the backgrounds and the additional objects is highly correlated with the label. While the test set consists of 8 environments based on combinations of the core and two spurious patterns, the training and validation set consist of four groups, based on combinations of the label and only one of the shortcuts.

1316

1317 E.4 TRAINING DETAILS 1318

1319 ERM For Waterbirds and CelebA, we utilize the ResNet50 checkpoints available in the GitHub repository of Kirichenko et al. (2023) as our base model. We 1320 use the ResNet-50 architecture provided by the torchvision package. In the 1321 case of CivilComments and MultiNLI, we adopt a similar approach to Kirichenko 1322 (2023),using BertForSequenceClassification.from\_pretrained et al. 1323 ('bert-base-uncased', ...) from the transformers package. The model is 1324 trained using the AdamW optimizer with a learning rate of  $10^{-5}$ , weight decay of  $10^{-4}$ , and a batch 1325 size of 16 for a total of 5 epochs. 1326

For the UrbanCars dataset, we adhere to the settings described in Li et al. (2023), which involves training a ResNet-50 model pretrained on ImageNet using the SGD optimizer with a learning rate of  $10^{-3}$ , momentum of 0.9, weight decay of  $10^{-4}$ , and a batch size of 128 for 300 epochs. For the Dominoes-CMF dataset, we train a ResNet18 model pretrained on ImageNet for 20 epochs with a batch size of 128 and an SGD optimizer with a learning rate of  $10^{-3}$ , momentum of 0.9, and weight decay of  $10^{-4}$ .

1332

1333 **EVaLS and EVaLS-GL** For every dataset, EIIL was utilized with a learning rate of 0.01, a total 1334 of 20000 steps, and a batch size of 128. The last layer of the model was trained on all datasets using the Adam optimizer. A batch size of 32 and a weight decay of  $10^{-4}$  were used for all datasets. Our 1335 1336 method was evaluated on the validation sets of each dataset, considering both fine-tuning and retrain-1337 ing of the last layer. For all datasets, with the exception of MultiNLI, retraining provided superior validation results. The specifics regarding the number of epochs and the ranges for hyperparameter 1338 search (including learning rate,  $\ell_1$ -regularization coefficient ( $\lambda$ ), and the number of selected samples 1339 (k)) for each dataset are as follows: 1340

# Waterbirds.

- epochs = 100,

- $lr = 5 \times 10^{-4}$
- 1344 1345

1341

1342

1349

 $- \lambda \in \{0, 0.01, 0.02, 0.03, 0.04, 0.05, 0.06, 0.07, 0.08, 0.09, 0.1, 0.2, 0.3, 0.4, 0.5\},\$ 

 $- k \in \{20, 25, 30, 35, 40, 45, 50, 55, 60\}.$ 

1347 • CelebA 1348

- epochs = 50,

-  $\ln = 5 \times 10^{-4}$ ,

1350 -  $\lambda \in \{0, 0.01, 0.02, 0.03, 0.04, 0.05, 0.06, 0.07, 0.08, 0.09, 0.1, 0.2, 0.3, 0.4, 0.5, 0.04, 0.05, 0.06, 0.07, 0.08, 0.09, 0.1, 0.2, 0.3, 0.4, 0.5, 0.04, 0.05, 0.06, 0.07, 0.08, 0.09, 0.1, 0.2, 0.3, 0.4, 0.5, 0.04, 0.05, 0.06, 0.07, 0.08, 0.09, 0.1, 0.2, 0.3, 0.4, 0.5, 0.04, 0.05, 0.06, 0.07, 0.08, 0.09, 0.1, 0.2, 0.3, 0.4, 0.5, 0.04, 0.05, 0.06, 0.07, 0.08, 0.09, 0.1, 0.2, 0.3, 0.4, 0.5, 0.04, 0.05, 0.06, 0.07, 0.08, 0.09, 0.1, 0.2, 0.3, 0.4, 0.5, 0.04, 0.05, 0.05,$ 1351  $0.6, 0.7, 0.8, 0.9, 1, 2\},\$ 1352  $- k \in \{50, 100, 150, 200, 250, 300\}.$ 1353 UrbanCars 1354 - epochs = 100, 1355 -  $\ln \in \{5 \times 10^{-4}, 10^{-3}\},\$ 1356 -  $\lambda \in \{0, 0.01, 0.02, 0.05, 0.1, 1\},\$ 1357  $- k \in \{10, 20, 30, 50, 63\}.$ 1358 1359 CivilComments 1360 - epochs = 50, 1361 -  $\ln \in \{10^{-4}, 5 \times 10^{-4}\},\$ 1362  $\lambda \in \{0, 0.01, 0.02, 0.03, 0.04, 0.05, 0.06, 0.07, 0.08, 0.09, 0.1, 0.2, 0.3, 0.4, 0.5, 0.09, 0.01, 0.02, 0.03, 0.04, 0.05, 0.01, 0.02, 0.02, 0.03, 0.04, 0.05, 0.01, 0.02, 0$ 1363  $0.6, 0.7, 0.8, 0.9, 1, 2\},\$ 1364  $- k \in \{500, 750, 1000, 1250, 1500\}.$ 1365 • MultiNLI 1367 - epochs = 200. -  $\ln \in \{10^{-3}, 10^{-2}\},\$ 1369  $-\lambda \in \{0, 0.01, 0.02, 0.03, 0.04, 0.05, 0.06, 0.07, 0.08, 0.09, 0.1, 0.2, 0.3, 0.4, 0.5\},\$ 1370  $- k \in \{20, 30, 40, 50, 60, 75, 100, 125, 150, 200, 250, 300\}.$ 1371 Dominoes-CMF 1372 - LogisticRegression (penalty="11", solver="liblinear") 1373  $-\lambda \in \{0.001, 0.003, 0.01, 0.02, 0.03, 0.05, 0.07, 0.1, 0.2, 0.3, 0.5, 0.7, 1.0, 3.0\},\$ 1374  $- k \in [10, 80].$ 1375 1376 CelebA-SHSG 1377 - LogisticRegression (penalty="11", solver="liblinear") 1378  $-\lambda \in \{0.001, 0.003, 0.01, 0.02, 0.03, 0.05, 0.07, 0.1, 0.2, 0.3, 0.5, 0.7, 1.0, 3.0\},\$ 1379  $- k \in [1, 100].$ 1380 1381 E.5 SENSITIVITY TO HYPERPARAMETERS 1382 1383 1384 1385 0.03 1386  $\ell_1$  $\ell_1$ 1387 1388 1389 1390  $\hat{k}^{30}$ 1391 1392 (a) Waterbirds (b) CelebA (c) UrbanCars 1393 Figure 9: WGA heatmap on  $D^{MS}$  for different hyperparameter settings across various datasets. 1394 1395

The parameters k (the number of selected samples from each loss tail) and  $\lambda$  (the  $\ell_1$  regularization 1396 factor) are automatically selected using the environment/group-based validation scheme proposed in our method. Sensitivity heatmaps demonstrate the impact of k and  $\lambda$  on the worst-group vali-1398 dation accuracy (WGA) across various datasets. Importantly, our results demonstrate that for most 1399 datasets, multiple hyperparameter combinations yield optimal or near-optimal performance, reduc-1400 ing the need for exhaustive searches. This suggests that the hyperparameter tuning process is not 1401 prohibitively difficult, and even relatively shallow or targeted hyperparameter searches suffice to identify optimal hyperparameter configurations. The difference in WGA between the best and worst 1402 hyperparameter settings for the Waterbirds, CelebA, and UrbanCars datasets is approximately 10%, 1403 16%, and 25%, respectively.

Table 10: Results of DFR and AFR with EIIL-inferred environment for model selection.

Method	Waterbirds	Celeba
DFR (with EIIL) AFR (with EIIL)	$\begin{array}{c} {\bf 92.21 \pm 0.02} \\ {82.6 \pm 0.04} \end{array}$	$85.55 \pm 1.0$ $72.5 \pm 0.01$

1410 Table 11: Performance comparison between misclassified sample selection and EVaLS on the Waterbirds, CelebA, and UrbanCars datasets. The mean and standard deviation values are calculated over three runs with different seeds.

Method	Wate	rbirds	Cel	ebA	Urba	nCars
	Worst	Average	Worst	Average	Worst	Average
Misclassified Selection EVaLS	$77.8_{\pm 5.2}\\88.4_{\pm 3.1}$	$94.0_{\pm 0.4}\\94.1_{\pm 0.1}$	$\begin{array}{c} 85.9_{\pm 1.0} \\ 85.3_{\pm 0.4} \end{array}$	$\begin{array}{c} 89.4_{\pm 0.8} \\ 89.4_{\pm 0.5} \end{array}$	$78.4_{\pm 4.5}\\82.1_{\pm 0.9}$	$\begin{array}{c} 86.9_{\pm 1.4} \\ 88.1_{\pm 0.9} \end{array}$

F ABLATION STUDY

#### 1423 F.1 USE OF EIIL WITH DFR AND AFR 1424

1425 We conducted an ablation study to investigate the impact of using environments inferred from EIIL on model selection. Specifically, we benchmarked the performance of DFR and AFR with EIIL-1426 inferred groups. The results, presented in Table 10, demonstrate the effectiveness of incorporating 1427 EIIL-inferred groups in model selection. The results show that while EIIL-inferred groups reduce the 1428 performance compared to ground-truth annotations for model selection, they still can be effective for 1429 robustness to an extent. Moreover, EVaLS outperforms these two methods when using EIIL inferred 1430 environments. 1431

1432

1404

1411

1412

1422

#### F.2 COMPARISON OF HIGH-LOSS AND MISCLASSIFIED-SAMPLE SELECTION 1433

1434 Several methods, such as JTT (Liu et al., 2021a), rely on misclassified points to address group 1435 imbalances by treating these points as belonging to a minority group. To verify the effectiveness of 1436 loss-based sampling in comparison with misclassification-based sample selection, we conducted an 1437 experiment by replacing loss-based sampling in in EVaLS with selecting misclassified samples and 1438 an equal number of randomly chosen correctly classified samples from each class. This results in 1439 degraded performance compared to EVaLS on the Waterbirds and UrbanCars datasets, and only a marginal improvement (with higher variance) on CelebA, as summarized in Table 11. 1440

1442 F.3 OTHER ENVIRONMENT INFERENCE METHODS

In addition to EIIL, other environment inference methods could be utilized for partitioning the model 1444 selection set into environments. 1445

1446 **Error Splitting** JTT Liu et al. (2021a) partitions data into two correctly classified and misclas-1447 sified sets based on the predictions of a model trained with ERM. We split each of these two sets 1448 based on labels of samples, obtaining  $|\mathcal{Y}| \times 2$  environments. 1449

1450

1441

1443

**Random Classifier Splitting** uses a random classifier to classify features obtained from a model 1451 trained with ERM into correctly classified and misclassified sets. Similar to error splitting, we split 1452 the sets based on class labels. The difference between error splitting and random classifier splitting 1453 is solely in the reinitialization of the classification layer. 1454

1455 The results for EVaLS-ES (EVaLS+Error Sampling) and EVaLS-RC (EVaLS+Random Classifier) are shown in Table 12. One limitation of error splitting is that in datasets with noisy labels or 1456 corrupted images, samples that an ERM model misclassifies may not always belong to minority 1457 groups. In these situations, choosing models based on their accuracy on corrupted data could lead 1458Table 12: The performances of three environment inference methods, when combined with loss-<br/>based sample selection, are evaluated on spurious correlation benchmarks. The mean and standard<br/>deviation values are calculated over three separate runs, each initiated with a different seed.

Method	Wate	erbirds	Cel	ebA	UrbanCars		
litetitet	Worst	Average	erage Worst A		Worst	Average	
EVaLS-ES EVaLS-RC	$\begin{array}{c} 82.1_{\pm 1.2} \\ 88.7_{\pm 1.0} \end{array}$	$94.3_{\pm 0.04}\\94.3_{\pm 1.1}$	$\begin{array}{c} 48.4_{\pm 11.6} \\ 78.1_{\pm 5.1} \end{array}$	$69.5_{\pm 6.5}$ 93.5 $_{\pm 0.2}$	$\begin{array}{c} 79.2_{\pm 2.9} \\ 82.4_{\pm 3.2} \end{array}$	$\begin{array}{c} 86.1_{\pm 0.9} \\ 88.2_{\pm 0.8} \end{array}$	
EVaLS	$88.4_{\pm 3.1}$	$94.1_{\pm 0.1}$	$85.3_{\pm 0.4}$	$89.4_{\pm 0.5}$	$82.1 \pm 0.9$	$88.1_{\pm 0.9}$	

1467 1468 1469

to the selection of models that are not robust to spurious correlations. This is demonstrated by the results of EVaLS-ES on the CelebA dataset.

This shortcoming of error splitting can be alleviated by employing a random classifier instead of the ERM-trained one. Due to the feature-level similarity between minority and majority samples in datasets affected by spurious correlation (Sohoni et al., 2020; Kirichenko et al., 2023; Lee et al., 2023), it is expected that the classifier can differentiate between the groups to some extent. As shown in Table 12, surprisingly, EVaLS-RC produces results that are generally comparable to EVaLS. However, the performance of this method may have high variance, depending on the different initializations of the classifier.

1479

# 1480 G CELEBA-SHSG DATASET FOR UNKNOWN SPURIOUS CORRELATIONS

1481

To further investigate the performance of EVaLS in scenarios with unknown spurious correlations, we propose the CelebA-SHSG (Straight Hair, Smiling, Gender) dataset. This dataset is a subset of the original CelebA (Liu et al., 2014), where the label "*Straight Hair*" is correlated with the attributes of smiling and being female.

The "Straight Hair" attribute is considered as the label, "Smiling" as the known spurious attribute, and gender as the unknown spurious attribute. Average accuracies and Worst group accuracies (WGA) are reported in Table 13 among 8 groups (all binary combinations of the label and spurious attributes). We set the spurious correlation of the known attribute to 80% and conduct experiments for various levels of unknown spurious correlation (similar to the Dominoes-CMF experiments). Spurious correlations are imposed by subsampling from the original CelebA dataset.

1492 The results demonstrate that methods that do not rely on group annotations for retraining or model 1493 selection achieve higher WGA among groups based on both known and unknown attributes. Specifi-1494 cally, EVaLS achieves higher WGA compared to EVaLS-GL, which uses loss-based sampling to create the retraining dataset and relies on group annotations for model selection. Furthermore, EVaLS-1495 GL outperforms DFR, which depends on group annotations for both retraining and model selection. 1496 EVaLS improves the WGA of ERM by 22.7% on average. The Oracle model uses group annota-1497 tions based on both known and unknown spurious attributes during retraining and model selection. 1498 Its WGA is 40.8% higher than that of DFR on average, which only uses annotations of the known 1499 spurious attribute. 1500

- 1501 These results further confirm the findings in Figure 4 for Dominoes-CMF.
- 1502 1503

1504

# H SOCIETAL IMPACTS

Real-world datasets often encapsulate social biases that stem from entrenched stereotypes and historical discrimination, affecting various groups such as genders and races. Machine learning methods,
which learn the correlation between patterns in input data and their targets (e.g., labels in a classification task) (Beery et al., 2018), inadvertently absorb this bias. This unintended consequence leads
to fairness issues in many applications. While strategies to mitigate such biases have been proposed
(as discussed comprehensively in Section A), societal biases are not always known and determined.
We believe that our work, as it addresses these unidentified biases, takes a significant step towards making machine learning fairer for our society.

1510	
1312	Table 13: A Comparison of ERM, DFR, DFR (Oracle), AFR, AFR+EIIL, EVaLS, and EVaLS-GL
1513	on the CelebA-SHSG with different spurious correlations for the unknown feature. Both the worst
1514	and average of test group accuracies are presented. The mean and standard deviation are calculated
1515	based on runs with three distinct seeds.

1516							
1517		85% Corr.		90% Corr.		95% Corr.	
1518	Method	Worst	Average	Worst	Average	Worst	Average
1519	ERM	$28.3_{\pm 0.6}$	$68.2_{\pm 0.6}$	$23.9_{\pm 1.5}$	$67.3_{\pm 0.3}$	$15.6_{\pm 2.6}$	$63.5_{\pm 0.8}$
1520	DFR (Oracle)	$63.1_{\pm 0.9}$	$71.7_{\pm 1.2}$	$59.2_{\pm 1.9}$	$70.0_{\pm 1.1}$	$58.4_{\pm 5.0}$	$67.7_{\pm 1.5}$
1521	DFR	$27.2_{\pm 2.2}$	$67.7_{\pm 0.2}^{-}$	$18.9_{\pm 0.7}$	$64.9_{\pm 0.3}^{-}$	$12.3_{\pm 1.6}$	$60.1_{\pm 0.3}$
1522	AFR	$28.1_{\pm 0.4}$	$68.0_{\pm 0.4}$	$24.3_{\pm 2.1}$	$65.7_{\pm 0.0}$	$15.7_{\pm 2.6}$	$63.1_{\pm 0.0}$
1523	AFR + EIIL	$41.3_{\pm 5.7}$	$63.2_{\pm 5.1}$	$36.3_{\pm 4.5}$	$69.8_{\pm 0.0}$	$45.0_{\pm 5.3}$	$63.2_{\pm 0.0}$
1524	EVaLS-GL	$30.5_{\pm 5.2}$	$68.6_{\pm 2.3}$	$26.3_{\pm 6.4}$	$67.4_{\pm 1.0}$	$19.3_{\pm 3.2}$	$61.6_{\pm 3.4}$
1525	EVaLS	$45.2_{\pm 2.9}$	$59.5_{\pm 2.7}$	$44.9_{\pm 3.1}$	$62.7_{\pm 1.8}$	$45.7_{\pm 2.2}$	$64.4_{\pm 1.8}$

# 1528 I COMPUTATIONAL RESOURCES

Each experiment was conducted on one of the following GPUs: NVIDIA H100 with 80G memory, NVIDIA A100 with 80G memory, NVIDIA Titan RTX with 24G memory, Nvidia GeForce RTX 3090 with 24G memory, and NVIDIA GeForce RTX 3080 Ti with 12G memory.