

Implicit Counterfactual Data Augmentation for Robust Learning

Xiaoling Zhou, *Member, IEEE*, Ou Wu, and Michael K. Ng, *Senior Member, IEEE*

Abstract—Machine learning models are prone to capturing the spurious correlations between non-causal attributes and classes, with counterfactual data augmentation being a promising direction for breaking these spurious associations. However, generating counterfactual data explicitly poses a challenge, and incorporating augmented data into the training process decreases training efficiency. This study proposes an Implicit Counterfactual Data Augmentation (ICDA) method to remove spurious correlations and make stable predictions. Specifically, first, a novel sample-wise augmentation strategy is developed that generates semantically and counterfactually meaningful deep features with distinct augmentation strength for each sample. Second, we derive an easy-to-compute surrogate loss on the augmented feature set when the number of augmented samples becomes infinite. Third, two concrete schemes are proposed, including direct quantification and meta-learning, to derive the key parameters for the robust loss. In addition, ICDA is explained from a regularization perspective, revealing its capacity to improve intra-class compactness and augment margins at both class and sample levels. Extensive experiments have been conducted across various biased learning scenarios covering both image and text datasets, demonstrating that ICDA consistently enhances the generalization and robustness performance of popular networks.

Index Terms—Counterfactual, implicit augmentation, spurious correlation, meta-learning, generalization.

I. INTRODUCTION

DEEP learning models are supposed to learn invariances and make stable predictions based on some right causes. However, models trained with empirical risk minimization are prone to learning spurious correlations and suffer from high generalization errors when the training and test distributions do not match [1]–[3]. For example, dogs are mostly on the grass in the training set. Thus, a dog in the water can easily be misclassified as a “drake” due to its rare scene context (“water”) in the “dog” class, which is illustrated in Fig. 1. A promising solution for improving the models’ generalization and robustness is to learn causal models [4] as if a model can concentrate more on the causal correlations but not the spurious associations between non-causal attributes and classes, stable and exact predictions are more likely.

Manuscript received XXX; revised XXX; accepted XXX. (Corresponding author: Ou Wu)

Xiaoling Zhou is with the Center for Applied Mathematics, Tianjin University, Tianjin 300072, China, and the National Engineering Research Center for Software Engineering, Peking University, Beijing 100091, China (e-mail: xiaolingzhou@stu.pku.edu.cn).

Ou Wu is with the Center for Applied Mathematics, Tianjin University, Tianjin 300072, China, and the Hangzhou Institute for Advanced Study, University of Chinese Academy of Sciences, Hangzhou 310024, China (e-mail: wuou@tju.edu.cn).

Michael K. Ng is with the Department of Mathematics, Hong Kong Baptist University, Hong Kong 999077, China (e-mail: michael-n@hkbu.edu.hk).

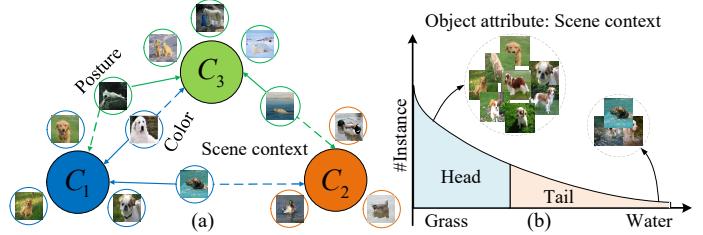


Fig. 1. (a): Illustration for images affected by spurious correlations due to rare attributes (e.g., posture, color, and scene context). C_1 , C_2 , and C_3 are the dog, drake, and polar bear classes, respectively. The solid line connects the sample’s ground-truth class, and the dotted line connects the class with a spurious correlation with the sample. (b): Illustration for attribute imbalance. Regarding the attribute of scene context, the majority of dogs in the training data are situated on grass, while only a small number are depicted in water. Imbalances in attributes generally lead to spurious correlations between non-causal attributes and labels in deep learning models.

Counterfactual augmentation has become popular for causal models because of its capacity to enhance model robustness and being model-agnostic. For instance, Lu et al. [5] and He et al. [6] augmented the data effectively by swapping identity pronouns in texts. Moreover, Chang et al. [7] introduced two new image generation procedures that included counterfactual and factual data augmentations to reduce spuriousness between backgrounds of images and labels, achieving higher accuracy in several challenging datasets. Mao et al. [2] utilized a novel strategy to learn robust representations that steered generative models to manufacture interventions on features caused by confounding factors. Nevertheless, the methods presented above suffer from several shortcomings. Specifically, it is not trivial to explicitly distinguish between causal and non-causal attributes, and the training efficiency will decline as excess augmented images are involved in training.

It should be mentioned that implicit data augmentation settles the inefficiency of explicit augmentation by avoiding the generation of excess samples. Implicit Semantic Data Augmentation (ISDA) [8] conducts a pioneering study on implicit data augmentation. It is inspired by the observation that deep features in a network are typically linearized, resulting in the existence of numerous semantic directions in the deep feature space. Then, it translates samples along the semantic directions in the feature space based on an assumed class-wise augmentation distribution. By deriving an upper bound on the expected cross-entropy (CE) loss, ISDA enables optimization of only the upper bound to achieve data augmentation in an efficient way. Subsequent studies on imbalance learning have expanded upon this approach. For instance, MetaSAug [9] optimizes the covariance matrix of the

tail classes on a balanced metadata set to mitigate the issue of inaccurate estimation arising from the insufficient number of samples in the tail classes, yielding good performance on imbalanced data. Besides, to generate more diverse samples for tail classes, Reasoning-based Implicit Semantic Data Augmentation (RISDA) [10] augments samples in tail classes using semantic vectors from not only the current class but also the relevant classes. However, these methods, specifically designed for imbalanced learning, may not effectively dismantle the spurious associations within deep learning models. Moreover, they adopt purely class-wise semantic augmentation strategies, and thus samples in the same class have identical augmentation distributions that are inaccurate. As illustrated in Fig. 1(a), samples in the same class may exhibit spurious correlations with different classes due to various attributes. Consequently, an ideal augmentation strategy should consider these sample-wise non-causal attributes.

This study proposes a novel sample-wise **Implicit Counterfactual Data Augmentation** (ICDA) method that facilitates both semantic and counterfactual augmentations. Semantic augmentation is accomplished by transforming samples along vectors drawn from the deep feature space of the ground-truth class. Moreover, counterfactual augmentation is realized by manipulating samples along vectors sourced from the deep feature spaces of non-target classes. The augmentation distribution and strength for each sample are determined based on class-wise statistical information and the degree of spurious correlation between the sample and each class. Then, we verify that ICDA is approximate to a novel robust surrogate loss (termed the ICDA loss) by considering the number of augmentations becoming infinite, making the process highly efficient. Furthermore, meta-learning is introduced to learn key parameters in this novel loss, which is analyzed and compared against existing methods in a unified regularization perspective, revealing that it enforces extra intra-class compactness by reducing the classes' mapped variances and encourages larger sample margins and class-boundary distances. Extensive experiments verify that ICDA consistently achieves state-of-the-art performance in several typical learning scenarios requiring the models to be robust and presenting a high generalization ability. Furthermore, the visualization results indicate that ICDA generates more diverse and meaningful counterfactual images with rare attributes, helping models break spurious correlations and affording stable predictions for the right reasons.

II. RELATED WORK

A. Data Augmentation

Counterfactual and implicit semantic augmentation strategies are reviewed here. Counterfactual augmentation generates hypothetical samples (i.e., counterfactuals) by making small changes to the original samples, which can be divided into hand-crafted [6], [7] and using causal generative models [11], [12], demonstrating competitive performance. However, explicitly finding the non-causal attributes is challenging and training models based on the augmented data is inefficient. Implicit semantic data augmentation [8]–[10] overcomes the

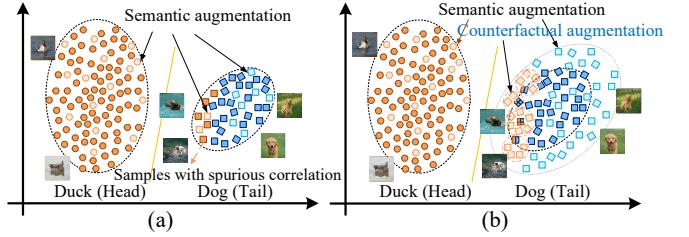


Fig. 2. (a): Diagram for ISDA, which only conducts semantic augmentation and treats all samples equally. (b): Diagram for ICDA, containing semantic and counterfactual augmentations. Samples in the tail class (Dog) and those with rare attributes (yellow ones in Dog class) are augmented most. The two axes mean the dimensions of the 2D feature space, in which each sample is represented by a dot or square. Solid and transparent samples are the original and augmented ones. Samples in the same red circle are augmented from the same sample. The augmentation strength is determined by the degree of spurious associations.

inefficiency of explicit data augmentation approaches as it does not generate excess samples and achieves the effect of augmentation only through optimizing the surrogate loss when the number of augmented samples becomes infinite, which is thus efficient. However, implicit data augmentation can not help overcome spurious associations. Besides, all samples or samples in the same class are treated equally in existing approaches, which is naturally not optimal.

B. Logit Adjustment

Logit vectors represent the outputs before the Softmax layer in the majority of deep classifiers. Logit adjustment approaches involve introducing perturbation terms to logits, aiming to bolster the robustness of models. This method is originally proposed in face recognition [13], [14], seeking to increase inter-class distance and intra-class compactness. Presently, logit adjustment is employed, implicitly or explicitly, in various contexts such as data augmentation [8], [10] and long-tailed classifications [15]–[17]. For example, the Logit-adjusted (LA) loss incorporates class proportion terms as perturbations, demonstrating effectiveness in imbalanced learning scenarios. Additionally, ISDA can be considered a logit adjustment approach, contributing significantly to the enhancement of models' generalization capabilities.

III. IMPLICIT COUNTERFACTUAL DATA AUGMENTATION

Notation. Consider training a network G with weights \mathbf{W} on a training set $D^{train} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$, where $y_i \in \{1, \dots, C\}$ is the label of the i -th sample \mathbf{x}_i over C classes. Let the H -dimensional vector $\mathbf{h}_i = G(\mathbf{x}_i, \mathbf{W})$ denote the deep feature of \mathbf{x}_i learned by G . Let $\mathbf{u}_i = f(\mathbf{h}_i) = \mathbf{w}\mathbf{h}_i + \mathbf{b}$ denote the logit vector, $\mathbf{w} = [\mathbf{w}_1, \dots, \mathbf{w}_C]^T \in \mathbb{R}^{C \times H}$, and $\mathbf{b} = [b_1, \dots, b_C]^T \in \mathbb{R}^C$. Let μ_c and Σ_c be the mean and covariance matrix of the deep features for class c . $\mathcal{N}(\mu, \Sigma)$ means a multivariate normal distribution with mean vector μ and covariance matrix Σ .

A. Counterfactual Data Augmentation

To mitigate spurious correlations between non-causal attributes and classes, we propose a counterfactual data augmentation strategy, which generates both meaningful semantic

and counterfactual samples. Considering that the spurious correlations between samples and classes are sample-wise, we devise and utilize sample-level augmentation distributions. To achieve semantic augmentation, perturbation vectors for the deep feature of each sample, \mathbf{h}_i , are sampled from a multivariate normal distribution $\mathcal{N}(\mathbf{0}, \Sigma_{y_i})$. To mitigate spurious correlations, we intervene on non-causal attributes that are spuriously correlated with other classes, while preserving the core object features to generate counterfactual instances. Specifically, the deep features of samples \mathbf{h}_i are transformed along the perturbation vectors extracted from the deep feature spaces of non-target classes, i.e., $\mathcal{N}(\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)$, where $c \neq y_i$. Consequently, when augmenting the deep feature \mathbf{h}_i to class c , the perturbation vectors are sampled from $\mathcal{N}(\mathbf{0} + \alpha_{i,c}\boldsymbol{\mu}_c, \Sigma_{y_i} + \alpha_{i,c}\boldsymbol{\Sigma}_c)$, where $\alpha_{i,c} (\geq 0)$ is determined by the degree of the spurious association between x_i and class c .

As for the augmentation strength, that is the number of augmented samples $\tilde{M}_{i,c}$ for sample x_i to class c , it is assumed to follow $\tilde{M}_{i,c} = (M\alpha_{i,c})/\pi_{y_i}$, where π_{y_i} denotes the proportion of class y_i and M is a constant. Consequently, the higher the degree of the spuriousness between x_i and class c and the smaller the π_{y_i} , the larger the number ($\tilde{M}_{i,c}$) of samples will be augmented from x_i to class c . Fig. 2(a) highlights that the existing augmentation approaches ignore the relationship between samples and the other class and all samples are treated equally, prohibiting a well-adjusted distribution. Fig. 2(b) presents our augmentation manner, in which samples in the tail class and the ones most spuriously correlated with attributes of the other class are augmented most, facilitating enhancing the generalization and robustness of models against spurious correlations.

During training, C feature means and covariance matrices are computed, one for each class. To enhance efficiency, the values of $\boldsymbol{\mu}_c$ and $\boldsymbol{\Sigma}_c$ are computed online by aggregating statistics from all mini-batches, which is given in Appendix A. Given that the estimated statistics information in the first few epochs is not quite informative, we add a scale parameter $\lambda = (t/\mathcal{T}) \times \lambda^0$ before the estimated $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$, where t and \mathcal{T} refer to the numbers of the current and total iterations. Additionally, λ_0 is a hyperparameter. The augmented feature $\mathbf{h}_{i,c}$ for \mathbf{h}_i to class c is obtained by translating \mathbf{h}_i along a random direction sampled from the above multivariate normal distribution. Equivalently, we have $\mathbf{h}_{i,c} \sim \mathcal{N}(\mathbf{h}_i + \lambda\alpha_{i,c}\boldsymbol{\mu}_c, \lambda(\Sigma_{y_i} + \alpha_{i,c}\boldsymbol{\Sigma}_c))$.

Notably, our augmentation strategy has distinct differences from current semantic augmentation methods:

- Their motivations are different. Our strategy aims to generate more counterfactual data for breaking spurious correlations, while the existing methods only generate diverse semantic data.
- Their granularities are different. Our augmentation strategy is sample-wise, which is fine-grained and pinpoint, while current schemes involve class-wise strategies.
- Our strategy highlights the augmentation strength, which is crucial in an augmentation strategy as inappropriate class and attribute distributions always cause spuriousness. However, it is overlooked by the existing methods.

B. New Loss under Implicit Augmentation

A naive method to implement ICDA is to explicitly augment the deep features of samples based on the designed augmentation distribution and strength. Specifically, for class c ($\neq y_i$), the deep features \mathbf{h}_i are augmented $\tilde{M}_{i,c}$ times utilizing perturbation vectors sampled from the corresponding distribution $\mathcal{N}(\mathbf{0} + \lambda\alpha_{i,c}\boldsymbol{\mu}_c, \lambda(\Sigma_{y_i} + \alpha_{i,c}\boldsymbol{\Sigma}_c))$. Consequently, an augmented feature set $\{\{\mathbf{h}_{i,c}^1, \dots, \mathbf{h}_{i,c}^{\tilde{M}_{i,c}}\}_{c=1, c \neq y_i}^C\}_{i=1}^N$ can be formed. Then, the corresponding CE loss for all augmented features is

$$\mathcal{L}_M(\mathbf{w}, \mathbf{b}, \mathbf{W}) = \frac{1}{\tilde{M}} \sum_{i=1}^N \sum_{c \neq y_i} \sum_{k=1}^{\tilde{M}_{i,c}} -\log \frac{\exp[f_{y_i}(\mathbf{h}_{i,c}^k)]}{\sum_{j=1}^C \exp[f_j(\mathbf{h}_{i,c}^k)]}, \quad (1)$$

where $\tilde{M} = \sum_{i=1}^N \sum_{c=1, c \neq y_i}^C \tilde{M}_{i,c}$ and $f_j(\mathbf{h}_{i,c}^k) = \mathbf{w}_j^T \mathbf{h}_{i,c}^k + b_j$. To augment more data while enhancing training efficiency, we let M in $\tilde{M}_{i,c}$ grow to infinity. Then, the expected CE loss for all augmented features is

$$\mathcal{L}_\infty(\mathbf{w}, \mathbf{b}, \mathbf{W}) = \frac{1}{\tilde{N}} \sum_{i=1}^N \sum_{c \neq y_i} \tilde{N}_{i,c} \mathbb{E}_{\mathbf{h}_{i,c}} [-\log \frac{\exp(f_{y_i}(\mathbf{h}_{i,c}))}{\sum_{j=1}^C \exp(f_j(\mathbf{h}_{i,c}))}], \quad (2)$$

where $\tilde{N}_{i,c} = \alpha_{i,c}/\pi_{y_i}$ and $\tilde{N} = \sum_{i=1}^N \sum_{c=1, c \neq y_i}^C \tilde{N}_{i,c}$. However, the above expected CE loss is hard to calculate. Then, we derive a more easy-to-compute surrogate loss for Eq. (2), which is as follows:

$$\begin{aligned} \mathcal{L}_s(\mathbf{w}, \mathbf{b}, \mathbf{W}) &= \frac{1}{\tilde{N}} \sum_{i=1}^N \frac{1}{\pi_{y_i}} \log(1 + \sum_{c \neq y_i} \exp[f_c(\mathbf{h}_i) - f_{y_i}(\mathbf{h}_i) + \phi_{i,c}]), \\ \phi_{i,c} &= (\lambda/2)P_{c,i} + \lambda Q_{c,i} + \beta\alpha_i, \end{aligned} \quad (3)$$

where $P_{c,i} = \Delta\mathbf{w}_{c,y_i}(\Sigma_{y_i} + \sum_{j=1, j \neq y_i}^C \hat{\alpha}_{i,j}\boldsymbol{\Sigma}_j)\Delta\mathbf{w}_{c,y_i}^T$ and $Q_{c,i} = \Delta\mathbf{w}_{c,y_i} \sum_{j=1, j \neq y_i}^C \hat{\alpha}_{i,j}\boldsymbol{\mu}_j$, in which $\Delta\mathbf{w}_{c,y_i} = \mathbf{w}_c^T - \mathbf{w}_{y_i}^T$ and $\hat{\alpha}_{i,j} = \alpha_{i,j}/(C-1)$. In addition, $\alpha_i = \sum_{j=1, j \neq y_i}^C \hat{\alpha}_{i,j}$. The inference details are presented in Appendix B. Consequently, instead of conducting the augmentation process explicitly, we can directly optimize this surrogate loss.

Although $\mathcal{L}_s(\mathbf{w}, \mathbf{b}, \mathbf{W})$ can be directly utilized during training, a more effective loss is leveraged after adopting the two following modifications: (1) Inspired by the manner in LA [16], the class-wise weight $1/\pi_{y_i}$ is replaced by a perturbation term on logits. (2) We only retain the term $\Delta\mathbf{w}_{c,y_i}\hat{\alpha}_{i,c}\boldsymbol{\mu}_c$ in $Q_{c,i}$. The reason for the proposed variation is detailed in Section V. Accordingly, the final ICDA training loss becomes

$$\begin{aligned} \bar{\mathcal{L}}_s(\mathbf{w}, \mathbf{b}, \mathbf{W}) &= \frac{1}{\tilde{N}} \sum_{i=1}^N \log(1 + \sum_{c \neq y_i} \exp[f_c(\mathbf{h}_i) - f_{y_i}(\mathbf{h}_i) + \hat{\phi}_{i,c}]), \\ \hat{\phi}_{i,c} &= (\lambda/2)P_{c,i} + \lambda\Delta\mathbf{w}_{c,y_i}\hat{\alpha}_{i,c}\boldsymbol{\mu}_c + \delta_{c,i} + \beta\alpha_i, \end{aligned} \quad (4)$$

where $\delta_{c,i} = \log(\pi_c/\pi_{y_i})$ and β is a hyperparameter which is fixed as 0.1 in our experiments. Notably, the ICDA loss can be considered as a generalization of several typical logit adjustment losses. For example, when $\lambda = \beta = 0$, our method can be reduced to LA. Additionally, for $\hat{\alpha}_{i,c} = \beta = 0$ ($c \neq y_i$) and balanced classes, ICDA degenerates to ISDA. Section V further demonstrates the superiority of the ICDA loss against current methods from a regularization perspective.

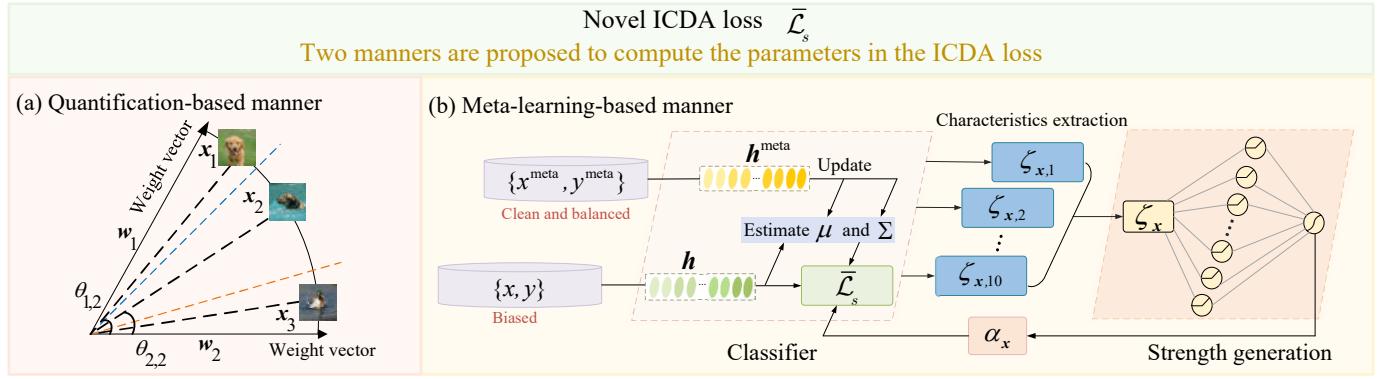


Fig. 3. Two manners for applying our proposed ICDA loss: quantification-based manner and meta-learning-based manner. (a): Illustration for the angle between sample feature and weight vector. The angle between the deep feature of x_2 and the classifier weights for class C_2 is smaller compared to the angle between the deep feature of x_1 and the classifier weights for class C_2 , attributed to the spurious correlation between x_2 and C_2 . (b): The overall structure of Meta-ICDA, which comprises three main components: the classifier, the characteristics extraction module, and the strength generation network.

IV. LEARNING WITH ICDA

In realization, parameters μ_c , Σ_c , and $\alpha_{i,c}$ in the ICDA training loss should be prefixed. Therefore, two approaches, including a direct quantification-based manner and a meta-learning-based manner, are proposed to optimize classifiers using the ICDA loss.

A. Direct Quantification-based Manner

The spurious correlation between sample x_i and class c can be directly quantified by the angle ($\theta_{i,c}$) between \mathbf{h}_i and the weight vector w_c of class c . Naturally, the larger the spurious correlation between \mathbf{h}_i and class c , the smaller the $\theta_{i,c}$ and the larger the $\cos \theta_{i,c}$. An illustration is presented in Fig. 3(a). Samples x_1 and x_2 both belong to class C_1 . Nevertheless, $\theta_{2,2}$ is smaller than $\theta_{1,2}$ as x_2 is more spuriously correlated with class C_2 .

Since $\alpha_{i,c}$ is determined by the degree of spurious correlation between x_i and class c , it should be positively correlated with $\cos \theta_{i,c}$. Moreover, only when the direction of \mathbf{h}_i is partially consistent with that of w_c (i.e., $\theta_{i,c} < 90^\circ$), the information of class c should be utilized to augment sample x_i . Thus, we denote $\alpha_{i,c} = \max(\cos \theta_{i,c}, 0)$, where a larger $\alpha_{i,c}$ value means a larger counterfactual augmentation strength. Then, we have $\alpha_i \propto \sum_{c \neq y_i} \max(\cos \theta_{i,c}, 0)$. Nevertheless, quantifying α_i through angle θ_{i,y_i} is more direct. If x_i is notably influenced by the spurious correlations with other classes, then θ_{i,y_i} will be large, and $\cos \theta_{i,y_i}$ will be small. Thus, α_i should be negatively correlated with $\cos \theta_{i,y_i}$. Meanwhile, the value range of α_i is restricted to $[0, 1]$. Therefore, we let $\alpha_i = (1 - \cos \theta_{i,y_i})/2$. This manner is empirically verified to be more effective.

B. Meta-learning-based Manner

If metadata are available, the extent a sample is affected by spurious correlation can be better determined by training a strength generation network. This manner is called Meta-ICDA. The input of the strength generation network involves ten training characteristics of samples ζ_i , including loss, margin, uncertainty, etc., denoted by $\zeta_{i,1}, \dots, \zeta_{i,10}$. Details for the

Algorithm 1 Meta-ICDA

Input: Training data D^{train} , metadata D^{meta} , batch sizes n and m , step sizes η_1 and η_2 , number of iterations \mathcal{T} .

Output: Learned $\tilde{\mathbf{W}}$ and Ω .

- 1: Initialize $\tilde{\mathbf{W}}^{(1)}$ and $\Omega^{(1)}$;
 - 2: **for** $t = 1$ to \mathcal{T} **do**
 - 3: Sample $\{(x_i, y_i)\}_{i=1}^n$ from D^{train} ;
 - 4: Sample $\{(x_i^{\text{meta}}, y_i^{\text{meta}})\}_{i=1}^m$ from D^{meta} ;
 - 5: Calculate current feature means $\mu^{(t)}$ and covariance matrices $\Sigma^{(t)}$;
 - 6: Formulate $\tilde{\mathbf{W}}^{(t)}(\Omega)$ by Eq. (5);
 - 7: Update $\Omega^{(t+1)}$ by Eq. (6);
 - 8: Update $\mu^{(t+1)}$ and $\Sigma^{(t+1)}$ by Eqs. (7) and (8);
 - 9: Update $\tilde{\mathbf{W}}^{(t+1)}$ by Eq. (9);
 - 10: **end for**
-

extracted training characteristics are presented in Appendix C. In this study, the network is a two-layer MLP and its output is the augmentation strength α_i . Thus, we have $\alpha_i = \text{MLP}(\zeta_i)$. Considering that the geodesic distance between the sample and other classes can well measure their correlation [18], $\alpha_{i,c}$ is still calculated by $\max(\cos \theta_{i,c}, 0)$. The estimated covariance matrices and the feature means are optimized on metadata because biased training data (e.g., imbalanced and noisy data) may not estimate the statistical information well. Fig. 3(b) illustrates the framework of Meta-ICDA, which includes three main parts: the classifier network, the strength generation network, and the characteristics extraction module. We utilize an online meta-learning-based learning strategy to alternatively update the parameters of the classifier and the strength generation network. The optimization process is detailed below.

To ease this paper's notation, the deep classifier network's parameters of \mathbf{W} and \mathbf{w} are denoted as $\tilde{\mathbf{W}}$. The deep classifier which includes both feature extractor G and classifier f is denoted as \tilde{f} . The parameters in the strength generation network are Ω . The small metadata set is denoted as $D^{\text{meta}} = \{(x_i^{\text{meta}}, y_i^{\text{meta}})\}_{i=1}^B$, where $B \leq N$.

TABLE I
REGULARIZATION TERMS AND REFLECTED GENERALIZATION FACTORS OF THE FOUR ALGORITHMS (LA, ISDA, RISDA, AND ICDA).

Method	Regularization term	Generalization factor
LA	$R_{LA} = \sum_{i=1}^N \sum_{c \neq y_i} q_{i,c} \delta_{c,i}$	✓ Class-wise margin
ISDA	$R_{ISDA} = \frac{\lambda}{2} \sum_{i=1}^N \sum_{c \neq y_i} q_{i,c} \Delta \mathbf{w}_{c,y_i} \Sigma_{y_i} \Delta \mathbf{w}_{c,y_i}^T$	✓ Intra-class compactness
RISDA	$R_{RISDA} = \sum_{i=1}^N \sum_{c \neq y_i} q_{i,c} [\alpha \Delta \mathbf{w}_{c,y_i} \sum_{j=1, j \neq y_i}^C \varepsilon_{y_j, j} \mu_j + \beta \Delta \mathbf{w}_{c,y_i} (\Sigma_{y_i} + \sum_{j=1, j \neq y_i}^C \varepsilon_{y_j, j} \Sigma_j) \Delta \mathbf{w}_{c,y_i}^T]$	✓ Intra-class compactness ✓ Class-boundary distance
ICDA	$R_{ICDA} = \sum_{i=1}^N \left\{ \sum_{c \neq y_i} q_{i,c} [\delta_{c,i} + \frac{\lambda}{2} \Delta \mathbf{w}_{c,y_i} (\Sigma_{y_i} + \sum_{j=1, j \neq y_i}^C \hat{\alpha}_{i,j} \Sigma_j) \Delta \mathbf{w}_{c,y_i}^T + \lambda \Delta \mathbf{w}_{c,y_i} \hat{\alpha}_{i,c} \mu_c] - \beta \alpha_i q_{i,y_i} \right\}$	✓ Sample-wise/class-wise margin ✓ Intra-class compactness ✓ Class-boundary distance

During this process, first, Ω is treated as the to-be-updated parameter, and the parameter of the deep classifier \tilde{f} , that is \tilde{W} , which is a function of Ω , is formulated. We utilize the stochastic gradient descent (SGD) optimizer to optimize the training loss on a minibatch of training samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ in each iteration, where n is the size of the mini-batch. Thus, \tilde{W} is formulated by the following equation:

$$\tilde{\mathbf{W}}^{(t)}(\Omega) = \tilde{\mathbf{W}}^{(t)} - \eta_1 \frac{1}{n} \sum_{i=1}^n \nabla_{\tilde{\mathbf{W}}} \ell_{ICDA}(\tilde{f}(\mathbf{x}_i), y_i; \alpha_i^{(t)})|_{\tilde{\mathbf{W}}^{(t)}}, \quad (5)$$

where η_1 is the step size. After extracting the training characteristics from the classifier, the parameter of the strength generation network Ω can be updated on a minibatch of metadata $\{(\mathbf{x}_i^{meta}, y_i^{meta})\}_{i=1}^m$ as follows:

$$\Omega^{(t+1)} = \Omega^{(t)} - \eta_2 \frac{1}{m} \sum_{i=1}^m \nabla_{\Omega} \ell_{CE}(\tilde{f}_{\tilde{\mathbf{W}}(\Omega^{(t)})}(\mathbf{x}_i^{meta}), y_i^{meta})|_{\Omega^{(t)}}, \quad (6)$$

where m and η_2 are the minibatch size of metadata and the step size, respectively. At the same time, the feature means and covariance matrices for all classes are optimized based on the metadata:

$$\Sigma^{(t+1)} = \Sigma^{(t)} - \eta_2 \frac{1}{m} \sum_{i=1}^m \nabla_{\Sigma} \ell_{CE}(\tilde{f}_{\tilde{\mathbf{W}}(\Omega^{(t)})}(\mathbf{x}_i^{meta}), y_i^{meta})|_{\Sigma^{(t)}}, \quad (7)$$

$$\mu^{(t+1)} = \mu^{(t)} - \eta_2 \frac{1}{m} \sum_{i=1}^m \nabla_{\mu} \ell_{CE}(\tilde{f}_{\tilde{\mathbf{W}}(\Omega^{(t)})}(\mathbf{x}_i^{meta}), y_i^{meta})|_{\mu^{(t)}}. \quad (8)$$

$\Sigma^{(t)}$ and $\mu^{(t)}$ refer to the covariance matrices and feature means of all classes in step t , respectively. Finally, the parameters of the classifier network can be updated with the obtained augmentation strengths $\alpha_i^{(t+1)}$:

$$\tilde{\mathbf{W}}^{(t+1)} = \tilde{\mathbf{W}}^{(t)} - \eta_1 \frac{1}{n} \sum_{i=1}^n \nabla_{\tilde{\mathbf{W}}} \ell_{ICDA}(\tilde{f}(\mathbf{x}_i), y_i; \alpha_i^{(t+1)})|_{\tilde{\mathbf{W}}^{(t)}}. \quad (9)$$

The steps of Meta-ICDA are presented in Algorithm 1.

V. EXPLANATION IN REGULARIZATION VIEW

This section conducts a deeper analysis considering regularization and reveals the ICDA's superiority against three advanced approaches: LA, ISDA, and RISDA. To our knowledge, this is the first time regularization has been used to explain these methods.

Using the first-order Taylor expansion of the loss, we have

$$\ell(\mathbf{u} + \Delta \mathbf{u}) \approx \ell(\mathbf{u}) + (\frac{\partial \ell}{\partial \mathbf{u}})^T \Delta \mathbf{u} = \ell(\mathbf{u}) + (\mathbf{q} - \mathbf{y})^T \Delta \mathbf{u}, \quad (10)$$

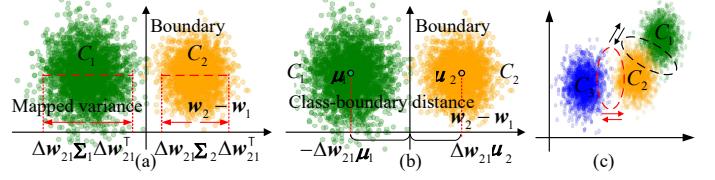


Fig. 4. Illustrations for the mapped variance (a) and class-boundary distance (b). (c) demonstrates that different samples in the same class should also have distinct augmentation directions.

where $\mathbf{q} = \text{softmax}(\mathbf{u})$ and \mathbf{y} is the one-hot label. Considering $R = (\mathbf{q} - \mathbf{y})^T \Delta \mathbf{u}$, the underlying regularizers of all approaches can be derived. The deviation process is presented in Appendix D. The regularizers and the factors affecting the generalization capability are summarized in Table I.

R_{LA} imposes greater punishment on the predictions $q_{i,c}$ ($c \neq y_i$) with a large $\delta_{c,i}$, improving their classification performance. Obviously, tail classes benefit more from LA. Thus, LA is prevalent in handling class imbalance.

R_{ISDA} contains $\Delta \mathbf{w}_{c,y_i} \Sigma_{y_i} \Delta \mathbf{w}_{c,y_i}^T$, and we prove that it is the mapped variance from samples of class y_i to the normal vector of the boundary between classes y_i and c (Fig. 4(a)).

Proof. If feature \mathbf{h} is on the boundary, we have

$$\mathbf{w}_c^T \mathbf{h} + b_c = \mathbf{w}_{y_i}^T \mathbf{h} + b_{y_i}. \quad (11)$$

Then, we know that the boundary between classes c and y_i is

$$\Delta \mathbf{w}_{c,y_i} \mathbf{h} + \Delta b_{c,y_i} = 0, \quad (12)$$

and $\Delta \mathbf{w}_{c,y_i} = \mathbf{w}_c^T - \mathbf{w}_{y_i}^T$ refers to the normal direction of the boundary between classes y_i and c . Thus, the value of the mapping $\Pi(\mathbf{h})$ of feature \mathbf{h} in class y to $\Delta \mathbf{w}_{c,y_i}$ is

$$\Pi(\mathbf{h}) = \Delta \mathbf{w}_{c,y_i} \mathbf{h} + \Delta b_{c,y_i}. \quad (13)$$

The (expected) variance of $\Pi(\mathbf{h})$ for $y = y_i$ denoted by Λ_{y_i} is as follows:

$$\begin{aligned} \Lambda_{c,y_i} &= E_{\Pi(\mathbf{h}):y=y_i} (\Pi(\mathbf{h}) - \bar{\Pi}(\mathbf{h})) (\Pi(\mathbf{h}) - \bar{\Pi}(\mathbf{h}))^T \\ &= E_{\mathbf{h}:y=y_i} [\Delta \mathbf{w}_{c,y_i} (\mathbf{h} - \bar{\mathbf{h}})(\mathbf{h} - \bar{\mathbf{h}})^T \Delta \mathbf{w}_{c,y_i}^T] \\ &= \Delta \mathbf{w}_{c,y_i} E_{\mathbf{h}:y=y_i} [(\mathbf{h} - \bar{\mathbf{h}})(\mathbf{h} - \bar{\mathbf{h}})^T] \Delta \mathbf{w}_{c,y_i}^T \\ &= \Delta \mathbf{w}_{c,y_i} \Sigma_{y_i} \Delta \mathbf{w}_{c,y_i}^T, \end{aligned} \quad (14)$$

where $\bar{\Pi}(\mathbf{h}) = E_{\Pi(\mathbf{h}):y=y_i} [\Pi(\mathbf{h})]$ and $\bar{\mathbf{h}} = E_{\mathbf{h}:y=y_i} [\mathbf{h}]$. \square

This term will force the model to decrease the mapped variances of class y_i towards the normal vectors of boundaries related to y_i , and thus increase intra-class compactness. This explains why ISDA performs well on standard datasets.

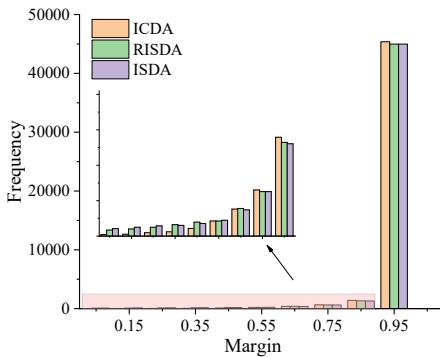


Fig. 5. Distributions of the margin values of models trained with the ISDA, RISDA, and ICDA losses.

However, as head classes have large training sizes, the punishment for their $q_{i,c}$ s and intra-class compactness is large. Xu et al. [19] revealed that classes with lower compactness, indicated by large variances, are more challenging and exhibit poorer performance. Consequently, ISDA further impairs the performance of tail classes and enlarges the performance gap between head and tail classes, as the variances of head classes decrease more than those of tail classes, which is undesirable in long-tailed classification.

The term $\Delta\mathbf{w}_{c,y_i}(\Sigma_{y_i} + \sum_{j=1, j \neq y_i}^C \varepsilon_{y_i,j} \Sigma_j) \Delta\mathbf{w}_{c,y_i}^T$ in R_{RISDA} is considered the mapped variances of more classes to the normal vector of the boundary between classes y_i and c . Therefore, it effectively decreases the intra-class compactnesses of more classes along each boundary and not just the ground-truth class. The term $\Delta\mathbf{w}_{c,y_i} \sum_{j=1, j \neq y_i}^C \varepsilon_{y_i,j} \mu_j$ can actually be divided into two parts: $\varepsilon_{y_i,c} \Delta\mathbf{w}_{c,y_i} \mu_c$ and $\varepsilon_{y_i,c'} \Delta\mathbf{w}_{c,y_i} \mu_{c'}$. Then, we prove that $\Delta\mathbf{w}_{c,y_i} \mu_c$ refers to the class-boundary distance between classes y_i and c , as illustrated in Fig. 4(b).

Proof. The boundary surface between classes y_i and c is

$$\Delta\mathbf{w}_{c,y_i} \mathbf{h} + \Delta b_{c,y_i} = 0. \quad (15)$$

Then, the distance from μ_c to the boundary is

$$d = \frac{|\Delta\mathbf{w}_{c,y_i} \mu_c + \Delta b_{c,y_i}|}{\|\Delta\mathbf{w}_{c,y_i}\|}. \quad (16)$$

As the feature mean μ_c must be classified correctly, we have $\Delta\mathbf{w}_{c,y_i} \mu_c + \Delta b_{c,y_i} > 0$. The bias term $\Delta b_{c,y_i}$ can be omitted. Thus, we have, when $\|\Delta\mathbf{w}_{c,y_i}\| = 1$, $\Delta\mathbf{w}_{c,y_i} \mu_c$ reflects the distance from μ_c to the boundary between classes y_i and c . Then, we explain why the term $-\Delta\mathbf{w}_{2,1} \mu_1$ in Fig. 4(b) is negative. As the feature mean μ_1 must be classified correctly, $\Delta\mathbf{w}_{1,2} \mu_1 > 0$ and $\Delta\mathbf{w}_{2,1} \mu_1 < 0$. Therefore, the distance between μ_1 and the boundary between classes C_1 and C_2 is the negative of $\Delta\mathbf{w}_{2,1} \mu_1$. \square

The regularization of $\varepsilon_{y_i,c} \Delta\mathbf{w}_{c,y_i} \mu_c$ will then force the boundary to move closer to μ_c and thus increase the class-boundary distance for class y_i , benefiting y_i . However, regularizing the second part seems unreasonable as $\mu_{c'}$ is supposed to have no bias towards both classes y_i and c . Ideally, this term keeps close to zero rather than having a negative value. Thus,

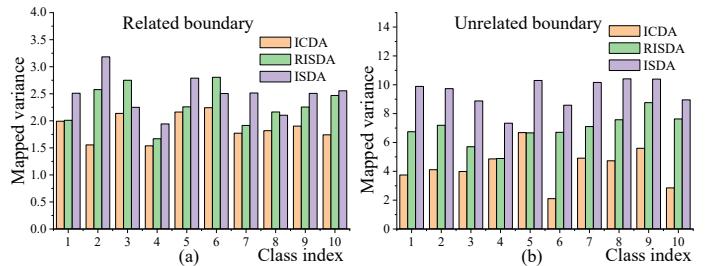


Fig. 6. The variation of mapped variances towards normal vectors of boundaries related to (a) and unrelated to (b) the ground-truth class on standard CIFAR10 using ResNet-32.

we removed this term from the derived ICDA loss, as stated in Section III.B.

Compared with other methods, R_{ICDA} can force models to simultaneously increase and decrease q_{i,y_i} and $q_{i,c}$, respectively. Thus, sample margins, especially those of hard ones will be enlarged because the harder the sample, the larger the α_i . Fig. 5 depicts the margin distributions of ISDA, RISDA, and ICDA, demonstrating that ICDA has fewer samples predicted correctly with small margins compared to the other two methods. In addition, the term $\Delta\mathbf{w}_{c,y_i} \mu_c$ in R_{ICDA} increases the class-boundary distance for class y_i . Like LA, the term $q_{i,c} \delta_{c,i}$ more increases the class-wise margins of the tail classes, manifesting that ICDA can deal well with imbalanced classification. Furthermore, $\Delta\mathbf{w}_{c,y_i} (\Sigma_{y_i} + \sum_{j=1, j \neq y_i}^C \hat{\alpha}_{i,j} \Sigma_j) \Delta\mathbf{w}_{c,y_i}^T$ is the mapped variance of all relevant classes to the normal vector of the boundary between classes y_i and c . Since this term is sample-wise, our punishment on the mapped variances is more refined and accurate than the class-wise approaches, enforcing better intra-class compactness. Fig. 4(c) reveals that although C_2 and C_1 are more confusing, the samples in the red circle are the most correlated with C_3 and cannot be taken seriously by the class-wise approaches. From Fig. 6, ICDA decreases the mapped variances on not only the boundaries related to the ground-truth class but also the unrelated ones to a higher degree. The β and λ parameters in R_{ICDA} can control the effect of each component.

VI. EXPERIMENTS

We empirically validate ICDA on several typical learning scenarios that require model generalization and robustness (i.e., biased datasets including both imbalanced and noisy data, subpopulation shifts datasets, generalized long-tailed datasets, and standard datasets) regarding performance and efficiency. Both image and text datasets are evaluated. For a fair comparison, Meta-ICDA is only compared when the competitor method utilizes meta-learning. We also visualize the augmented samples in the original input space and the attention of the trained model on several images. Finally, we conduct ablation studies and sensitivity tests. Regarding the hyperparameter settings in ICDA, λ^0 is selected in $\{0.1, 0.25, 0.5, 0.75, 1\}$, and β is set to 0.1 in all subsections.

A. Experiments on Biased Datasets

Experiments on Long-tailed CIFAR Datasets

TABLE II

TOP-1 ACCURACY ON LONG-TAILED CIFAR DATASETS. BOLD AND UNDERLINED NUMBERS ARE THE BEST AND SECOND-BEST RESULTS.

Dataset	CIFAR10		CIFAR100	
	100:1	10:1	100:1	10:1
Class-balanced CE [21]	72.68%	86.90%	38.77%	57.57%
Class-balanced Focal [21]	74.57%	87.48%	39.60%	57.99%
LDAM [15]	73.55%	87.32%	40.60%	57.29%
LDAM-DRW [15]	78.12%	88.37%	42.89%	58.78%
LA [16]	77.67%	88.93%	43.89%	58.34%
LPL [24]	77.95%	89.41%	44.25%	60.97%
ALA [23]	77.65%	88.32%	43.67%	58.92%
De-confound-TDE [26]	<u>80.60%</u>	88.50%	44.10%	59.60%
MixUp [25]	73.10%	87.10%	39.50%	58.00%
ISDA [8]	72.55%	87.02%	37.40%	55.51%
RISDA [10]	79.89%	89.36%	50.16%	62.38%
ICDA (Ours)	<u>81.69%</u>	<u>90.62%</u>	<u>50.18%</u>	<u>63.45%</u>
Meta-Weight-Net [27]	73.57%	87.55%	41.61%	58.91%
MetaSAug [9]	80.54%	89.44%	46.87%	61.73%
Meta-ICDA (Ours)	<u>82.47%</u>	<u>91.13%</u>	<u>50.96%</u>	<u>63.97%</u>

Settings. Long-tailed CIFAR is the long-tailed version of the CIFAR [20] data. The original CIFAR10 (CIFAR100) dataset consists of 50,000 images for 10 (100) classes with a balanced class distribution. Following Cui et al. [21], we discard some training samples to construct imbalanced datasets. Two training sets with imbalance ratios of 100:1 and 10:1 are built. We train ResNet-32 [22] with an initial learning rate of 0.1 and the standard SGD with the momentum of 0.9 and a weight decay of 5×10^{-4} . The learning rate is decayed by 0.1 at the 120-th and 160-th epochs. As for the meta-learning-based algorithms, the initial learning rate is 0.1 and it is decayed by 0.01 at the 160-th and 180-th epochs following MetaSAug [9]. We randomly select ten images per class from training data to construct metadata.

Several classical and advanced robust loss functions and data augmentation approaches that are mainly designed for long-tailed classifications are compared, including Class-balanced CE loss [21], Class-balanced Focal loss, LDAM [15], LDAM-DRW [15], ISDA [8], LA [16], ALA [23], LPL [24], MixUp [25], and RISDA [10]. Besides, De-confound-TDE [26], which uses causal intervention in training and counterfactual reasoning in inference, is also involved in our comparison. Two meta-learning-based methods including Meta-Weight-Net [27] and MetaSAug [9] are also compared.

Results. Table II reports the results on long-tailed CIFAR data, which are divided into two groups according to the usage of

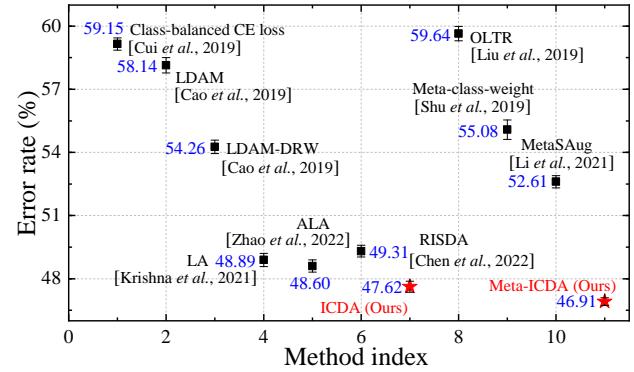


Fig. 7. Top-1 error rate on ImageNet-LT.

meta-learning. The results reveal that ICDA significantly outperforms other reweighting, solely logit adjustment, and implicit semantic augmentation methods, demonstrating that our sample-wise counterfactual augmentation strategy deals well with long-tailed classification. Although ICDA and RISDA achieve comparable performance on CIFAR100 with an imbalance ratio of 100:1, ICDA outperforms RISDA in other cases. Additionally, ICDA consistently surpasses De-confound-TDE, which uses causal intervention in training and counterfactual reasoning in inference. Our Meta-ICDA achieves state-of-the-art performance compared to all approaches.

To evaluate efficiency, we record the additional training time for ICDA and compare it against CE and ISDA. Table III reports the additional training time introduced by ICDA loss compared with CE loss on various backbones. The additional time introduced by ISDA loss compared with CE loss can be seen in the ISDA paper [8]. The results reveal that only a little time is increased by ICDA loss, and the values of training time for ICDA and ISDA are nearly equivalent.

Experiments on Long-tailed ImageNet Dataset

Settings. ImageNet [36] is a benchmark visual recognition dataset, which contains 1,281,167 training images and 50,000 validation images. Liu et al. [37] built the long-tailed version of ImageNet, which is denoted as ImageNet-LT. After discarding some training samples, ImageNet-LT remains 115,846 training examples in 1,000 classes. The imbalance ratio of ImageNet-LT is 256:1. Following MetaSAug [9], we adopt the original validation set to test methods. Ten images per class which are selected from the balanced validation set compiled by Liu et al. [37] are utilized to construct our metadata. ResNet-50 [22] is used as the backbone network. The learning rate is decayed by 0.1 at the 60-th and the 80-th epochs. The batch size is set to 64. Only the last fully connected layer is finetuned for training efficiency.

Methods designed for long-tailed classification including Class-balanced CE loss [21], OLTR [37], LDAM [15], LDAM-DRW [15], LA [16], ALA [23], RISDA [10], Meta-class-weight [27], and MetaSAug [9] are compared.

Results. Fig. 7 highlights that ICDA achieves good performance among the robust losses. Meta-ICDA significantly outperforms all competitor methods, including the meta semantic augmentation approach, proving that our proposed approach is more effective on long-tailed data.

TABLE III

ADDITIONAL TRAINING TIME INCREASED BY ICDA LOSS COMPARED WITH CE LOSS.

Networks	Params	Additional cost	
		CIFAR10	CIFAR100
ResNet-32	0.5M	6.9%	6.9%
ResNet-56	0.9M	7.1%	7.2%
ResNet-110	1.7M	6.8%	6.7%
DenseNet-BC-121	8M	5.6%	5.3%
DenseNet-BC-265	33.3M	5.3%	5.1%
Wide ResNet-16-8	11.0M	6.8%	7.0%
Wide ResNet-28-10	36.5M	6.7%	6.8%

TABLE IV
TOP-1 ACCURACY ON CIFAR DATASETS WITH UNIFORM AND FLIP NOISES.

Dataset	CIFAR10		CIFAR100		CIFAR10		CIFAR100	
	Flip				Uniform			
Noise type	20%	40%	20%	40%	40%	60%	40%	60%
Noise ratio								
CE loss	76.83%	70.77%	50.86%	43.01%	68.07%	53.12%	51.11%	30.92%
<i>LDMI</i> [28]	86.70%	84.00%	62.26%	57.23%	85.90%	79.60%	63.16%	55.37%
JoCoR [29]	90.78%	83.67%	65.21%	45.44%	89.15%	64.54%	65.45%	44.43%
D2L [30]	87.66%	83.89%	63.48%	51.83%	85.60%	68.02%	52.10%	41.11%
Co-teaching [31]	82.83%	75.41%	54.13%	44.85%	74.81%	73.06%	46.20%	35.67%
APL [32]	87.23%	80.08%	59.37%	52.98%	86.49%	79.22%	57.84%	49.13%
ISDA [8]	88.90%	86.14%	64.36%	59.48%	88.11%	83.12%	65.15%	58.19%
RISDA [10]	85.48%	81.12%	61.81%	54.60%	83.25%	76.31%	54.09%	45.57%
ICDA (Ours)	91.81%	88.76%	66.85%	61.57%	90.23%	84.91%	67.24%	60.26%
GLC [33]	89.68%	88.92%	63.07%	62.22%	88.28%	83.49%	61.31%	50.81%
MentorNet [34]	86.36%	81.76%	61.97%	52.66%	87.33%	82.80%	61.39%	36.87%
L2RW [35]	87.86%	85.66%	57.47%	50.98%	86.92%	82.24%	60.79%	48.15%
Meta-Weight-Net [27]	90.33%	87.54%	64.22%	58.64%	89.27%	84.07%	67.73%	58.75%
MetaSAug [9]	90.42%	87.73%	66.47%	61.43%	89.32%	84.65%	66.50%	59.84%
Meta-ICDA (Ours)	92.46%	90.21%	67.54%	63.26%	91.14%	85.86%	68.92%	61.80%

Experiments on Noisy CIFAR Datasets

Settings. Following Shu et al [27], two settings of corrupted labels are adopted, namely, uniform and pair-flip noise labels; 1,000 images with clean labels in the validation set are selected as the metadata. Wide ResNet-28-10 (WRN-28-10) [38] and ResNet-32 [22] are adopted as the classifiers for the uniform and pair-flip noises, respectively. The initial learning rate and batch size are set to 0.1 and 128, respectively. For ResNet, standard SGD with the momentum of 0.9 and a weight decay of 1×10^{-4} is utilized. For Wide ResNet, standard SGD with the momentum of 0.9 and a weight decay of 5×10^{-4} is utilized. As for the meta-learning-based algorithms, the initial learning rate is 0.1 and it is decayed by 0.01 at the 160-th and 180-th epochs following MetaSAug [9].

Several robust loss functions including Information-theoretic Loss (*LDMI*) [28], JoCoR [29], Co-teaching [31], D2L [30], and APL [32] are compared. The meta-learning-based methods, including MentorNet [34], L2RW [35], GLC [33], and Meta-Weight-Net [27] are also involved in comparison. We also compared our proposed ICDA with three implicit data augmentation methods, including ISDA [8], RISDA [10], and MetaSAug [9].

Results. Table IV reports the results of CIFAR data with flip and uniform noise, respectively. ICDA notably surpasses all competitor approaches including robust loss functions and the

class-level implicit data augmentation approaches. Besides, Meta-ICDA achieves state-of-the-art performance compared with other meta-learning-based manners, manifesting that our proposed method can effectively improve the generalization and robustness of models on noisy data.

B. Experiments on Subpopulation Shifts Datasets

Settings. Four subpopulation shifts datasets are evaluated, including CMNIST, Waterbirds [44], CelebA [12], and CivilComments [52], in which the domain information is highly spuriously correlated with the labels. Detailed descriptions of the datasets are shown in Appendix E. In the subsequent trials, ResNet-50 is utilized as the backbone network for the first three image datasets, while DistilBert [53] is adopted for the text set CivilComments. The initial learning rates for CMNIST and Waterbirds are 1×10^{-3} , while those for CelebA and CivilComments are 1×10^{-4} and 1×10^{-5} , respectively. The values of weight decay are 1×10^{-4} for CMNIST, Waterbirds, and CelebA, and 0 for CivilComments. The values of batch size for CMNIST, Waterbirds, and CelebA are 16, and that for CivilComments is 8. For the three image classification datasets, SGD optimizer is utilized, while Adam is utilized for CivilComments.

Robust methods, including IRM [40], IB-IRM [41], V-REx [42], CORAL [43], GroupDRO [44], DomainMix [45],

TABLE V
AVERAGE AND WORST-GROUP ACCURACY ON SUBPOPULATION SHIFTS DATASETS.

Dataset	CelebA		CMNIST		Waterbirds		CivilComments	
	Avg.	Worst	Avg.	Worst	Avg.	Worst	Avg.	Worst
Method								
UW [39]	92.9%	83.3%	72.2%	66.0%	95.1%	88.0%	89.8%	69.2%
IRM [40]	94.0%	77.8%	72.1%	70.3%	87.5%	75.6%	88.8%	66.3%
IB-IRM [41]	93.6%	85.0%	72.2%	70.7%	88.5%	76.5%	89.1%	65.3%
V-REx [42]	92.2%	86.7%	71.7%	70.2%	88.0%	73.6%	90.2%	64.9%
CORAL [43]	93.8%	76.9%	71.8%	69.5%	90.3%	79.8%	88.7%	65.6%
GroupDRO [44]	92.1%	87.2%	72.3%	68.6%	91.8%	90.6%	89.9%	70.0%
DomainMix [45]	93.4%	65.6%	51.4%	48.0%	76.4%	53.0%	90.9%	63.6%
Fish [46]	93.1%	61.2%	46.9%	35.6%	85.6%	64.0%	89.8%	71.1%
LISA [39]	92.4%	89.3%	74.0%	73.3%	91.8%	89.2%	89.2%	72.6%
ICDA (Ours)	93.3%	90.7%	76.1%	75.3%	92.9%	90.7%	91.1%	73.5%

TABLE VI
TOP-1 ACCURACY AND PRECISION OF CLT, GLT, AND ALT PROTOCOLS ON IMAGENET-GLT.

Protocol	CLT		GLT		ALT	
Method	Acc.	Prec.	Acc.	Prec.	Acc.	Prec.
CE loss	42.52%	47.92%	34.75%	40.65%	41.73%	41.74%
cRT [47]	45.92%	45.34%	37.57%	37.51%	41.59%	41.43%
LWS [47]	46.43%	45.90%	37.94%	38.01%	41.70%	41.71%
De-confound-TDE [26]	45.70%	44.48%	37.56%	37.00%	41.40%	42.36%
BLSsoftmax [48]	45.79%	46.27%	37.09%	38.08%	41.32%	41.37%
LA [16]	46.53%	45.56%	37.80%	37.56%	41.73%	41.74%
BBN [49]	46.46%	49.86%	37.91%	41.77%	43.26%	43.86%
LDAM [15]	46.74%	46.86%	38.54%	39.08%	42.66%	41.80%
IFL [50]	45.97%	52.06%	37.96%	44.47%	45.89%	46.42%
MixUp [25]	38.81%	45.41%	31.55%	37.44%	42.11%	42.42%
RandAug [51]	46.40%	52.13%	38.24%	44.74%	46.29%	46.32%
ISDA [8]	47.66%	51.98%	39.44%	44.26%	47.62%	47.46%
RISDA [10]	49.31%	50.64%	38.45%	42.77%	47.33%	46.33%
ICDA (Ours)	52.11%	55.05%	42.73%	47.49%	50.52%	49.68%
MetaSAug [9]	50.53%	55.21%	41.27%	47.38%	49.12%	48.56%
Meta-ICDA (Ours)	52.76%	56.71%	44.15%	49.32%	51.74%	51.43%

Fish [46], and LISA [39], are compared. Upweighting (UW) is suitable for subpopulation shifts, so we also use it for comparison. We only compare ICDA with other methods for fair comparisons as all these approaches do not rely on meta-learning. Following Yao et al. [39], the worst-group accuracy is used to compare the performance of all methods.

Results. Table V reports the results of the four subpopulation shifts datasets. The performance of methods that learn invariant predictors with explicit regularizers, e.g., IRM, IB-IRM, and V-REx, is not consistent across datasets. For example, V-REx outperforms IRM on CelebA, but it fails to achieve better performance than IRM on CMNIST, Waterbirds, and CivilComments. Opposing, ICDA consistently achieves an appealing performance on all datasets, demonstrating ICDA’s effectiveness in breaking spurious correlations and achieving invariant feature learning. Although ICDA and GroupDRO achieve similar performance on Waterbirds, ICDA far exceeds GroupDRO on the other three datasets.

C. Experiments on Generalized Long-tailed Datasets

Settings. Tang et al. [50] proposed a novel learning problem, namely, generalized long-tailed classification, in which two

new benchmarks, including MSCOCO-GLT and ImageNet-GLT, were proposed. Each benchmark has three protocols, i.e., CLT, ALT, and GLT, in which class distribution, attribute distribution, and both class and attribute distributions are changed from training to testing, respectively. More details of the two benchmarks can be seen in [50]. The training and testing configurations follow those in the IFL [50] paper. ResNeXt-50 [54] is used as the backbone network for all methods except for BBN [49]. Both Top-1 accuracy and precision are presented. All models are trained with a batch size of 256 and an initial learning rate of 0.1. SGD optimizer is utilized with a weight decay of 5×10^{-4} and the momentum of 0.9. Here, Meta-ICDA is exclusively evaluated on the ImageNet-GLT benchmark. To collect the attribute-wise balanced metadata, images from each class in a balanced validation set compiled by Liu et al. [37] are clustered into 6 groups by KMeans using a pre-trained ResNet-50 model. From each group and class, 10 images are sampled to construct the metadata.

As for the compared methods, we studied two-stage resampling methods, including cRT [47] and LWS [47], posthoc distribution adjustment methods including De-confound-TDE [26] and LA [16], multi-branch models with diverse

TABLE VII
TOP-1 ACCURACY AND PRECISION OF CLT, GLT, AND ALT PROTOCOLS ON MSCOCO-GLT.

Protocol	CLT		GLT		ALT	
Method	Acc.	Prec.	Acc.	Prec.	Acc.	Prec.
CE loss	72.34%	76.61%	63.79%	70.52%	50.17%	50.94%
cRT [47]	73.64%	75.84%	64.69%	68.33%	49.97%	50.37%
LWS [47]	72.60%	75.66%	63.60%	68.81%	50.14%	50.61%
De-confound-TDE [26]	73.79%	74.90%	66.07%	68.20%	50.76%	51.68%
BLSsoftmax [48]	72.64%	75.25%	64.07%	68.59%	49.72%	50.65%
LA [16]	75.50%	76.88%	66.17%	68.35%	50.17%	50.94%
BBN [49]	73.69%	77.35%	64.48%	70.20%	51.83%	51.77%
LDAM [15]	75.57%	77.70%	67.26%	70.70%	55.52%	56.21%
IFL [50]	74.31%	78.90%	65.31%	72.24%	52.86%	53.49%
MixUp [25]	74.22%	78.61%	64.45%	71.13%	48.90%	49.53%
RandAug [51]	76.81%	79.88%	67.71%	72.73%	53.69%	54.71%
ISDA [8]	77.32%	79.23%	67.57%	72.89%	54.43%	54.62%
RISDA [10]	76.34%	79.27%	66.85%	72.66%	54.58%	53.98%
ICDA (Ours)	78.82%	81.33%	68.78%	74.29%	56.48%	57.81%

TABLE VIII
TOP-1 ERROR RATE ON STANDARD CIFAR DATASETS.

Backbone	ResNet-110		WRN-28-10	
	CIFAR10	CIFAR100	CIFAR10	CIFAR100
Large Margin [55]	6.46%	28.00%	3.69%	18.48%
Disturb Label [56]	6.61%	28.46%	3.91%	18.56%
Focal loss [57]	6.68%	28.28%	3.62%	18.22%
Center loss [58]	6.38%	27.85%	3.76%	18.50%
Lq loss [59]	6.69%	28.78%	3.78%	18.43%
WGAN [60]	6.63%	-	3.81%	-
CGAN [61]	6.56%	28.25%	3.84%	18.79%
ACGAN [62]	6.32%	28.48%	3.81%	18.54%
infoGAN [63]	6.59%	27.64%	3.81%	18.44%
ISDA [8]	<u>5.98%</u>	<u>26.35%</u>	<u>3.58%</u>	<u>17.98%</u>
RISDA [10]	6.47%	28.42%	3.79%	18.46%
ICDA (Ours)	4.89%	25.21%	3.01%	17.03%

sampling strategies like BBN [49], invariant feature learning methods like IFL [50], and reweighting loss functions like BLSoftmax [48] and LDAM [15]. We also compare some data augmentation methods, including MixUp [25], RandAug [51], ISDA [8], RISDA [10], and MetaSAug [9].

Results. Tables VI and VII report the results of the three protocols for ImageNet-GLT and MSCOCO-GLT, respectively, some of which are from the IFL [50] paper. ICDA notably improves model performance in all three protocols, demonstrating that it can well break the spurious associations caused by imbalanced attribute and class distributions, while the majority of previous LT algorithms using rebalancing strategies fail to improve the robustness against the attribute-wise bias. Additionally, we found that augmentation methods generally perform better than other long-tailed transfer learning approaches on GLT protocols.

D. Experiments on Standard CIFAR Datasets

Settings. To verify that ICDA has a good augmentation effect, it is compared with a number of advanced methods ranging from robust loss functions (i.e., Large Margin [55], Disturb Label [56], Focal loss [57], Center loss [58], and Lq loss [59]) to explicit (i.e., WGAN [60], CGAN [61], ACGAN [62], and infoGAN [63]) and implicit (i.e., ISDA [8] and RISDA [10]) augmentation methods on standard CIFAR data. ResNet-110 and WRN-28-10 models are utilized. Regarding the hyperparameter settings, the initial learning rate and the batch size are set to 0.1 and 128, respectively. For ResNet, standard SGD with the momentum of 0.9 and a weight decay of 1×10^{-4} is utilized. For Wide ResNet, standard SGD with the momentum

of 0.9 and a weight decay of 5×10^{-4} is utilized. The learning rate is decayed by 0.1 at the 120-th and 160-th epochs.

Results. The results are reported in Table VIII. ICDA achieves the best performance compared with other explicit and implicit augmentation approaches. Moreover, GAN-based methods perform poorly on CIFAR100 due to a limited training size. Additionally, these methods impose excess calculations and decrease training efficiency. Although ISDA affords lower error and is more efficient than GAN-based schemes, it can not surpass ICDA as ICDA assists the models in breaking spurious correlations of models.

E. Visualization Results

Following ISDA’s visualization manner, we map the augmented features back into the pixel space. The corresponding results are presented in Fig. 8, highlighting that ICDA can generate more diverse and meaningful counterfactual images and notably alter the non-intrinsic attributes, e.g., scene contexts and viewpoints, compared with ISDA and RISDA.

Additionally, Grad-CAM [64] is utilized to visualize the regions that models used for making predictions. Fig. 9(a) manifests that ISDA focuses on the background or other nuisances for false predictions, while ICDA focuses tightly on the causal regions corresponding to the object, assisting models in making correct classifications. For example, for the image of “Brittany spaniel”, ISDA utilizes spurious context “Water”, so its prediction is “Drake”, while the model trained with ICDA attends more to the dog, contributing to a correct prediction. Therefore, in addition to performance gains, the ICDA predictions are made for the right reasons. Fig. 9(b) presents some images that are wrongly classified by both ISDA and ICDA. Although the two methods make false predictions, ICDA still helps the model to concentrate more on the causal attributes by breaking the spurious associations between the nuisances and classes. More visualization results are shown in Appendix F.

F. Ablation and Sensitivity Studies

To get a better understanding of the effect of varying components, we evaluate the following three settings of ICDA. Setting I: Without the covariance matrices and feature means of the other classes, i.e., $\alpha_{i,c} = 0$. Setting II: Without the class-level logit perturbation term, i.e., removing $\delta_{c,i}$. Setting III: Without the sample-level logit perturbation term, i.e., removing $\beta\alpha_i$. Since the proportion of each class is the



Fig. 8. Visualization of images augmented by ISDA, RISDA, and ICDA.

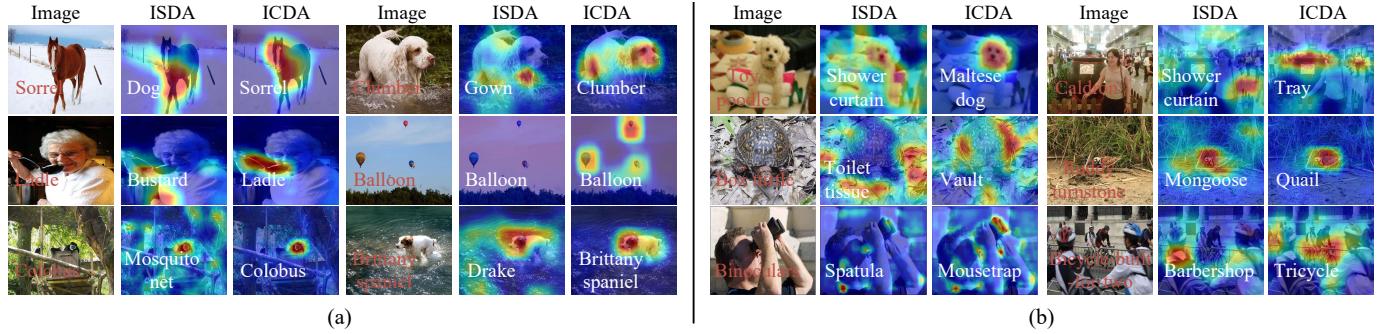


Fig. 9. (a) Visualization of the regions that the model used for making predictions. Blue and red imply that the region is indecisive and very discriminative, respectively. White texts are the predicted labels. Red texts are the ground-truth labels of the images. (b) More results of model visualization, in which both ISDA and ICDA make false predictions.

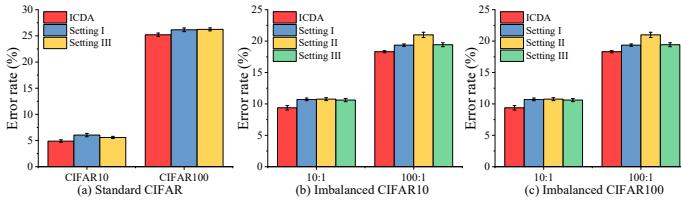


Fig. 10. Results of ablation studies on standard and imbalanced CIFAR data.

same on standard data, we only evaluate Settings I and III on the standard data. The ablation results are presented in Fig. 10, revealing that all three components are crucial and necessary for imbalanced data. Additionally, the statistics information of the other classes and the sample-level perturbation term are critical for standard data. Without each of them, the performance of ICDA will be weakened.

To study how the hyperparameters in ICDA (i.e., λ^0 and β) affect our method's performance, several sensitivity tests are conducted, where ResNet-110 is used as the backbone network. The corresponding results on standard CIFAR data are shown in Fig. 11, revealing that ICDA achieves superior performance for $0.01 \leq \beta \leq 0.5$ and $0.25 \leq \lambda^0 \leq 1$. When β and λ^0 are too large, the model is easier to overfit and underfit, respectively. Empirically, we recommend $\beta = 0.1$ and $\lambda^0 = 0.5$ for a naive implementation or a starting point of hyperparameter searching.

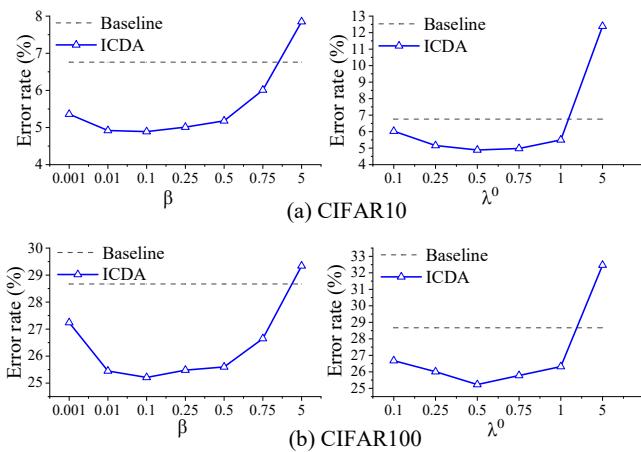


Fig. 11. Results of the sensitivity tests on standard CIFAR datasets.

VII. CONCLUSION

This study proposes a sample-wise implicit counterfactual data augmentation (ICDA) method to break spurious correlations and make stable predictions. Our method can be formulated as a novel robust loss, easily adopted by any classifier, and is considerably more efficient than explicit augmentation approaches. Two manners, including direct quantification and meta-learning, are introduced to learn the key parameters in the robust loss. Furthermore, the regularization analysis demonstrates that ICDA improves intra-class compactness, class and sample-wise margins, and class-boundary distances. Extensive experimental comparison and visualization results on several typical learning scenarios demonstrate the proposed method's effectiveness and efficiency.

REFERENCES

- [1] J. Kaddour, A. Lynch, Q. Liu, M. J. Kusner, and R. Silva, "Causal machine learning: A survey and open problems," *arXiv:2206.15475*, 2022.
- [2] C. Mao, A. Cha, A. Gupta, H. Wang, J. Yang, and C. Vondrick, "Generative interventions for causal learning," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2021, pp. 3946–3955.
- [3] K. Liu, F. Xue, D. Guo, P. Sun, S. Qian, and R. Hong, "Multimodal graph contrastive learning for multimedia-based recommendation," *IEEE Trans. Multimedia*, vol. 25, pp. 9343–9355, 2023.
- [4] B. Schölkopf, "Causality for machine learning," in *Proc. Int. Conf. Learn. Representations*, 2017.
- [5] K. Lu, P. Mardziel, F. Wu, P. Amancharla, and A. Datta, *Logic, language, and security*. Berlin, Germany: Springer, 2020.
- [6] J. He, M. Xia, C. Fellbaum, and D. Chen, "Mabel: Attenuating gender bias using textual entailment data," in *Proc. Conf. Empir. Methods Nat. Lang. Process.*, 2022, pp. 9681–9702.
- [7] C.-H. Chang, G. A. Adam, and A. Goldenberg, "Towards robust classification model by counterfactual and invariant data generation," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2021, pp. 15207–15216.
- [8] Y. Wang, G. Huang, S. Song, X. Pan, Y. Xia, and C. Wu, "Regularizing deep networks with semantic data augmentation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 44, no. 7, pp. 3733–3748, 2021.
- [9] S. Li, K. Gong, C.-H. Liu, Y. Wang, F. Qiao, and X. Cheng, "Metasaug: Meta semantic augmentation for long-tailed visual recognition," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2021, pp. 5208–5217.
- [10] X. Chen, Y. Zhou, D. Wu, W. Zhang, Y. Zhou, B. Li, and W. Wang, "Imagine by reasoning: A reasoning-based implicit semantic data augmentation for long-tailed classification," in *Proc. Assoc. Adv. Artif. Intell.*, 2022, pp. 356–364.
- [11] A. Sauer and A. Geiger, "Counterfactual generative networks," in *Proc. Int. Conf. Learn. Representations*, 2021.
- [12] Z. Liu, P. Luo, X. Wang, and X. Tang, "Deep learning face attributes in the wild," in *Proc. Int. Conf. Conf. Comput. Vis.*, 2016, pp. 3730–3738.

- [13] W. Liu, Y. Wen, Z. Yu, and M. Yang, "Large-margin softmax loss for convolutional neural networks," in *Proc. Int. Conf. Mach. Learn.*, 2016, pp. 507–516.
- [14] H. Wang, Y. Wang, Z. Zhou, X. Ji, D. Gong, J. Zhou, Z. Li, and W. Liu, "Costface: Large margin cosine loss for deep face recognition," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 5265–5274.
- [15] K. Cao, C. Wei, A. Gaidon, N. Arechiga, and T. Ma, "Learning imbalanced datasets with label-distribution-aware margin loss," in *Proc. Int. Conf. Neural Inf. Process. Syst.*, 2019, pp. 1567–1578.
- [16] A. Krishna, M. Sadeep, J. Ankit, S. Rawat, H. Jain, A. Veit, and S. Kumar, "Long-tail learning via logit adjustment," in *Proc. Int. Conf. Learn. Representations*, 2021.
- [17] Y. Hu, J. Gao, and C. Xu, "Learning multi-expert distribution calibration for long-tailed video classification," *IEEE Trans. Multimedia*, vol. 26, pp. 555–567, 2023.
- [18] S. Wu and X. Gong, "Boundaryface: A mining framework with noise label self-correction for face recognition," in *Proc. Eur. Conf. Comput. Vis.*, 2022, pp. 91–106.
- [19] H. Xu, X. Liu, Y. Li, A. K. Jain, and J. Tang, "To be robust or to be fair: Towards fairness in adversarial training," in *Proc. Int. Conf. Mach. Learn.*, 2021, pp. 11492–11501.
- [20] A. Krizhevsky, "Learning multiple layers of features from tiny images," *Tech. Rep.*, pp. 1–60, 2009.
- [21] Y. Cui, M. Jia, T.-Y. Lin, Y. Song, and S. Belongie, "Class-balanced loss based on effective number of samples," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 9260–9269.
- [22] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 770–778.
- [23] Y. Zhao, W. Chen, X. Tan, K. Huang, and J. Zhu, "Adaptive logit adjustment loss for long-tailed visual recognition," in *Proc. Assoc. Adv. Artif. Intell.*, 2022, pp. 3472–3480.
- [24] M. Li, F. Su, O. Wu, and J. Zhang, "Logit perturbation," in *Proc. Assoc. Adv. Artif. Intell.*, 2022, pp. 1359–1366.
- [25] H. Zhang, M. Cisse, Y. N. Dauphin *et al.*, "Mixup: Beyond empirical risk minimization," in *Proc. Int. Conf. Learn. Representations*, 2018.
- [26] K. Tang, J. Huang, and H. Zhang, "Long-tailed classification by keeping the good and removing the bad momentum causal effect," in *Proc. Int. Conf. Neural Inf. Process. Syst.*, 2020, pp. 1513–1524.
- [27] J. Shu, Q. Xie, L. Yi, Q. Zhao, S. Zhou, Z. Xu, and D. Meng, "Meta-weight-net: Learning an explicit mapping for sample weighting," in *Proc. Int. Conf. Neural Inf. Process. Syst.*, 2019, pp. 1919–1930.
- [28] Y. Xu, P. Cao, Y. Kong, and Y. Wang, "L_dmi: A novel information-theoretic loss function for training deep nets robust to label noise," in *Proc. Int. Conf. Neural Inf. Process. Syst.*, 2019, pp. 6225–6236.
- [29] H. Wei, L. Feng, X. Chen *et al.*, "Combating noisy labels by agreement: A joint training method with co-regularization," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 13723–13732.
- [30] X. Ma, Y. Wang, M. E. Houle, S. Zhou, S. Erfani, S. Xia, S. Wijewickrema, and J. Bailey, "Dimensionality-driven learning with noisy labels," in *Proc. Int. Conf. Mach. Learn.*, 2018, pp. 3355–3364.
- [31] B. Han, Q. Yao, X. Yu, G. Niu, M. Xu, W. Hu, I. W. Tsang, and M. Sugiyama, "Co-teaching: Robust training of deep neural networks with extremely noisy labels," in *Proc. Int. Conf. Neural Inf. Process. Syst.*, 2018, pp. 8536–8546.
- [32] X. Ma, H. Huang, Y. Wang, S. Romano, and J. B. Sarah Erfani, "Normalized loss functions for deep learning with noisy labels," in *Proc. Int. Conf. Mach. Learn.*, 2020, pp. 6543–6553.
- [33] D. Hendrycks, M. Mazeika, D. Wilson, and K. Gimpel, "Using trusted data to train deep networks on labels corrupted by severe noise," in *Proc. Int. Conf. Neural Inf. Process. Syst.*, 2018, pp. 10477–10486.
- [34] L. Jiang, Z. Zhou, T. Leung, L.-J. Li, and F.-F. Li, "Mentornet: Learning data-driven curriculum for very deep neural networks on corrupted labels," in *Proc. Int. Conf. Mach. Learn.*, 2018, pp. 2304–2313.
- [35] M. Ren, W. Zeng, B. Yang, and R. Urtasun, "Learning to reweight examples for robust deep learning," in *Proc. Int. Conf. Mach. Learn.*, 2018, pp. 4334–4343.
- [36] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy *et al.*, "Imagenet large scale visual recognition challenge," *Int. J. Comput. Vis.*, vol. 115, pp. 211–252, 2015.
- [37] Z. Liu, Z. Miao, X. Zhan, J. Wang, B. Gong, and S. X. Yu, "Large-scale long-tailed recognition in an open world," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 2532–2541.
- [38] S. Zagoruyko and N. Komodakis, "Wide residual networks," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2017.
- [39] H. Yao, Y. Wang, S. Li, L. Zhang, W. Liang, J. Zou, and C. Finn, "Improving out-of-distribution robustness via selective augmentation," in *Proc. Int. Conf. Mach. Learn.*, 2022, pp. 25407–25437.
- [40] M. Arjovsky, L. Bottou, I. Gulrajani, and D. Lopez-Paz, "Invariant risk minimization," *arXiv:1907.02893*, 2019.
- [41] K. Ahuja, E. Caballero, D. Zhang, Y. Bengio, I. Mitliagkas *et al.*, "Invariance principle meets information bottleneck for out-of-distribution generalization," in *Proc. Int. Conf. Neural Inf. Process. Syst.*, 2021.
- [42] D. Krueger, E. Caballero, J.-H. Jacobsen, A. Zhang, J. Binas, D. Zhang, R. L. Priol, and A. Courville, "Out-of-distribution generalization via risk extrapolation (rex)," in *Proc. Int. Conf. Mach. Learn.*, 2021, pp. 5815–5826.
- [43] B. Sun and K. Saenko, "Deep coral: Correlation alignment for deep domain adaptation," in *Proc. Eur. Conf. Comput. Vis.*, 2016, pp. 443–450.
- [44] S. Sagawa, P. W. Koh, T. B. Hashimoto, and P. Liang, "Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization," in *Proc. Int. Conf. Learn. Representations*, 2020.
- [45] M. Xu, J. Zhang, B. Ni, T. Li, C. Wang, Q. Tian, and W. Zhang, "Adversarial domain adaptation with domain mixup," in *Proc. Assoc. Adv. Artif. Intell.*, 2020, pp. 6502–6509.
- [46] Y. Shi, J. Seely, P. H. Torr, N. Siddharth, A. Hannun, N. Usunier, and G. Synnaeve, "Gradient matching for domain generalization," in *Proc. Int. Conf. Learn. Representations*, 2022.
- [47] B. Kang, S. Xie, M. Rohrbach, Z. Yan, A. Gordo, J. Feng, and etc, "Decoupling representation and classifier for long-tailed recognition," in *Proc. Int. Conf. Learn. Representations*, 2020.
- [48] J. Ren, C. Yu, shunan sheng, X. Ma, H. Zhao, S. Yi, and hongsheng Li, "Balanced meta-softmax for long-tailed visual recognition," in *Proc. Int. Conf. Neural Inf. Process. Syst.*, 2020, pp. 4175–4186.
- [49] B. Zhou, Q. Cui, X.-S. Wei *et al.*, "Bbn: Bilateral-branch network with cumulative learning for long-tailed visual recognition," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 9716–9725.
- [50] K. Tang, M. Tao, J. Qi, Z. Liu, and H. Zhang, "Invariant feature learning for generalized long-tailed classification," in *Proc. Eur. Conf. Comput. Vis.*, 2022, pp. 709–726.
- [51] E. D. Cubuk, B. Zoph, J. Shlens, and Q. V. Le, "Randaugment: Practical automated data augmentation with a reduced search space," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 3008–3017.
- [52] D. Borkan, L. Dixon, J. Sorensen, N. Thain *et al.*, "Nuanced metrics for measuring unintended bias with real data for text classification," in *Companion World Wide Web Conf.*, 2019, pp. 491–500.
- [53] V. Sanh, L. Debut, J. Chaumond *et al.*, "Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter," *arXiv:1910.01108*, 2019.
- [54] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, "Aggregated residual transformations for deep neural networks," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 5987–5995.
- [55] W. Liu, Y. Wen, Z. Yu, and M. Yang, "Large-margin softmax loss for convolutional neural networks," in *Proc. Int. Conf. Mach. Learn.*, 2016, pp. 507–516.
- [56] L. Xie, J. Wang, Z. Wei, M. Wang, and Q. Tian, "Disturblabel: Regularizing cnn on the loss layer," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 4753–4762.
- [57] T.-Y. Lin, P. Goyal, R. Girshick, K. He *et al.*, "Focal loss for dense object detection," in *Proc. Int. Conf. Conf. Comput. Vis.*, 2017, pp. 2999–3007.
- [58] Y. Wen, K. Zhang, Z. Li1, and Y. Qiao, "A discriminative feature learning approach for deep face recognition," in *Proc. Eur. Conf. Comput. Vis.*, 2016, pp. 499–515.
- [59] Z. Zhang and M. R. Sabuncu, "Generalized cross entropy loss for training deep neural networks with noisy labels," in *Proc. Int. Conf. Neural Inf. Process. Syst.*, 2018, pp. 8792–8802.
- [60] M. Arjovsky, S. Chintala, and L. Bottou, "Wasserstein gan," in *Proc. Int. Conf. Mach. Learn.*, 2017, pp. 214–223.
- [61] M. Mirza and S. Osindero, "Conditional generative adversarial nets," *arXiv:1411.1784*, 2014.
- [62] A. Odena, C. Olah, and J. Shlens, "Conditional image synthesis with auxiliary classifier gans," in *Proc. Int. Conf. Mach. Learn.*, 2017, pp. 2642–2651.
- [63] X. Chen, Y. Duan, R. Houthooft, J. Schulman, I. Sutskever, and P. Abbeel, "Infogan: Interpretable representation learning by information maximizing generative adversarial nets," in *Proc. Int. Conf. Neural Inf. Process. Syst.*, 2016, pp. 2180–2188.
- [64] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-cam: Visual explanations from deep networks via gradient-based localization," in *Proc. Int. Conf. Conf. Comput. Vis.*, 2016, pp. 618–626.