

Training Debiased Subnetworks with Contrastive Weight Pruning

Geon Yeong Park¹ Sangmin Lee² Sang Wan Lee^{1*} Jong Chul Ye^{1,2,3*}
¹Bio and Brain Engineering, ²Mathematical Sciences, ³Kim Jaechul Graduate School of AI
Korea Advanced Institute of Science and Technology (KAIST), Daejeon, Korea
{pky3436, leeleesang, sangwan, jong.ye}@kaist.ac.kr

Abstract

Neural networks are often biased to spuriously correlated features that provide misleading statistical evidence that does not generalize. This raises an interesting question: “Does an optimal unbiased functional subnetwork exist in a severely biased network? If so, how to extract such subnetwork?” While empirical evidence has been accumulated about the existence of such unbiased subnetworks, these observations are mainly based on the guidance of ground-truth unbiased samples. Thus, it is unexplored how to discover the optimal subnetworks with biased training datasets in practice. To address this, here we first present our theoretical insight that alerts potential limitations of existing algorithms in exploring unbiased subnetworks in the presence of strong spurious correlations. We then further elucidate the importance of bias-conflicting samples on structure learning. Motivated by these observations, we propose a Debiased Contrastive Weight Pruning (DCWP) algorithm, which probes unbiased subnetworks without expensive group annotations. Experimental results demonstrate that our approach significantly outperforms state-of-the-art debiasing methods despite its considerable reduction in the number of parameters.

1. Introduction

While deep neural networks have made substantial progress in solving challenging tasks, they often undesirably rely on spuriously correlated features or dataset bias, if present, which is considered one of the major hurdles in deploying models in real-world applications. For example, consider recognizing desert foxes and cats from natural images. If the background scene (e.g., a desert) is spuriously correlated to the type of animal, the neural networks might use the background information as a shortcut to classification, resulting in performance degradation in different backgrounds (e.g., a desert fox in the house).

To investigate the origin of the spurious correlations, this paper considers shortcut learning as a fundamental architec-

tural design issue of neural networks. Specifically, if any available information channels in deep networks’ structure could transmit the information of spuriously correlated features (*spurious features* from now on), networks would exploit those features as long as they are sufficiently predictive. It naturally follows that pruning weights on spurious features can purify the biased latent representations, thereby improving performances on bias-conflicting samples¹. We conjecture that this neural pruning may improve the generalization of the network in a way that reduces the effective dimension of spurious features, considering that the failure of Out-of-Distribution (OOD) generalization may arise due to high-dimensional spurious features [26, 34].

Recently, Zhang et al. [37] has empirically demonstrated the existence of subnetworks that are less susceptible to spurious features. Based on the modular property of neural networks [5], they prune out weights that are closely related to the spurious attributes. While [37] affords us valuable insights on the importance of neural architectures, the study has limitation in that such neural pruning requires sufficient number of ground-truth bias-conflicting samples. Thus, *how to discover the optimal subnetworks in practice when the dataset is highly biased?*

To address this, we first present a simple theoretical observation that reveals the limitations of existing substructure probing methods in searching unbiased subnetworks. Specifically, we reveal that there exists an unavoidable generalization gap in the subnetworks obtained by standard pruning algorithms in the presence of strong spurious correlations. Our analysis also shows that trained models may inevitably rely on the spuriously correlated features in a practical training setting with finite training time and a number of samples.

In addition, we show that sampling more bias-conflicting data makes it possible to identify spurious weights. Specifically, bias-conflicting samples require that the weights as-

¹The *bias-aligned* samples refer to data with a strong correlation between (potentially latent) spurious features and target labels (e.g., cat in the house). The *bias-conflicting* samples refer to the opposite cases where spurious correlations do not exist (e.g., cat in the desert).

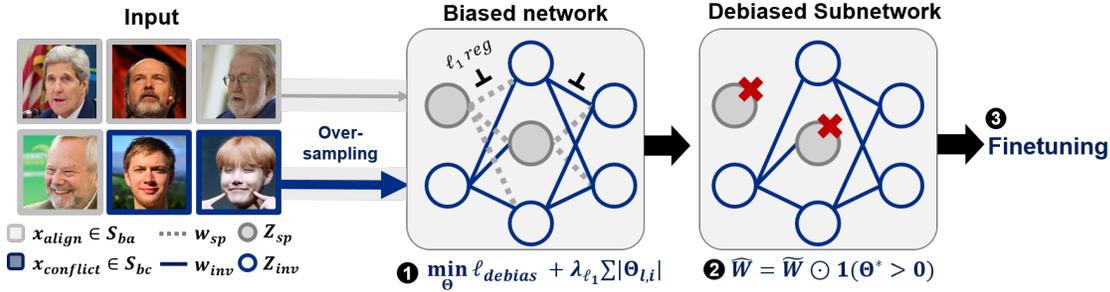


Figure 1. **Concept:** We demonstrate an inevitable generalization gap of subnetworks obtained by standard pruning methods including [37]. Based on these observations, we design a novel subnetwork probing framework by fully exploiting unbiased samples.

sociated with spurious features should be pruned out as the spurious features do not help predict bias-conflicting samples. Our theoretical observations suggest that balancing the ratio between the number of bias-aligned and bias-conflicting samples is crucial in finding the optimal unbiased subnetworks.

In practice, the dataset may severely lack diversity for bias-conflicting samples due to the potential pitfalls in data collection protocols or human prejudice. Since it is often highly laborious to supplement enough bias-conflicting samples, we propose a novel debiasing scheme called Debiased Contrastive Weight Pruning (DCWP) that uses the oversampled bias-conflicting data to search unbiased subnetworks.

As shown in Fig. 1, DCWP is comprised of two stages: (1) identifying the bias-conflicting samples without expensive annotations on spuriously correlated attributes, and (2) training the pruning parameters to obtain weight pruning masks with the sparsity constraint and *debiased* loss function. Here, the debiased loss includes a weighted cross-entropy loss for the identified bias-conflicting samples and an alignment loss to further reduce the geometrical alignment gap between bias-aligned and bias-conflicting samples within each class.

We demonstrate that DCWP consistently outperforms state-of-the-art debiasing methods across various biased datasets, including the Color-MNIST [23, 27], Corrupted CIFAR-10 [13], Biased FFHQ [21] and CelebA [25], even without direct supervision on the bias type. Our approach improves the accuracy on the unbiased evaluation dataset by 86.74% \rightarrow 93.41%, 27.86% \rightarrow 35.90% on Colored-MNIST and Corrupted CIFAR-10 compared to the second best model, respectively, even when 99.5% of samples are bias-aligned.

2. Related works

Spurious correlations. A series of empirical works have shown that the deep networks often find shortcut solutions relying on spuriously correlated attributes, such as the tex-

ture of image [10], language biases [12], or sensitive variables such as ethnicity or gender [7, 28]. Such behavior is of practical concern because it deteriorates the reliability of deep networks in sensitive applications like healthcare, finance, and legal services [4].

Debiasing frameworks. Recent studies to train a debiased network robust to spurious correlations can be roughly categorized into approaches (1) leveraging annotations of spurious attributes, i.e., bias label [29, 36], (2) presuming specific type of bias, e.g., texture [1, 9] or (3) without using explicit kinds of supervisions on dataset bias [22, 27]. The authors in [15, 29] optimize the worst-group error by using training group information. For practical implementation, reweighting or subsampling protocols are often used with increased model regularization [30]. Liu et al.; Sohoni et al. [24, 32] extend these approaches to the settings without expensive group annotations. Goel et al.; Kim et al. [11, 21] provide bias-tailored augmentations to balance the majority and minority groups. In particular, these approaches have mainly focused on better approximation and regularization of worst-group error combined with advanced data sampling, augmentation, or retraining strategies.

Studying impacts of neural architectures. Recently, the effects of deep neural network architecture on generalization performance have been explored. Diffenderfer et al. [6] employ recently advanced lottery-ticket-style pruning algorithms [8] to design the compact and robust network architecture. Bai et al. [2] directly optimize the neural architecture in terms of accuracy on OOD samples. Zhang et al. [37] demonstrate the effectiveness of pruning weights on spurious attributes, but the solution for discriminating such spurious weights lacks robust theoretical justifications, resulting in marginal performance gains. To fully resolve the above issues, we carry out a theoretical case study, and build a novel pruning algorithm that distills the representations to be independent of the spurious attributes.

3. Theoretical insights

3.1. Problem setup

Consider a supervised setting of predicting labels $Y \in \mathcal{Y}$ from input samples $X \in \mathcal{X}$ by a classifier $f_\theta : \mathcal{X} \rightarrow \mathcal{Y}$ parameterized by $\theta \in \Theta$. Following [37], let $(X^e, Y^e) \sim P^e$, where $X^e \in \mathcal{X}$ and $Y^e \in \mathcal{Y}$ refer to the input random variable and the corresponding label, respectively, and $e \in \mathcal{E} = \{1, 2, \dots, E\}$ denotes the index of environment, P^e is the corresponding distribution, and the set \mathcal{E} corresponds to every possible environments. We further assume that \mathcal{E} is divided into training environments \mathcal{E}_{train} and unseen test environments \mathcal{E}_{test} , i.e. $\mathcal{E} = \mathcal{E}_{train} \cup \mathcal{E}_{test}$.

For a given a loss function $\ell : \mathcal{X} \times \mathcal{Y} \times \Theta \rightarrow \mathbb{R}^+$, the standard training protocol for the empirical risk minimization (ERM) is to minimize the expected loss with a training environment $e \in \mathcal{E}_{train}$:

$$\hat{\theta}_{ERM} = \arg \min_{\theta} \mathbb{E}_{(X^e, Y^e) \sim \hat{P}^e} [\ell(X^e, Y^e; \theta)], \quad (1)$$

where \hat{P}^e is the empirical distribution over the training data. Our goal is to learn a model with good performance on OOD samples of $e \in \mathcal{E}_{test}$.

3.2. Motivating example

We conjecture that neural networks trained by ERM indiscriminately rely on predictive features, including those spuriously correlated ones [34].

To verify this conjecture, we present a simple binary classification example $(\mathbf{X}^e, Y^e) \sim P^e$, where $Y^e \in \mathcal{Y} = \{-1, 1\}$ represents the corresponding target label, and a sample $\mathbf{X}^e \in \mathcal{X} = \{-1, 1\}^{D+1} \in \mathbb{R}^{D+1}$ is constituted with both the invariant feature $Z_{inv}^e \in \{-1, 1\}$ and spurious features $\mathbf{Z}_{sp}^e \in \{-1, 1\}^D$, i.e. $\mathbf{X}^e = (Z_{inv}^e, \mathbf{Z}_{sp}^e)$. Suppose, furthermore, $Z_{sp,i}^e$ denote the i -th spurious feature component of \mathbf{Z}_{sp}^e . Note that we assume $D \gg 1$ to simulate the model heavily relies on spurious features [26, 37].

We consider the setting where the training environment $e \in \mathcal{E}_{train}$ is highly biased. In other words, we suppose that $Z_{inv}^e = Y^e$, and each of the i -th spurious feature component $Z_{sp,i}^e$ is independent and identically distributed (i.i.d) Bernoulli variable: i.e. $Z_{sp,i}^e$ independently takes a value equal to Y^e with a probability p^e and $-Y^e$ with a probability $1 - p^e$, where $p^e \in (0.5, 1], \forall e \in \mathcal{E}_{train}$. Note that $p^e \rightarrow 1$ as the environment is severely biased. A test environment $e \in \mathcal{E}_{test}$ is assumed to have $p^e = 0.5$, which implies that the spurious feature is totally independent with Y^e . Then we introduce a linear classifier f parameterized by a weight vector $\mathbf{w} = (w_{inv}, \mathbf{w}_{sp}) \in \mathbb{R}^{D+1}$, where $w_{inv} \in \mathbb{R}$ and $\mathbf{w}_{sp} \in \mathbb{R}^D$. In this example, we consider a class of pretrained classifiers parameterized by $\tilde{\mathbf{w}}(t) = (\tilde{w}_{inv}(t), \tilde{w}_{sp,1}(t), \dots, \tilde{w}_{sp,D}(t))$, where $t < T$ is a finite pretraining time for some sufficiently large T . Time t will be often omitted in notations for simplicity.

Our goal is to obtain the optimal sparse classifier with a highly biased training dataset. To achieve this, we introduce a binary weight pruning mask \mathbf{m} as $\mathbf{m} = (m_{inv}, \mathbf{m}_{sp}) \in \{0, 1\}^{D+1}$ for the pretrained weights, which is a significant departure from the theoretical setting in [37]. Specifically, let $m_{inv} \sim \text{Bern}(\pi_{inv})$, where π_{inv} and $1 - \pi_{inv}$ represents the probability of preserving (i.e. $m_{inv} = 1$) and pruning out (i.e. $m_{inv} = 0$), respectively. Similarly, let $m_{sp,i} \sim \text{Bern}(\pi_{sp,i}), \forall i$. Then, our optimization goal is to estimate the pruning probability parameter $\boldsymbol{\pi} = (\pi_1, \dots, \pi_{D+1}) = (\pi_{inv}, \pi_{sp,1}, \dots, \pi_{sp,D})$, where $\mathbf{m} \sim P(\boldsymbol{\pi})$ is a mask sampled with probability parameters $\boldsymbol{\pi}$. Accordingly, our main loss function for the pruning parameters given the environment e can be defined as follows:

$$\begin{aligned} \ell^e(\boldsymbol{\pi}) &= \frac{1}{2} \mathbb{E}_{\mathbf{X}^e, Y^e, \mathbf{m}} [1 - Y^e \hat{Y}^e] \\ &= \frac{1}{2} \mathbb{E}_{\mathbf{X}^e, Y^e, \mathbf{m}} \left[1 - Y^e \cdot \text{sgn} \left(\tilde{\mathbf{w}}^T(\mathbf{X}^e \odot \mathbf{m}) \right) \right], \end{aligned} \quad (2)$$

where \hat{Y}^e is the prediction of binary classifier, $\tilde{\mathbf{w}}$ is the pretrained weight vector, $\text{sgn}(\cdot)$ represents the sign function, and \odot represents element-wise product.

We first derive the upper-bound of the training loss $\ell^e(\boldsymbol{\pi})$ to illustrate the difficulty of learning optimal pruning parameters in a biased data setting. The proof can be found in Supplementary Material.

Theorem 1. (Training and test bound) Assume that $p^e > 1/2$ in the biased training environment $e \in \mathcal{E}_{train}$. Define $\tilde{\mathbf{w}}(t)$ as weights pretrained for a finite time $t < T$. Then the upper bound of the error of training environment w.r.t. pruning parameters $\boldsymbol{\pi}$ is given as:

$$\ell^e(\boldsymbol{\pi}) \leq 2 \exp \left(- \frac{2(\pi_{inv} + (2p^e - 1) \sum_{i=1}^D \alpha_i(t) \pi_{sp,i})^2}{4 \sum_{i=1}^D \alpha_i(t)^2 + 1} \right), \quad (3)$$

where the weight ratio $\alpha_i(t) = \tilde{w}_{sp,i}(t) / \tilde{w}_{inv}(t)$ is bounded below some positive constant. Given a test environment $e \in \mathcal{E}_{test}$ with $p^e = 1/2$, the upper bound of the error of test environment w.r.t. $\boldsymbol{\pi}$ is given as:

$$\ell^e(\boldsymbol{\pi}) \leq 2 \exp \left(- \frac{2\pi_{inv}^2}{4 \sum_{i=1}^D \alpha_i(t)^2 + 1} \right), \quad (4)$$

which implies that there is an unavoidable gap between training bound and test bound.

The detailed proof of Theorem 1 is provided in the supplementary material. This mismatch of the bounds is attributed to the contribution of $\pi_{sp,i}$ on the training bound (3). Intuitively, the networks prefer to preserve both \tilde{w}_{inv} and $\tilde{w}_{sp,i}$ in the presence of strong spurious correlations due

to the inherent sensitivity of ERM to all kinds of predictive features [17, 34]. This behavior is directly reflected in the training bound, where increasing either π_{inv} or $\pi_{sp,i}$, i.e., the probability of preserving weights, decreases the training bound. This inertia of spurious weights may prevent themselves from being primarily pruned against the sparsity constraint.

We note that the unintended reliance on spurious features is fundamentally rooted to the positivity of the weight ratio $\alpha_i(t)$. In the proof of Theorem 1 in Supplementary Material, we show some intriguing properties of $\alpha_i(t)$: (1) If infinitely many data and sufficient training time is provided, the gradient flow converges to the optimal solution which is invariant to \mathbf{Z}_{sp}^e , i.e., $\alpha_i(t) \rightarrow 0$. In this ideal situation, the gap between training and test bound is closed, thereby guaranteeing generalizations of obtained subnetworks. (2) However, given a finite time $t < T$ with a strongly biased dataset in practice, $\alpha_i(t)$ is bounded below by some positive constant, resulting in an inevitable generalization gap.

Theorem 1 implies that the classifier may preserve spurious weights due to the lack of bias-conflicting samples, which serve as counterexamples that spurious features themselves fail to explain. It motivates us to analyze the training bound in another environment η where we can systematically augment bias-conflicting samples. Specifically, consider $\mathbf{X}^\eta = (\mathbf{Z}_{inv}^\eta, \mathbf{Z}_{sp}^\eta)$, where $Z_{inv}^\eta = Y^\eta$ and mixture distribution of \mathbf{Z}_{sp}^η given $Y^\eta = y$ is defined in an element wise as follows:

$$P_{mix}^\eta(Z_{sp,i}^\eta | Y^\eta = y) = \phi P_{debias}^\eta(Z_{sp,i}^\eta | Y^\eta = y) + (1 - \phi) P_{bias}^\eta(Z_{sp,i}^\eta | Y^\eta = y), \quad (5)$$

where ϕ is a scalar mixture weight,

$$P_{debias}^\eta(Z_{sp,i}^\eta | Y^\eta = y) = \begin{cases} 1, & \text{if } Z_{sp,i}^\eta = -y \\ 0, & \text{if } Z_{sp,i}^\eta = y \end{cases} \quad (6)$$

is a debiasing distribution to weaken the correlation between Y^η and $Z_{sp,i}^\eta$ by setting the value of $Z_{sp,i}^\eta$ as $-Y^\eta$, and

$$P_{bias}^\eta(Z_{sp,i}^\eta | Y^\eta = y) = \begin{cases} p^\eta, & \text{if } Z_{sp,i}^\eta = y \\ 1 - p^\eta, & \text{if } Z_{sp,i}^\eta = -y \end{cases} \quad (7)$$

is a biased distribution similarly defined in the previous environment $e \in \mathcal{E}_{train}$. Given this new environment η , the degree of spurious correlations can be controlled by ϕ . This leads to a training bound as follow:

Theorem 2. (Training bound with the mixture distribution) Assume that the defined mixture distribution P_{mix}^η is biased, i.e., for all $i \in \{1, \dots, D\}$,

$$P_{mix}^\eta(Z_{sp,i}^\eta = -y | Y^e = y) \leq P_{mix}^\eta(Z_{sp,i}^\eta = y | Y^\eta = y). \quad (8)$$

Then, ϕ satisfies $0 \leq \phi \leq 1 - \frac{1}{2p^\eta}$. Then the upper bound of the error of training environment η w.r.t. the pruning parameters is given by

$$\ell^\eta(\boldsymbol{\pi}) \leq 2 \exp\left(-\frac{2(\pi_{inv} + (2p^\eta(1 - \phi) - 1) \sum_{i=1}^D \alpha_i(t) \pi_{sp,i})^2}{4 \sum_{i=1}^D \alpha_i(t)^2 + 1}\right). \quad (9)$$

Furthermore, when $\phi = 1 - \frac{1}{2p^\eta}$, the mixture distribution is perfectly debiased, and we have

$$\ell^\eta(\boldsymbol{\pi}) \leq 2 \exp\left(-\frac{2\pi_{inv}^2}{4 \sum_{i=1}^D \alpha_i(t)^2 + 1}\right), \quad (10)$$

which is equivalent to the test bound in (4).

The detailed proof is provided in the supplementary material. Our new training bound (31) suggests that the significance of $\pi_{sp,i}$ on training bound decreases as ϕ progressively increases, and at the extreme end with $\phi = 1 - \frac{1}{2p^\eta}$, it can be easily shown that $P_{mix}^\eta(Z_{sp,i}^\eta | Y^\eta = y) = \frac{1}{2}$ for both $y = 1$ and $y = -1$ so that $Z_{sp,i}^\eta$ turns out to be random. In other words, by plugging $\phi = 1 - \frac{1}{2p^\eta}$ into (31), we can minimize the gap between training and test error bound, which guarantees the improved OOD generalization.

4. Debiased Contrastive Weight Pruning

Our theoretical observations elucidate the importance of balancing between the bias-aligned and bias-conflicting samples in discovering the optimal unbiased subnetworks structure. While the true analytical form of the debiasing distribution is unknown in practice, we aim to approximate such unknown distribution with existing bias-conflicting samples and simulate the mixture distribution P_{mix}^η with modifying sampling strategy. To this end, we propose a Debiased Contrastive Weight Pruning (DCWP) algorithms that learn the unbiased subnetworks structure from the original full-size network.

Consider a L layer neural networks as a function $f_{\mathbf{W}} : \mathcal{X} \rightarrow \mathbb{R}^C$ parameterized by weights $\mathbf{W} = \{\mathbf{W}_1, \dots, \mathbf{W}_L\}$, where $C = |\mathcal{Y}|$ is the number of classes. Analogous to the earlier works on pruning, we introduce binary weight pruning masks $\mathbf{m} = \{\mathbf{m}_1, \dots, \mathbf{m}_L\}$ to model the subnetworks as $f(\cdot; \mathbf{m}_1 \odot \mathbf{W}_1, \dots, \mathbf{m}_L \odot \mathbf{W}_L)$. We denote such subnetworks as $f_{\mathbf{m} \odot \mathbf{W}}$ for the notational simplicity. We treat each entry of \mathbf{m}_l as an independent Bernoulli variable, and model their logits as our new pruning parameters $\boldsymbol{\Theta} = \{\boldsymbol{\Theta}_1, \dots, \boldsymbol{\Theta}_L\}$ where $\boldsymbol{\Theta}_l \in \mathbb{R}^{n_l}$ and n_l represents the dimensionality of the l -th layer weights \mathbf{W}_l . Then $\pi_{l,i} = \sigma(\boldsymbol{\Theta}_{l,i})$ denotes the probability of preserving the i -th weight of l -th layer $\mathbf{W}_{l,i}$ where σ refers to a sigmoid

function. To enable the end-to-end training, the Gumbel-softmax trick [18] for sampling masks together with ℓ_1 regularization term of Θ is adopted as a sparsity constraint. With a slight abuse of notations, $\mathbf{m} \sim G(\Theta)$ denotes a set of masks sampled with logits Θ by applying Gumbel-softmax trick.

Then our main optimization problem is defined as follows:

$$\min_{\Theta} \ell_{debias}(\{(x_i, y_i)\}_{i=1}^{|S|}; \tilde{\mathbf{W}}, \Theta) + \lambda_{\ell_1} \sum_{l,i} |\Theta_{l,i}|, \quad (11)$$

where S denotes the index set of whole training samples, $\lambda_{\ell_1} > 0$ is a Lagrangian multiplier, $\tilde{\mathbf{W}}$ represents the pretrained weights and ℓ_{debias} is our main objective which will be illustrated later. Note that we freeze the pretrained weights $\tilde{\mathbf{W}}$ during training pruning parameters Θ . We interchangeably use $\ell_{debias}(\{(x_i, y_i)\}_{i=1}^{|S|}; \Theta)$ and $\ell_{debias}(S; \Theta)$ in the rest of the paper. For comparison with our formulation, we recast the optimization problem of [37] with our notations as follows:

$$\min_{\Theta} \ell(\{(x_i, y_i)\}_{i=1}^{|S|}; \tilde{\mathbf{W}}, \Theta) + \lambda_{\ell_1} \sum_{l,i} |\Theta_{l,i}|, \quad (12)$$

where [37] uses the cross entropy (CE) loss function for ℓ .

Bias-conflicting sample mining In the first stage, we identify bias-conflicting training samples which empower functional modular probing. Specifically, we train a bias-capturing model and treat an error set S_{bc} of the index of misclassified training samples as bias-conflicting sample proxies. Our framework is broadly compatible with various bias-capturing models, where we mainly leverage the ERM model trained with generalized cross entropy (GCE) loss [39]:

$$\ell_{GCE}(x_i, y_i; \mathbf{W}_B) = \frac{1 - p_{y_i}(x_i; \mathbf{W}_B)^q}{q}, \quad (13)$$

where $q \in (0, 1]$ is a hyperparameter controlling the degree of bias amplification, \mathbf{W}_B is the parameters of the bias-capturing model, and $p_{y_i}(x_i; \mathbf{W}_B)$ is a softmax output value of the bias-capturing model assigned to the target label y_i . Compared to the CE loss, the gradient of the GCE loss up-weights the samples with a high probability of predicting the correct target, amplifying the network bias by putting more emphasis on easy-to-predict samples [27].

To preclude the possibility that the generalization performance of DCWP is highly dependent on the behavior of the bias-capturing model, we demonstrate in Section 5 that DCWP is reasonably robust to the degradation of accuracy on capturing bias-conflicting samples. Details about the bias-capturing model and simulation settings are presented in the supplementary material.

Upweighting Bias-conflicting samples After mining the index set of bias-conflicting sample proxies S_{bc} , we treat $S_{ba} = S \setminus S_{bc}$ as the index set of majority bias-aligned samples. Then we calculate the weighted cross entropy (WCE) loss $\ell_{WCE}(\{(x_i, y_i)\}_{i=1}^{|S|}; \tilde{\mathbf{W}}, \Theta)$ as follows:

$$\ell_{WCE}(S; \tilde{\mathbf{W}}, \Theta) := \mathbb{E}_{\mathbf{m} \sim G(\Theta)} [\lambda_{up} \ell_{bc}(S_{bc}; \mathbf{m}, \tilde{\mathbf{W}}) + \ell_{ba}(S_{ba}; \mathbf{m}, \tilde{\mathbf{W}})], \quad (14)$$

where $\lambda_{up} \geq 1$ is an upweighting hyperparameter, and

$$\ell_{bc}(S_{bc}; \mathbf{m}, \tilde{\mathbf{W}}) = \frac{1}{|S_{bc}|} \sum_{i \in S_{bc}} \ell_{CE}(x_i, y_i; \mathbf{m} \odot \tilde{\mathbf{W}}), \quad (15)$$

where ℓ_{CE} denotes the cross entropy loss. ℓ_{ba} is defined as similar to ℓ_{bc} .

The expectation is approximated with Monte Carlo estimates, where the number of mask \mathbf{m} sampled per iteration is set to 1 in practice. To implement (14), we oversample the samples in S_{bc} for λ_{up} times more than the samples in S_{ba} . This sampling strategy is aimed at increasing the mixture weight ϕ of the proposed mixture distribution P_{mix}^η in (5), while we empirically approximate the unknown bias-conflicting group distribution with the sample set S_{bc} .

Note that although simple oversampling of bias-conflicting samples may not lead to the OOD generalization due to the inductive bias towards memorizing a few counterexamples in overparameterized neural networks [30], such failure is unlikely reproduced in learning *pruning* parameters under the strong sparsity constraint. We sample new weight masks \mathbf{m} for each training iteration in a stochastic manner, effectively precluding the overparameterized networks from potentially memorizing the minority samples. As a result, DCWP exhibits reasonable performance even with few bias-conflicting samples.

Bridging the alignment gap by pruning To fully utilize the bias-conflicting samples, we consider the sample-wise relation between bias-conflicting samples and majority bias-aligned samples. Zhang et al. [38] demonstrates that the deteriorated OOD generalization is potentially attributed to the distance gap between same-class representations; bias-aligned representations are more closely aligned than bias-conflicting representations, although they are generated from the same-class samples. We hypothesized that well-designed pruning masks could alleviate such geometrical misalignment. Specifically, ideal weight sparsification may guide each latent dimension to be independent of spurious attributes, thereby preventing representations from being misaligned with spuriously correlated latent dimensions. This motivates us to explore pruning masks by contrastive learning. (Related illustrative example in appendix)

Following the conventional notations of contrastive learning, we denote $f_{\tilde{\mathbf{W}}}^{enc} : \mathcal{X} \rightarrow \mathbb{R}^{n_L-1}$ as an encoder

parameterized by $\mathbf{W} = (\mathbf{W}_1, \dots, \mathbf{W}_{L-1})$ which maps samples into the representations at penultimate layer. Let $f_{\mathbf{W}_L}^{cls} : \mathbb{R}^{n_L} \rightarrow \mathbb{R}^C$ be the classification layer parameterized by \mathbf{W}_L . Then $f_{\mathbf{W}}(\mathbf{x}) = f_{\mathbf{W}_L}^{cls}(f_{\mathbf{W}}^{enc}(\mathbf{x}))$, $\forall \mathbf{x} \in \mathcal{X}$. We similarly define $f_{\mathbf{m} \odot \mathbf{W}}^{enc}$ and $f_{\mathbf{m}_L \odot \mathbf{W}_L}^{cls}$. For the i -th sample x_i , let $\mathbf{z}_i(\mathbf{W}) = \text{norm}(f_{\mathbf{W}}^{enc}(x_i))$ be the normalized representations lies on the unit hypersphere, and similarly define $\mathbf{z}_i(\mathbf{m} \odot \mathbf{W})$. We did not consider projection networks [3, 20] for architectural simplicity. Given index subsets of training samples $\mathcal{V}, \mathcal{V}^+$, the supervised contrastive loss [20] function is defined as follows:

$$\ell_{con}(\mathcal{V}, \mathcal{V}^+; \mathbf{W}) = \sum_{i \in \mathcal{V}} \frac{-1}{|\mathcal{V}^+(y_i)|} \sum_{j \in \mathcal{V}^+(y_i)} \log \frac{\exp(\mathbf{z}_i(\mathbf{W}) \cdot \mathbf{z}_j(\mathbf{W})/\tau)}{\sum_a \exp(\mathbf{z}_i(\mathbf{W}) \cdot \mathbf{z}_a(\mathbf{W})/\tau)}, \quad (16)$$

where $a \in \mathcal{V} \setminus \{i\}$, $\tau > 0$ is a temperature hyperparameter, and $\mathcal{V}^+(y_i) = \{k \in \mathcal{V}^+ : y_k = y_i, k \neq i\}$ indicates the index set of samples with target label y_i . Then, we define the debiased alignment loss as follows:

$$\ell_{align}(\{x_i, y_i\}_{i=1}^{|\mathcal{S}|}; \tilde{\mathbf{W}}, \Theta) = \mathbb{E}_{\mathbf{m} \sim G(\Theta)} \left[\ell_{con}(S_{bc}, S; \mathbf{m} \odot \tilde{\mathbf{W}}) + \ell_{con}(S_{ba}, S_{bc}; \mathbf{m} \odot \tilde{\mathbf{W}}) \right], \quad (17)$$

where the expectation is approximated with Monte Carlo estimates as in (14). Intuitively, (17) reduces the gap between bias-conflicting samples and others (first term), while preventing bias-aligned samples from being aligned too close each other (second term, more discussions in appendix).

Finally, our debiased loss in (11) is defined as follows:

$$\ell_{debias}(S; \tilde{\mathbf{W}}, \Theta) = \ell_{WCE}(S; \tilde{\mathbf{W}}, \Theta) + \lambda_{align} \ell_{align}(S; \tilde{\mathbf{W}}, \Theta), \quad (18)$$

where $\lambda_{align} > 0$ is a balancing hyperparameter.

Fine-tuning after pruning After solving (11) by gradient-descent optimization, we can obtain the pruning parameters Θ^* . This allows us to uncover the structure of unbiased subnetworks with binary weight masks $\mathbf{m}^* = \{\mathbf{m}_1^*, \dots, \mathbf{m}_L^*\}$, where $\mathbf{m}_l^* = \{\mathbb{1}(\sigma(\Theta_{l,i}^*) > 1/2) | 1 \leq i \leq n_l\}$, $\forall l \in \{1, \dots, L\}$, and n_l is a dimensionality of the l -th weight. After pruning, we finetune the survived weights $\hat{\mathbf{W}} = \mathbf{m}^* \odot \tilde{\mathbf{W}}$ using ℓ_{WCE} in (14) and $\lambda_{align} \ell_{align}$ in (17). Interestingly, we empirically found that the proposed approach works well without the reset [8] (Related experiments in Section 5). Accordingly, we resume the training while fixing the unpruned pretrained weights. The pseudocode of DCWP is provided in Algorithm 1.

Algorithm 1 Debiased Contrastive Weight Pruning (DCWP)

- 1: **Input:** Dataset $D = \{(x_i, y_i)_{i=1}^{|\mathcal{S}|}\}$, pruning parameters Θ , Training iterations T_1, T_2, T_3 .
 - 2: **Output:** Trained pruning parameters Θ^* and finetuned weights \mathbf{W}^*
 - 3:
 - 4: **Stage 1. Mining debiased samples**
 - 5: Update the weights of bias-capturing network \mathbf{W}_b on D for T_1 iterations.
 - 6: Identify S_{bc} and S_{ba} .
 - 7:
 - 8: **Stage 2. Debiased Contrastive Weight Pruning**
 - 9: Pretrain the main network on D . Denote the pretrained weights as $\tilde{\mathbf{W}}$.
 - 10: **for** $t = 1$ **to** T_2 **do**
 - 11: Update Θ with $\ell_{debias}(S; \tilde{\mathbf{W}}, \Theta) + \lambda_{\ell_1} \sum_{l,i} |\Theta_{l,i}|$ as in (11).
 - 12: **end for**
 - 13: Prune out weight as $\hat{\mathbf{W}} = \tilde{\mathbf{W}} \odot \mathbb{1}(\Theta^* > 0)$.
 - 14: Update $\hat{\mathbf{W}}$ with ℓ_{WCE} and $\lambda_{align} \ell_{align}$ on D for T_3 iterations.
-

Table 1. Unbiased test accuracy evaluated on CMNIST, CIFAR10-C and bias-conflict test accuracy evaluated on BFFHQ. Models requiring supervisions on dataset bias are denoted with \checkmark , while others are denoted with \times . Results are averaged on 4 different random seeds.

Dataset	Ratio (%)	ERM	EnD	Rebias	MRM	LfF	DisEnt	DCWP
		\times	\checkmark	\checkmark	\times	\times	\times	\times
CMNIST	0.5	62.36	84.32	69.12	60.98	83.73	86.74	93.41
	1.0	81.73	94.98	84.65	80.42	88.44	93.15	95.98
	2.0	89.33	97.01	91.96	89.31	92.67	95.15	97.16
	5.0	95.22	98.00	96.74	95.23	94.90	96.76	98.02
CIFAR10-C	0.5	22.02	23.93	21.73	23.92	27.02	27.86	35.90
	1.0	28.00	27.61	28.09	27.77	31.44	34.62	41.56
	2.0	34.63	36.62	35.57	33.53	38.49	41.95	49.01
	5.0	45.66	43.67	48.22	47.00	46.16	49.15	56.17
BFFHQ	0.5	52.25	59.80	54.90	54.75	56.50	55.50	60.35

Table 2. Worst-group and average test accuracies on CelebA (Blonde). (\checkmark , \times) here represents $\text{Idx} = (6, 4)$ (w/ and w/o pruning) in Table 3, respectively, which shows the impacts of pruning.

Models	ERM	DisEnt	JTT [24]	DCWP (\times)	DCWP (\checkmark)
Worst-group	47.02	65.26	76.80	67.85	79.30
Average	97.80	67.88	93.98	95.89	94.50

5. Experimental results

5.1. Methods

Datasets To show the effectiveness of the proposed pruning algorithms, we evaluate the generalization performance of several debiasing approaches on Colored MNIST (CM-

NIST), Corrupted CIFAR-10 (CIFAR10-C), Biased FFHQ (BFFHQ) with varying ratio of bias-conflicting samples, i.e., bias ratio. We report unbiased accuracy [22, 27] on the test set, which includes a balanced number of samples from each data group. We also report bias-conflict accuracy for some experiments, which is the average accuracy on bias-conflicting samples included in an unbiased test set. Specifically, we report the bias-conflict accuracy on BFFHQ in which half of the unbiased test samples are bias-aligned, while the model with the best-unbiased accuracy is selected (Unbiased accuracy in Table 4). For CelebA (blonde) [14, 29], we report worst-group and average accuracy following [29] considering that abundant samples are included in (Blonde Hair=0, Male=0) bias-conflicting group. We use the same data splits from [14].

Baselines We compare DCWP with vanilla network trained by ERM, and the following state-of-the-art debiasing approaches: EnD [33], Rebias [1], MRM [37], LfF [27], JTT [24] and DisEnt [22]. EnD relies on the annotations on the spurious attribute of training samples, i.e., bias labels. Rebias relies on prior knowledge about the type of dataset bias (e.g., texture). MRM, LfF, JTT and DisEnt do not presume such bias labels or prior knowledge about dataset bias. Notably, MRM is closely related to DCWP where it probes the unbiased functional subnetwork with standard cross entropy. Details about other simulation settings are provided in Supplementary Material.

5.2. Evaluation results

As shown in Table 1, we found that DCWP outperforms other state-of-the-art debiasing methods by a large margin. Moreover, the catastrophic pitfalls of the existing pruning method become evident, where MRM fails to search for unbiased subnetworks. It underlines that the proposed approach for utilizing bias-conflicting samples plays a pivotal role in discovering unbiased subnetworks.

5.3. Quantitative analyses

Ablation studies To quantify the extent of performance improvement achieved by each introduced module, we analyzed the dependency of model performance on: (a) pruning out spurious weights following the trained parameters, (b) using alignment loss or (c) oversampling identified bias-conflicting samples when training Θ and \tilde{W} . To emphasize the contribution of each module, we intentionally use an SGD optimizer which results in lower baseline accuracy (and for other CMNIST experiments in this subsection as well). Table 3 shows that every module plays an important role in OOD generalization, while (a) pruning contributes significantly comparing (1→2, +7.19%), (3→5, +11.59%) or (4→6, +8.68%).

Dependency on bias-capturing models To evaluate the reliability of DCWP, we compare different version of DCWP

Table 3. Ablation study on CMNIST (Bias ratio=1%). Unbiased accuracy is reported. $\text{Idx} = 2$ uses ℓ_{WCE} only for training pruning parameters Θ while using ℓ_{CE} for retraining. $\text{Idx} = 3, 4$ does not conduct pruning and finetune the pretrained weights \tilde{W} by oversampling minorities or using alignment loss.

Idx	(a) Pruning	(b) ℓ_{align}	(c) ℓ_{WCE}	Accuracy (%)
1	-	-	-	43.10
2	✓	-	-	50.29
3	-	-	✓	73.20
4	-	✓	✓	79.28
5	✓	-	✓	84.79
6	✓	✓	✓	87.96

which does not rely on the dataset-tailored mining algorithms. We posit that early stopping [24] is an easy plug-and-play method to train the bias-capturing model in general. Thus we newly train DCWP_{ERM} which collects bias-conflicting samples by using the early-stopped ERM model. Table 4 shows that DCWP_{ERM} outperforms other baselines even though the precision, the fraction of samples in S_{bc} that are indeed bias-conflicting, or recall, the fraction of the bias-conflicting samples that are included in S_{bc} , were significantly dropped. It implies that DCWP may perform reasonably well with the limited number and quality of bias-conflicting samples.

Table 4. Robustness dependency of DCWP on the performance of bias-capturing models. We set the bias ratio as 1% for CIFAR10-C. Results are averaged on 4 different random seeds.

Dataset	Model	Accuracy			Mining metrics	
		bias-align	bias-conflict	unbiased	precision	recall
CIFAR10-C	DisEnt	80.04	26.51	34.62	-	-
	DCWP_{ERM}	94.33	<u>29.75</u>	<u>36.21</u>	19.71	79.53
	DCWP	<u>91.68</u>	35.99	41.56	85.97	74.89
BFFHQ	DisEnt	89.80	55.55	72.68	-	-
	LfF	96.05	56.50	76.30	-	-
	DCWP_{ERM}	99.45	<u>56.90</u>	<u>78.20</u>	20.18	28.39
	DCWP	<u>98.85</u>	60.35	79.60	30.61	31.25

Do we need to reset weights? While it becomes widespread wisdom that remaining weights should be reset to their initial ones from the original network after pruning [8], we analyze whether such reset is also required for the proposed pruning framework. We compared the training dynamics of different models such as: (1) ERM model, (2) MRM_{debias} which solves (11) instead of (12) to obtain the weight pruning masks, and (3) DCWP. Note that MRM_{debias} reset the unpruned weights to its initialization after pruning. Figure 2a shows that although MRM_{debias} makes a considerable advance, weight reset inevitably limits the performance gain. Moreover, finetuning the biased model significantly improves the generalization performance within only a few iterations, which implies that the proposed neural pruning can further boost the accuracy

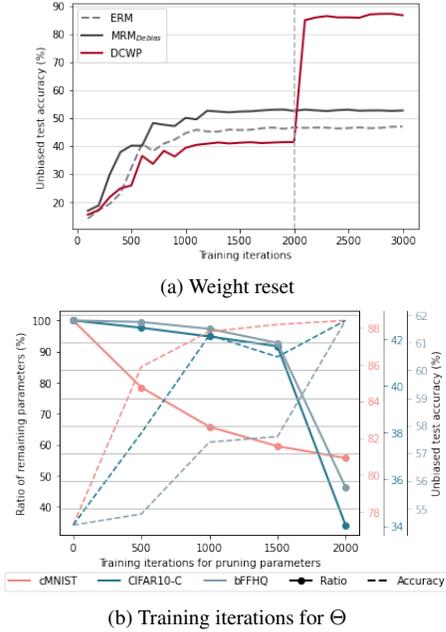


Figure 2. (a) Comparison study on finetuning and weight resetting (CMNIST, bias ratio=1%). For DCWP, after pretraining weights for 2000 iterations, we pause and start training pruning parameters (vertical dotted line in the figure). After convergence, we mask out and finetune weights for another 1000 iterations. For MRM_{debias}, we reset the unpruned weight to its initialization and retrain for 3000 iterations. (b) Sensitivity analysis on the training iterations for pruning parameter Θ . Bias ratio=1% for both CMNIST and CIFAR10-C. Bias-conflict accuracy is reported for BFFHQ.

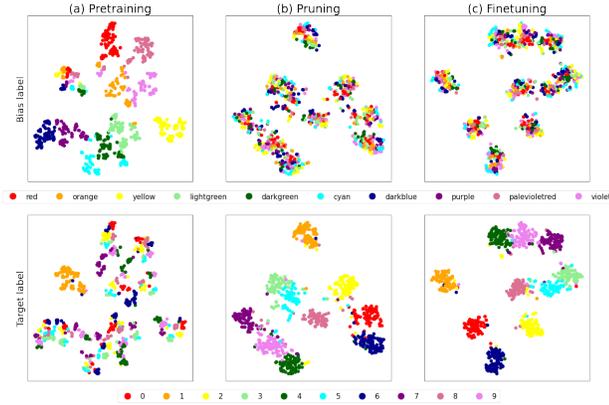


Figure 3. t-SNE visualization of representations encoded from unbiased test samples after (a) pretraining, (b) pruning and (c) finetuning (CMNIST, bias ratio=0.5%). Each point is painted following its label (i.e., bias label in first row, and target label in second row).

without weight reset. This finding allows us to debias large-scale pretrained models *without* retraining by simple pruning and finetuning.

Sensitivity analysis on training iterations We also analyzed the hyperparameter sensitivity on the training iterations of the pruning parameter Θ . The unbiased test accuracy is evaluated with weight pruning masks generated by Θ trained for {500, 1000, 1500, 2000} iterations on each dataset. Figure 2b shows that the accuracy increases as more (potentially biased) weights are pruned out. It implies that the proposed method can compress the networks to a substantial extent while significantly improving the OOD generalization performance.

Visualization of learned latent representations. We visualized latent representations of unbiased test samples in CMNIST after (a) pretraining, (b) pruning, and (c) finetuning. Note that we did not reset or finetune the weights in (b). As reported in Figure 3, biased representations in (a) are misaligned along with bias labels as discussed in section 4. However, after pruning, the representations were well-aligned with respect to the class of digits even without modifying the values of pretrained weights. It implies that the geometrical misalignment of representations can be addressed by pruning spurious weights while finetuning with ℓ_{debias} can further improve the generalizations.

6. Conclusion

This paper presented a novel functional subnetwork probing method for OOD generalization. Our goal was to find a winning functional lottery ticket [37], which can achieve better OOD performance compared to its counterpart full network, given a highly biased dataset in practice. We provided theoretical insights and empirical evidence to show that the minority samples provide an important clue for probing the optimal unbiased subnetworks. Simulations on various benchmark datasets demonstrated that our model significantly outperforms state-of-the-art debiasing methods. The proposed method is memory efficient and potentially compatible with many other debiasing methods.

Acknowledgments

This work was supported by the National Research Foundation of Korea (NRF) grant funded by the Korea government (MSIT) (NRF-2020R1A2B5B03001980), Institute for Information & Communications Technology Planning & Evaluation (IITP) grant funded by the Korea government (No. 2017-0-00451, No. RS-2023-00233251, System3 reinforcement learning with high-level brain functions), KAIST Key Research Institute (Interdisciplinary Research Group) Project and Field-oriented Technology Development Project for Customs Administration through National Research Foundation of Korea (NRF) funded by the Ministry of Science & ICT and Korea Customs Service(**NRF-2021M3I1A1097938**).

Appendix: Training Debiased Subnetworks with Contrastive Weight Pruning

The supplementary material is organized as follows. We first present the proof for Theorem 1 and 2. In section 8, we extend the presented theoretical example in the main paper to illustrate the risks of geometrical misalignment of embeddings arising from strong spurious correlations. Additional results are reported in section 9. Optimization setting, hyperparameter configuration, and other experimental details are provided in section 10.

7. Proofs

In this section, we present the detailed proofs for Theorems 1 and 2 explained in the main paper, followed by an illustration about the dynamics of weight ratio $\alpha_i(t) = \tilde{w}_{sp,i}(t)/\tilde{w}_{inv}(t)$.

7.1. Proof of Theorem 1

Theorem 1. (Training and test bound) Assume that $p^e > 1/2$ in the biased training environment $e \in \mathcal{E}_{train}$. Define $\tilde{\mathbf{w}}(t)$ as weights pretrained for a finite time $t < T$. Then the upper bound of the error of training environment w.r.t. pruning parameters $\boldsymbol{\pi}$ is given as:

$$\ell^e(\boldsymbol{\pi}) \leq 2 \exp \left(- \frac{2(\pi_{inv} + (2p^e - 1) \sum_{i=1}^D \alpha_i(t) \pi_{sp,i})^2}{4 \sum_{i=1}^D \alpha_i(t)^2 + 1} \right), \quad (19)$$

where the weight ratio $\alpha_i(t) = \tilde{w}_{sp,i}(t)/\tilde{w}_{inv}(t)$ is bounded below some positive constant. Given a test environment $e \in \mathcal{E}_{test}$ with $p^e = \frac{1}{2}$, the upper bound of the error of test environment w.r.t. $\boldsymbol{\pi}$ is given as:

$$\ell^e(\boldsymbol{\pi}) \leq 2 \exp \left(- \frac{2\pi_{inv}^2}{4 \sum_{i=1}^D \alpha_i(t)^2 + 1} \right), \quad (20)$$

which implies that there is an unavoidable gap between training bound and test bound.

Proof. We omit time t in $\tilde{\mathbf{w}}(t)$ and $\alpha_i(t)$ for notational simplicity throughout the proof of Theorem 1 and 2.

We recall the loss function defined in the main paper for convenience.

$$\begin{aligned} \ell^e(\boldsymbol{\pi}) &= \frac{1}{2} \mathbb{E}_{\mathbf{X}^e, Y^e, \mathbf{m}} [1 - Y^e \hat{Y}^e] \\ &= \frac{1}{2} \mathbb{E}_{\mathbf{X}^e, Y^e, \mathbf{m}} \left[1 - Y^e \cdot \text{sgn} \left(\tilde{\mathbf{w}}^T (\mathbf{X}^e \odot \mathbf{m}) \right) \right], \end{aligned} \quad (21)$$

where \hat{Y}^e is the prediction of binary classifier, $\tilde{\mathbf{w}}$ is the pretrained weight vector, $\text{sgn}(\cdot)$ represents the sign function, and \odot represents element-wise product.

The prediction from the classifier \hat{Y}^e is defined as

$$\begin{aligned} \hat{Y}^e &= \text{sgn} \left(\tilde{\mathbf{w}}^T (\mathbf{X}^e \odot \mathbf{m}) \right) \\ &= \text{sgn} \left(\mathcal{O}^e \right), \end{aligned} \quad (22)$$

where

$$\mathcal{O}^e := \tilde{w}_{inv} m_{inv} Z_{inv}^e + \sum_{i=1}^D \tilde{w}_{sp,i} m_{sp,i} Z_{sp,i}^e. \quad (23)$$

Assume that Y^e is uniformly distributed binary random variable. Then,

$$\mathbb{E}_{\mathbf{X}^e, Y^e, \mathbf{m}} [Y^e \hat{Y}^e] = \frac{1}{2} \mathbb{E}_{\mathbf{X}^e, \mathbf{m}} \left[\hat{Y}^e | Y^e = 1 \right] - \frac{1}{2} \mathbb{E}_{\mathbf{X}^e, \mathbf{m}} \left[\hat{Y}^e | Y^e = -1 \right], \quad (24)$$

where

$$\begin{aligned} \mathbb{E}_{\mathbf{X}^e, \mathbf{m}} \left[\hat{Y}^e | Y^e = 1 \right] &= \mathbb{E}_{\mathbf{X}^e, \mathbf{m}} \left[\text{sgn} \left(\mathcal{O}^e \right) \mid Y^e = 1 \right] \\ &= P(\mathcal{O}^e > 0 \mid Y^e = 1) - P(\mathcal{O}^e < 0 \mid Y^e = 1) \\ &= 1 - 2P(\mathcal{O}^e < 0 \mid Y^e = 1), \end{aligned} \quad (25)$$

and

$$\begin{aligned}\mathbb{E}_{\mathbf{X}^e, \mathbf{m}}[\hat{Y}^e | Y^e = -1] &= P(\mathcal{O}^e > 0 | Y^e = -1) - P(\mathcal{O}^e < 0 | Y^e = -1) \\ &= -\mathbb{E}_{\mathbf{X}^e, \mathbf{m}}[\hat{Y}^e | Y^e = 1],\end{aligned}\tag{26}$$

where we use $P(\mathcal{O}^e < 0 | Y^e = 1) = P(\mathcal{O}^e > 0 | Y^e = -1)$ and $P(\mathcal{O}^e > 0 | Y^e = 1) = P(\mathcal{O}^e < 0 | Y^e = -1)$ thanks to the symmetry. Therefore, we have

$$\begin{aligned}\ell^e(\boldsymbol{\pi}) &= \frac{1}{2} \mathbb{E}_{\mathbf{X}^e, Y^e, \mathbf{m}}[1 - Y^e \hat{Y}^e] \\ &= \frac{1}{2} - \frac{1}{2} \mathbb{E}_{\mathbf{X}^e, \mathbf{m}}[\hat{Y}^e | Y^e = 1] \\ &= P(\mathcal{O}^e < 0 | Y^e = 1).\end{aligned}\tag{27}$$

In order to derive a concentration inequality of $\ell^e(\boldsymbol{\pi})$, we compute a conditional expectation as follows:

$$\begin{aligned}\mathbb{E}_{\mathbf{X}^e, \mathbf{m}}[\mathcal{O}^e | Y^e = 1] &= \mathbb{E}_{\mathbf{X}^e, \mathbf{m}}\left[\tilde{w}_{inv} m_{inv} Z_{inv}^e + \sum_{i=1}^D \tilde{w}_{sp,i} m_{sp,i} Z_{sp,i}^e \mid Y^e = 1\right] \\ &= \mathbb{E}_{\mathbf{X}^e, \mathbf{m}}\left[\tilde{w}_{inv} m_{inv} + \sum_{i=1}^D \tilde{w}_{sp,i} m_{sp,i} Z_{sp,i}^e \mid Y^e = 1\right] \\ &= \tilde{w}_{inv} \pi_{inv} + \mathbb{E}_{\mathbf{X}^e, \mathbf{m}}\left[\sum_{i=1}^D \tilde{w}_{sp,i} m_{sp,i} Z_{sp,i}^e \mid Y^e = 1\right] \\ &= \tilde{w}_{inv} \pi_{inv} + \sum_{i=1}^D (2p^e - 1) \tilde{w}_{sp,i} \pi_{sp,i},\end{aligned}\tag{28}$$

where the last equality follows from the independence of $Z_{sp,\cdot}$ and $m_{sp,\cdot}$ as assumed in the main paper. Then,

$$\begin{aligned}P(\mathcal{O}^e < 0 | Y^e = 1) &= P\left(\mathcal{O}^e - \mathbb{E}_{\mathbf{X}^e, \mathbf{m}}[\mathcal{O}^e] < -\mathbb{E}_{\mathbf{X}^e, \mathbf{m}}[\mathcal{O}^e] \mid Y^e = 1\right) \\ &\leq P\left(\left|\mathcal{O}^e - \mathbb{E}_{\mathbf{X}^e, \mathbf{m}}[\mathcal{O}^e]\right| > \mathbb{E}_{\mathbf{X}^e, \mathbf{m}}[\mathcal{O}^e] \mid Y^e = 1\right) \\ &\leq 2 \exp\left(-\frac{2\mathbb{E}_{\mathbf{X}^e, \mathbf{m}}[\mathcal{O}^e | Y^e = 1]^2}{\tilde{w}_{inv}^2 + \sum_{i=1}^D 4\tilde{w}_{sp,i}^2}\right) \\ &\leq 2 \exp\left(-\frac{2(\tilde{w}_{inv} \pi_{inv} + \sum_{i=1}^D (2p^e - 1) \tilde{w}_{sp,i} \pi_{sp,i})^2}{\tilde{w}_{inv}^2 + \sum_{i=1}^D 4\tilde{w}_{sp,i}^2}\right) \\ &\leq 2 \exp\left(-\frac{2(\pi_{inv} + \sum_{i=1}^D (2p^e - 1) \alpha_i \pi_{sp,i})^2}{1 + \sum_{i=1}^D 4\alpha_i^2}\right),\end{aligned}\tag{29}$$

where the second inequality is obtained using Hoeffding's inequality, third inequality is from (28), and last inequality is obtained by dividing both denominator and numerator with \tilde{w}_{inv}^2 . We use the definition of weight ratio $\alpha_i = \tilde{w}_{sp,i} / \tilde{w}_{inv}$. For the second inequality, we use that $\tilde{w}_{inv} m_{inv} Z_{inv}^e \in \{0, \tilde{w}_{inv}\}$ and $\tilde{w}_{sp,i} m_{sp,i} Z_{sp,i}^e \in \{-\tilde{w}_{sp,i}, 0, \tilde{w}_{sp,i}\} \forall i$ in (23) to obtain the denominator.

Finally, the proof for the positivity of $\alpha_i(t)$ comes from Proposition 1 in section 1.3 in this appendix. This concludes the proof. \square

7.2. Proof of Theorem 2

Theorem 2. (Training bound with the mixture distribution) Assume that the defined mixture distribution P_{mix}^η is biased, i.e., for all $i \in \{1, \dots, D\}$,

$$P_{mix}^\eta(Z_{sp,i}^\eta = -y | Y^\eta = y) \leq P_{mix}^\eta(Z_{sp,i}^\eta = y | Y^\eta = y).\tag{30}$$

Then, ϕ satisfies $0 \leq \phi \leq 1 - \frac{1}{2p^\eta}$. Then the upper bound of the error of training environment η w.r.t. the pruning parameters is given by

$$\ell^\eta(\boldsymbol{\pi}) \leq 2 \exp \left(- \frac{2(\pi_{inv} + (2p^\eta(1 - \phi) - 1) \sum_{i=1}^D \alpha_i(t) \pi_{sp,i})^2}{4 \sum_{i=1}^D \alpha_i(t)^2 + 1} \right). \quad (31)$$

Furthermore, when $\phi = 1 - \frac{1}{2p^\eta}$, the mixture distribution is perfectly debiased, and we have

$$\ell^\eta(\boldsymbol{\pi}) \leq 2 \exp \left(- \frac{2\pi_{inv}^2}{4 \sum_{i=1}^D \alpha_i(t)^2 + 1} \right), \quad (32)$$

which is equivalent to the test bound in (20).

Proof. Recall that $Z_{sp,i}^\eta$ follows the mixture distribution P_{mix}^η :

$$P_{mix}^\eta(Z_{sp,i}^\eta | Y^\eta = y) = \phi P_{debias}^\eta(Z_{sp,i}^\eta | Y^\eta = y) + (1 - \phi) P_{bias}^\eta(Z_{sp,i}^\eta | Y^\eta = y), \quad (33)$$

where

$$P_{debias}^\eta(Z_{sp,i}^\eta | Y^\eta = y) = \begin{cases} 1, & \text{if } Z_{sp,i}^\eta = -y \\ 0, & \text{if } Z_{sp,i}^\eta = y \end{cases} \quad (34)$$

is a debiasing distribution to weaken the correlation between Y^η and $Z_{sp,i}^\eta$ by setting the value of $Z_{sp,i}^\eta$ as $-Y^\eta$, and

$$P_{bias}^\eta(Z_{sp,i}^\eta | Y^\eta = y) = \begin{cases} p^\eta, & \text{if } Z_{sp,i}^\eta = y \\ 1 - p^\eta, & \text{if } Z_{sp,i}^\eta = -y. \end{cases} \quad (35)$$

Then, with definition in (34) and (35),

$$\begin{aligned} P_{mix}^\eta(Z_{sp,i}^\eta = -y | Y^\eta = y) &= \phi + (1 - \phi)(1 - p^\eta) \\ P_{mix}^\eta(Z_{sp,i}^\eta = y | Y^\eta = y) &= (1 - \phi)p^\eta, \end{aligned} \quad (36)$$

for $y \in \{-1, 1\}$. Then, based on the assumption, $\phi + (1 - \phi)(1 - p^\eta) \leq (1 - \phi)p^\eta$, which gives $\phi \leq 1 - \frac{1}{2p^\eta}$. Specifically, if $\phi = 1 - \frac{1}{2p^\eta}$, it turns out that $P_{mix}^\eta(Z_{sp,i}^\eta = -y | Y^\eta = y) = P_{mix}^\eta(Z_{sp,i}^\eta = y | Y^\eta = y) = \frac{1}{2}$, which implies that spurious features turns out to be random and the mixture distribution becomes perfectly debiased. If $\phi = 0$, the mixture distribution boils down into a biased distribution as similarly defined in the environment $e \in \mathcal{E}_{train}$.

The prediction from the classifier \mathcal{O}^η is defined as similar to \mathcal{O}^e in (23). Then in order to derive a concentration inequality of $\ell^\eta(\boldsymbol{\pi})$, we derive a conditional expectation of \mathcal{O}^η as done in (28):

$$\begin{aligned} \mathbb{E}_{\mathbf{X}^\eta, \mathbf{m}}[\mathcal{O}^\eta | Y^\eta = 1] &= \mathbb{E}_{\mathbf{X}^\eta, \mathbf{m}} \left[\tilde{w}_{inv} m_{inv} Z_{inv}^\eta + \sum_{i=1}^D \tilde{w}_{sp,i} m_{sp,i} Z_{sp,i}^\eta \mid Y^\eta = 1 \right] \\ &= \mathbb{E}_{\mathbf{X}^\eta, \mathbf{m}} \left[\tilde{w}_{inv} m_{inv} + \sum_{i=1}^D \tilde{w}_{sp,i} m_{sp,i} Z_{sp,i}^\eta \mid Y^\eta = 1 \right]. \end{aligned} \quad (37)$$

Then, with the definition in (33), the second term in the above conditional expectation of (37) is defined as follows:

$$\begin{aligned} \mathbb{E}_{\mathbf{X}^\eta, \mathbf{m}} \left[\sum_{i=1}^D \tilde{w}_{sp,i} m_{sp,i} Z_{sp,i}^\eta \mid Y^\eta = 1 \right] &= \sum_{i=1}^D \tilde{w}_{sp,i} \pi_{sp,i} \left(\phi \mathbb{E}_{debias}[Z_{sp,i}^\eta | Y^\eta = 1] + (1 - \phi) \mathbb{E}_{bias}[Z_{sp,i}^\eta | Y^\eta = 1] \right) \\ &= \sum_{i=1}^D \tilde{w}_{sp,i} \pi_{sp,i} \left(\phi \cdot (-1) + (1 - \phi)(2p^\eta - 1) \right) \\ &= \sum_{i=1}^D \tilde{w}_{sp,i} \pi_{sp,i} (2p^\eta(1 - \phi) - 1), \end{aligned} \quad (38)$$

where \mathbb{E}_{debias} and \mathbb{E}_{bias} in the first equality denote the conditional expectation with respect to distribution P_{debias}^η and P_{bias}^η in (34) and (35), respectively. Plugging (38) into (37), we get

$$\mathbb{E}_{\mathbf{X}^\eta, \mathbf{m}}[\mathcal{O}^\eta | Y^\eta = 1] = \tilde{w}_{inv}\pi_{inv} + \sum_{i=1}^D (2p^\eta(1-\phi) - 1)\tilde{w}_{sp,i}\pi_{sp,i}. \quad (39)$$

Then we can derive the upper bound of $\ell^\eta(\boldsymbol{\pi}) = P(\mathcal{O}^\eta < 0 | Y^\eta = 1)$ similarly to (29):

$$\begin{aligned} P(\mathcal{O}^\eta < 0 | Y^\eta = 1) &\leq P\left(\left|\mathcal{O}^\eta - \mathbb{E}_{\mathbf{X}^\eta, \mathbf{m}}[\mathcal{O}^\eta] \right| > \mathbb{E}_{\mathbf{X}^\eta, \mathbf{m}}[\mathcal{O}^\eta] | Y^\eta = 1\right) \\ &\leq 2 \exp\left(-\frac{2\mathbb{E}_{\mathbf{X}^\eta, \mathbf{m}}[\mathcal{O}^\eta | Y^\eta = 1]^2}{\tilde{w}_{inv}^2 + 4\sum_{i=1}^D \tilde{w}_{sp,i}^2}\right) \\ &\leq 2 \exp\left(-\frac{2(\tilde{w}_{inv}\pi_{inv} + \sum_{i=1}^D (2p^\eta(1-\phi) - 1)\tilde{w}_{sp,i}\pi_{sp,i})^2}{\tilde{w}_{inv}^2 + 4\sum_{i=1}^D \tilde{w}_{sp,i}^2}\right) \\ &\leq 2 \exp\left(-\frac{2(\pi_{inv} + \sum_{i=1}^D (2p^\eta(1-\phi) - 1)\alpha_i\pi_{sp,i})^2}{1 + \sum_{i=1}^D 4\alpha_i^2}\right), \end{aligned} \quad (40)$$

where the first inequality is obtained by Hoeffding's inequality, and second inequality is from (39). The denominator is obtained as same as in (29), since $\tilde{w}_{inv}m_{inv}Z_{inv}^\eta \in \{0, \tilde{w}_{inv}\}$ and $\tilde{w}_{sp,i}m_{sp,i}Z_{sp,i}^\eta \in \{-\tilde{w}_{sp,i}, 0, \tilde{w}_{sp,i}\} \forall i$ as-is. If we plug-in the upper bound value of $\phi = 1 - \frac{1}{2p^\eta}$ obtained from (36) into (40), it boils down into the test bound in (20). \square

7.3. Dynamics of the weight ratio

We omit an index of environment e in the proposition below for notational simplicity.

Proposition 1. Consider a binary classification problem of linear classifier $f_{\mathbf{w}}$ under exponential loss. Let $(\mathbf{X}, Y) \sim P$, where each input random variable \mathbf{X} and the corresponding label Y is generated by

$$\mathbf{X} = \begin{pmatrix} Z_{inv} \\ \mathbf{Z}_{sp} \end{pmatrix}, Y = Z_{inv},$$

where $\mathbf{Z}_{sp} = (2\mathbf{z} - 1)Z_{inv}$ for a random variable $\mathbf{z} \in \{0, 1\}^D$ which is chosen from multivariate Bernoulli distribution ($z_i \sim \text{Bern}(p)$) with $p > \frac{1}{2}$, i.e., p denotes p^e in the main paper. Let $\mathbf{w} = \begin{pmatrix} w_{inv} \\ \mathbf{w}_{sp} \end{pmatrix} \in \mathbb{R}^{D+1}$ be the weight of the linear classifier $f_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$. Assume that $0 < w_{inv}(0)$, i.e., w_{inv} is initialized with a positive value, and $0 < w_{sp,i}(0) < \frac{1}{2} \log \frac{p}{1-p}$. Then, after sufficient time of training, w_{inv} diverges to $+\infty$ and $w_{sp,i}$ converges to $\frac{1}{2} \log \frac{p}{1-p}$, which means $\alpha_i := \frac{w_{sp,i}}{w_{inv}}$ converges to 0 for all $i \in \{1, 2, \dots, D\}$. More precisely,

$$\log\left(e^{w_{inv}(0)} + [4p(1-p)]^{\frac{D}{2}t}\right) \leq w_{inv}(t) \leq \log\left(e^{w_{inv}(0)} + t \prod_{i=1}^D \left(pe^{-w_{sp,i}(0)} + \sqrt{p(1-p)}\right)\right).$$

However, for a fixed $t < T$, each α_i is positive and its lower bound converges to some positive value.

Proof. In this proof, $w_{inv}(t)$ denotes the invariant weight at time t , while we often omit the time t and interchangeably use w_{inv} for notational simplicity, and likewise for $w_{sp,i}(t)$.

Note that the network output is given by

$$\begin{aligned} f_{\mathbf{w}}(\mathbf{x}) &= \mathbf{w}^T \mathbf{x} \\ &= Z_{inv}w_{inv} + \mathbf{Z}_{sp}^T \mathbf{w}_{sp} \\ &= Z_{inv}w_{inv} + \sum_{i=1}^D Z_{sp,i}w_{sp,i}. \end{aligned}$$

The exponential loss is defined by

$$\begin{aligned}
L(\mathbf{w}) &= \mathbb{E}_{(\mathbf{X}, Y)}[e^{-f_{\mathbf{w}}(\mathbf{X})Y}] \\
&= \mathbb{E}_{\mathbf{z}} \left[\exp \left(-(Z_{inv} w_{inv} + \sum_{i=1}^D Z_{sp,i} w_{sp,i}) Z_{inv} \right) \right] \\
&= \mathbb{E}_{\mathbf{z}} \left[\exp \left(-w_{inv} - (2z_1 - 1)w_{sp,1} - \dots - (2z_D - 1)w_{sp,D} \right) \right] \\
&= e^{-w_{inv}} \prod_{i=1}^D \mathbb{E}_{\mathbf{z}} [e^{-(2z_i - 1)w_{sp,i}}] \\
&= e^{-w_{inv}} \prod_{i=1}^D (pe^{-w_{sp,i}} + (1-p)e^{w_{sp,i}}).
\end{aligned}$$

Then, thanks to symmetry of w_{sp} , it is enough to consider $\alpha := \frac{w_{sp,1}}{w_{inv}}$. We first compute the gradient:

$$\begin{aligned}
\frac{\partial L}{\partial w_{inv}} &= -e^{-w_{inv}} \prod_{i=1}^D (pe^{-w_{sp,i}} + (1-p)e^{w_{sp,i}}) \\
\frac{\partial L}{\partial w_{sp,1}} &= -e^{-w_{inv}} (pe^{-w_{sp,1}} - (1-p)e^{w_{sp,1}}) \prod_{i=2}^D (pe^{-w_{sp,i}} + (1-p)e^{w_{sp,i}}).
\end{aligned}$$

Since $\frac{d}{dt} w_{inv} = -\frac{\partial L}{\partial w_{inv}}$, the dynamics is given by the following differential equations.

$$\begin{aligned}
\frac{d}{dt} w_{inv} &= e^{-w_{inv}} \prod_{i=1}^D (pe^{-w_{sp,i}} + (1-p)e^{w_{sp,i}}) \\
\frac{d}{dt} w_{sp,1} &= e^{-w_{inv}} (pe^{-w_{sp,1}} - (1-p)e^{w_{sp,1}}) \prod_{i=2}^D (pe^{-w_{sp,i}} + (1-p)e^{w_{sp,i}}).
\end{aligned}$$

First we show that $w_{inv}(t)$ diverges to $+\infty$ as t goes ∞ . We show this by computing its lower bound.

$$\begin{aligned}
\frac{d}{dt} w_{inv} &= e^{-w_{inv}} \prod_{i=1}^D (pe^{-w_{sp,i}} + (1-p)e^{w_{sp,i}}) \\
&\geq e^{-w_{inv}} \prod_{i=1}^D (2\sqrt{p(1-p)}) \\
&= e^{-w_{inv}} [4p(1-p)]^{\frac{D}{2}},
\end{aligned}$$

where the inequality is obtained by AM-GM inequality. This implies $e^{w_{inv}} dw_{inv} \geq [4p(1-p)]^{\frac{D}{2}} dt$. Integrating both sides from 0 to t , we get

$$e^{w_{inv}(t)} - e^{w_{inv}(0)} \geq [4p(1-p)]^{\frac{D}{2}} t$$

or

$$w_{inv}(t) \geq \log \left(e^{w_{inv}(0)} + [4p(1-p)]^{\frac{D}{2}} t \right), \quad (41)$$

which shows that $w_{inv}(t)$ diverges to $+\infty$ as $t \rightarrow \infty$. Note also that w_{inv} strictly increases since $\frac{d}{dt} w_{inv} > 0$.

For $w_{sp,i}$, $\frac{d}{dt} w_{sp,i} = 0$ implies $w_{sp,i}$ converges to $w_{sp,i}^*$ such that

$$pe^{-w_{sp,i}^*} - (1-p)e^{w_{sp,i}^*} = 0,$$

namely, $w_{sp,i}^* = \frac{1}{2} \log \frac{p}{1-p}$.

As similar to w_{inv} , $w_{sp,1}$ strictly increases if and only if $w_{sp,1} < \frac{1}{2} \log \frac{p}{1-p}$. Based on the assumptions that $0 < w_{sp,i}(0) < \frac{1}{2} \log \frac{p}{1-p}$, we conclude that $w_{sp,1}$ monotonically converges to $\frac{1}{2} \log \frac{p}{1-p}$. As p goes to 1, $\frac{1}{2} \log \frac{p}{1-p}$ is sufficiently large and we can assume $w_{sp,i}(0) < \frac{1}{2} \log \frac{p}{1-p}$.

Now, we fix $0 < t < T$ for given T and compute an upper bound of w_{inv} . Using $w_{sp,i}(t) < \frac{1}{2} \log \frac{p}{1-p}$, we get

$$\begin{aligned} \frac{d}{dt} w_{inv} &= e^{-w_{inv}} \prod_{i=1}^D (pe^{-w_{sp,i}} + (1-p)e^{w_{sp,i}}) \\ &< e^{-w_{inv}} \prod_{i=1}^D \left(pe^{-w_{sp,i}(0)} + (1-p)\sqrt{\frac{p}{1-p}} \right) \\ &= e^{-w_{inv}} \prod_{i=1}^D \left(pe^{-w_{sp,i}(0)} + \sqrt{p(1-p)} \right) \end{aligned}$$

which implies

$$e^{w_{inv}} dw_{inv} < \prod_{i=1}^D \left(pe^{-w_{sp,i}(0)} + \sqrt{p(1-p)} \right) dt.$$

Integrating both sides from 0 to t , we get

$$w_{inv}(t) < \log \left(e^{w_{inv}(0)} + \prod_{i=1}^D \left(pe^{-w_{sp,i}(0)} + \sqrt{p(1-p)} \right) t \right). \quad (42)$$

Similarly, we compute a lower bound of $w_{sp,1}$ on $0 < t < T$. Before we start, note that $w_{inv}(t) < w_{inv}(T) =: M$ from monotonicity.

$$\begin{aligned} \frac{d}{dt} w_{sp,1} &= e^{-w_{inv}} (pe^{-w_{sp,1}} - (1-p)e^{w_{sp,1}}) \prod_{i=2}^D (pe^{-w_{sp,i}} + (1-p)e^{w_{sp,i}}) \\ &> e^{-M} (pe^{-w_{sp,1}} - (1-p)e^{w_{sp,1}}) \prod_{i=2}^D (2\sqrt{p(1-p)}) \\ &= e^{-M} [4p(1-p)]^{\frac{D-1}{2}} (pe^{-w_{sp,1}} - (1-p)e^{w_{sp,1}}) \end{aligned}$$

induces

$$\frac{1}{pe^{-w_{sp,1}} - (1-p)e^{w_{sp,1}}} dw_{sp,1} > e^{-M} [4p(1-p)]^{\frac{D-1}{2}} dt.$$

Integrating both sides from 0 to $t < T$, we get

$$\left[\frac{1}{\sqrt{p(1-p)}} \tanh^{-1} \left(\sqrt{\frac{1-p}{p}} e^{w_{sp,1}} \right) \right]_0^t > e^{-M} [4p(1-p)]^{\frac{D-1}{2}} t$$

or

$$w_{sp,1}(t) > \frac{1}{2} \log \frac{p}{1-p} + \log \tanh \left(\tanh^{-1} \left(\sqrt{\frac{1-p}{p}} e^{w_{sp,1}(0)} \right) + e^{-M} 2^{D-1} [p(1-p)]^{\frac{D-1}{2}} t \right). \quad (43)$$

Combining (42) and (43), we conclude that

$$\alpha_p(t) = \frac{w_{sp,1}(t)}{w_{inv}(t)} \quad (44)$$

$$> \frac{\frac{1}{2} \log \frac{p}{1-p} + \log \tanh \left(\tanh^{-1} \left(\sqrt{\frac{1-p}{p}} e^{w_{sp,1}(0)} \right) + e^{-M} 2^{D-1} [p(1-p)]^{\frac{D-1}{2}} t \right)}{\log \left(e^{w_{inv}(0)} + t \prod_{i=1}^D \left(pe^{-w_{sp,i}(0)} + \sqrt{p(1-p)} \right) \right)} \quad (45)$$

for $0 < t < T$. Note that $\alpha_p(t)$ is positive in $0 < t < T$, since both $w_{sp,1}(t)$ and $w_{inv}(t)$ is monotonically increasing in $0 < t < T$, and $0 < w_{sp,1}(0), w_{inv}(0)$ by assumptions.

The numerator becomes

$$\begin{aligned} & \frac{1}{2} \log \frac{p}{1-p} + \log \tanh \left(\tanh^{-1} \left(\sqrt{\frac{1-p}{p}} e^{w_{sp,1}(0)} \right) + e^{-M} 2^{D-1} [p(1-p)]^{\frac{D}{2}} t \right) \\ &= \log \left[\sqrt{\frac{p}{1-p}} \tanh \left(\tanh^{-1} \left(\sqrt{\frac{1-p}{p}} e^{w_{sp,1}(0)} \right) + e^{-M} 2^{D-1} [p(1-p)]^{\frac{D}{2}} t \right) \right] \\ &= \log \left[\sqrt{\frac{p}{1-p}} \left(\sqrt{\frac{1-p}{p}} e^{w_{sp,1}(0)} + e^{-M} 2^{D-1} [p(1-p)]^{\frac{D}{2}} t \operatorname{sech}^2 c \right) \right] \end{aligned}$$

for some c such that

$$\tanh^{-1} \left(\sqrt{\frac{1-p}{p}} e^{w_{sp,1}(0)} \right) < c < \tanh^{-1} \left(\sqrt{\frac{1-p}{p}} e^{w_{sp,1}(0)} + e^{-M} 2^{D-1} [p(1-p)]^{\frac{D}{2}} t \right).$$

We use $f(x+y) = f(x) + yf'(c)$ by the Mean Value Theorem (MVT) at the last line.

Notably, if we take a limit $p \rightarrow 1$, the numerator becomes

$$\lim_{p \rightarrow 1} \log \left[e^{w_{sp,1}(0)} + e^{-M} 2^{D-1} p^{\frac{D+1}{2}} (1-p)^{\frac{D-1}{2}} t \operatorname{sech}^2 c \right] = w_{sp,1}(0).$$

Similarly, the denominator becomes

$$\begin{aligned} & \lim_{p \rightarrow 1} \log \left(e^{w_{inv}(0)} + t \prod_{i=1}^D \left(p e^{-w_{sp,i}(0)} + \sqrt{p(1-p)} \right) \right) \\ &= \log \left(e^{w_{inv}(0)} + t \prod_{i=1}^D e^{-w_{sp,i}(0)} \right) \\ &= \log \left(e^{w_{inv}(0)} + t \exp \left(- \sum_{i=1}^D w_{sp,i}(0) \right) \right) \end{aligned}$$

Therefore, for a fixed $0 < t < T$, we conclude that

$$\begin{aligned} \lim_{p \rightarrow 1} \alpha_p(t) &= \lim_{p \rightarrow 1} \frac{w_{sp,1}(t)}{w_{inv}(t)} \\ &\geq \frac{w_{sp,1}(0)}{\log \left(e^{w_{inv}(0)} + t \exp \left(- \sum_{i=1}^D w_{sp,i}(0) \right) \right)} \\ &> \frac{w_{sp,1}(0)}{\log \left(e^{w_{inv}(0)} + T \exp \left(- \sum_{i=1}^D w_{sp,i}(0) \right) \right)} \\ &\geq \frac{w_{sp,1}(0)}{\log T + \frac{1}{T} \exp \left(w_{inv}(0) + \sum_{i=1}^D w_{sp,i}(0) \right) - \sum_{i=1}^D w_{sp,i}(0)} \end{aligned} \tag{46}$$

where we use the inequality $\log(x+y) \leq \log x + \frac{y}{x}$ in the last line. \square

The key insights from Proposition 1 can be summarized as follows:

- (1) Weight ratio $\alpha_i(t)$ converges to 0 as $t \rightarrow \infty$.
- (2) However, for a fixed $t < T$, $\alpha_i(t) > 0$.
- (3) When $t < T$ and $p \rightarrow 1$, i.e., the environment is almost perfectly biased, the convergence rate of (1) is remarkably slow as in (46). In other words, there exists $c > 0$ such that $\frac{c}{\log t} < \alpha_p(t)$ over $0 < t < T$ if p is sufficiently close to 1.

This results afford us intriguing perspective on the fundamental factors behind the biased classifiers. If we situate the presented theoretical example in an ideal scenario in which infinitely many data and sufficient training time is provided, our result (1) shows that the pretrained classifier becomes fully invariant to the spurious correlations. However, in practical setting with finite training time and number of samples, our result (2) shows that the pretrained model inevitably rely on the spuriously correlated features.

Beyond theoretical results, we empirically observe that the weight ratio α_i of pretrained classifiers indeed increases as $p^e \rightarrow 1$. We simulate the example presented in section 3.2 of the main paper, where the dimensionality D is set to 15, and probability p^e varies from 0.6 (weakly biased) to 0.99 (severely biased). We train a linear classifier for 500 epochs with batch size of 1024, and measure the unbiased accuracy on test samples generated from environment $e \in \mathcal{E}_{test}$. We also measure weight ratio $\text{mean}(\tilde{w}_{sp})/\tilde{w}_{inv}$, where $\text{mean}(\tilde{w}_{sp})$ denotes the average of pretrained spurious weights $\{w_{sp,i}\}_{i=1}^D$. To enable the end-to-end training, we use binary cross entropy loss instead of exponential loss, with setting $\mathcal{Y} = \{0, 1\}$ instead of $\mathcal{Y} = \{-1, 1\}$. We do not consider pruning process in this implementation. Figure 4 shows that the weight ratio increases to 1 in average as $p^e \rightarrow 1$. It implies that the spurious features \mathbf{Z}_{sp}^e participate almost equally to the invariant feature Z_{inv}^e in the presence of strong spurious correlations. In this worst case, it is frustratingly difficult to discriminate weights necessary for OOD generalization in biased environment, resulting in the failure of learning optimal pruning parameters. Simulation results are averaged on 15 different random seeds.

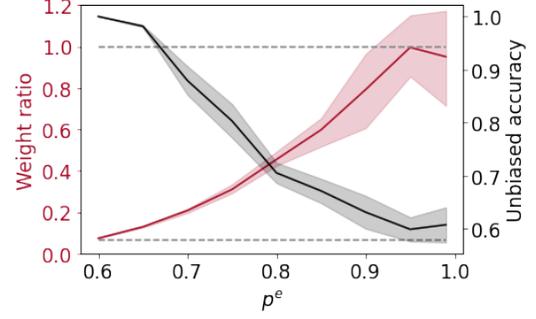


Figure 4. Implemented results of presented example.

8. Example of geometrical misalignment

In this section, we present a simple example illustrating the potential adverse effect of spurious correlations on latent representations. Consider independent arbitrary samples within the same class $\mathbf{X}_i^b, \mathbf{X}_j^b \sim P_{\mathbf{X}^b|Y^b=y}^b$ and $\mathbf{X}^d \sim P_{\mathbf{X}^d|Y^d=y}^d$ for a common $y \in \{-1, 1\}$ and environments b, d where $b \in \mathcal{E}_{train}$ and $d \in \mathcal{E}_{test}$. Let $\mathbf{W} \in \mathbb{R}^{Q \times (D+1)}$ be a weight matrix representation of a linear mapping $T : \{-1, 1\}^{D+1} \rightarrow \mathbb{R}^Q$ which encodes the embedding vector of a given sample. We denote such embedding as $\mathbf{h}^e = \mathbf{W}\mathbf{X}^e$ for some $e \in \mathcal{E}$. We assume that \mathbf{W} is initialized as to be semi-orthogonal [16, 31] for simplicity. Then the following lemma reveals the geometrical misalignment of embeddings in the presence of strong spurious correlations:

Lemma 1. *Given $y \in \{-1, 1\}$, let $\mathbf{h}_i^b, \mathbf{h}_j^b, \mathbf{h}^d$ be embeddings of $\mathbf{X}_i^b, \mathbf{X}_j^b, \mathbf{X}^d$ respectively. Then, the expected cosine similarity between \mathbf{h}_i^b and \mathbf{h}^d is derived as:*

$$\mathbb{E} \left[\frac{\langle \mathbf{h}_i^b, \mathbf{h}^d \rangle}{\|\mathbf{h}_i^b\| \cdot \|\mathbf{h}^d\|} \mid Y^b = y, Y^d = y \right] = \frac{1}{D+1}, \quad (47)$$

while the expected cosine similarity between \mathbf{h}_i^b and \mathbf{h}_j^b is derived as:

$$\mathbb{E} \left[\frac{\langle \mathbf{h}_i^b, \mathbf{h}_j^b \rangle}{\|\mathbf{h}_i^b\| \cdot \|\mathbf{h}_j^b\|} \mid Y^b = y \right] = \frac{1 + D(2p^b - 1)^2}{D+1}, \quad (48)$$

where p^b is a probability parameter of Bernoulli distribution of i.i.d variable $Z_{sp,i}^b$, similar to p^e in the main paper.

Proof. Let $\mathbf{X}^e = \mathbf{V}_{inv}^e + \mathbf{V}_{sp}^e$ for the sample from an arbitrary environment e in general, where $\mathbf{v}_{inv}^e, \mathbf{v}_{sp}^e \in \{-1, 1\}^{D+1}$ are invariant and spurious component vector, respectively:

$$\mathbf{V}_{inv,j}^e = \begin{cases} Z_{inv}^e, & \text{if } j = 1 \\ 0, & \text{otherwise,} \end{cases} \quad (49)$$

$$\mathbf{V}_{sp,j}^e = \begin{cases} Z_{sp,j}^e, & \text{if } j = 2, \dots, D+1 \\ 0, & \text{otherwise.} \end{cases} \quad (50)$$

Thus, \mathbf{V}_{inv}^e and \mathbf{V}_{sp}^e are orthogonal. Given $Y^b = y$ and $Y^d = y$ for some $y \in \{-1, 1\}$, the cosine similarity between \mathbf{h}_i^b and \mathbf{h}^d is expressed as follows:

$$\begin{aligned}
\mathbb{E} \left[\frac{\langle \mathbf{h}_i^b, \mathbf{h}^d \rangle}{\|\mathbf{h}_i^b\| \|\mathbf{h}^d\|} \mid Y^b = y, Y^d = y \right] &= \mathbb{E} \left[\frac{\langle \mathbf{X}_i^b, \mathbf{W}^T \mathbf{W} \mathbf{X}^d \rangle}{\|\mathbf{h}_i^b\| \|\mathbf{h}^d\|} \mid Y^b = y, Y^d = y \right] \\
&= \mathbb{E} \left[\frac{\langle \mathbf{X}_i^b, \mathbf{X}^d \rangle}{D+1} \mid Y^b = y, Y^d = y \right] \\
&= \mathbb{E} \left[\frac{\langle \mathbf{V}_{i,inv}^b + \mathbf{V}_{i,sp}^b, \mathbf{V}_{inv}^d + \mathbf{V}_{sp}^d \rangle}{D+1} \mid Y^b = y, Y^d = y \right] \\
&= \frac{1}{D+1},
\end{aligned} \tag{51}$$

where $\mathbf{V}_{i,inv}^b$ and $\mathbf{V}_{i,sp}^b$ represent the invariant and spurious component vector of \mathbf{X}_i^b , respectively, and the second equality comes from the semi-orthogonality of \mathbf{W} . The last equality comes from the orthogonality of spurious component vector from different environment $b \in \mathcal{E}_{train}$ and $d \in \mathcal{E}_{test}$.

On the other hand, the expected cosine similarity between two arbitrary embeddings \mathbf{h}_i^b and \mathbf{h}_j^b from the biased environment b is expressed as follows:

$$\begin{aligned}
\mathbb{E} \left[\frac{\langle \mathbf{h}_i^b, \mathbf{h}_j^b \rangle}{\|\mathbf{h}_i^b\| \|\mathbf{h}_j^b\|} \mid Y^b = y \right] &= \mathbb{E} \left[\frac{\langle \mathbf{V}_{i,inv}^b + \mathbf{V}_{i,sp}^b, \mathbf{V}_{j,inv}^b + \mathbf{V}_{j,sp}^b \rangle}{D+1} \mid Y^e = y \right] \\
&= \frac{1 + D(2p^b - 1)^2}{D+1},
\end{aligned} \tag{52}$$

where the last equality comes from the expectation of product of independent Bernoulli variables. \square

The gap between (47) and (48) unveils the imbalance of distance between same-class embeddings from different environments on the unit hypersphere; embeddings from the training environment are more closely aligned to other embeddings from the same environment than embeddings from test environment at initial even when all samples are generated within the same class. While the Lemma 1 is only applicable to the initialized \mathbf{W} before training, such imbalance may be worsened if \mathbf{W} learns to project the samples on the high-dimensional subspace where most of its basis are independent to the invariant features. This sparks interests in designing weight pruning masks to aggregate the representations from same-class samples all together. Indeed, in this simple example, we can address this misalignment by masking out every weight in \mathbf{W} except the first column, which is associated with the invariant feature.

From this point of view, we revisit the proposed alignment loss in main paper:

$$\ell_{align} \left(\{x_i, y_i\}_{i=1}^{|S|}; \tilde{\mathbf{W}}, \Theta \right) = \mathbb{E}_{\mathbf{m} \sim G(\Theta)} \left[\ell_{con}(S_{bc}, S; \mathbf{m} \odot \tilde{\mathbf{W}}) + \ell_{con}(S_{ba}, S_{bc}; \mathbf{m} \odot \tilde{\mathbf{W}}) \right], \tag{53}$$

where the first term reduces the gap between bias-conflicting samples and others, while the second term prevents bias-aligned samples from being aligned too close each other. In other words, the first term is aimed at increasing the cosine similarity between representations of same-class samples with different spurious attributes, as \mathbf{h}_i^b and \mathbf{h}^d in this example. The second term serves as a regularizer that pulls apart same-class bias-aligned representations, as \mathbf{h}_i^b and \mathbf{h}_j^b in this example. Thus we can leverage abundant bias-aligned samples as negatives regardless of their class in second term, while [38] limits the negatives to samples with different target label but same bias label, which are often highly scarce in a biased dataset.

9. Additional results

Comparisons to the pruning baselines. Pruning (debiasing) appears to suffer from the generalization-efficiency tradeoff; improving computational efficiency (OOD generalization) does not always guarantee improvement in OOD generalization (efficiency). Unlike this, our framework reliably improves both generalization and efficiency as shown in Table 5. Note that the standard pruning algorithms [35] fail to improve the unbiased accuracy in CIFAR10-C.

Analysis of sparsity level. One may concern that a trade-off between performance and sparsity may exist. For example, networks with mild sparsity may still be over-parameterized and thus not fully debiased, whereas networks with high sparsity do not have enough capacity to preserve the averaged accuracy. In order to investigate the trade-off, we measure the unbiased accuracy by explicitly controlling the pruning ratio with varying λ_{ℓ_1} . Figure 5 shows that (1) the trade-off between

Table 5. Test (unbiased) accuracy (%) on standard and corrupted CIFAR10 (Bias ratio=5%). Pruning ratio=90.0% for GraSP and (92.4%, 90.4%) for DCWP on (CIFAR10, CIFAR10-C).

Dataset	Full-size	GraSP [35]	DCWP
CIFAR10	86.76	85.64	86.32
CIFAR10-C	45.66	44.21	60.24

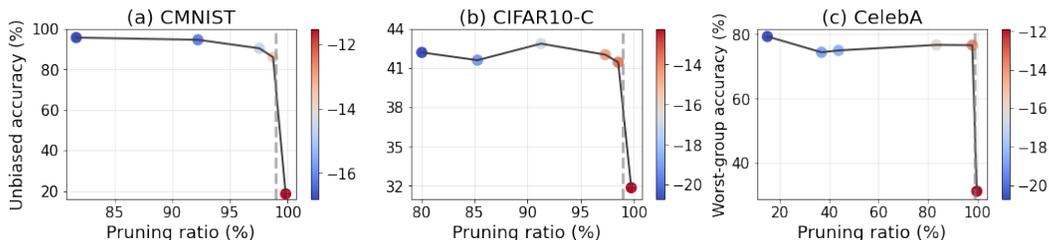


Figure 5. Analysis on the sparsity level. Bias ratio=1% for (a), (b). Color bar: log-scaled λ_{ℓ_1} . Dotted line: ratio = 99%.

performance and sparsity does exist, while (2) the proposed framework is reasonably tolerant to high sparsity in terms of generalization. We conjecture that such tolerance is owing to the *prioritized* elimination of spurious weights; the networks can be compressed to a significant extent without hurting the generalization after pruning out the spurious weights.

10. Experimental setup

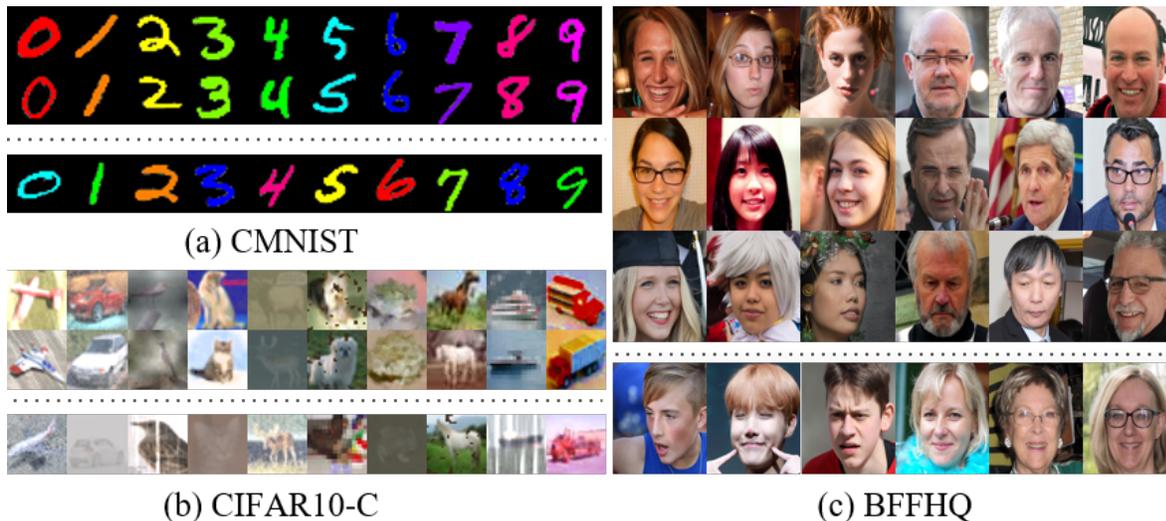


Figure 6. Example images of datasets. The images above the dotted line denote the bias-aligned samples, while the ones below the dotted line are the bias-conflicting samples. For CMNIST and CIFAR10-C, each column indicates each class. For BFFHQ, the group of three columns indicates each class.

10.1. Datasets

We mainly follow [22, 27] to evaluate our framework on Color-MNIST (CMNIST), Corrupted CIFAR-10 (CIFAR10-C) and Biased FFHQ (BFFHQ) as presented in Figure 6.

CMNIST. We first consider the prediction task of digit class which is spuriously correlated to the pre-assigned color, following the existing works [1, 22, 27, 33]. Each digit is colored with certain type of color, following [22, 27]. The ratio

of bias-conflicting samples, i.e., bias ratio, is varied in range of $\{0.5\%, 1.0\%, 2.0\%, 5.0\%\}$, where the exact number of (bias-aligned, bias-conflicting) samples is set to: (54,751, 249)-0.5%, (54,509, 491)-1%, (54,014, 986)-2%, and (52,551, 2,449)-5%.

CIFAR10-C. Each sample in this dataset is generated by corrupting original samples in CIFAR-10 with certain types of corruption. Among 15 different corruptions introduced in the original paper [13], we select 10 types which are Brightness, Contrast, Gaussian Noise, Frost, Elastic Transform, Gaussian Blur, Defocus Blur, Impulse Noise, Saturate, and Pixelate, following [22]. Each of these corruption is spuriously correlated to the object classes of CIFAR-10, which are Plane, Car, Bird, Cat, Deer, Dog, Frog, Horse, Ship, and Truck. We use the samples corrupted in most severe level among five different severity, following [22]. The exact number of (bias-aligned, bias-conflicting) samples is set to: (44,832, 228)-0.5%, (44,527, 442)-1%, (44,145, 887)-2%, and (42,820, 2,242)-5%.

BFFHQ. Each sample in this biased dataset are selected from Flickr-Faces-HQ (FFHQ) Dataset [19], where we conduct binary classifications with considering (Age, Gender) as target and spuriously correlated attribute pair following [21,22]. Specifically, majority of training images correspond to either young women (i.e., aged 10-29) or old men (i.e., aged 40-59). This dataset consists of 19,104 number of such bias-aligned samples and 96 number of bias-conflicting samples, i.e., old women and young men.

CelebA. For CelebA, we consider (Blonde Hair, Male) as (target, spurious) attribute pair, following [14,27,29]. Pixel resolutions and batch size are 256×256 and 128, respectively. The exact number of samples for the prediction task follows that from [14].

10.2. Simulation settings

Architecture details. We use a simple convolutional network with three convolution layers for CMNIST, with feature map dimensions of 64, 128 and 256, each followed by a ReLU activation and a batch normalization layer following [37]. For CIFAR10-C and BFFHQ, we use ResNet-18 with pretrained weights provided in PyTorch `torchvision` implementations. Each convolutional network and ResNet-18 includes 1.3×10^6 and 2.2×10^7 number of parameters, respectively. We assign a pruning parameter for each weight parameter except bias in deep networks. Each of pruning parameter is initialized with value 1.5 so that the initial probability of preserving the corresponding weight is set to $\sigma(1.5) \approx 0.8$ in default.

Training details. We first train bias-capturing networks using GCE loss ($q=0.7$) for (CMNIST, BFFHQ, CelebA), with (2000, 10000, 10000) iterations, respectively. For CIFAR10-C, we use epoch-ensemble-based mining algorithms presented in [40], which select samples cooperated with an ensemble of predictions at each epoch to prevent overfitting. We use b-score threshold $\tau = 0.8$ and the confidence threshold $\eta = 0.05$ as suggested in the original paper.

Then, main networks are pretrained for 10000 iterations using an Adam optimizer with learning rates 0.01, 0.001, 0.001 and 0.0001 for CMNIST, CIFAR10-C, BFFHQ, and CelebA, respectively.

We train pruning parameters for 2000 iterations using a learning rate 0.01, upweighting hyperparameter $\lambda_{up} = 80$ and a balancing hyperparameter $\lambda_{align} = 0.05$ for each dataset. We use a Lagrangian multiplier $\lambda_{\ell_1} = 10^{-8}$ for CMNIST, and $\lambda_{\ell_1} = 10^{-9}$ for CIFAR10-C, BFFHQ and CelebA. Specifically, we set λ_{ℓ_1} by considering the size of deep networks, where we found that the value within range $\mathcal{O}(0.1 * n^{-1})$ serves as a good starting point where n is the number of parameters.

After pruning, we finetune the networks with decaying learning rate to 0.001 for CMNIST and 0.0005 for others. We use $\lambda_{align} = 0.05$ consistently. Then, we use $\lambda_{up} = 80$ for BFFHQ, $\lambda_{up} = 20$ for CelebA, and $\lambda_{up} = \{10, 30, 50, 80\}$ for CMNIST and CIFAR10-C with $\{0.5\%, 1.0\%, 2.0\%, 5.0\%\}$ of bias ratio, respectively.

Considering the pruning as a strong regularization, we did not use additional capacity control techniques such as early stopping or strong ℓ_2 regularization presented in [24,30].

Data augmentations. We did not use any kinds of data augmentations which may implicitly enforce networks to encode invariances. For the BFFHQ and CelebA dataset, we only apply random horizontal flip. For the CIFAR10-C dataset, we take 32×32 random crops from image padded by 4 pixels followed by random horizontal flip, following [27]. We do not use any kinds of augmentations in CMNIST.

Baselines. We use the official implementations of Rebias, LfF, DisEnt and JTT released by authors, and reproduce EnD and MRM by ourselves. For DisEnt, we use the official hyperparameter configurations provided in the original paper. We use $q = 0.7$ for LfF as suggested by authors on every experiment. For Rebias, we use the official hyperparameter configurations for CMNIST, and train for 200 epochs using Adam optimizer with learning rate 0.001 and RBF kernel radius of 1 for other datasets. For MRM, we use λ_{ℓ_1} of 10^{-8} for CMNIST following the original paper, and 10^{-9} for the others. For EnD, we set the multipliers α for disentangling and β for entangling to 1.

References

- [1] Hyojin Bahng, Sanghyuk Chun, Sangdoon Yun, Jaegul Choo, and Seong Joon Oh. Learning de-biased representations with biased representations. In *International Conference on Machine Learning*, pages 528–539. PMLR, 2020. [2](#), [7](#), [18](#)
- [2] Haoyue Bai, Fengwei Zhou, Lanqing Hong, Nanyang Ye, S-H Gary Chan, and Zhenguo Li. Nas-ood: Neural architecture search for out-of-distribution generalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8320–8329, 2021. [2](#)
- [3] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020. [6](#)
- [4] Sam Corbett-Davies and Sharad Goel. The measure and mismeasure of fairness: A critical review of fair machine learning. *arXiv preprint arXiv:1808.00023*, 2018. [2](#)
- [5] Róbert Csordás, Sjoerd van Steenkiste, and Jürgen Schmidhuber. Are neural nets modular? inspecting functional modularity through differentiable weight masks. *arXiv preprint arXiv:2010.02066*, 2020. [1](#)
- [6] James Diffenderfer, Brian Bartoldson, Shreya Chaganti, Jize Zhang, and Bhavya Kailkhura. A winning hand: Compressing deep networks can improve out-of-distribution robustness. *Advances in Neural Information Processing Systems*, 34:664–676, 2021. [2](#)
- [7] Michael Feldman, Sorelle A Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkatasubramanian. Certifying and removing disparate impact. In *proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*, pages 259–268, 2015. [2](#)
- [8] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. *arXiv preprint arXiv:1803.03635*, 2018. [2](#), [6](#), [7](#)
- [9] Songwei Ge, Shlok Mishra, Chun-Liang Li, Haohan Wang, and David Jacobs. Robust contrastive learning using negative samples with diminished semantics. *Advances in Neural Information Processing Systems*, 34:27356–27368, 2021. [2](#)
- [10] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. *arXiv preprint arXiv:1811.12231*, 2018. [2](#)
- [11] Karan Goel, Albert Gu, Yixuan Li, and Christopher Ré. Model patching: Closing the subgroup performance gap with data augmentation. *arXiv preprint arXiv:2008.06775*, 2020. [2](#)
- [12] Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel R Bowman, and Noah A Smith. Annotation artifacts in natural language inference data. *arXiv preprint arXiv:1803.02324*, 2018. [2](#)
- [13] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *arXiv preprint arXiv:1903.12261*, 2019. [2](#), [19](#)
- [14] Youngkyu Hong and Eunho Yang. Unbiased classification through bias-contrastive and bias-balanced learning. *Advances in Neural Information Processing Systems*, 34:26449–26461, 2021. [7](#), [19](#)
- [15] Weihua Hu, Gang Niu, Issei Sato, and Masashi Sugiyama. Does distributionally robust supervised learning give robust classifiers? In *International Conference on Machine Learning*, pages 2029–2037. PMLR, 2018. [2](#)
- [16] Wei Hu, Lechao Xiao, and Jeffrey Pennington. Provable benefit of orthogonal initialization in optimizing deep linear networks. *arXiv preprint arXiv:2001.05992*, 2020. [16](#)
- [17] Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander Madry. Adversarial examples are not bugs, they are features. *Advances in neural information processing systems*, 32, 2019. [4](#)
- [18] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*, 2016. [5](#)
- [19] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019. [19](#)
- [20] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *Advances in Neural Information Processing Systems*, 33:18661–18673, 2020. [6](#)
- [21] Eungyeup Kim, Jihyeon Lee, and Jaegul Choo. Biaswap: Removing dataset bias with bias-tailored swapping augmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14992–15001, 2021. [2](#), [19](#)
- [22] Jungsoo Lee, Eungyeup Kim, Juyoung Lee, Jihyeon Lee, and Jaegul Choo. Learning debiased representation via disentangled feature augmentation. *Advances in Neural Information Processing Systems*, 34:25123–25133, 2021. [2](#), [7](#), [18](#), [19](#)
- [23] Yi Li and Nuno Vasconcelos. Repair: Removing representation bias by dataset resampling. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9572–9581, 2019. [2](#)
- [24] Evan Z Liu, Behzad Haghgoo, Annie S Chen, Aditi Raghunathan, Pang Wei Koh, Shiori Sagawa, Percy Liang, and Chelsea Finn. Just train twice: Improving group robustness without training group information. In *International Conference on Machine Learning*, pages 6781–6792. PMLR, 2021. [2](#), [6](#), [7](#), [19](#)
- [25] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of the IEEE international conference on computer vision*, pages 3730–3738, 2015. [2](#)
- [26] Vaishnavh Nagarajan, Anders Andreassen, and Behnam Neyshabur. Understanding the failure modes of out-of-distribution generalization. *arXiv preprint arXiv:2010.15775*, 2020. [1](#), [3](#)

- [27] Junhyun Nam, Hyuntak Cha, Sungsoo Ahn, Jaeho Lee, and Jinwoo Shin. Learning from failure: De-biasing classifier from biased classifier. *Advances in Neural Information Processing Systems*, 33:20673–20684, 2020. [2](#), [5](#), [7](#), [18](#), [19](#)
- [28] Arvind Narayanan. Translation tutorial: 21 fairness definitions and their politics. In *Proc. Conf. Fairness Accountability Transp., New York, USA*, volume 1170, page 3, 2018. [2](#)
- [29] Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. *arXiv preprint arXiv:1911.08731*, 2019. [2](#), [7](#), [19](#)
- [30] Shiori Sagawa, Aditi Raghunathan, Pang Wei Koh, and Percy Liang. An investigation of why overparameterization exacerbates spurious correlations. In *International Conference on Machine Learning*, pages 8346–8356. PMLR, 2020. [2](#), [5](#), [19](#)
- [31] Andrew M Saxe, James L McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *arXiv preprint arXiv:1312.6120*, 2013. [16](#)
- [32] Nimit Sohoni, Jared Dunmon, Geoffrey Angus, Albert Gu, and Christopher Ré. No subclass left behind: Fine-grained robustness in coarse-grained classification problems. *Advances in Neural Information Processing Systems*, 33:19339–19352, 2020. [2](#)
- [33] Enzo Tartaglione, Carlo Alberto Barbano, and Marco Grangetto. End: Entangling and disentangling deep representations for bias correction. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 13508–13517, 2021. [7](#), [18](#)
- [34] Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry. Robustness may be at odds with accuracy. *arXiv preprint arXiv:1805.12152*, 2018. [1](#), [3](#), [4](#)
- [35] Chaoqi Wang, Guodong Zhang, and Roger Grosse. Picking winning tickets before training by preserving gradient flow. *arXiv preprint arXiv:2002.07376*, 2020. [17](#), [18](#)
- [36] Zeyu Wang, Klint Qinami, Ioannis Christos Karakozis, Kyle Genova, Prem Nair, Kenji Hata, and Olga Russakovsky. Towards fairness in visual recognition: Effective strategies for bias mitigation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8919–8928, 2020. [2](#)
- [37] Dinghui Zhang, Kartik Ahuja, Yilun Xu, Yisen Wang, and Aaron Courville. Can subnetwork structure be the key to out-of-distribution generalization? In *International Conference on Machine Learning*, pages 12356–12367. PMLR, 2021. [1](#), [2](#), [3](#), [5](#), [7](#), [8](#), [19](#)
- [38] Michael Zhang, Nimit S Sohoni, Hongyang R Zhang, Chelsea Finn, and Christopher Ré. Correct-n-contrast: A contrastive approach for improving robustness to spurious correlations. *arXiv preprint arXiv:2203.01517*, 2022. [5](#), [17](#)
- [39] Zhilu Zhang and Mert Sabuncu. Generalized cross entropy loss for training deep neural networks with noisy labels. *Advances in neural information processing systems*, 31, 2018. [5](#)
- [40] Bowen Zhao, Chen Chen, Qi Ju, and Shutao Xia. Learning debiased models with dynamic gradient alignment and bias-conflicting sample mining. *arXiv preprint arXiv:2111.13108*, 2021. [19](#)