
Trained Models Tell Us How to Make Them Robust to Spurious Correlation without Group Annotation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Classifiers trained with Empirical Risk Minimization (ERM) tend to rely on at-
2 tributes that have high spurious correlation with the target. This can degrade the
3 performance on underrepresented (or *minority*) groups that lack these attributes,
4 posing significant challenges for both out-of-distribution generalization and fair-
5 ness objectives. Many studies aim to improve robustness to spurious correlation,
6 yet nearly all require group annotation for training and/or model selection. This
7 constrains their applicability in situations where the nature of the spurious correla-
8 tion is not known, or when group labels for certain spurious attributes are either
9 insufficient or completely absent. To meet the demand for effectively enhancing the
10 model robustness under minimal assumptions about group annotation, we propose
11 Environment-based Validation and Loss-based Sampling (EVALS). It uses the losses
12 from a trained model to construct a balanced dataset of high-loss and low-loss
13 samples in which the training data group imbalance is mitigated. This results in
14 a significant robustness to group shifts when equipped with a simple mechanism
15 of last layer retraining. Furthermore, by utilizing environment inference methods
16 for creating diverse environments with correlation shifts, EVALS can potentially
17 eliminate the need for group annotation in the validation data. In such a context, the
18 worst environment accuracy acts as a reliable surrogate throughout the retraining
19 process for tuning hyperparameters and finding a model that performs well across
20 diverse group shifts. EVALS effectively achieves group robustness, showing that
21 group annotation is not necessary even for validation. It is a fast, straightforward,
22 and effective approach that reaches near-optimal worst group accuracy without
23 needing group annotations, marking a new chapter in the robustness of trained
24 models against spurious correlation.

25 1 Introduction

26 Training deep learning models using Empirical Risk Minimization (ERM) on a dataset, poses the
27 risk of relying on *spurious correlation*. These are correlations between certain patterns in the
28 training dataset and the target (e.g., the class label in a classification task) despite lacking any causal
29 relationship. Learning such correlations as shortcuts can negatively impact the models' accuracy on
30 *minority groups* that do not contain the spurious patterns associated with the target [1, 2]. This problem
31 leads to concerns regarding fairness [3], and can also cause a marked reduction in the performance.
32 This occurs particularly when minority groups, which are underrepresented during training, become
33 overrepresented at the time of testing, as a result of shifts within the subpopulations [4]. Hence,
34 ensuring robustness to group shifts and developing methods that improve *worst group accuracy*
35 (WGA) is crucial for achieving both fairness and robustness in the realm of deep learning.

36 Many studies have proposed solutions to address this challenge. A promising line of research
37 focuses on increasing the contribution of minority groups in the model’s training [1, 5–7]. A strong
38 assumption that is considered by some previous works is having access to group annotations for
39 training or fully/partially fine-tuning a pretrained model [8, 7, 1]. The study by Kirichenko et al. [1]
40 proposes that retraining the last layer of a model on a dataset which is balanced in terms of group
41 annotation can effectively enhance the model’s robustness against shifts in spurious correlation. While
42 these works have shown tremendous robustness performance, their assumption for the availability of
43 group annotation restricts their usage.

44 In many real-world applications, the process of labeling samples according to their respective groups
45 can be prohibitively expensive, and sometimes impractical, especially when all minority groups
46 may not be identifiable beforehand. A widely adopted strategy in these situations involves the
47 indirect inference of various groups, followed by the training of models using a loss function that is
48 balanced across groups [5, 9, 10, 4]. The loss value of the model or its similar metrics is a popular
49 signal for recognizing minority groups [5, 9–11]. While most of these techniques necessitate full
50 training of a model, Qiu et al. [9] attempt to adapt the DFR method [1] with the aim of preserving
51 computational efficiency while simultaneously improving robustness to group shift. However, this
52 method still requires group annotations of the validation set for model selection and hyperparameter
53 tuning. Consequently, this constitutes a restrictive assumption when adequate annotations for certain
54 groups are not supplied. It also applies to situations where some shortcut attributes are completely
55 unidentified.

56 In this study, we present a novel strategy that effectively mitigates reliance on spurious correlation,
57 completely eliminating the need for group annotations during both training and retraining. More
58 interestingly, we provide empirical evidence indicating that group annotations are not necessary,
59 even for model selection. We show that assembling a diverse collection of environments for model
60 selection, which reflect group shifts can serve as an effective alternative approach. Our proposed
61 method, Environment-based Validation and Loss-based Sampling (EVALS), is a technique that
62 strengthens the robustness of trained models against spurious correlation, all without relying on group
63 annotations. EVALS is pioneering in its ability to eliminate the need for group annotations at *every*
64 *phase*, including the model selection step. EVALS posits that in the absence of group annotations, a
65 set of *environments* showcasing group shifts is sufficient. Worst Environment Accuracy (WEA) could
66 then be utilized for model selection. Our findings demonstrate that utilizing environment inference
67 methods [12], or even dividing the validation data based on the predictions of a random linear layer
68 atop a trained model’s feature space can markedly enhance group robustness. Figure 1 demonstrates
69 the overall procedure of the main parts of EVALS.

70 Our empirical observations support prior research which suggests that high-loss data points in a
71 trained model may signal the presence of minority groups [5, 9, 10]. Our method, EVALS, evenly
72 selects from both high-loss and low-loss data to form a balanced dataset that is used for last-layer
73 retraining. We offer theoretical explanations for the effectiveness of this approach in addressing group
74 imbalances, and experimentally show the superiority of our efficient solution to the previous strategies.
75 Comprehensive experiments conducted on spurious correlation benchmarks such as CelebA [13],
76 Waterbirds [7], and UrbanCars [14], demonstrate that EVALS achieves optimal accuracy. Moreover,
77 when group annotations are accessible solely for model selection, our approach, EVALS-GL, exhibits
78 enhanced performance against various distribution shifts, including attribute imbalance, as seen in
79 MultiNLI [15], and class imbalance, exemplified by CivilComments [16]. We further present a
80 new dataset, *Dominoes Colored-MNIST-FashionMNIST*, which depicts a situation featuring multiple
81 independent shortcuts, that group annotations are only available for part of them (see Section 2.2). In
82 this setting, we show that strategies with lower levels of group supervision are paradoxically more
83 effective in mitigating the reliance on both known and unknown shortcuts.

84 The main contributions of this paper are summarized as follows:

- 85 • We present EVALS, a simple yet effective approach that enhances model robustness against
86 spurious correlation without relying on ground-truth group annotations.
- 87 • We offer both theoretical and practical insights on how balanced sampling from high-loss and
88 low-loss samples can result in a dataset in which the group imbalance is notably mitigated.
- 89 • Using simple environment inference techniques, EVALS leverages worst environment accu-
90 racy as a reliable indicator for model selection.

- 91 • EVaLS attains near-optimal worst group accuracies or even exceeds them in spurious
92 correlation benchmarks, all with zero group annotations.
- 93 • When group annotations are available for model selection, EVaLS delivers state-of-the-art
94 performance across a variety of subpopulation shift benchmarks.
- 95 • We introduce a new dataset consisting of two spurious features in which partial supervision
96 may negatively impact the performance of the underrepresented groups.

97 2 Preliminaries

98 2.1 Problem Setting

99 We assume a general setting of a supervised learning problem with distinct data partitions \mathcal{D}^{tr} for
100 training, \mathcal{D}^{val} for validation, and $\mathcal{D}^{\text{test}}$ for final evaluation. Each dataset comprises a set of paired
101 samples (x, y) , where $x \in \mathcal{X}$ represents the data and $y \in \mathcal{Y}$ denotes the corresponding labels.
102 Conventionally, \mathcal{D}^{tr} , \mathcal{D}^{val} , and $\mathcal{D}^{\text{test}}$ are assumed to be uniformly sampled from the same distribution.
103 However, this idealized assumption does not hold in many real-world problems where distribution
104 shift is inevitable. In this context, we consider the sub-population shift problem [4]. In a general
105 form of this setting, it is assumed that data samples consist of different groups \mathcal{G}_i , where each
106 group comprises samples that share a property. More specifically, the overall data distribution
107 $p(x, y) = \sum_i \alpha_i p_i(x, y)$ is a composition of individual group distributions $p_i(x, y)$ weighted by their
108 respective proportions α_i , where $\sum_i \alpha_i = 1$. In this work, we assume that \mathcal{D}^{tr} , \mathcal{D}^{val} , and $\mathcal{D}^{\text{test}}$ are
109 composed of identical groups but with a different set of mixing coefficients $\{\alpha_i\}$. It is noteworthy
110 that the validation set may have approximately identical coefficients to those of the training or testing
111 sets, or it may have entirely different coefficients.

112 Several kinds of subpopulation shifts are defined in the literature, including class imbalance, attribute
113 imbalance, and spurious correlation [4]. Class imbalance refers to the cases where there is a difference
114 between the proportion of samples from each class, while attribute imbalance occurs when instances
115 with a certain attribute are underrepresented in the training data, even though this attribute may not
116 necessarily be a reliable predictor of the label. On the other hand, spurious correlation occurs when
117 various groups are differentiated by spurious attributes that are partially predictive and correlated with
118 class labels but are causally irrelevant. More precisely, we can consider a set of spurious attributes
119 \mathcal{S} that partition the data into $|\mathcal{S}| \times |\mathcal{Y}|$ groups. When the concurrence of a spurious attribute with a
120 label is significantly higher than its correlation with other labels, that spurious attribute could become
121 predictive of the label, resulting in deep models relying on the spurious attributes as shortcuts instead
122 of the core ones. This is followed by a decrease in the model’s performance on groups that do not
123 have this attribute.

124 Given a class, the group containing samples with correlated spurious attributes is referred to as
125 *majority* group of that class, while the other groups are called the *minority* groups. As an example,
126 in the Waterbirds dataset [7], for which the task is to classify images of birds into landbird and
127 waterbird, there are spurious attributes $\{\textit{water background}, \textit{land background}\}$. Each background is
128 spuriously correlated with its associated label, decompose the data into two majority groups *waterbird*
129 *on water background*, and *landbird on land background*, and two minority groups *waterbird on land*
130 *background* and *landbird on water background*. Our goal is to make the classifier robust to spurious
131 attributes by increasing performance for all groups.

132 2.2 Robustness of a Trained Model to an Unknown Shortcut

133 In scenarios where group annotations are absent, traditional methods that depend on these annotations
134 for training or model selection become infeasible. Moreover, as previously discussed by [14], when
135 data contains multiple spurious attributes and annotations are only available for some of them, such
136 methods would make the model robust only to the known spurious attributes. To explore such complex
137 scenarios, we introduce the *Dominoes Colored-MNIST-FashionMNIST (Dominoes CMF)* dataset
138 (Figure 3(d)). Drawing inspiration from Pagliardini et al. [17] and Arjovsky et al. [18], Dominoes
139 CMF merges an image from CIFAR10 [19] at the top with a colored (red or green) MNIST [20] or
140 FashionMNIST [21] image at the bottom. The primary label is derived from the CIFAR10 image,
141 while the bottom part introduces two independent spurious attributes: color and shape. Although

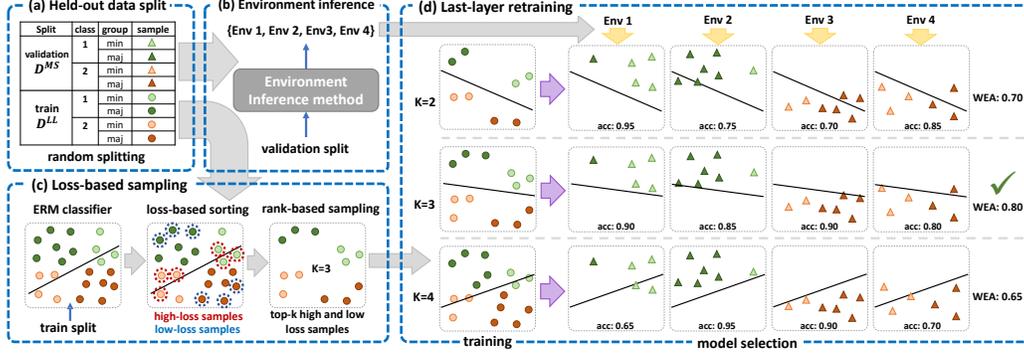


Figure 1: Overview of the proposed method. Given an ERM-trained model (similar to DFR [1]), the following steps are performed: (a) we randomly split the held-out dataset into train and validation splits. (b) An environment inference method is utilized to infer diverse environments from the validation split. (c) We evaluate train split 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. 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.

142 annotations for shape are provided for training and model selection, color remains an unknown
 143 variable until testing. For more details on the dataset refer to the Appendix.

144 The illustrations in Figure 3(a-c) depict the outlined scenario. A classifier trained using ERM is
 145 dependent on both spurious features (Figure 3(b)). Yet, achieving robustness against one spurious
 146 correlation (Figure 3(c)), does not ensure robustness against both (Figure 3(a)). In Section 4 we
 147 show that our method, which does not rely on the group annotations of the identified group, achieves
 148 enhanced robustness against both spurious correlations, outperforming strategies that depend on the
 149 known group’s information.

150 3 Environment-based Validation and Loss-based Sampling

151 Our method, EVaLS, is designed to improve the robustness of deep learning models to group shifts
 152 without the need for group annotation. In line with the DFR [1] approach, we utilize a classifier
 153 defined as $f = h_\phi \circ g_\theta$, where g_θ represents a deep neural network serving as a feature extractor, and
 154 h_ϕ denotes a linear classifier. The classifier is initially trained with the ERM objective on the training
 155 dataset \mathcal{D}^{tr} . Subsequently, we freeze the feature extractor g_θ and focus solely on retraining the last
 156 linear layer h_ϕ using the validation dataset \mathcal{D}^{val} as a held-out dataset.

157 We randomly divide the validation set \mathcal{D}^{val} into two subsets, \mathcal{D}^{LL} and \mathcal{D}^{MS} which are used for last
 158 layer training and model selection, respectively. In Section 3.1 we explain how to sample a subset
 159 of \mathcal{D}^{LL} that statistically handles the group shifts inherent in the dataset. In Section 3.2 we describe
 160 how \mathcal{D}^{MS} is divided into different environments that are later used for model selection. The optimal
 161 number of selected samples from \mathcal{D}^{LL} and other hyperparameters is determined based on the worst
 162 environment accuracies among environments that are obtained from \mathcal{D}^{MS} . By combining our novel
 163 sampling and validation strategy, we aim to provide a robust linear classifier h_{ϕ^*} that significantly
 164 improves the accuracy of underrepresented groups without requiring group annotations of training
 165 or validation sets. Figure 1 illustrates the comprehensive workflow of the EVaLS methodology.
 166 Finally in Section 3.3, we provide theoretical support for the loss-based sampling procedure and its
 167 effectiveness.

168 3.1 Loss-Based Instance Sampling

169 Following previous works [5, 10, 9], we use the loss value as an indicator for identifying minority
 170 groups. We first evaluate classifier f on samples within \mathcal{D}^{LL} and choose k samples with the highest
 171 and lowest loss values for a given k . By combining these $2k$ samples from each class, we construct a

172 balanced set $\mathcal{D}^{balanced}$, consisting of high-loss and low-loss samples (see Figure 1(c)). $\mathcal{D}^{balanced}$ is
173 then used for the training of the last layer of the model.

174 As depicted in Figure 2, the proportion of minority samples among various percentiles of samples
175 with the highest loss values increases as we select a smaller subset of samples with the highest loss.
176 This suggests that high and low-loss samples could serve as effective representatives of minority
177 and majority groups, respectively. In Section 3.3, we offer theoretical insights explaining why this
178 approach could lead to the creation of group-balanced data.

179 3.2 Partitioning Validation Set into Environments

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

186 The concept of an *environment*, as frequently discussed in the invariant learning literature, denotes
187 partitions of data that exhibit different distributions. A model that consistently excels across these
188 varied environments, achieving impressive worst environment accuracy (WEA), is likely to perform
189 equally well across different groups in the test set. Several methods for inferring environments
190 with notable distribution shifts have been introduced [12, 22]. Environment Inference for Invariant
191 Learning (EIL) [12], leverages the predictions from an earlier trained ERM model to divide the data
192 into two distinct environments that significantly deviate from the invariant learning principle proposed
193 by Arjovsky et al. [18], thus creating environments with distribution shifts. Initially, EIL is employed
194 to split \mathcal{D}^{MS} into two environments. Subsequently, each environment is further divided based on
195 sample labels, resulting in $2 \times |\mathcal{Y}|$ environments. To measure the difference between the distribution
196 of environments, we define *group shift* of a class as the absolute difference in the proportion of
197 a minority group between two environments of that class. A higher group shift suggests a more
198 distinct separation between environments. As detailed in the Appendix, environments inferred by
199 EIL demonstrate an average group shift of 28.7% over datasets with spurious correlation. Further
200 information about EIL and the group shift quantities for each dataset can be found in the Appendix.

201 We demonstrate that even more straightforward techniques, such as applying a random linear layer
202 over the feature embedding space and distinguishing environments based on correctly and incorrectly
203 classified samples of each class, can be effective to an extent in several cases (See Appendix E.2).
204 It underscores that the feature space of a trained model is a valuable resource of information for
205 identifying groups affected by spurious correlations. This supports the logic of previous research that
206 employs clustering [23] or contrastive methods [24] in this space to differentiate between groups.

207 3.3 Theoretical Analysis

208 In this subsection, we provide theoretical insights into why loss-based sampling in a class can be
209 utilized to create a balanced dataset of each group under sufficient conditions. We will show the close
210 relationship between the existence of a balanced dataset and the difference between the minority vs.
211 majority group means, calculated based on the logits of an ERM-trained classifier. Such logits are
212 known to depend on spurious features. Hence the mentioned group mean difference is expected to be
213 high if spurious features are present in the dataset.

214 Consider a binary classification problem with a cross-entropy loss function. Let logits be denoted as
215 L . Because loss is a monotonic function of logits, the tails of the distribution of loss across samples
216 are equivalent to that of the logits in each class.

217 We assume that in feature space (output of g_θ) samples from the minority and majority of a class are
218 derived from Gaussian distributions. So, we can consider $\mathcal{N}(\mu_{min}, \sigma_{min}^2)$ and $\mathcal{N}(\mu_{maj}, \sigma_{maj}^2)$ as the
219 distribution of minority and majority samples in logits space.

220 **Proposition 3.1** (Feasibility Of Loss-based Group Balancing). *Suppose that L is derived from the*
221 *mixture of two distributions $\mathcal{N}(\mu_{min}, \sigma_{min}^2)$ and $\mathcal{N}(\mu_{maj}, \sigma_{maj}^2)$ with proportion of ε and $1 - \varepsilon$,*
222 *respectively, where $\varepsilon \leq \frac{1}{2}$. Under sufficient (see App.C) and necessary conditions on μ_{min} , μ_{maj} ,*
223 *σ_{min} and σ_{maj} including inequality 1, there exists α and β such that restricting L to the α -left and*

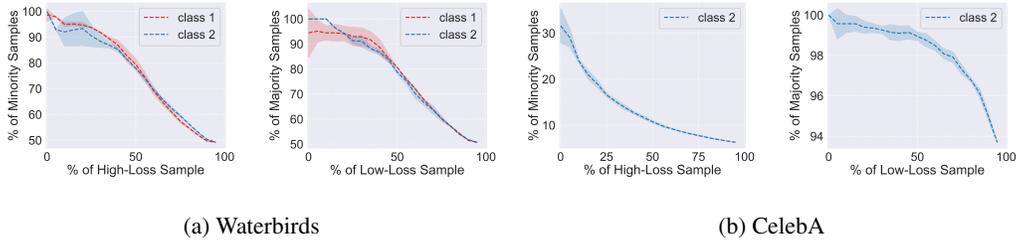


Figure 2: The proportion of minority(majority) samples across different classes within various percentages of \mathcal{D}^{LL} samples with highest (lowest) loss for the Waterbirds (a) and CelebA (b) datasets. Minority group samples are more prevalent among high-loss samples, while majority group samples dominate the low-loss areas. The error bars are calculated across three ERM models.¹

224 β -right tails of its distribution results in a group-balanced distribution; in which both components
 225 are equally represented.

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

226 We provide an outline for proof of Proposition 3.1 here and leave the complete and formal proof and
 227 also exact bounds to Appendix C. We also analyze the conditions and effects of spurious correlation
 228 in satisfying these conditions. To proceed with the outline we first define a key concept to outline our
 229 proof.

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

$$\Delta_\epsilon P_{\text{mixture}}(I) \triangleq \epsilon P_1(x \in I) - (1 - \epsilon)P_2(x \in I)$$

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

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

243 Although the parameters α and β are theoretically established under certain conditions, their actual
 244 values are undetermined. Therefore, validation data is necessary to ascertain them. For practicality
 245 and simplicity in this study, we consider that $\alpha = \beta$ and explore its corresponding sample number
 246 (the count of high- and low-loss samples) from a predefined set of possibilities. By leveraging the
 247 worst environment accuracy, as elaborated in Section 3.2, we identify the optimal candidate that
 248 ensures uniform accuracy across all environments.

249 4 Experiments

250 In this section, we evaluate the effectiveness of our proposed method through comprehensive experi-
 251 ments on multiple datasets and compare it with various methods and baselines. We begin by briefly

¹Note that in the CelebA dataset, only the "blond hair" class includes a minority group.

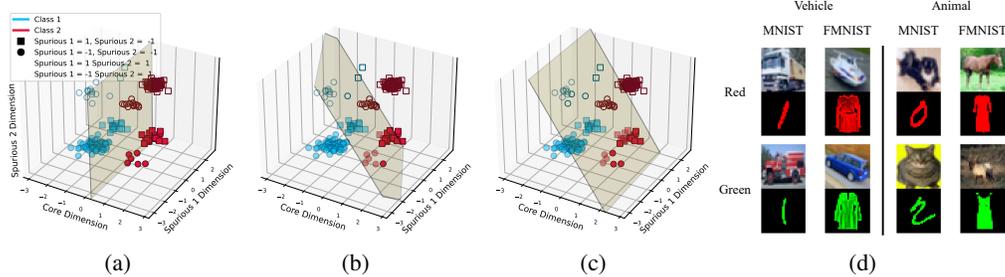


Figure 3: Two spurious correlations in a dataset. (a) If both spurious attributes are known, they can be utilized to fit a classifier that captures the essential attributes. (b) In the absence of knowledge about both spurious attributes, the model would depend on them for classification, leading to incorrect classification of minority samples. (c) If one spurious attribute is unknown (Spurious 2), the model becomes robust only to the known spurious correlation (Spurious 1), but it still underperforms on minority samples. (d) The Dominoes-CMF dataset, which contains two spurious attributes.

252 describing evaluation datasets and then introduce baselines and comparative methods. Finally, we
 253 report and fully explain the results.

254 **Datasets** Our method, along with other baselines, is evaluated on Waterbirds [7], CelebA [13],
 255 UrbanCars [14], CivilComments [16], and MultiNLI [15]. As per the study by Yang et al. [4],
 256 Waterbirds, CelebA, and UrbanCars among these datasets exhibit spurious correlation. Among the
 257 rest, CivilComments has class and attribute imbalance, whereas MultiNLI exhibits attribute imbalance.
 258 For additional details on the datasets, please refer to the Appendix.

259 **Baselines** We compare our method with four baselines in addition to standard ERM. **GroupDRO** [7]
 260 trains a model on the data with the objective of minimizing its average loss on the minority samples.
 261 This method requires group labels of both the training and validation sets. **DFR** [1] argues that
 262 models trained with ERM are capable of extracting the core features of images. Thus, it first trains a
 263 model with ERM, and retrains only the last linear classifier layer on a group-balanced subset of the
 264 validation or the held-out training data. While DFR reduces the number of group-annotated samples,
 265 it still requires group labels in the training phase. **GroupDRO + EIIL** [12] infers environments of
 266 the training set and trains a model with GroupDRO on the inferred environments. **JTT** [5] first trains
 267 a model with ERM on the dataset, and then retrains it on the dataset by upweighting the samples that
 268 were misclassified by the initial ERM model. **AFR** [9] trains a model with ERM on a portion of the
 269 training set, and retrains the classifier on the weighted held-out training data. The weights assigned to
 270 retraining samples are based on the loss of the ERM model, upweighting samples from the minority
 271 groups. Group DRO + EIIL, JTT and AFR remove the reliance on group annotation in the training
 272 phase. However, unlike our method, they all require group labels for model selection.

273 **Setup** Similar to all the works mentioned in Section 4, we use ResNet-50 [25] pretrained on
 274 ImageNet [26] for image classification tasks. We used random crop and random horizontal flip
 275 as data augmentation, similar to [1]. For a fair comparison with the baselines, we did not employ
 276 any data augmentation techniques in the process of retraining the last layer of the model. For the
 277 CivilComments and MultiNLI, we use pretrained BERT [27] and crop sentences to 220 tokens length.
 278 In EvaLS, we use the implementation of EIIL by spuco package [28] for environments inference on
 279 the model selection set with 20000 steps, SGD optimizer, and learning rate 10^{-2} for all datasets.

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

288 **4.1 Results**

289 The results of our experiments along with the reported results on GroupDRO [7], DFR [1], JTT [5],
290 and AFR [9] on five datasets are shown in Table 1. The reported results for GroupDRO, DFR, JTT,
291 and AFR except those for the UrbanCars are taken from Qiu et al. [9]. For EIL+Group DRO, the
292 results are reported from Zhang et al. [24]. We report only the worst group accuracy of methods in
293 Table 1. The average group accuracies are documented in the Appendix. The Group Info column
294 shows whether group annotation is required for training or model selection entry for each method.

295 When compared to other methods with the same level of supervision, EVaLS-GL outperforms on four
296 of the five datasets, achieving near-optimal worst group accuracy on Waterbirds, demonstrating the
297 effectiveness of loss-based sample selection compared to the weighting scheme in AFR [9]. Given
298 that AFR employs exponential weights with a temperature parameter to assign a positive weight
299 to all samples, proportional to the model’s assigned probability of the correct class, an increase
300 in the number of low-loss samples will lead to a corresponding rise in their cumulative weight.
301 Consequently, in situations where spurious correlation is high and an uptick in majority samples leads
302 to a greater proportion of low-loss over high-loss samples, determining an appropriate parameter
303 becomes challenging.

304 The comparison between EVaLS and Group DRO + EIL indicates that when environments are
305 available instead of groups, our method, which uses environments solely for model selection and
306 utilizes loss-based sampling, is more effective than GroupDRO, a potent invariant learning method,
307 which uses this annotation for training.

308 Regarding the UrbanCars, which contains an un-annotated spurious attribute, Li et al. [14] has shown
309 that shortcut mitigation methods often struggle to address multiple shortcuts simultaneously. Notably,
310 techniques such as DFR [1] which are designed to reduce reliance on a specific shortcut feature,
311 cannot make the model robust to an unknown shortcut. In contrast, our experiments suggest that
312 loss-based methods can mitigate the impact of both labeled and unlabeled shortcut features more
313 effectively. Also, in the case of CivilComments, which is viewed as a benchmark for class imbalance,
314 EVaLS-GL exceeds all prior methods, even those with complete group annotation, thanks to the class
315 balancing for the training of the last layer.

316 Our evaluation of EVaLS is based on the spurious correlation benchmarks. This is because, in
317 other instances of subpopulation shift, the attributes that differ across groups are not predictive of
318 the label, thereby reducing the visibility of these attributes’ effects in the model’s final layers [29].
319 Consequently, EIL, which depends on output logits for prediction, might not effectively separate
320 the groups. This observation is further supported by our findings related to the degree of group
321 shift between the environments inferred by EIL for each class in the CivilComments and MultiNLI
322 datasets. The average group shift (defined in the Section 3.2) in the environments of the minority
323 class of CivilComments is only $5.6_{\pm 0.8}\%$. Also, environments associated with Classes 1 and 2 in
324 MultiNLI show only $1.1_{\pm 0.3}\%$ and $1.9_{\pm 1.0}\%$ group shift respectively. More results and ablation
325 studies can be found in the Appendix.

326 **Mitigating Multiple Shortcut Attributes** To evaluate the performance of our method in the case
327 of unknown spurious correlations, we train a ResNet-18 [25] model on the *Dominoes-CMF* dataset.
328 We apply DFR [1], EVaLS-GL, and EVaLS on top of the trained ERMs to assess their ability to
329 mitigate multiple shortcuts. For the last layer training set, we consider the MNIST/Fashion-MNIST
330 feature as the known group label, and the color as the unknown attribute. The results are shown in
331 Table 2. To clarify, we calculate the worst-group accuracy on the validation set considering only the
332 label of one shortcut, *i.e.*, the lowest accuracy among the four groups based on the combination of the
333 target label and the single known shortcut label. Note that EVaLS does not require group annotations.

334 Our results confirm findings by Li et al. [14], suggesting that methods using group labels mitigate
335 reliance on the known shortcut but not necessarily on the unknown one. EVaLS-GL mitigates this
336 phenomenon using its loss-based sampling approach, but surprisingly EVaLS even outperforms
337 EVaLS-GL. Combining a loss-based sampling approach for last layer training and environment-based
338 model selection, results in a completely group-annotation-free method in a multi-shortcut setting and
339 successfully re-weights features to perform well with respect to both spurious attributes.

Table 1: A comparison of the worst group accuracy across various methods, ours included, on five datasets. The Group Info column indicates if each method utilizes group labels of the training/validation data, with \checkmark denoting that group information is employed during both the training and validation stages. Bold numbers are the highest results overall, while underlined ones are the best among methods that may require group annotation only for model selection. CivilComments is class imbalanced, MultiNLI has imbalanced attributes, and the other three datasets have spurious correlations. The \times sign indicates that the dataset is out of the scope of the method. The mean and standard deviation are calculated over three runs with different seeds.

Method	Group Info Train/Val	Datasets				
		Waterbirds	CelebA	UrbanCars	CivilComments	MultiNLI
GDRO [7]	\checkmark/\checkmark	91.4	88.9	-	69.9	77.7
DFR [1]	\times/\checkmark	92.9\pm0.2	88.3 \pm 1.1	79.6 \pm 2.22	70.1 \pm 0.8	74.7 \pm 0.7
GDRO + EIIL [12]	\times/\checkmark	77.2 \pm 1	81.7 \pm 0.8	-	67.0 \pm 2.4	-
JTT [5]	\times/\checkmark	86.7	81.1	-	69.3	72.6
AFR [9]	\times/\checkmark	<u>90.4\pm1.1</u>	82.0 \pm 0.5	80.2 \pm 2.0	68.7 \pm 0.6	73.4 \pm 0.6
EVaLS-GL (Ours)	\times/\checkmark	89.4 \pm 0.3	84.6 \pm 1.6	82.27\pm1.16	80.5\pm0.4	<u>75.1\pm1.2</u>
ERM	\times/\times	66.4 \pm 2.3	47.4 \pm 2.3	18.67 \pm 2.01	61.2 \pm 3.6	64.8 \pm 1.9
EVaLS (Ours)	\times/\times	88.4 \pm 3.1	<u>85.3\pm0.4</u>	82.13 \pm 0.92	\times	\times

Table 2: Worst test group accuracy of ERM, DFR, EVaLS, and EVaLS-GL on the Dominoes-CMF Dataset. The mean and standard deviation are calculated based on runs with three distinct seeds.

	ERM	DFR	EVaLS-GL	EVaLS
Worst Group Accuracy	50.6 \pm 1.0	60.2 \pm 1.2	63.6 \pm 1.3	67.1\pm4.2

340 5 Discussion

341 This study presents EVaLS, a novel approach to improve robustness to spurious correlations with
342 zero group annotation. EVaLS uses loss-based sampling to create a balanced training dataset that
343 effectively disrupts spurious correlations and employs EIIL to infer environments for model selection.
344 We also explore situations with multiple spurious correlations where not all spurious factors are
345 known. In this context, we introduce Dominoes-CMF, a dataset in which two factors are spuriously
346 correlated with the label, but only one is identified. Our findings suggest that EVaLS attains near-
347 optimal worst test group accuracy on spurious correlation datasets. We also present EVaLS-GL, which
348 needs group labels only for model selection. Our empirical tests on various datasets demonstrate
349 EVaLS-GL outperforms state-of-the-art methods requiring group data during evaluation or training.

350 Note that this paper remains consistent with the findings of Lin et al. [30]. Our approach does not
351 involve identifying spurious attributes without auxiliary information. Instead, the objective is to make
352 a trained model robust against its reliance on shortcuts. Specifically, conditioning on what a trained
353 model learns, we ascertain that both the loss value and the model’s feature space are instrumental in
354 mitigating shortcuts and effectuating notable shifts among groups.

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

360 EVaLS represents a significant advancement in the development of methods for enhancing model
361 fairness and robustness without prior knowledge about group annotations. Future work could explore
362 developing environment inference methods effective for other types of subpopulation shift, such as
363 attribute and class imbalance.

364 **References**

- 365 [1] Polina Kirichenko, Pavel Izmailov, and Andrew Gordon Wilson. Last layer re-training is suffi-
366 cient for robustness to spurious correlations. In *The Eleventh International Conference on Learn-*
367 *ing Representations*, 2023. URL <https://openreview.net/forum?id=Zb6c8A-Fghk>.
- 368 [2] Tyler LaBonte, Vidya Muthukumar, and Abhishek Kumar. Towards last-layer retraining for
369 group robustness with fewer annotations. In *Thirty-seventh Conference on Neural Information*
370 *Processing Systems*, 2023. URL <https://openreview.net/forum?id=kshC3NOP6h>.
- 371 [3] Tatsunori Hashimoto, Megha Srivastava, Hongseok Namkoong, and Percy Liang. Fairness
372 without demographics in repeated loss minimization. In *International Conference on Machine*
373 *Learning*, pages 1929–1938. PMLR, 2018.
- 374 [4] Yuzhe Yang, Haoran Zhang, Dina Katabi, and Marzyeh Ghassemi. Change is hard: A closer
375 look at subpopulation shift. In *International Conference on Machine Learning*, 2023.
- 376 [5] Evan Z Liu, Behzad Haghgoo, Annie S Chen, Aditi Raghunathan, Pang Wei Koh, Shiori
377 Sagawa, Percy Liang, and Chelsea Finn. Just train twice: Improving group robustness without
378 training group information. In Marina Meila and Tong Zhang, editors, *Proceedings of the*
379 *38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine*
380 *Learning Research*, pages 6781–6792. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/liu21f.html>.
- 382 [6] Yu Yang, Eric Gan, Gintare Karolina Dziugaite, and Baharan Mirzasoleiman. Identifying
383 spurious biases early in training through the lens of simplicity bias. *ArXiv*, abs/2305.18761,
384 2023. URL <https://api.semanticscholar.org/CorpusID:258967752>.
- 385 [7] Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust
386 neural networks. In *International Conference on Learning Representations*, 2019.
- 387 [8] Junhyun Nam, Jaehyung Kim, Jaeho Lee, and Jinwoo Shin. Spread spurious attribute: Improving
388 worst-group accuracy with spurious attribute estimation. In *International Conference on*
389 *Learning Representations*, 2021.
- 390 [9] Shikai Qiu, Andres Potapczynski, Pavel Izmailov, and Andrew Gordon Wilson. Simple and fast
391 group robustness by automatic feature reweighting. In *International Conference on Machine*
392 *Learning*, pages 28448–28467. PMLR, 2023.
- 393 [10] Junhyun Nam, Hyuntak Cha, Sungsoo Ahn, Jaeho Lee, and Jinwoo Shin. Learning from failure:
394 De-biasing classifier from biased classifier. *Advances in Neural Information Processing Systems*,
395 33:20673–20684, 2020.
- 396 [11] Fahimeh Hosseini Noohdani, Parsa Hosseini, Aryan Yazdan Parast, HamidReza Yaghoubi
397 Araghi, and Mahdieh Soleymani Baghshah. Decompose-and-compose: A compositional
398 approach to mitigating spurious correlation. *CoRR*, abs/2402.18919, 2024. doi: 10.48550/
399 ARXIV.2402.18919. URL <https://doi.org/10.48550/arXiv.2402.18919>.
- 400 [12] Elliot Creager, Joern-Henrik Jacobsen, and Richard Zemel. Environment inference for invariant
401 learning. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International*
402 *Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*,
403 pages 2189–2200. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/creager21a.html>.
- 404 [13] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the
405 wild. *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 3730–3738,
406 2014. URL <https://api.semanticscholar.org/CorpusID:459456>.
- 407 [14] Zhiheng Li, Ivan Evtimov, Albert Gordo, Caner Hazirbas, Tal Hassner, Cristian Canton Ferrer,
408 Chenliang Xu, and Mark Ibrahim. A whac-a-mole dilemma: Shortcuts come in multiples where
409 mitigating one amplifies others. In *Proceedings of the IEEE/CVF Conference on Computer*
410 *Vision and Pattern Recognition (CVPR)*, pages 20071–20082, June 2023.
- 411

- 412 [15] Adina Williams, Nikita Nangia, and Samuel R Bowman. A broad-coverage challenge corpus
413 for sentence understanding through inference. *arXiv preprint arXiv:1704.05426*, 2017.
- 414 [16] Daniel Borkan, Lucas Dixon, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. Nuanced
415 metrics for measuring unintended bias with real data for text classification. In *Companion*
416 *proceedings of the 2019 world wide web conference*, pages 491–500, 2019.
- 417 [17] Matteo Pagliardini, Martin Jaggi, François Fleuret, and Sai Praneeth Karimireddy. Agree to
418 disagree: Diversity through disagreement for better transferability. In *The Eleventh International*
419 *Conference on Learning Representations*, 2022.
- 420 [18] Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk mini-
421 mization, 2020.
- 422 [19] Alex Krizhevsky. Learning multiple layers of features from tiny images. pages 32–33, 2009.
423 URL <https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf>.
- 424 [20] Yann LeCun and Corinna Cortes. MNIST handwritten digit database. 2010. URL <http://yann.lecun.com/exdb/mnist/>.
425
- 426 [21] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image
427 dataset for benchmarking machine learning algorithms, 2017. URL <http://arxiv.org/abs/1708.07747>. cite arxiv:1708.07747Comment: Dataset is freely available at
428 <https://github.com/zalando-research/fashion-mnist> Benchmark is available at [http://fashion-](http://fashion-mnist.s3-website.eu-central-1.amazonaws.com/)
429 [mnist.s3-website.eu-central-1.amazonaws.com/](http://fashion-mnist.s3-website.eu-central-1.amazonaws.com/).
430
- 431 [22] Jiashuo Liu, Zheyuan Hu, Peng Cui, Bo Li, and Zheyuan Shen. Heterogeneous risk minimization.
432 In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on*
433 *Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 6804–
434 6814. PMLR, 18–24 Jul 2021. URL [https://proceedings.mlr.press/v139/liu21h.](https://proceedings.mlr.press/v139/liu21h.html)
435 [html](https://proceedings.mlr.press/v139/liu21h.html).
- 436 [23] Nimit Sohoni, Jared Dunnmon, Geoffrey Angus, Albert Gu, and Christopher Ré. No subclass
437 left behind: Fine-grained robustness in coarse-grained classification problems. *Advances in*
438 *Neural Information Processing Systems*, 33:19339–19352, 2020.
- 439 [24] Michael Zhang, Nimit Sharad Sohoni, Hongyang R. Zhang, Chelsea Finn, and Christopher Ré.
440 Correct-n-contrast: A contrastive approach for improving robustness to spurious correlations.
441 In *NeurIPS 2021 Workshop on Distribution Shifts: Connecting Methods and Applications*, 2021.
442 URL https://openreview.net/forum?id=Q41k1_DwS3Y.
- 443 [25] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
444 recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*,
445 pages 770–778, 2016. doi: 10.1109/CVPR.2016.90.
- 446 [26] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng
447 Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual
448 recognition challenge. *International journal of computer vision*, 115:211–252, 2015.
- 449 [27] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of
450 deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and
451 Tamar Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter*
452 *of the Association for Computational Linguistics: Human Language Technologies, Volume 1*
453 *(Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association
454 for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL [https://aclanthology.](https://aclanthology.org/N19-1423)
455 [org/N19-1423](https://aclanthology.org/N19-1423).
- 456 [28] Siddharth Joshi, Yu Yang, Yihao Xue, Wenhan Yang, and Baharan Mirzsoleiman. To-
457 wards mitigating spurious correlations in the wild: A benchmark & a more realistic dataset.
458 *ArXiv*, abs/2306.11957, 2023. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:259211935)
459 [259211935](https://api.semanticscholar.org/CorpusID:259211935).

- 460 [29] Yoonho Lee, Annie S Chen, Fahim Tajwar, Ananya Kumar, Huaxiu Yao, Percy Liang, and
461 Chelsea Finn. Surgical fine-tuning improves adaptation to distribution shifts. In *The Eleventh
462 International Conference on Learning Representations*, 2023. URL [https://openreview.
463 net/forum?id=APuPRxjHvZ](https://openreview.net/forum?id=APuPRxjHvZ).
- 464 [30] Yong Lin, Shengyu Zhu, Lu Tan, and Peng Cui. Zin: When and how to learn invariance without
465 environment partition? In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh,
466 editors, *Advances in Neural Information Processing Systems*, volume 35, pages 24529–24542.
467 Curran Associates, Inc., 2022. URL [https://proceedings.neurips.cc/paper_files/
468 paper/2022/file/9b77f07301b1ef1fe810aae96c12cb7b-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/9b77f07301b1ef1fe810aae96c12cb7b-Paper-Conference.pdf).
- 469 [31] Zheyang Shen, Jiashuo Liu, Yue He, Xingxuan Zhang, Renzhe Xu, Han Yu, and Peng Cui.
470 Towards out-of-distribution generalization: A survey. *ArXiv*, abs/2108.13624, 2021. URL
471 <https://api.semanticscholar.org/CorpusID:237364121>.
- 472 [32] Seonguk Seo, Joon-Young Lee, and Bohyung Han. Information-theoretic bias reduction via
473 causal view of spurious correlation. In *Proceedings of the AAAI Conference on Artificial
474 Intelligence*, volume 36, pages 2180–2188, 2022.
- 475 [33] Yuzhen Mao, Zhun Deng, Huaxiu Yao, Ting Ye, Kenji Kawaguchi, and James Zou. Last-layer
476 fairness fine-tuning is simple and effective for neural networks. *arXiv preprint arXiv:2304.03935*,
477 2023.
- 478 [34] David Krueger, Ethan Caballero, Joern-Henrik Jacobsen, Amy Zhang, Jonathan Binas, Dinghui
479 Zhang, Remi Le Priol, and Aaron Courville. Out-of-distribution generalization via risk extrapo-
480 lation (rex). In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International
481 Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*,
482 pages 5815–5826. PMLR, 18–24 Jul 2021. URL [https://proceedings.mlr.press/v139/
483 krueger21a.html](https://proceedings.mlr.press/v139/krueger21a.html).
- 484 [35] Alexandre Rame, Corentin Dancette, and Matthieu Cord. Fishr: Invariant gradient variances for
485 out-of-distribution generalization. In *International Conference on Machine Learning*, pages
486 18347–18377. PMLR, 2022.
- 487 [36] Faruk Ahmed, Yoshua Bengio, Harm van Seijen, and Aaron Courville. Systematic generalisation
488 with group invariant predictions. In *International Conference on Learning Representations*,
489 2021. URL <https://openreview.net/forum?id=b9P0imzZFJ>.
- 490 [37] Matteo Pagliardini, Martin Jaggi, François Fleuret, and Sai Praneeth Karimireddy. Agree
491 to disagree: Diversity through disagreement for better transferability. *arXiv preprint*,
492 arXiv:2202.04414, 2022.
- 493 [38] Harshay Shah, Kaustav Tamuly, Aditi Raghunathan, Prateek Jain, and Praneeth Netrapalli.
494 The pitfalls of simplicity bias in neural networks. *Advances in Neural Information Processing
495 Systems*, 33:9573–9585, 2020.
- 496 [39] Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay
497 Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, Tony Lee,
498 Etienne David, Ian Stavness, Wei Guo, Berton Earnshaw, Imran Haque, Sara M Beery, Jure
499 Leskovec, Anshul Kundaje, Emma Pierson, Sergey Levine, Chelsea Finn, and Percy Liang.
500 Wilds: A benchmark of in-the-wild distribution shifts. In Marina Meila and Tong Zhang,
501 editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of
502 *Proceedings of Machine Learning Research*, pages 5637–5664. PMLR, 18–24 Jul 2021. URL
503 <https://proceedings.mlr.press/v139/koh21a.html>.
- 504 [40] Sara Beery, Grant Van Horn, and Pietro Perona. Recognition in terra incognita. In *Proceedings
505 of the European Conference on Computer Vision (ECCV)*, September 2018.

506 A Related Work

507 Robustness to spurious correlation is a critical concern across various machine learning subfields.
508 It is a form of out-of-distribution generalization [31] where the distribution shift arises from the
509 disproportionate representation of minority groups—those instances that are devoid of the correlated
510 spurious patterns associated with their labels [4]. The issue of spurious correlation also intersects
511 with the discourse on fairness in machine learning. [32, 33].

512 Past studies have proposed a range of strategies to mitigate the models’ reliance on spurious correla-
513 tion. Broadly speaking, these methods can be categorized according to the degree of supervision they
514 require regarding group labels.

515 Invariant learning (IL) methods [18, 34, 35] operate under the assumption of having access to a
516 collection of environments that comprise group shift. By imposing invariant conditions on these envi-
517 ronments, IL methods strive to create classifiers robust against group-sensitive features. IRM [18] is
518 designed to learn a feature extractor, which, when utilized, guarantees the existence of a classifier that
519 would be optimal in all training environments. VREx [34] aims to decrease the risk variance among
520 different training environments. PGI [36] works by minimizing the distance between the expected
521 softmax distribution of labels, conditioned on inputs across both majority and minority environments.
522 Lastly, Fishr [35] focuses on bringing the variance of risk gradients closer together across different
523 training environments. For scenarios which environments are not available, environment inference
524 methods [12, 22] are used to obtain a set of environments. Creager et al. [12] introduce environment
525 inference for invariant learning (EIIL), which tries to partition samples into two groups such that the
526 objective of IRM [18] is maximized. HRM [22] aims to optimize both an environment inference
527 module and an invariant prediction module jointly, with the goal of achieving an invariant predictor.

528 When group annotations are accessible, various methods leverage this information to equalize the
529 impact of different groups on the model’s loss. The Group Distributionally Robust Optimization
530 (GDRO) approach [7], for instance, focuses on optimizing the loss for the worst-performing group
531 during training. Kirichenko et al. [1] has shown that models can still learn and extract core data
532 features even in the presence high spurious correlation. Consequently, They suggest that retraining
533 just the last layer of a model initially trained with Empirical Risk Minimization (ERM) can effectively
534 reduce reliance on spurious correlation for predicting class labels. This method, termed Deep Feature
535 Re-weighting (DFR), has been validated as not only highly effective but also significantly more
536 efficient than earlier techniques that necessitated retraining the full model [8, 7]. However, availability
537 of group annotations is considered a serious restrictive assumption.

538 Several recent studies have endeavored to enhance model robustness against spurious correlation,
539 even in the absence of group annotations [5, 24, 9, 2, 6]. Liu et al. [5] introduce a two-stage method
540 that involves training a model using ERM for a number of epochs before retraining it to give more
541 weight to misclassified samples. The study by Zhang et al. [24] employs the same two-stage training
542 process, but with a twist for the second stage: they utilize contrastive methods. The goal is to
543 bring samples from the same class but with divergent predictions closer in the feature space, while
544 simultaneously increasing the separation between samples from different classes that have similar
545 predictions. Another method, known as automatic feature reweighting (AFR) [9], reweights the last
546 layer of an ERM-pretrained model to favor samples that the original model was less accurate on.
547 LaBonte et al. [2] refine the last layer of an ERM-trained model through class-balanced finetuning,
548 identifying challenging data points by comparing the classifier’s predictions with those of an early-
549 stopped version. While these methods have significantly reduced the reliance on group annotations,
550 some are still required for validation and model selection. This remains a constraint, particularly
551 when the spurious correlation is completely unknown.

552 For making a trained model robust to spurious correlation with zero group annotations, recently,
553 LaBonte et al. [2] have empirically demonstrated that the class-balanced retraining of a model
554 pretrained with ERM can effectively improve the WGA for certain datasets. However, this approach
555 fails in datasets with a high degree of spurious correlation.

Table 3: The average and variation percentage (%) (across 3 seeds) of group shift between the inferred environments using EIIL [12] 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		
	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

556 B Environment Inference for Invariant Learning

557 Consider the training dataset $\mathcal{D}^{tr} = \{(x^{(i)}, y^{(i)}) | x^{(i)} \in \mathcal{X}, y^{(i)} \in \mathcal{Y}\}$, where \mathcal{X} and \mathcal{Y} represent the
558 input and output spaces, respectively. This dataset can be partitioned into different environments
559 $\mathcal{E}^{tr} = \{e_1, \dots, e_n\}$, such that for any $i \neq j$, the data distribution in e_i and e_j differs. The objective
560 of invariant learning is to train a predictor that performs consistently across all environments in \mathcal{E}^{tr} .
561 Under certain conditions, this predictor is also expected to perform well on e^{tst} , a test environment
562 with a distribution distinct from the training data. Invariant Risk Minimization (IRM) [18] approaches
563 this problem by learning a feature extractor $\Phi(\cdot)$ such that a classifier $\omega(\cdot)$ exists, where $\omega \circ \Phi(\cdot)$
564 performs consistently across all training environments. The practical implementation of the IRM
565 objective is to minimize

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

566 where $\bar{\omega}$ is a constant scalar with a value of 1.0, λ is a hyperparameter, and $R^e(f) =$
567 $\mathbf{E}_{(x,y) \sim p_e} [l(f(x), y)]$ is referred to as the risk on environment e .

568 In real-world scenarios, training environments might not always be available. To address this,
569 Environment Inference for Invariant Learning (EIIL) [12] partitions samples into two environments
570 in a way that maximizes the objective in Eq 2.

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

$$\mathbf{q}^* = \arg \max_{\mathbf{q}} \|\nabla_{\bar{\omega}} \tilde{R}^e(\bar{\omega} \circ \Phi)\|. \quad (3)$$

575 As discussed in Creager et al. [12], using a biased base model Φ could lead to environments exhibiting
576 varying degrees of spurious correlation. During the inference phase, the soft assignment is converted
577 to a hard assignment. The average group shift between the inferred environments using EIIL is
578 illustrated in Table 3.

579 **C Theoretical Analysis**

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

583 Consider a binary classification problem with a cross-entropy loss function. Let logits be denoted
 584 as L . Because loss is a monotonic function of logits, the tails of the distribution of loss across
 585 samples are equivalent to that of the logits in each class. We assume that in feature space (output
 586 of g_θ) samples from the minority and majority of a class are derived from Gaussian distributions
 587 $\mathcal{N}(h_{\min}, t_{\min}^2 I_d)$ and $\mathcal{N}(h_{\text{maj}}, t_{\text{maj}}^2 I_d)$, respectively. Before diving into the group balance problem we
 588 initially show that the distribution of minority and majority samples in the logit space (output of h_ϕ)
 589 are Gaussian too.

590 **Lemma C.1** (Gaussian Distribution of Logits). *if $Z \sim \mathcal{N}(h, t^2 I_d)$ in feature space and $W \in \mathbb{R}^d$
 591 then logits $L = \langle W, Z \rangle \sim \mathcal{N}(Wh, t^2 \|W\|^2)$*

592 *Proof.* Let $Z \sim \mathcal{N}(h, t^2 I_d)$.

593 Consider the linear combination $L = \langle W, Z \rangle = W^T Z$, where $W \in \mathbb{R}^d$ which is a univariate
 594 gaussian.

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

596 **1. Mean of L**

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

597 Therefore, the mean of L is Wh .

598 **2. Variance of L :**

599 The variance of L can be computed using the properties of covariance. Recall that if $Z \sim \mathcal{N}(h, t^2 I_d)$,
 600 then the covariance matrix of Z is $t^2 I_d$.

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

$$\text{Var}(L) = \text{Var}(W^T Z) = W^T \text{Cov}(Z) W.$$

602 Given $\text{Cov}(Z) = t^2 I_d$, we have:

$$\text{Var}(L) = W^T (t^2 I_d) W = t^2 W^T I_d W = t^2 \|W\|^2,$$

603 where $\|W\|$ denotes the Euclidean norm of W .

604 Combining the mean and variance results, we conclude that L is normally distributed with mean Wh
 605 and variance $t^2 \|W\|^2$:

$$L = \langle W, Z \rangle \sim \mathcal{N}(Wh, t^2 \|W\|^2).$$

606 Thus, we have proved that if $Z \sim \mathcal{N}(h, t^2 I_d)$, then the logits $L = \langle W, Z \rangle$ follow the distribution
 607 $\mathcal{N}(Wh, t^2 \|W\|^2)$. \square

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

610 Next, we prove the more formal version of the main proposition 3.1 which describes the existence of
 611 a balanced dataset, only after we define a key concept, *proportional density difference* (illustrated in
 612 figure 4) to outline our proof.

Definition C.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 $\Delta_\varepsilon P_{mixture}(I)$.

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

613 **Definition C.2** (Tail Proportional Density Difference). For a mixture distribution $\varepsilon P_1(x) + (1 -$
614 $\varepsilon)P_2(x)$, we define $tail_L(\alpha)$ as $\Delta_\varepsilon P_{mixture}((-\infty, \alpha])$ and $tail_R(\beta)$ as $-\Delta_\varepsilon P_{mixture}((\beta, +\infty))$.

Corollary C.1.

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

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

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

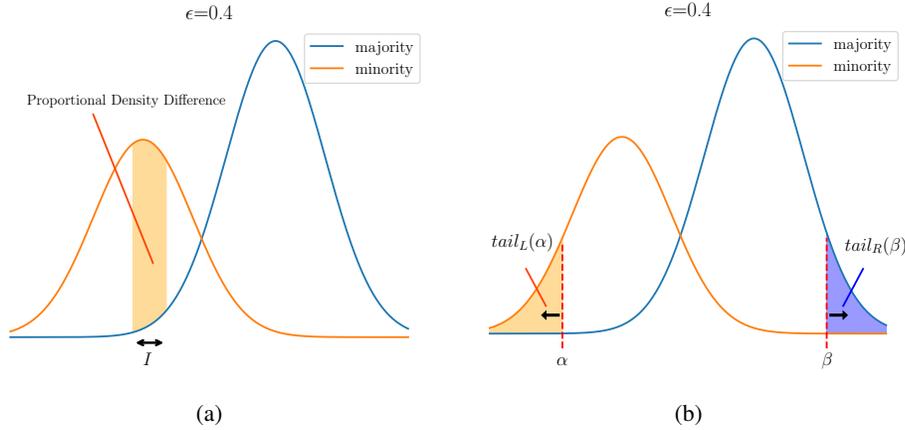


Figure 4: (a) Illustration of proportion density difference C.1, (b) equation of $tail_L(\alpha) = tail_R(\beta)$ at C.2.

616 **Proposition C.1** (Feasibility Of Loss-based Group Balancing). Suppose that L is derived from
617 the mixture of two distributions $\mathcal{N}(\mu_{min}, \sigma_{min}^2)$ and $\mathcal{N}(\mu_{maj}, \sigma_{maj}^2)$ with proportion of ε and $1 - \varepsilon$,
618 respectively, where $\varepsilon \leq \frac{1}{2}$. There exists α and β such that restricting L to the α -left and β -right tails
619 of its distribution results in a group-balanced distribution if and only if $\sigma_{min} \geq \sigma_{maj}$ or

$$tail_L\left(\frac{-B + \sqrt{\Delta}}{2A}\right) > 0 \quad (4)$$

620 and

$$\varepsilon \geq \text{sigmoid}\left(-\frac{(\mu_{maj} - \mu_{min})^2}{2(\sigma_{maj}^2 - \sigma_{min}^2)} - \log\left(\frac{\sigma_{maj}}{\sigma_{min}}\right)\right) \quad (5)$$

621 where $A = \left(\frac{1}{2\sigma_{maj}^2} - \frac{1}{2\sigma_{min}^2}\right)$, $B = \left(\frac{\mu_{min}}{\sigma_{min}^2} - \frac{\mu_{maj}}{\sigma_{maj}^2}\right)$ and $\Delta = \frac{(\mu_{min} - \mu_{maj})^2}{\sigma_{min}^2 \sigma_{maj}^2} - 4\left[\log\left(\frac{\sigma_{maj}}{\sigma_{min}}\right) +$
622 $\log\left(\frac{\varepsilon}{1-\varepsilon}\right)\right]\left[\frac{1}{2\sigma_{maj}^2} - \frac{1}{2\sigma_{min}^2}\right]$.

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

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

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

$$T_{\alpha,\beta} = \begin{cases} High & \text{if } L \geq \beta \\ Mid & \text{if } \alpha < L < \beta \\ Low & \text{if } L \leq \alpha \end{cases} \quad (6)$$

635 We can rewrite the problem in formal form as finding an α and β which satisfies the following
636 equation:

$$P(g = \min | T_{\alpha,\beta} \neq Mid) = P(g = \max | T_{\alpha,\beta} \neq Mid) \quad (7)$$

637 Equation 5 now can be rewritten to a more suitable form:

$$P(g = \min | T_{\alpha,\beta} \neq Mid) = P(g = \max | T_{\alpha,\beta} \neq Mid) \quad (8)$$

$$\Leftrightarrow \frac{P(T_{\alpha,\beta} \neq Mid | g = \min) P(g = \min)}{P(T_{\alpha,\beta} \neq Mid)} = \frac{P(T_{\alpha,\beta} \neq Mid | g = \max) P(g = \max)}{P(T_{\alpha,\beta} \neq Mid)} \quad (9)$$

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

$$\Leftrightarrow \varepsilon P(T_{\alpha,\beta} \neq Mid | g = \min) = (1 - \varepsilon) P(T_{\alpha,\beta} \neq Mid | g = \max) \quad (11)$$

$$\Leftrightarrow \varepsilon \left[P(T_{\alpha,\beta} = Low | g = \min) + P(T_{\alpha,\beta} = High | g = \min) \right] = \quad (12)$$

$$(1 - \varepsilon) \left[P(T_{\alpha,\beta} = Low | g = \max) + P(T_{\alpha,\beta} = High | g = \max) \right] \quad (13)$$

$$\Leftrightarrow \varepsilon \left[P(L \leq \alpha | g = \min) + P(L \geq \beta | g = \min) \right] = \quad (14)$$

$$(1 - \varepsilon) \left[P(L \leq \alpha | g = \max) + P(L \geq \beta | g = \max) \right] \quad (15)$$

$$\Leftrightarrow \varepsilon \left[F^{\min}(\alpha) + (1 - F^{\min}(\beta)) \right] = (1 - \varepsilon) \left[F^{\max}(\alpha) + (1 - F^{\max}(\beta)) \right] \quad (16)$$

$$\Leftrightarrow \varepsilon F^{\min}(\alpha) - (1 - \varepsilon) F^{\max}(\alpha) = (1 - \varepsilon) \left[1 - F^{\max}(\beta) \right] - \varepsilon \left[1 - F^{\min}(\beta) \right] \quad (17)$$

638 We can see the left side of equation 17 is just a function of α . The same goes for the right side of
639 the equation which is a function of β .

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

$$tail_L(\alpha) = tail_R(\beta) \quad (18)$$

642 which is shown in figure 4.

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

644 **Lemma C.2.** $tail_L(\alpha)$ and $tail_R(\beta)$ are continuous functions and $\lim_{\alpha \rightarrow -\infty} tail_L(\alpha) = 0$,
645 $\lim_{\alpha \rightarrow +\infty} tail_L(\alpha) = 2\varepsilon - 1 < 0$, $\lim_{\beta \rightarrow +\infty} tail_R(\beta) = 0$ and $\lim_{\beta \rightarrow -\infty} tail_R(\beta) = 1 - 2\varepsilon > 0$.
646

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

648 **Step 2** Solving the equation 18 for simple cases.

649 **Lemma C.3.** $tail_R(\mu_{maj}) > \frac{1}{2} - \varepsilon \geq 0$

Proof.

$$tail_R(\mu_{maj}) = (1 - \varepsilon) \left[1 - F^{maj}(\mu_{maj}) \right] - \varepsilon \left[1 - F^{min}(\mu_{maj}) \right] \quad (19)$$

$$= (1 - \varepsilon) \left[1 - \phi(0) \right] - \varepsilon \left[1 - \phi\left(\frac{\mu_{maj} - \mu_{min}}{\sigma_{min}}\right) \right] \quad (20)$$

$$> \frac{(1 - \varepsilon)}{2} - \varepsilon \left(1 - \frac{1}{2}\right) = \frac{1 - 2\varepsilon}{2} = \frac{1}{2} - \varepsilon \quad (21)$$

650

□

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

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

$$\frac{dtail_L(\alpha)}{d\alpha} = \varepsilon f^{min}(\alpha) - (1 - \varepsilon) f^{maj}(\alpha) = \varepsilon f^{maj}(\alpha) \left[\frac{f^{min}(\alpha)}{f^{maj}(\alpha)} - \frac{1 - \varepsilon}{\varepsilon} \right] \quad (22)$$

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

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

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

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

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

663 **Step 3** *Solving equation 18 for special case $\sigma_{min} < \sigma_{maj}$* In case of $\sigma_{min} \leq \sigma_{maj}$, having $\Delta > 0$
 664 is a necessary condition, also derivative of $tail_L(\alpha)$ is only positive in $\left(\frac{-b - \sqrt{\Delta}}{2a}, \frac{-b + \sqrt{\Delta}}{2a} \right)$ so the
 665 maximum of $tail_L$ is either in $-\infty$ or in $\frac{-b + \sqrt{\Delta}}{2a}$. Having $tail_L\left(\frac{-b + \sqrt{\Delta}}{2a}\right) > 0$ next to $\Delta > 0$
 666 condition, would be the necessary and also sufficient in this case.

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

$$4AC = \frac{\mu_{min}^2}{\sigma_{min}^4} - \frac{\mu_{min}^2}{\sigma_{maj}^2 \sigma_{min}^2} - \frac{\mu_{maj}^2}{\sigma_{maj}^2 \sigma_{min}^2} + \frac{\mu_{maj}^2}{\sigma_{maj}^4} + 4 \left[\log \left(\frac{\sigma_{maj}}{\sigma_{min}} \right) + \log \left(\frac{\varepsilon}{1 - \varepsilon} \right) \right] \left[\frac{1}{2\sigma_{maj}^2} - \frac{1}{2\sigma_{min}^2} \right]$$

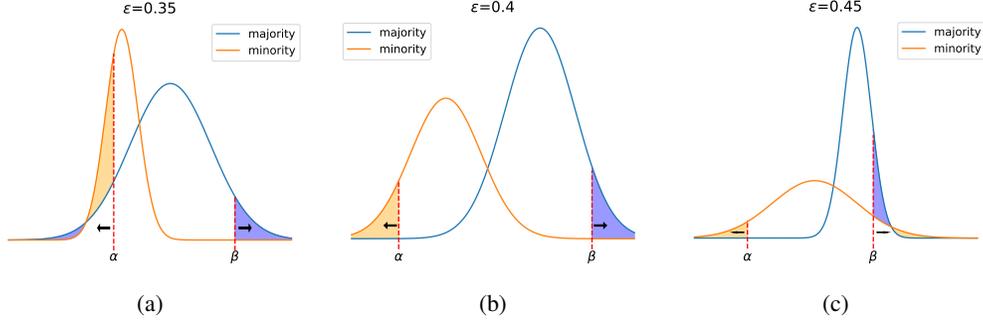


Figure 5: Tail thresholds for three cases: (a) minority group variance is less than majority ($\sigma_{\min} < \sigma_{\text{maj}}$), (b) the variance of two groups are equal ($\sigma_{\min} = \sigma_{\text{maj}}$) and (c) the variance of the minority group is more than majority ($\sigma_{\min} > \sigma_{\text{maj}}$).

$$\begin{aligned} \Delta &= \frac{(\mu_{\min} - \mu_{\text{maj}})^2}{\sigma_{\min}^2 \sigma_{\text{maj}}^2} - 4 \left[\log \left(\frac{\sigma_{\text{maj}}}{\sigma_{\min}} \right) + \log \left(\frac{\epsilon}{1 - \epsilon} \right) \right] \left[\frac{1}{2\sigma_{\text{maj}}^2} - \frac{1}{2\sigma_{\min}^2} \right] \geq 0 \\ &\iff (\mu_{\min} - \mu_{\text{maj}})^2 \geq 2 \left[\log \left(\frac{1 - \epsilon}{\epsilon} \right) - \log \left(\frac{\sigma_{\text{maj}}}{\sigma_{\min}} \right) \right] [\sigma_{\text{maj}}^2 - \sigma_{\min}^2] \\ &\iff \epsilon \geq \text{sigmoid} \left(- \frac{(\mu_{\text{maj}} - \mu_{\min})^2}{2(\sigma_{\text{maj}}^2 - \sigma_{\min}^2)} - \log \left(\frac{\sigma_{\text{maj}}}{\sigma_{\min}} \right) \right) \end{aligned}$$

667 Next, we investigate properties of the conditions of the proposition C.1 in case of $\sigma_{\text{maj}} < \sigma_{\min}$.
 668 Schematic interpretation of these conditions is presented in figure 6.

- 669 • As equation 5 indicates, the minority group is not allowed to be too underrepresented. This
 670 especially has a direct relation with the difference of means. The more mean values of
 671 groups are different, the more imbalance can be mitigated through loss-based sampling.
 672 Mean value difference is especially affected by the spurious correlation, it escalates as the
 673 model relies on spurious correlation and also when the spurious features between groups are
 674 too different.
- 675 • On the other hand condition 4 is more complex and doesn't have a simple closed form, we
 676 analytically describe its behaviors by fixating the means and calculating the valid values for
 677 ϵ . As the results show in figure 6, most of ϵ are feasible in for $\sigma_{\min} < \Delta\mu$ as we can see the
 678 possible region declines with an increase of σ_{\min} and valid ϵ values cease to exist.

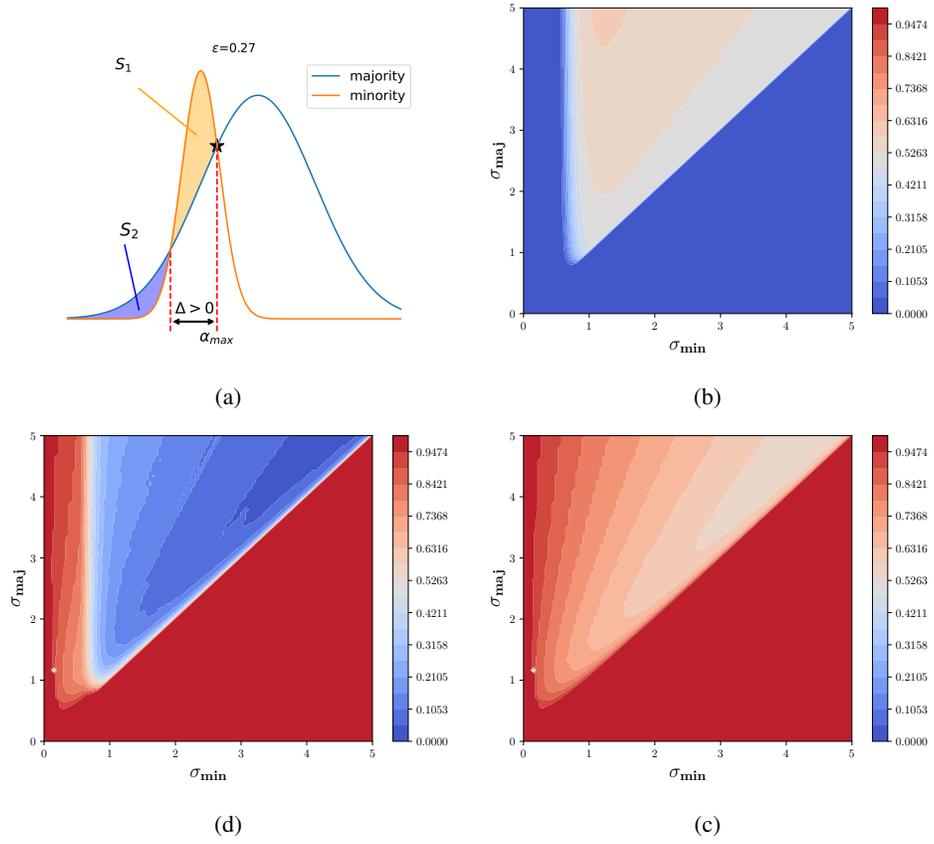


Figure 6: (a) Conditions if $\sigma_{min} > \sigma_{maj}$, (b), (c), (d) Minimum, maximum and interval length of feasible ϵ values across $(\sigma_{min}, \sigma_{maj})$ field for $\mu_{min} = 0, \mu_{maj} = 1$.

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 \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 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 do not require group annotation in the training phase.

Method	Group Info Train/Val	Waterbirds		CelebA		UrbanCars	
		Worst	Best	Worst	Best	Worst	Best
GDRO [7]	\checkmark/\checkmark	91.4	93.5	88.9	92.9	-	-
DFR [1]	\times/\checkmark	<u>92.9\pm0.2</u>	94.2 \pm 0.4	88.3 \pm 1.1	91.3 \pm 0.3	79.6 \pm 2.22	87.5 \pm 0.6
GDRO + EIIL [12]	\times/\checkmark	77.2 \pm 1	<u>96.5\pm0.2</u>	81.7 \pm 0.8	85.7 \pm 0.1	-	-
JTT [5]	\times/\checkmark	86.7	93.3	81.1	88.0	-	-
AFR [9]	\times/\checkmark	90.4 \pm 1.1	94.2 \pm 1.2	82.0 \pm 0.5	91.3 \pm 0.3	80.2 \pm 2.0	87.1 \pm 1.2
EVaLS-GL (Ours)	\times/\checkmark	89.4 \pm 0.3	95.1 \pm 0.3	84.6 \pm 1.6	91.1 \pm 0.6	<u>82.27\pm1.16</u>	<u>88.2\pm0.6</u>
ERM	\times/\times	66.4 \pm 2.3	90.3 \pm 0.5	47.4 \pm 2.3	<u>95.5\pm0.0</u>	18.67 \pm 2.01	76.5 \pm 4.6
EVaLS (Ours)	\times/\times	88.4 \pm 3.1	94.1 \pm 0.1	<u>85.3\pm0.4</u>	89.4 \pm 0.5	82.13 \pm 0.92	88.1 \pm 0.9

Table 5: A comparison of the various methods, ours included, on CivilComments and MultiNLI. The Group Info column indicates if each method utilizes group labels of the training/validation data, 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 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 do not require group annotation in the training phase.

Method	Group Info Train/Val	CivilComments		MultiNLI	
		Worst	Best	Worst	Best
GDRO [7]	\checkmark/\checkmark	69.9	88.9	77.7	81.4
DFR [1]	\times/\checkmark	70.1 \pm 0.8	87.2 \pm 0.3	74.7 \pm 0.7	82.1 \pm 0.2
GDRO + EIIL [12]	\times/\checkmark	67.0 \pm 2.4	90.5 \pm 0.2	-	-
JTT [5]	\times/\checkmark	69.3	91.1	72.6	78.6
AFR [9]	\times/\checkmark	68.7 \pm 0.6	89.8 \pm 0.6	73.4 \pm 0.6	81.4 \pm 0.2
EVaLS-GL (Ours)	\times/\checkmark	<u>80.5\pm0.4</u>	88.0 \pm 0.4	<u>75.1\pm1.2</u>	81.6 \pm 0.2
ERM	\times/\times	61.2 \pm 3.6	<u>92.0\pm0.0</u>	64.8 \pm 1.9	<u>82.6\pm0.0</u>

679 D Experimental Details

680 D.1 Complete Results

681 The complete results on Waterbirds, CelebA, and UrbanCars, in addition to complete results on
 682 CivilComments and MultiNLI are reported in Tables 4 and 5 respectively. The results for all methods
 683 except Group DRO + EIIL on all datasets except UrbanCars are reported by Qiu et al. [9]. The
 684 results for Group DRO + EIIL are taken from Zhang et al. [24]. Also, the results of our method and
 685 DFR are shown in Table 6

686 D.2 Dominoes-Colored-MNIST-FashionMNIST

687 **Dominoes-Colored-MNIST-FashionMNIST (Dominoes-CMF)** is a synthetic dataset. We adopt
 688 a similar approach to previous works [37, 38, 1] using a modified version of the *Dominoes* binary
 689 classification dataset. This dataset consists of images with the top half showing CIFAR-10 images
 690 [19], divided into two meaningful classes: vehicles (airplane, car, ship, truck) and animals (cat,
 691 dog, horse, deer). The bottom half displays either MNIST [20] images from classes $\{0 - 3\}$ or
 692 Fashion-MNIST [21] images from classes $\{T\text{-shirt, Dress, Coat, Shirt}\}$. The complex feature (top

Table 6: A Comparison of ERM, DFR, EVaLS, and EVaLS-GL on the Dominoes-CMF Dataset. Both the worst and average of test group accuracies are presented. The mean and standard deviation are calculated based on runs with three distinct seeds.

Method	Worst	Average
ERM	50.6 \pm 1.0	84.1 \pm 0.0
DFR	60.2 \pm 1.2	84.6 \pm 0.4
EVaLS-GL	63.6 \pm 1.3	78.7 \pm 1.5
EVaLS	67.1\pm4.2	78.6 \pm 2.0

693 half) serves as the core feature and the simple feature (bottom half) is linearly separable and correlated
 694 with the class label at 75%. Furthermore, inspired by the approaches in Zhang et al. [24], Arjovsky
 695 et al. [18], we intentionally introduce an additional spurious attribute by artificially coloring a subset
 696 of images in the following manner: 90% of the bottom half images in class c_1 are randomly assigned
 697 a red color, while 10% are assigned a green color, and vice versa for class c_2 . See Table 7 for more
 details about the dataset statistics.

Table 7: *Dominoes-CMF* Dataset Statistics

CIFAR-10 Class	Top part	Bottom part	
	Color	MNIST	Fashion-MNIST
c_1 (Vehicle)	Red	13,500	4,500
	Green	1,500	500
c_2 (Animal)	Red	500	1,500
	Green	4,500	13,500
Total		40,000	

698

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

CIFAR-10 Class	Top part	Bottom part	
	Color	MNIST	Fashion-MNIST
c_1 (Vehicle)	Red	99.2 \pm 0.01%	95.2 \pm 1.1%
	Green	84.5 \pm 2.4%	54.7 \pm 0.5%
c_2 (Animal)	Red	56.8 \pm 5.6%	86.7 \pm 2.4%
	Green	96.2 \pm 0.5%	99.3 \pm 0.2%

699 D.3 Datasets

700 **Waterbirds [7]** The dataset comprises images of diverse bird species, classified into two categories:
 701 waterbirds and landbirds. Each image features a bird set against a backdrop of either water or land.
 702 Interestingly, the background scene acts as a spurious feature in this classification task. Waterbirds are
 703 primarily shown against water backgrounds, and landbirds against land backgrounds. Consequently,
 704 waterbirds on water and landbirds on land form the minority groups in the training data. It’s important
 705 to note that the validation dataset for waterbirds is group-balanced, meaning birds from each class are
 706 equally represented against both water and land backgrounds. This dataset is mainly categorized as a
 707 spurious correlation dataset.

708 **CelebA [13]** is a widely used dataset in image classification tasks, featuring annotations for 40
 709 binary facial attributes such as hair color, gender, and age. Hair color classification is particularly
 710 prominent in literature focusing on spurious correlation robustness. Notably, gender serves as a
 711 spurious attribute within this dataset, where a significant majority 94% of individuals with blond hair

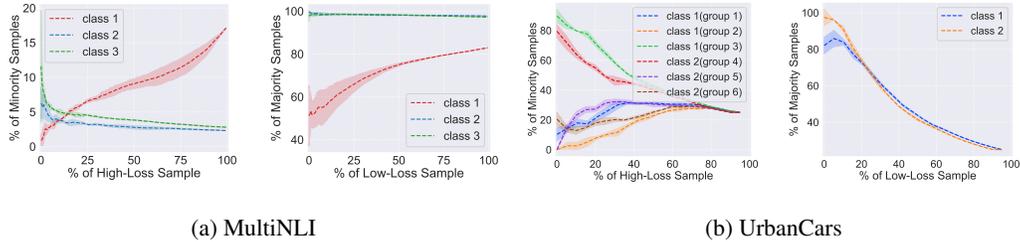


Figure 7: The proportion of minority and majority samples across different classes within various percentages of \mathcal{D}^{LL} samples with highest and lowest loss for the MultiNLI (a) and UrbanCars (b) datasets. MultiNLI exhibits attribute imbalance rather than spurious correlation, which explains its different behavior compared to Waterbirds and CelebA.

712 are women, while men with blond hair represent a minority group. In addition to spurious correlation
 713 in the class of blond hair, this dataset also exhibits class imbalance.

714 **MultiNLI [15]** dataset involves a text classification task focused on determining the relationship
 715 between pairs of sentences: contradiction, entailment, or neutral. Sentences containing negation
 716 words such as "no" or "never" are under-represented in all three classes, inducing attribute imbalance
 717 in the dataset. Figure 7 illustrates the distinct behavior of this dataset compared to other datasets that
 718 contain spurious attributes.

719 **CivilComments [16]** dataset, as part of the WILDS benchmark, involves a text classification task
 720 focused on labeling online comments as either "toxic" or "not toxic". Each comment is associated
 721 with 8 attributes, including gender (male, female), sexual orientation (LGBTQ), race (black, white),
 722 and religion (Christian, Muslim, or other), based on whether these characteristics are mentioned
 723 in the comment. While there is a small attribute imbalance in the dataset, it can be categorized into
 724 datasets with class imbalance. In this paper, we use the implementation of the dataset by the WILDS
 725 package [39].

726 **UrbanCars [14]** is an image classification dataset with multiple shortcuts. Each image in the
 727 dataset consists of a car in the center of the image on a natural scene background, with another object
 728 to the right of the image. Images are labeled *Urban* or *City* according to the type of car present in
 729 the center. However, each of the backgrounds and the additional objects is highly correlated with
 730 the label. While the test set consists of 8 environments based on combinations of the core and two
 731 spurious patterns, the training and validation set consist of four groups, based on combinations of the
 732 label and only one of the shortcuts.

733 D.4 Training Details

734 **ERM** For Waterbirds and CelebA, we utilize the ResNet50 checkpoints available in
 735 the GitHub repository of Kirichenko et al. [1] as our base model. We use the
 736 ResNet-50 architecture provided by the torchvision package. In the case of Civil-
 737 Comments and MultiNLI, we adopt a similar approach to Kirichenko et al. [1], using
 738 `BertForSequenceClassification.from_pretrained('bert-base-uncased', ...)` from
 739 the transformers package. The model is trained using the AdamW optimizer with a learning
 740 rate of 10^{-5} , weight decay of 10^{-4} , and a batch size of 16 for a total of 5 epochs.

741 For the UrbanCars dataset, we adhere to the settings described in Li et al. [14], which involves
 742 training a ResNet-50 model pretrained on ImageNet using the SGD optimizer with a learning rate
 743 of 10^{-3} , momentum of 0.9, weight decay of 10^{-4} , and a batch size of 128 for 300 epochs. For the
 744 Dominoes-CMF dataset, we train a ResNet18 model pretrained on ImageNet for 20 epochs with a
 745 batch size of 128 and an SGD optimizer with a learning rate of 10^{-3} , momentum of 0.9, and weight
 746 decay of 10^{-4} .

747 **EValS and EValS-GL** For every dataset, EIL was utilized with a learning rate of 0.01, a total of
 748 20000 steps, and a batch size of 128. The last layer of the model was trained on all datasets using the

749 Adam optimizer. A batch size of 32 and a weight decay of 10^{-4} were used for all datasets. Our method
750 was evaluated on the validation sets of each dataset, considering both fine-tuning and retraining of the
751 last layer. For all datasets, with the exception of MultiNLI, retraining provided superior validation
752 results. The specifics regarding the number of epochs and the ranges for hyperparameter search
753 (including learning rate, l_1 -regularization coefficient (λ), and the number of selected samples (k)) for
754 each dataset are as follows:

755 • **Waterbirds.**

- 756 – epochs = 100,
- 757 – lr = 5×10^{-4} ,
- 758 – $\lambda \in \{0, 0.01, 0.02, 0.03, 0.04, 0.05, 0.06, 0.07, 0.08, 0.09, 0.1, 0.2, 0.3, 0.4, 0.5\}$,
- 759 – $k \in \{20, 25, 30, 35, 40, 45, 50, 55, 60\}$.

760 • **CelebA**

- 761 – epochs = 50,
- 762 – lr = 5×10^{-4} ,
- 763 – $\lambda \in \{0, 0.01, 0.02, 0.03, 0.04, 0.05, 0.06, 0.07, 0.08, 0.09, 0.1, 0.2, 0.3, 0.4, 0.5,$
764 $0.6, 0.7, 0.8, 0.9, 1, 2\}$,
- 765 – $k \in \{50, 100, 150, 200, 250, 300\}$.

766 • **UrbanCars**

- 767 – epochs = 100,
- 768 – lr $\in \{5 \times 10^{-4}, 10^{-3}\}$,
- 769 – $\lambda \in \{0, 0.01, 0.02, 0.05, 0.1, 1\}$,
- 770 – $k \in \{10, 20, 30, 50, 63\}$.

771 • **CivilComments**

- 772 – epochs = 50,
- 773 – lr = 5×10^{-4} ,
- 774 – $\lambda \in \{0, 0.01, 0.02, 0.03, 0.04, 0.05, 0.06, 0.07, 0.08, 0.09, 0.1, 0.2, 0.3, 0.4, 0.5,$
775 $0.6, 0.7, 0.8, 0.9, 1, 2\}$,
- 776 – $k \in \{500, 750, 1000, 1250, 1500\}$.

777 • **MultiNLI**

- 778 – epochs = 200,
- 779 – lr $\in \{10^{-2}, 10^{-3}\}$,
- 780 – $\lambda \in \{0, 0.01, 0.02, 0.03, 0.04, 0.05, 0.06, 0.07, 0.08, 0.09, 0.1, 0.2, 0.3, 0.4, 0.5\}$,
- 781 – $k \in \{20, 30, 40, 50, 60, 75, 100, 125, 150, 200, 250, 300\}$.

782 **E Ablation Study**

783 **E.1 Use of EIIL with DFR and AFR**

784 We conducted an ablation study to investigate the impact of using environments inferred from EIIL on
 785 model selection. Specifically, we benchmarked the performance of DFR and AFR with EIIL-inferred
 786 groups. The results, presented in Table 9, demonstrate the effectiveness of incorporating EIIL-inferred
 787 groups in model selection. The results show that while EIIL-inferred groups reduce the performance
 788 compared to ground-truth annotations for model selection, they still can be effective for robustness to
 789 an extent. Moreover, EVaLS outperforms these two methods when using EIIL inferred environments.

Table 9: Results of DFR and AFR with EIIL-inferred environment for model selection.

Method	Waterbirds	Celeba
DFR (with EIIL)	92.21 ± 0.02	85.55 ± 1.0
AFR (with EIIL)	82.6 ± 0.04	72.5 ± 0.01

790 **E.2 Other Group Inference Methods**

791 In addition to EIIL, other group inference methods could be utilized for partitioning the model
 792 selection set into environments.

793 **Error Splitting** JTT [5] partitions data into two correctly classified and misclassified sets based
 794 on the predictions of a model trained with ERM. We split each of these two sets based on labels of
 795 samples, obtaining $|\mathcal{D}| \times 2$ environments.

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

800 The results for EVaLS-ES (EVaLS+Error Sampling) and EVaLS-RC (EVaLS+Random Classifier) are
 801 shown in Table 10. One limitation of error splitting is that in datasets with noisy labels or corrupted
 802 images, samples that an ERM model misclassifies may not always belong to minority groups. In these
 803 situations, choosing models based on their accuracy on corrupted data could lead to the selection of
 804 models that are not robust to spurious correlations. This is demonstrated by the results of EVaLS-ES
 805 on the CelebA dataset.

806 This shortcoming of error splitting can be alleviated by employing a random classifier instead of
 807 the ERM-trained one. Due to the feature-level similarity between minority and majority samples in
 808 datasets affected by spurious correlation [23, 1, 29], it is expected that the classifier can differentiate
 809 between the groups to some extent. As shown in Table 10, surprisingly, EVaLS-RC produces results
 810 that are generally comparable to EVaLS. However, the performance of this method may have high
 811 variance, depending on the different initializations of the classifier.

Table 10: The performances of three environment inference methods, when combined with loss-based sample selection, are evaluated on spurious correlation benchmarks. The mean and standard deviation values are calculated over three separate runs, each initiated with a different seed.

Method	Waterbirds		CelebA		UrbanCars	
	Worst	Average	Worst	Average	Worst	Average
EVaLS-ES	82.1 \pm 1.2	94.3\pm0.04	48.4 \pm 11.6	69.5 \pm 6.5	79.2 \pm 2.9	86.1 \pm 0.9
EVaLS-RC	88.7\pm1.0	94.3 \pm 1.1	78.1 \pm 5.1	93.5\pm0.2	82.4\pm3.2	88.2\pm0.8
EVaLS	88.4 \pm 3.1	94.1 \pm 0.1	85.3\pm0.4	89.4 \pm 0.5	79.4 \pm 3.1	86.5 \pm 1.5

812 **F Societal Impacts**

813 Real-world datasets often encapsulate social biases that stem from entrenched stereotypes and
814 historical discrimination, affecting various groups such as genders and races. Machine learning
815 methods, which learn the correlation between patterns in input data and their targets (e.g., labels
816 in a classification task) [40], inadvertently absorb this bias. This unintended consequence leads to
817 fairness issues in many applications. While strategies to mitigate such biases have been proposed
818 (as discussed comprehensively in Section A), societal biases are not always known and determined.
819 We believe that our work, as it addresses these unidentified biases, takes a significant step towards
820 making machine learning fairer for our society.

821 **G Computational Resources**

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

825 **NeurIPS Paper Checklist**

826 **1. Claims**

827 Question: Do the main claims made in the abstract and introduction accurately reflect the
828 paper's contributions and scope?

829 Answer: [\[Yes\]](#)

830 Justification: The scope of the effectiveness and main claims are clearly demonstrated in the
831 abstract and introduction.

832 Guidelines:

- 833 • The answer NA means that the abstract and introduction do not include the claims
834 made in the paper.
- 835 • The abstract and/or introduction should clearly state the claims made, including the
836 contributions made in the paper and important assumptions and limitations. A No or
837 NA answer to this question will not be perceived well by the reviewers.
- 838 • The claims made should match theoretical and experimental results, and reflect how
839 much the results can be expected to generalize to other settings.
- 840 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
841 are not attained by the paper.

842 **2. Limitations**

843 Question: Does the paper discuss the limitations of the work performed by the authors?

844 Answer: [\[Yes\]](#)

845 Justification: The limitations are discussed in the discussion section.

846 Guidelines:

- 847 • The answer NA means that the paper has no limitation while the answer No means that
848 the paper has limitations, but those are not discussed in the paper.
- 849 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 850 • The paper should point out any strong assumptions and how robust the results are to
851 violations of these assumptions (e.g., independence assumptions, noiseless settings,
852 model well-specification, asymptotic approximations only holding locally). The authors
853 should reflect on how these assumptions might be violated in practice and what the
854 implications would be.
- 855 • The authors should reflect on the scope of the claims made, e.g., if the approach was
856 only tested on a few datasets or with a few runs. In general, empirical results often
857 depend on implicit assumptions, which should be articulated.
- 858 • The authors should reflect on the factors that influence the performance of the approach.
859 For example, a facial recognition algorithm may perform poorly when image resolution
860 is low or images are taken in low lighting. Or a speech-to-text system might not be
861 used reliably to provide closed captions for online lectures because it fails to handle
862 technical jargon.
- 863 • The authors should discuss the computational efficiency of the proposed algorithms
864 and how they scale with dataset size.
- 865 • If applicable, the authors should discuss possible limitations of their approach to
866 address problems of privacy and fairness.
- 867 • While the authors might fear that complete honesty about limitations might be used by
868 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
869 limitations that aren't acknowledged in the paper. The authors should use their best
870 judgment and recognize that individual actions in favor of transparency play an impor-
871 tant role in developing norms that preserve the integrity of the community. Reviewers
872 will be specifically instructed to not penalize honesty concerning limitations.

873 **3. Theory Assumptions and Proofs**

874 Question: For each theoretical result, does the paper provide the full set of assumptions and
875 a complete (and correct) proof?

876 Answer: [\[Yes\]](#)

877 Justification: All the lemmas and propositions are stated upon exact definitions, assumptions
878 and conditions. All the theorems, formulas, and proofs in the paper are numbered and
879 cross-referenced.

880 Guidelines:

- 881 • The answer NA means that the paper does not include theoretical results.
- 882 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
883 referenced.
- 884 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 885 • The proofs can either appear in the main paper or the supplemental material, but if
886 they appear in the supplemental material, the authors are encouraged to provide a short
887 proof sketch to provide intuition.
- 888 • Inversely, any informal proof provided in the core of the paper should be complemented
889 by formal proofs provided in appendix or supplemental material.
- 890 • Theorems and Lemmas that the proof relies upon should be properly referenced.

891 4. Experimental Result Reproducibility

892 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
893 perimental results of the paper to the extent that it affects the main claims and/or conclusions
894 of the paper (regardless of whether the code and data are provided or not)?

895 Answer: [Yes]

896 Justification: The training procedure is described accurately and all the training details and
897 hyperparameters required for reproducing the results are provided.

898 Guidelines:

- 899 • The answer NA means that the paper does not include experiments.
- 900 • If the paper includes experiments, a No answer to this question will not be perceived
901 well by the reviewers: Making the paper reproducible is important, regardless of
902 whether the code and data are provided or not.
- 903 • If the contribution is a dataset and/or model, the authors should describe the steps taken
904 to make their results reproducible or verifiable.
- 905 • Depending on the contribution, reproducibility can be accomplished in various ways.
906 For example, if the contribution is a novel architecture, describing the architecture fully
907 might suffice, or if the contribution is a specific model and empirical evaluation, it may
908 be necessary to either make it possible for others to replicate the model with the same
909 dataset, or provide access to the model. In general, releasing code and data is often
910 one good way to accomplish this, but reproducibility can also be provided via detailed
911 instructions for how to replicate the results, access to a hosted model (e.g., in the case
912 of a large language model), releasing of a model checkpoint, or other means that are
913 appropriate to the research performed.
- 914 • While NeurIPS does not require releasing code, the conference does require all submis-
915 sions to provide some reasonable avenue for reproducibility, which may depend on the
916 nature of the contribution. For example
 - 917 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
918 to reproduce that algorithm.
 - 919 (b) If the contribution is primarily a new model architecture, the paper should describe
920 the architecture clearly and fully.
 - 921 (c) If the contribution is a new model (e.g., a large language model), then there should
922 either be a way to access this model for reproducing the results or a way to reproduce
923 the model (e.g., with an open-source dataset or instructions for how to construct
924 the dataset).
 - 925 (d) We recognize that reproducibility may be tricky in some cases, in which case
926 authors are welcome to describe the particular way they provide for reproducibility.
927 In the case of closed-source models, it may be that access to the model is limited in
928 some way (e.g., to registered users), but it should be possible for other researchers
929 to have some path to reproducing or verifying the results.

930 5. Open access to data and code

931 Question: Does the paper provide open access to the data and code, with sufficient instruc-
932 tions to faithfully reproduce the main experimental results, as described in supplemental
933 material?

934 Answer: [Yes]

935 Justification: Codes and information of datasets that are constructed or reused in the paper
936 are anonymized and included in the main paper and supplementary material.

937 Guidelines:

- 938 • The answer NA means that paper does not include experiments requiring code.
- 939 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
940 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 941 • While we encourage the release of code and data, we understand that this might not be
942 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
943 including code, unless this is central to the contribution (e.g., for a new open-source
944 benchmark).
- 945 • The instructions should contain the exact command and environment needed to run to
946 reproduce the results. See the NeurIPS code and data submission guidelines ([https://
947 nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 948 • The authors should provide instructions on data access and preparation, including how
949 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 950 • The authors should provide scripts to reproduce all experimental results for the new
951 proposed method and baselines. If only a subset of experiments are reproducible, they
952 should state which ones are omitted from the script and why.
- 953 • At submission time, to preserve anonymity, the authors should release anonymized
954 versions (if applicable).
- 955 • Providing as much information as possible in supplemental material (appended to the
956 paper) is recommended, but including URLs to data and code is permitted.

957 6. Experimental Setting/Details

958 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
959 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
960 results?

961 Answer: [Yes]

962 Justification: The training details, hyperparameters, model selection criteria, etc. have are
963 written in the Appendix and the data and metadata have been provided in our code.

964 Guidelines:

- 965 • The answer NA means that the paper does not include experiments.
- 966 • The experimental setting should be presented in the core of the paper to a level of detail
967 that is necessary to appreciate the results and make sense of them.
- 968 • The full details can be provided either with the code, in appendix, or as supplemental
969 material.

970 7. Experiment Statistical Significance

971 Question: Does the paper report error bars suitably and correctly defined or other appropriate
972 information about the statistical significance of the experiments?

973 Answer: [Yes]

974 Justification: All tables report standard deviation and how it was computed and the plot
975 contains error bar (also by standard deviation).

976 Guidelines:

- 977 • The answer NA means that the paper does not include experiments.
- 978 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
979 dence intervals, or statistical significance tests, at least for the experiments that support
980 the main claims of the paper.

- 981 • The factors of variability that the error bars are capturing should be clearly stated (for
982 example, train/test split, initialization, random drawing of some parameter, or overall
983 run with given experimental conditions).
- 984 • The method for calculating the error bars should be explained (closed form formula,
985 call to a library function, bootstrap, etc.)
- 986 • The assumptions made should be given (e.g., Normally distributed errors).
- 987 • It should be clear whether the error bar is the standard deviation or the standard error
988 of the mean.
- 989 • It is OK to report 1-sigma error bars, but one should state it. The authors should
990 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
991 of Normality of errors is not verified.
- 992 • For asymmetric distributions, the authors should be careful not to show in tables or
993 figures symmetric error bars that would yield results that are out of range (e.g. negative
994 error rates).
- 995 • If error bars are reported in tables or plots, The authors should explain in the text how
996 they were calculated and reference the corresponding figures or tables in the text.

987 8. Experiments Compute Resources

988 Question: For each experiment, does the paper provide sufficient information on the com-
989 puter resources (type of compute workers, memory, time of execution) needed to reproduce
1000 the experiments?

1001 Answer: [No]

1002 Justification: The paper does provide details about the hardware used for the experiments.
1003 However, since experiments were done on different hardwares, the computational resources
1004 needed for each individual experiment are not documented.

1005 Guidelines:

- 1006 • The answer NA means that the paper does not include experiments.
- 1007 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
1008 or cloud provider, including relevant memory and storage.
- 1009 • The paper should provide the amount of compute required for each of the individual
1010 experimental runs as well as estimate the total compute.
- 1011 • The paper should disclose whether the full research project required more compute
1012 than the experiments reported in the paper (e.g., preliminary or failed experiments that
1013 didn't make it into the paper).

1014 9. Code Of Ethics

1015 Question: Does the research conducted in the paper conform, in every respect, with the
1016 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

1017 Answer: [Yes]

1018 Justification: All codes and rules have been thoroughly reviewed and checked, with no
1019 instances of non-compliance found.

1020 Guidelines:

- 1021 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 1022 • If the authors answer No, they should explain the special circumstances that require a
1023 deviation from the Code of Ethics.
- 1024 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
1025 eration due to laws or regulations in their jurisdiction).

1026 10. Broader Impacts

1027 Question: Does the paper discuss both potential positive societal impacts and negative
1028 societal impacts of the work performed?

1029 Answer: [Yes]

1030 Justification: In the Social Impacts section we discuss that our work can significantly
1031 contribute to fairness in machine learning. We did not find any negative social impacts of
1032 our work.

1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The paper poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: Every asset that we utilized for our implementations have been appropriately referenced, both within the paper itself and in the code (if needed). Although we did not specify the names of their respective licenses, you can find these details on the webpages we've cited.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.

- 1085 • The authors should state which version of the asset is used and, if possible, include a
1086 URL.
- 1087 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 1088 • For scraped data from a particular source (e.g., website), the copyright and terms of
1089 service of that source should be provided.
- 1090 • If assets are released, the license, copyright information, and terms of use in the
1091 package should be provided. For popular datasets, paperswithcode.com/datasets
1092 has curated licenses for some datasets. Their licensing guide can help determine the
1093 license of a dataset.
- 1094 • For existing datasets that are re-packaged, both the original license and the license of
1095 the derived asset (if it has changed) should be provided.
- 1096 • If this information is not available online, the authors are encouraged to reach out to
1097 the asset's creators.

13. New Assets

1099 Question: Are new assets introduced in the paper well documented and is the documentation
1100 provided alongside the assets?

1101 Answer: [NA]

1102 Justification: The paper does not release new assets.

1103 Guidelines:

- 1104 • The answer NA means that the paper does not release new assets.
- 1105 • Researchers should communicate the details of the dataset/code/model as part of their
1106 submissions via structured templates. This includes details about training, license,
1107 limitations, etc.
- 1108 • The paper should discuss whether and how consent was obtained from people whose
1109 asset is used.
- 1110 • At submission time, remember to anonymize your assets (if applicable). You can either
1111 create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and Research with Human Subjects

1113 Question: For crowdsourcing experiments and research with human subjects, does the paper
1114 include the full text of instructions given to participants and screenshots, if applicable, as
1115 well as details about compensation (if any)?

1116 Answer: [NA]

1117 Justification: The paper does not involve crowdsourcing nor research with human subjects.

1118 Guidelines:

- 1119 • The answer NA means that the paper does not involve crowdsourcing nor research with
1120 human subjects.
- 1121 • Including this information in the supplemental material is fine, but if the main contribu-
1122 tion of the paper involves human subjects, then as much detail as possible should be
1123 included in the main paper.
- 1124 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
1125 or other labor should be paid at least the minimum wage in the country of the data
1126 collector.

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

1129 Question: Does the paper describe potential risks incurred by study participants, whether
1130 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
1131 approvals (or an equivalent approval/review based on the requirements of your country or
1132 institution) were obtained?

1133 Answer: [NA]

1134 Justification: The paper does not involve crowdsourcing nor research with human subjects.

1135 Guidelines:

- 1136 • The answer NA means that the paper does not involve crowdsourcing nor research with
1137 human subjects.
- 1138 • Depending on the country in which research is conducted, IRB approval (or equivalent)
1139 may be required for any human subjects research. If you obtained IRB approval, you
1140 should clearly state this in the paper.
- 1141 • We recognize that the procedures for this may vary significantly between institutions
1142 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
1143 guidelines for their institution.
- 1144 • For initial submissions, do not include any information that would break anonymity (if
1145 applicable), such as the institution conducting the review.