

Orthogonal Uncertainty Representation of Data Manifold for Robust Long-Tailed Learning

Yanbiao Ma
Xidian University
Xi'an, China
ybmamail@stu.xidian.edu.cn

Licheng Jiao*
Xidian University
Xi'an, China
lchjiao@mail.xidian.edu.cn

Fang Liu
Xidian University
Xi'an, China
f63liu@163.com

Shuyuan Yang
Xidian University
Xi'an, China
syyang@xidian.edu.cn

Xu Liu
Xidian University
Xi'an, China
xuliu361@163.com

Lingling Li
Xidian University
Xi'an, China
lli@xidian.edu.cn

ABSTRACT

In scenarios with long-tailed distributions, the model's ability to identify tail classes is limited due to the under-representation of tail samples. Class rebalancing, information augmentation, and other techniques have been proposed to facilitate models to learn the potential distribution of tail classes. The disadvantage is that these methods generally pursue models with balanced class accuracy on the data manifold, while ignoring the ability of the model to resist interference. By constructing noisy data manifold, we found that the robustness of models trained on unbalanced data has a long-tail phenomenon. That is, even if the class accuracy is balanced on the data domain, it still has bias on the noisy data manifold. However, existing methods cannot effectively mitigate the above phenomenon, which makes the model vulnerable in long-tailed scenarios. In this work, we propose an Orthogonal Uncertainty Representation (OUR) of feature embedding and an end-to-end training strategy to improve the long-tail phenomenon of model robustness. As a general enhancement tool, OUR has excellent compatibility with other methods and does not require additional data generation, ensuring fast and efficient training. Comprehensive evaluations on long-tailed datasets show that our method significantly improves the long-tail phenomenon of robustness, bringing consistent performance gains to other long-tailed learning methods.

CCS CONCEPTS

• **Computing methodologies** → **Computer vision**; *Image representations*; Supervised learning by classification.

KEYWORDS

Long-tailed distribution, Imbalanced Learning, Model bias

*Corresponding author(s).

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

MM '23, October 29–November 3, 2023, Ottawa, ON, Canada

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 979-8-4007-0108-5/23/10...\$15.00

<https://doi.org/10.1145/3581783.3611698>

ACM Reference Format:

Yanbiao Ma, Licheng Jiao, Fang Liu, Shuyuan Yang, Xu Liu, and Lingling Li. 2023. Orthogonal Uncertainty Representation of Data Manifold for Robust Long-Tailed Learning. In *Proceedings of the 31st ACM International Conference on Multimedia (MM '23)*, October 29–November 3, 2023, Ottawa, ON, Canada. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3581783.3611698>

1 INTRODUCTION

Long-tailed recognition is an important challenge in computer vision, manifested by models trained on long-tailed data that tend to perform poorly for classes with few samples. Previous research has attributed this phenomenon to the fact that the few samples in the tail classes do not well represent their true distribution, resulting in a shift between the test and training domains [4, 31, 40]. Numerous methods have been proposed to mitigate model bias. Class rebalancing [2, 3, 7, 14, 17, 18, 20, 23, 26, 28, 33, 34, 37, 39, 41, 42, 45], for example, aims to boost the weight of losses arising from tail classes, thereby pushing the decision boundary away from the tail class and improving the probability of correctly classifying the underlying distribution. Information augmentation [4, 8, 12, 15, 19, 21, 22, 24, 25, 32, 36, 38, 39, 43], on the other hand, expands the observed distribution of the tail classes by introducing prior knowledge to facilitate the model learning of the underlying distribution. It is important to note that these methods default to the model being able to learn adequately and fairly at least for the samples in the training domain. However, we find that even if a model has balanced class accuracy over the training domain, its robustness still exhibits a long-tailed distribution, and existing methods do not improve the phenomenon well.

Recent study [10] indicates that moving in the direction orthogonal to the data manifold produces a series of noisy data manifolds, and the samples on these noisy data manifolds are noisy versions of the real samples. We construct the noisy data manifold corresponding to the training samples (i.e. the data manifold) on sample-balanced MNIST and CIFAR-10. It is found that ResNet-18 trained on the data manifold can correctly recognize noisy samples from all classes with high confidence (Fig.1A, even images that are meaningless to the human eye) and that the class accuracy on the noisy data manifold is balanced (Fig.1B). The same experiments are then performed on the CIFAR-10-LT. Unexpectedly, we find that although ResNet-18 performs well and fairly for each class on

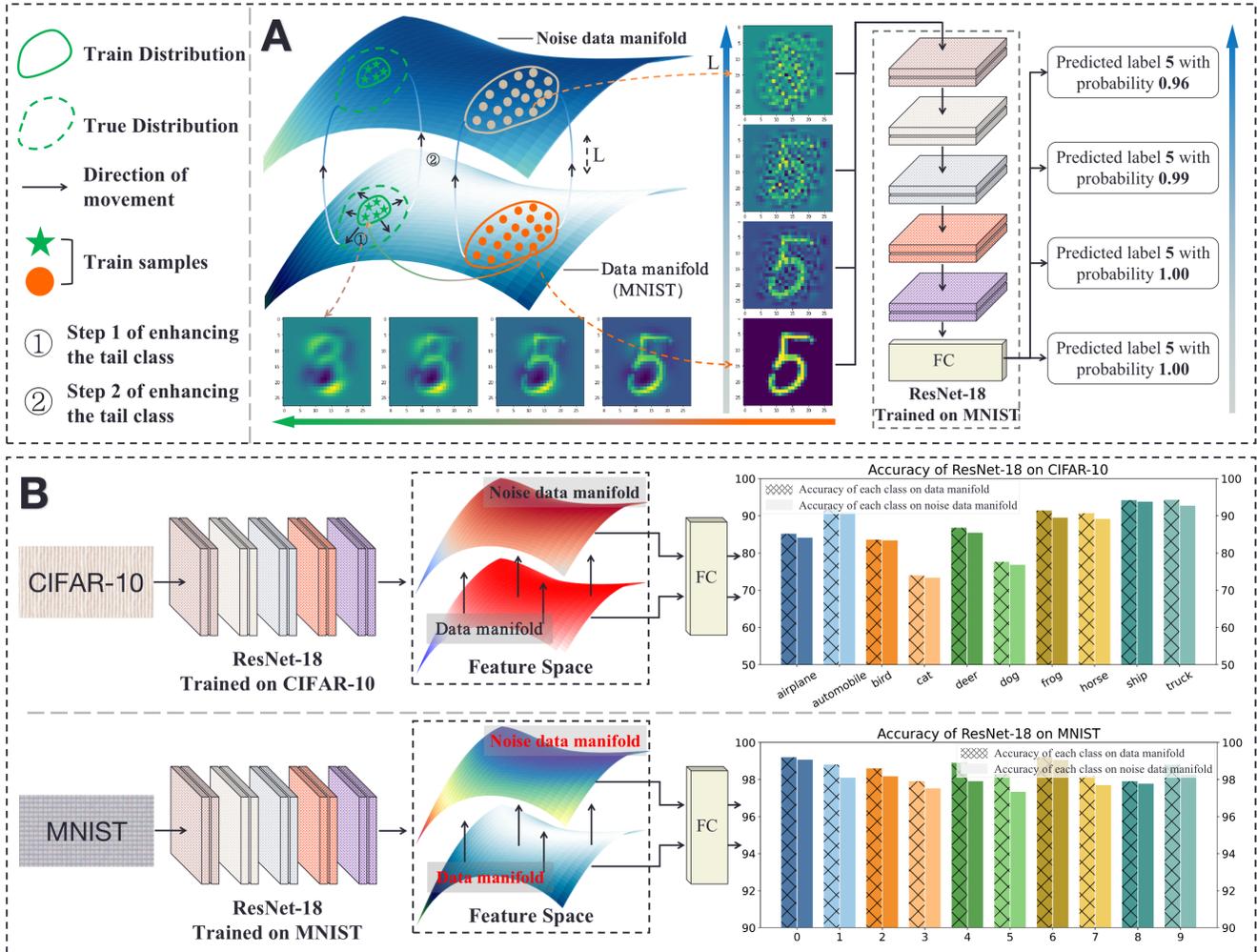


Figure 1: A: Moving in the orthogonal direction along the data manifold produces a series of images where the level of noise increases with distance. Moving on the data manifold corresponds to successive changes in the different classes of samples. A trained deep neural network can predict samples on noisy data stream shapes with a high confidence level. B: The network trained on the balanced dataset has excellent robustness. ResNet-18 was trained on CIFAR-10 and MNIST respectively and then tested for its performance on noisy data manifolds.

the CIFAR-10-LT training set, the accuracy of the model for the tail class decreases rapidly as the distance between the noisy data manifold and the data manifold increases, leading to a long-tailed distribution of model robustness (Fig.3A). This suggests that the decision surface is not only biased towards the tail class in terms of data manifold, but that the degree of bias towards the tail class increases with increasing distance between the noisy data manifold and the data manifold (Fig.3B).

The long-tailed phenomenon of robustness makes the model more vulnerable in test scenarios. To alleviate this phenomenon, we propose the orthogonal uncertainty representation of tail classes in the feature space. It is well-compatible with existing methods to improve the performance of the model on both the underlying distribution and the noisy data manifold. The main contributions of this work are summarised as follows.

- 1) We discover and define the long-tailed phenomenon of model robustness on unbalanced datasets and propose a corresponding measure of unbalance, *RIF*. (Section 2)
- 2) We propose the orthogonal uncertainty representation (*OUR*) of feature embedding. *OUR* is simple and efficient, plug and-play, and does not affect the speed of inference. (Section 3.1)
- 3) We solve the problem that calculating the orthogonal direction of feature manifolds interrupts training and consumes time and video memory, enabling end-to-end and low-cost applications *OUR*. (Section 3.2)
- 4) Comprehensive experiments show that our method has excellent compatibility and generality, demonstrating superior performance on multiple long-tailed datasets and effectively improving the long-tailed phenomenon of model robustness.

2 MOTIVATION: THE LONE-TAILED PHENOMENON OF MODEL ROBUSTNESS

In this section, we first introduce the method of constructing noisy data manifolds, and then discover and define the long-tail phenomenon of model robustness and the measure of imbalance factor. Finally, we analyze the factors that affect the performance of the tail class, thus pointing out the directions and goals of the research.

2.1 Data manifold and noise data manifold

2.1.1 Constructing noisy data manifold. The manifold distribution law [16] considers that natural images distribute around a low-dimensional manifold in a high-dimensional space, called a data manifold. As shown in Fig.1A, [10] found that moving the sample points (i.e., images) along the direction orthogonal to the data manifold, the noise of the images continues to increase, and these sample points in the orthogonal direction constitute the noisy data manifold. The following describes how to generate the noisy data manifold.

Given an image dataset with the number of samples N , assume that the size of the image is $l \times w \times h = d$ and all samples are denoted as $X = [x_1, \dots, x_N] \in \mathbb{R}^{d \times N}$. In the d -dimensional sample space, each image is considered a point, and the set of points corresponding to all images constitutes the data manifold. The intrinsic dimension of a data manifold is usually smaller than the dimension d of the linear space, so a direction vector $U \in \mathbb{R}^d$ orthogonal to the data manifold can be found to construct the noisy data manifold. If U is strictly orthogonal to the data manifold, then its inner product with any vector $(x_i - c) \in \mathbb{R}^d, i = 1, \dots, N$ is 0, where $c = \frac{1}{N} \sum_{i=1}^N x_i$. Therefore, we solve for U by optimizing the following objective.

$$\min \sum_{i=1}^N ((x_i - c)^T U)^2. \quad (1)$$

Let $y_i = x_i - c \in \mathbb{R}^d$, then the optimization objective is transformed into

$$\min \sum_{i=1}^N (y_i^T U)^2 = \min \sum_{i=1}^N U^T y_i y_i^T U = \min (U^T (\sum_{i=1}^N y_i y_i^T) U). \quad (2)$$

Let $Y = [y_1, \dots, y_N] \in \mathbb{R}^{d \times N}$ and $\sum_{i=1}^N y_i y_i^T = YY^T \in \mathbb{R}^{d \times d}$. The optimization objective can be equivalent to

$$\begin{aligned} \min (U^T YY^T U), \\ \text{s.t. } U^T U = 1. \end{aligned} \quad (3)$$

Construct the Lagrangian function $L(U, \lambda) = U^T YY^T U - \lambda(U^T U - 1)$, where λ is a coefficient. By making $\frac{\partial L(U, \lambda)}{\partial U}$ and $\frac{\partial L(U, \lambda)}{\partial \lambda}$ equal to 0, respectively, we get

$$\begin{aligned} YY^T U &= \lambda U, \\ U^T U &= 1. \end{aligned} \quad (4)$$

Obviously, YY^T is the covariance matrix of X and U is the eigenvector of YY^T . Further, from $(YY^T U, U) = (\lambda U, U)$ we can get

$$\lambda = (YY^T U, U) = U^T (YY^T)^T U = U^T YY^T U. \quad (5)$$

Therefore, in combination with equation (3), the optimization objective is ultimately equivalent to $\min_U(\lambda)$. The above results

show that U can take the eigenvector corresponding to the smallest eigenvalue of YY^T .

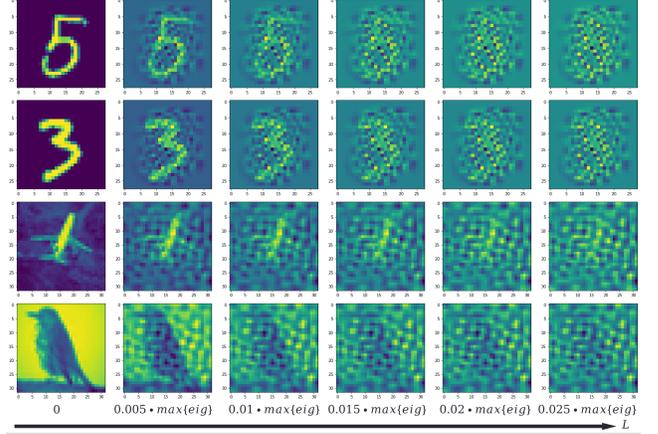


Figure 2: The above two rows show the variation process of the samples in the MNIST dataset moving along U . The lower two rows show the variation process of the airplane and bird images in CIFAR-10 moving along U . $\max\{eig\}$ denotes the maximum eigenvalue of YY^T . It can be observed that when $L = 0.02 \times \max\{eig\}$, it is already difficult for the human eye to distinguish these images.

As shown in Fig.1A, the data manifold formed by the sample set $X = [x_1, \dots, x_N] \in \mathbb{R}^{d \times N}$ is shifted by a distance L along the direction U to obtain the sample set $X' = X + LU$, and X' forms a noise data manifold. The farther the noisy data manifold is from the data manifold, the noisier the image in X' is. We correlate the value of L with the maximum eigenvalue of YY^T . Studies on MNIST and CIFAR-10 show that images on noisy data manifold are almost unrecognizable by the human eye when L is about 2% of the maximum eigenvalue of YY^T (Fig.2). In the following, L defaults to 2% of the largest eigenvalue of the sample covariance matrix if not otherwise specified.

2.1.2 Why the orthogonal direction? Moving a sample in a non-orthogonal direction may imply changes in the main features of that sample, **resulting in a conflict between the moved sample and the label**. For example, in Fig.1A, moving from the orange sample to the green sample on the data manifold corresponds to a gradual change in the image from the number 5 to the number 3, while the label of the orange sample is kept constant. Of course, shifts in non-orthogonal directions may be a potential method for generating adversarial samples, which we do not discuss in depth in this work. Since the direction of random noise is uncontrollable, it cannot be used to produce the noisy data manifold. In Section 4 we compare noise-based data augmentation with our method.

2.2 Equity in model robustness on balanced data

Surprisingly, Fig.1A shows that ResNet-18 trained on MNIST can correctly classify samples on noisy data manifolds with high confidence even in the face of heavily noisy images that are meaningless to the human eye. ResNet-18 has the same setup as ResNet-32 in section 4.1.

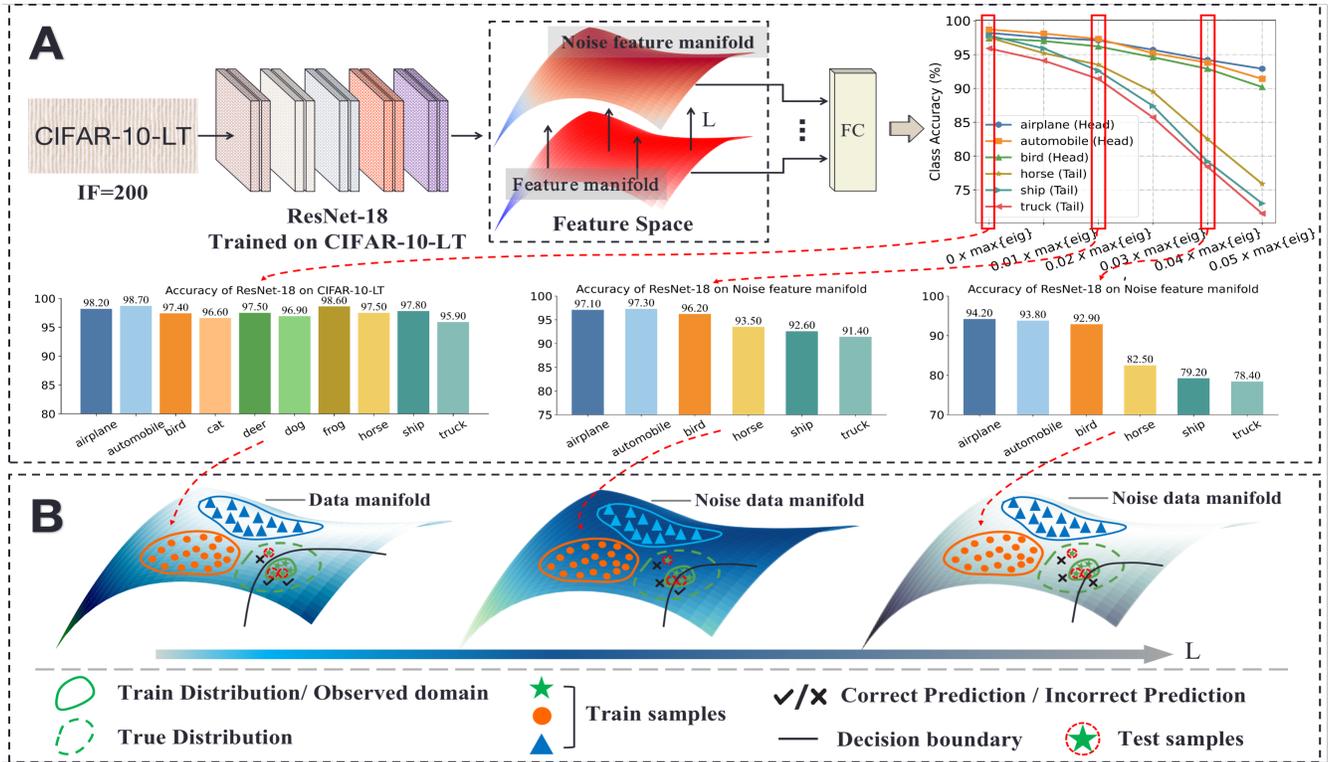


Figure 3: A: Long-tailed phenomenon for robustness of ResNet-18 trained on CIFAR-10-LT. The RIF increases as L increases. B: When the test sample is outside the observed domain, the model may predict incorrectly. Combined with the trend in RIF , we speculate that as L increases, the decision surface on the noisy data manifold becomes more biased towards the tail class.

We explore and find that the above phenomenon also exists in the feature space. First, ResNet-18 is trained on sample-balanced MNIST and CIFAR-10 respectively, and extracts the features of all samples, which constitute the feature data manifolds (referred to as feature manifolds below). We test the classification accuracy of ResNet-18 on the feature manifolds as well as the noisy feature manifolds respectively, and the experimental results are shown in Fig.1B. ResNet-18 has excellent generalization performance for each class on noisy feature manifolds, and the overall test accuracy on noisy feature manifolds for CIFAR-10 and MNIST is only 1.04% and 0.49% lower than the overall test accuracy on feature manifolds, respectively.

The above experiments combined show that a model trained on a sample-balanced data manifold, whether in image space or feature space, typically has fair and well robustness on noisy data manifolds, and is far superior to humans. What makes us curious is whether the above phenomenon also exists in long-tailed data. From the sample level, the noise generalization ability of the model trained on the balanced dataset is relatively balanced for each sample. Does the same phenomenon exist in the long-tailed data?

2.3 Inequities in model robustness on long-tailed data

First think about a question: when training a classification model on a long-tailed dataset, without considering the performance of

the model outside the training domain, does the model learn equally well on the training set for both head and tail class samples? We trained ResNet-18 on CIFAR-10-LT (imbalance factor = 200) and show the training accuracy of ResNet-18 on all classes in Fig.3A. As can be seen, ResNet-18 can recognize training samples for both head and tail classes with high accuracy, which indicates that ResNet-18 adequately fits the distribution of the training domain. However, is this sufficient to prove that the model learns fairly for each class over the training domain?

We further explore the ability of the model to generalize over noisy data manifolds in the feature space. First, training ResNet-18 on CIFAR-10-LT and extracting the features of all training samples, the maximum eigenvalue of the sample covariance matrix is denoted as $\max\{eig\}$. Then the shift distance was gradually increased along the orthogonal direction of the feature manifold until $L = 0.05 \times \max\{eig\}$, and multiple noisy feature manifolds were generated in this process. Testing the classification performance of ResNet-18 on each noisy feature manifold, the experimental results are shown in Fig.3A. We find that the performance of the tail classes decreases earlier and faster than the head classes as the distance L increases, resulting in class accuracies on the noisy feature manifold that exhibit increasingly extreme long-tailed distributions. **This suggests that the robustness of models trained on long-tailed datasets to classes also exhibits a long-tailed distribution.** Not only is the decision surface skewed toward the tail class on the feature manifold, but it is also skewed even more heavily toward

the tail class in the noisy feature manifold (see Fig.3B). In the following, we formally define the long-tailed phenomenon of model robustness and its imbalance metric.

Definition 2.1 (The long-tailed phenomenon of model robustness). Given a dataset containing C classes (data manifold) and training a classification model on it, test the class accuracy A_1, \dots, A_C of the model on the data manifold. Then construct a noisy data manifold and test the class accuracy A'_1, \dots, A'_C of the classification model on it. The difference in class accuracy, $A_i - A'_i, i = 1, \dots, C$, reflects the robustness of the model to the classes. The long-tailed phenomenon of model robustness arises when the difference in the accuracy of the classes is unbalanced.

Definition 2.2 (Imbalance factor for model robustness). The imbalance factor for model robustness is defined as

$$RIF = \max\{A_i - A'_i\} - \min\{A_i - A'_i\} (i = 1, \dots, C).$$

A larger RIF indicates that the model is more imbalanced in its robustness to the class. When the robustness is balanced, RIF = 0.

Even when fully fitting the training domain, the model is not as robust to the tail classes as the head classes, and this unfair performance may be even worse on the test set. The above results lead us to consider more comprehensively the reasons affecting the performance of the tail class.

2.4 Rethinking the factors affecting tail class performance in long-tailed recognition

First, we define the meaning of two concepts. For a class, the observed distribution is the distribution consisting of the available samples, and the underlying distribution is the true distribution in addition to the observed distribution. Combining with the discoveries in Section 2.3, we believe that the reasons for the performance limitations of the tail class are as follows.

- 1) As shown in Fig.3B, the few samples of the tail classes do not well represent its true distribution, so the performance of the model is limited outside the training domain [4, 31]. How to recover the underlying distribution of the tail classes is the key to the study.
- 2) The robustness of the model to the classes shows a long-tailed distribution. As illustrated in Fig.3B, the decision boundary of the model trained on long-tailed data is not only biased towards the tail class in the data manifold but also biased more severely in the noisy data manifold. This limits the noise invariance of the model on tail classes.

Previous studies have not taken into account the second cause of damage to the tail class, and simply expanding the data on the data manifold is not sufficient to mitigate the severe bias of the decision surface on the noisy data manifold. In this work, we propose a simple and efficient method with good compatibility to compensate for the shortcomings of existing methods.

3 ROBUST DEEP LONG-TAILED LEARNING

Consistent with [1, 12, 16, 19], we are interested in ways to enhance tail classes in the feature space. In the following, the data manifold in the feature space is called the feature manifold, and the corresponding noisy data manifold is called the noisy feature

manifold. To mitigate the long tail phenomenon of model robustness, we propose the orthogonal uncertainty representation (*OUR*) of features in Section 3.1. In Section 3.2, we propose an end-to-end training scheme that substantially reduces the time cost and memory consumption of applying *OUR*.

3.1 Orthogonal uncertainty representation

The orthogonal uncertainty representation aims to augment the tail class samples along the orthogonal direction of the feature manifold. Given a long-tailed dataset X containing C classes and a deep neural network $Model = \{f(x, \theta_1), g(z, \theta_2)\}$, where $f(x, \theta_1)$ denotes a feature sub-network with parameter θ_1 and $g(z, \theta_2)$ denotes a classifier with parameter θ_2 . The feature embedding corresponding to X is assumed to be $Z = f(X, \theta_1) = [z_1, \dots, z_N] \in \mathbb{R}^{p \times N}$, where $N = \sum_{i=1}^C N_i$, p is the sample dimension and N_i denotes the sample number of class i . Z forms a feature manifold and calculates the sample covariance matrix $\Sigma_Z = \frac{1}{N} ZZ^T \in \mathbb{R}^{p \times p}$. The maximum and minimum eigenvalues of Σ_Z are denoted by λ_{\max} and λ_{\min} , respectively. The orthogonal direction $U \in \mathbb{R}^p$ of the feature manifold is the eigenvector corresponding to λ_{\min} .

Suppose a batch of samples is encoded by $f(X_B, \theta_1)$ as $Z_B \in \mathbb{R}^{p \times bs}$ and bs as the batch size, where the feature embedding belonging to tail class t is $Z_{B,t} = [z_t^1, \dots, z_t^{n_t}] \in \mathbb{R}^{p \times n_t}$. We model the uncertainty representation of features by applying perturbations along the direction U for $Z_{B,t}$, thereby enhancing the noise invariance of the model to the tail class t on noisy data manifolds. The specific form can be formulated as

$$\begin{aligned} & \text{Orthogonal Uncertainty Representation of } Z_{B,t} \\ OUR(Z_{B,t}) &= \overbrace{Z_{B,t} + \mu \lambda_{mean} [\varepsilon_1 U, \dots, \varepsilon_{n_t} U]} \in \mathbb{R}^{p \times n_t}, \quad (6) \\ & \varepsilon_i \sim N(0, 1), i = 1, \dots, n_t. \end{aligned}$$

$\varepsilon_1, \dots, \varepsilon_{n_t}$ all follow a standard Gaussian distribution and are independent of each other, they increase the uncertainty of each feature embedding in $Z_{B,t}$. λ_{mean} denotes the average of the top 10 eigenvalues. Our study in Section 2 shows that images on noisy feature manifolds are no longer recognizable to the human eye when the distance to the feature manifold is $0.02\lambda_{max}$. To improve the stability of *OUR*, we use λ_{mean} instead of λ_{max} to perturb the features. The $\mu\lambda_{mean}$ term therefore guarantees a reasonable range of perturbations and the reasonable choice of μ will be discussed in Section 4.2. However, an additional consideration is that calculating the orthogonal direction U of the feature manifold interrupts training. We detail the difficulties faced in practice in the next subsection.

3.2 End-to-end training with *OUR*

The parameters of the model change continuously with training, resulting in an offset between features extracted from the same sample at different periods. Therefore, when applying the orthogonal uncertainty representation (*OUR*) in the feature space, the orthogonal direction of the feature manifold needs to be continuously updated. However, calculating the orthogonal direction of the feature manifold requires re-extracting features from the entire dataset, which significantly increases the time cost and interrupts the training, complicating the training process.

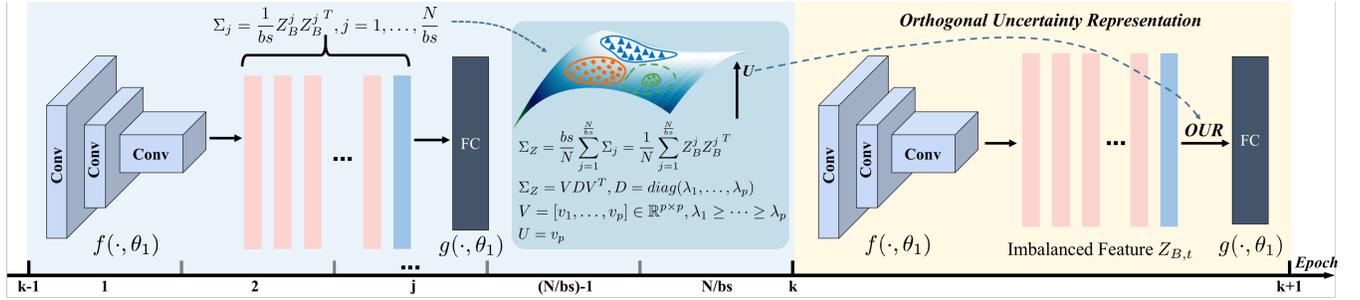


Figure 4: End-to-end training process with OUR. When applying OUR in the k -th training epoch, the covariance matrix corresponding to each batch feature needs to be calculated and saved in the $k-1$ -th epoch. At the end of the $k-1$ -th epoch, the saved N/bs covariance matrices are used to calculate the corresponding feature covariance matrix of the entire dataset, which further yields the orthogonal direction U of the feature manifold and is utilized in the k -th epoch.

The feature slow shift phenomenon [39] indicates that as the training epoch increases, the shift between the historical and latest features of the same sample decreases to the point where the historical features can be used to approximate the latest features. Assume that OUR is applied from the k -th epoch. Given that the entire dataset can be traversed in a single training epoch, the intuitive solution is to save all batches of features from the $k-1$ -th training epoch in place of the latest features Z for the entire dataset. Before the k -th epoch, the covariance matrix is calculated with the saved features, and then the orthogonal direction U of the feature manifold is further obtained. And so on, the historical features saved in the k -th epoch will be used to calculate the orthogonal direction used in the $k+1$ -th epoch. [6, 39, 40] shows that in the training of classification models, only 5 epochs are usually needed to make the shift of features small enough. And k is significantly larger than 5, so there is no need to be concerned about the feature shift not being small enough.

Although the above approach avoids extracting features from the entire dataset, additional storage space is still required to hold all the features generated in an epoch. To further reduce memory consumption, we propose to calculate and save the covariance matrix $\Sigma_j = \frac{1}{bs} Z_B^j Z_B^{jT}$, $j = 1, \dots, \frac{N}{bs}$ for each batch of features in the k -th epoch, where $Z_B^j \in \mathbb{R}^{p \times bs}$ denotes the j -th batch of features. The sum of $\Sigma_1, \dots, \Sigma_{\frac{N}{bs}}$ is then used to approximate the covariance matrix of all features from the dataset. Specifically, when the k -th epoch ends, the feature covariance matrix of the entire dataset can be approximated as

$$\Sigma_Z = \frac{bs}{N} \sum_{j=1}^{\frac{N}{bs}} \Sigma_j = \frac{bs}{N} \left(\frac{1}{bs} Z_B^1 Z_B^{1T} + \dots + \frac{1}{bs} Z_B^{\frac{N}{bs}} Z_B^{\frac{N}{bs}T} \right) = \frac{1}{N} \sum_{j=1}^{\frac{N}{bs}} Z_B^j Z_B^{jT} \in \mathbb{R}^{p \times p}$$

Calculate the orthogonal direction U of the feature manifold based on Σ_Z and apply it to the next training epoch. Fig.4 illustrates the process of applying OUR end-to-end. At the cost of negligible memory consumption (only $\frac{N}{bs}$ covariance matrices need to be stored), we solve the problem of interrupted training and time consumption when calculating the orthogonal direction of the feature manifold.

Since OUR is dedicated to mitigating the long-tailed phenomenon of robustness found for the first time, it has good compatibility and generality as it does not overlap with the research aims of other methods. Our experiments also show that existing methods are not effective in mitigating the long-tailed phenomenon of robustness (Fig.6). Listing 1 demonstrates a simple implementation of OUR that can easily be combined with existing methods. Fig.1A illustrates the two steps in enhancing the tail classes. We hope to co-train to improve the model's performance on both underlying distribution and noisy data manifold.

Listing 1: End-to-end training with OUR

```

1 for epoch in range(M):
2     Q = np.empty([N/bs, p, p]) # bs is batch size
3     # X_B: data, y_B: labels
4     for X_B, y_B, j in loader(N/bs):
5         Z_B = f(X_B, theta_1)
6         if epoch == k-1:
7             Sigma = np.matmul(Z_B, Z_B.T)
8             Q[j] = Sigma
9         elif epoch >= k:
10            Sigma = np.matmul(Z_B, Z_B.T)
11            Q[j] = Sigma
12            for i in range(C):
13                # Execute OUR on the tail category
14                if Z_B_i is a tail category:
15                    Z_B_i = OUR(Z_B_i, U, mu)
16            y^ = g(Z_B, theta_2)
17            loss = loss function(y_B, y^)
18            loss.backward()
19            optimizer.step()
20            Sigma_Z = np.sum(Q, axis = 0)/N
21            vals, vecs = np.linalg.eig(Sigma_Z)
22            U = vecs[:, p-1]
23            lambda_mean = np.mean(vecs[:, 0:10])

```

4 EXPERIMENTS

4.1 Datasets and Experimental Setting

Datasets. We evaluated our method on four long-tail benchmark datasets CIFAR-10-LT, CIFAR-100 LT [7], ImageNet-LT [27], and iNaturalist 2018 [30]. Long-tailed CIFAR is the artificially produced imbalance dataset using its balanced version. We chose three long-tailed versions with imbalance factors (IF) of 10, 50, and 100 for

training. ImageNet-LT contains a total of 1000 classes with an imbalance factor of 256. iNaturalist 2018 is a large-scale species classification dataset with a long-tailed distribution and imbalance factor of 500. In this work, the official training and testing splits of all datasets are used for a fair comparison.

Experimental Setting. In accordance with the previous setup [5, 25, 27], we adopt ResNet-32 [11] as the backbone network on the CIFAR-10/100-LT and adopt an SGD optimizer with momentum 0.9 for all experiments. The batch size is set to 128, the initial learning rate is 0.1, and a total of 200 epochs are trained. Linear warm-up of the learning rate is used in the first five epochs, with the learning rate decaying by 0.1 times at 160 and 180 epochs respectively. We employ ResNeXt-50 [35] on ImageNet-LT and ResNet-50 on iNaturalist 2018 as the backbone network, training 200 epochs. In all experiments, the batch size is set to 256 (for ImageNet-LT) / 512 (for iNaturalist 2018), the initial learning rate is 0.1 (linear LR decay), and the SGD optimizer with a momentum of 0.9 is used to train all models.

4.2 Effect of hyper-parameter μ

μ determines the degree of uncertainty of the feature embedding, and when $\mu = 0$, *OUR* does not perform a transformation on the feature embedding. Since we observe that the human eye can barely recognize samples on noisy data manifolds when $\mu = 0.02$. Therefore, we explore the effect of μ on *OUR* in the interval $[0, 0.1]$. Experimental results on CIFAR-10-LT, CIFAR-100-LT and ImageNet-LT are shown in Fig.5. It can be seen that the performance of the model increases and then decreases as μ increases. When μ is too small, the perturbation of the feature embedding is weak and the model does not learn sufficiently from the noise. Due to the scarcity of tail class samples, the model’s learning of the original feature distribution is easily disturbed when μ is too large. Specifically, optimal performance is achieved on CIFAR-10/100-LT when μ is taken as 0.02 and 0.03 and on ImageNet when μ is taken as 0.01 and 0.02 for *OUR*.

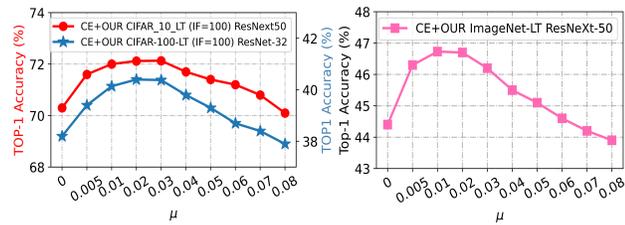


Figure 5: Ablation on CIFAR-10-LT (IF=100), CIFAR-100-LT (IF=100), and ImageNet-LT for select hyperparameter μ .

4.3 Evaluation on CIFAR-10/100-LT

Table 1 shows the experimental results on CIFAR-10/100-LT. The pink rows show improvements to the cross-entropy loss as well as cost-sensitive learning methods, and the blue rows show improvements to the information augmentation methods. It can be observed that our proposed *OUR* significantly improves the existing methods on all datasets. In particular, it also shows effectiveness on several state-of-the-art methods (RIDE+CMO, GCL, ResLT), which indicates that existing methods lack attention to the long-tailed

Table 1: Comparison on CIFAR-10-LT and CIFAR-100-LT. The accuracy (%) of Top-1 is reported. The best and second-best results are shown in underlined bold and **bold**, respectively.

Dataset	CIFAR-10-LT			CIFAR-100-LT		
Backbone Net	ResNet-32					
Imbalance Factor	100	50	10	100	50	10
BBN [44]	79.8	82.1	88.3	42.5	47.0	59.1
De-c-TDE [29]	80.6	83.6	88.5	44.1	50.3	59.6
Cross Entropy	70.3	74.8	86.3	38.2	43.8	55.7
+ OUR	72.1	76.1	87.5	39.7	45.2	56.8
CB-Focal [7]	74.5	79.2	87.4	39.6	45.3	57.9
+ OUR	75.9	80.1	88.2	40.8	46.5	58.6
LDAM-DRW [3]	77.0	81.0	88.2	42.0	46.6	58.7
+ OUR	78.1	81.8	88.9	42.8	47.4	59.3
MiSLAS [43]	82.1	85.7	90.0	47.0	52.3	63.2
+ OUR	83.4	86.7	90.8	48.1	53.1	63.9
RIDE + CMO [25]	-	-	-	50.0	53.0	60.2
+ OUR	-	-	-	50.8	53.9	60.7
GCL [18]	82.7	85.5	-	48.7	53.6	-
+ OPeN [38]	83.1	85.8	-	49.2	53.9	-
+ OUR	83.7	86.3	-	49.8	54.5	-
ResLT [5]	80.4	83.5	89.1	45.3	50.0	60.8
+ OPeN [38]	80.8	83.8	89.6	45.8	50.4	61.2
+ OUR	81.6	84.3	90.0	46.5	50.9	61.7

phenomenon of robustness. It is important to note that CE is not designed for long-tailed classification, but our approach improves the performance of CE by 1.8% and 1.5% on CIFAR-10-LT (IF = 100) and CIFAR-100-LT (IF = 100) by improving the long-tailed phenomenon of model robustness alone. This further indicates that the long-tailed distribution of robustness in the long-tailed scenario limits the performance of the classification model. *OUR* performs better on datasets with larger IF, which is in line with our expectation since the long tail of model robustness becomes more severe when the data are more unbalanced. For example, on the CIFAR-10-LT (IF = 100, 50, 10) and CIFAR-100-LT (IF = 100, 50, 10), *OUR* achieved performance gains of 1.4%, 0.9%, 0.8%, 1.2%, 1.2% and 0.7% for CB-Focal.

We also compare *OUR* with the latest noise-based augmentation method OPeN [38], and it can be observed that OPeN provides very weak improvements on GCL and ResLT. Random noise may cause changes in the main features of the samples, resulting in the ambiguity between samples and labels (refer to the analysis in Section 2.1.2). Compared with OPeN, our method not only does not add additional training samples but also has better performance.

4.4 Evaluation on ImageNet-LT and iNat 2018

We report in Table 2 not only the overall performance of *OUR* but also add the evaluation results of *OUR* on three subsets (Head, Middle, Tail) of the two datasets. It can be observed that *OUR* improves the overall performance of the other methods by at least about 1% on both datasets. For the ImageNets dataset, *OUR* performs superiorly, delivering performance gains of 1.6%, 1.4%, 1.2% and 1.5% for CE, Focal, LADE and OFA, respectively. For the iNaturalist dataset, *OUR* improves the overall performance of CE, Focal and OFA by 1.5%, 1.4% and 1.3%, respectively. *OUR* also maintains consistent superiority in the face of the latest state-of-the-art methods. *OUR* improves the overall performance of RIDE+CMO by 1% and 0.9% on the two datasets and improves the overall performance of ResLT by

Table 2: Comparison on ImageNet-LT and iNaturalist2018. The Top-1 Acc (%) is reported. The best and the second-best results are shown in underline bold and **bold, respectively.**

Methods	ImageNet-LT				iNaturalist 2018			
	ResNext-50				ResNet-50			
	H	M	T	Overall	H	M	T	Overall
DisAlign [10]	59.9	49.9	31.8	52.9	68.0	71.3	69.4	70.2
MiSLAS [1]	65.3	50.6	33.0	53.4	73.2	72.4	70.4	71.6
DiVe [11]	64.0	50.4	31.4	53.1	70.6	70.0	67.5	69.1
PaCo [12]	63.2	51.6	39.2	54.4	69.5	72.3	73.1	72.3
RIDE (3*) [12]	66.2	51.7	34.9	54.9	70.2	72.2	72.7	72.2
GCL [8]	-	-	-	54.9	-	-	-	72.0
CE	65.9	37.5	7.70	44.4	67.2	63.0	56.2	61.7
+ OUR	65.0	38.4	14.5	46.0	67.3	63.9	60.5	63.2
Focal Loss [3]	67.0	41.0	13.1	47.2	-	-	-	61.1
+ OUR	67.2	42.5	19.7	48.6	68.6	63.4	57.8	62.5
LDAM [2]	60.0	49.2	31.9	51.1	-	-	-	64.6
+ OUR	60.6	50.0	33.5	52.2	69.0	66.9	62.1	65.5
LADE [13]	62.3	49.3	31.2	51.9	-	-	-	69.7
+ OUR	62.4	50.5	34.4	53.1	72.2	70.6	65.9	70.7
OFA [9]	47.3	31.6	14.7	35.2	-	-	-	65.9
+ OUR	47.2	32.8	18.6	36.7	69.7	68.2	64.8	67.2
RIDE + CMO [7]	66.4	54.9	35.8	56.2	70.7	72.6	73.4	72.8
+ OPeN	66.7	55.1	37.0	56.8	70.4	73.4	74.1	73.2
+ CR	66.5	55.7	37.9	57.2	70.5	73.9	74.8	73.7
ResLT [14]	63.0	50.5	35.5	53.0	68.5	69.9	70.4	70.2
+ OPeN	63.3	51.3	36.2	53.6	68.6	70.5	71.2	70.7
+ OUR	63.5	51.7	37.3	54.3	68.8	71.1	72.0	71.3

RIDE (3*) denotes the RIDE model with 3 experts. RIDE in RIDE+CMO comes with 3 experts. H, M, and T denote the Head (more than 100 images), Middle (20-100 images), and Tail (less than 20 images) subsets of the dataset, respectively.

1.3% and 1.1% on the two datasets, respectively. It should be noted that CMO already expands the richness of tail classes by pasting the foreground of tail classes into the background of head classes, so the samples are balanced. Even in such a case, *OUR* still improves the performance of RIDE+CMO, which is a solid indication that *OUR* has excellent compatibility and does not conflict with the goals pursued by existing methods.

OUR delivers the most significant improvement for the tail subset. On ImageNet-LT, *OUR* improves the performance of CE, Focal, and OFA on the tail subset by 7.5%, 6.6%, and 4.6%, respectively. Compared to OPeN, our method results in better performance of RIDE+CMO and ResLT on both datasets. Specifically, RIDE+CMO+*OUR* outperforms RIDE+CMO with OPeN by 0.9% and 0.7%, respectively, on the tail subset of both datasets, and ResLT+*OUR* outperforms ResLT+OPeN by 1.1% and 0.8%, respectively. Comprehensive experiments on CIFAR-10/100-LT, ImageNet-LT, and iNaturalist 2018 show that our method is stable and excellent, outperforming the recently advanced noise-based augmentation method OPEN. The experiments and analysis in Section A further demonstrate the performance gap between OPEN and *OUR*.

4.5 Evaluation with multiple backbones

To fully evaluate the generality of *OUR*, we adopt different backbone networks on ImageNet-LT to demonstrate the effective improvement of *OUR* for other methods. The experimental results are shown in Table 3, with the **cyan rows** illustrating the method with ResNet-18 as the backbone network and the **violet rows** illustrating the method with ResNeXt-101-32×4d as the backbone network. *OUR* brings about a 1% performance gain in overall accuracy for all

Table 3: Top-1 Accuracy (%) with Various ResNet Backbones.

Methods	Backbone Net	Head	Middle	Tail	Overall
CE	ResNet-10	59.7	29.4	5.7	37.3
+OUR	ResNet-10	60.0	30.3	8.1(+2.4)	38.5(+1.2)
cRT	ResNet-10	53.8	41.3	25.4	43.2
+OUR	ResNet-10	54.0	42.2	26.7(+1.3)	44.0(+0.8)
LWS	ResNet-10	51.8	42.2	28.1	43.4
+OUR	ResNet-10	52.4	42.5	29.2(+1.1)	44.1(+0.7)
ResLT	ResNet-10	52.3	41.6	27.6	43.0
+OUR	ResNet-10	52.7	42.5	29.0(+1.4)	43.9(+0.9)
CE	ResNeXt-101	69.6	44.6	15.6	49.6
+OUR	ResNeXt-101	69.9	45.7	17.1(+1.5)	50.6(+1.0)
cRT	ResNeXt-101	66.2	50.4	30.8	53.3
+OUR	ResNeXt-101	66.4	51.7	32.6(+1.8)	54.4(+1.1)
LWS	ResNeXt-101	65.7	51.4	34.7	54.0
+OUR	ResNeXt-101	66.0	52.0	36.3(+1.6)	54.8(+0.8)
ResLT	ResNeXt-101	63.3	53.3	40.3	55.1
+OUR	ResNeXt-101	63.8	54.1	41.7(+1.4)	56.0(+0.9)

methods, with CE+*OUR* outperforming CE by 1.2% overall when ResNet-18 is used as the backbone network. Consistent with Table 2, *OUR* has the most significant performance in tail classes. For example, when ResNet-10 is used as the backbone, CE+*OUR* outperforms CE by 2.4% on the tail subset. When ResNeXt-101-32×4d is adopted as the backbone, cRT+*OUR* outperforms cRT by 1.8% on the tail subset. The evaluation of various backbones solidly demonstrates that our method can work in a wide range of scenarios.

4.6 CONCLUSION

This work finds and defines the long-tail phenomenon of model robustness in data imbalance scenarios and proposes a corresponding calculation of the imbalance factor (i.e., *RIF*). Then, we propose the orthogonal uncertainty representation (*OUR*) of features to mitigate the model bias on data manifolds and noisy data manifolds. Also, an end-to-end training scheme is proposed for efficient and fast application of *OUR*. Although our approach strongly mitigates the degree of imbalance in model robustness, it still needs to be improved. In the future, we hope more work will focus on the long-tail phenomenon of model robustness to make the model more robust in data imbalance scenarios.

ACKNOWLEDGMENTS

This work was supported in part by the Key Scientific Technological Innovation Research Project by Ministry of Education, the State Key Program and the Foundation for Innovative Research Groups of the National Natural Science Foundation of China (61836009), the Major Research Plan of the National Natural Science Foundation of China (91438201, 91438103, and 91838303), the National Natural Science Foundation of China (U22B2054, U1701267, 62076192, 62006177, 61902298, 61573267, 61906150, and 62276199), the 111 Project, the Program for Cheung Kong Scholars and Innovative Research Team in University (IRT 15R53), the ST Innovation Project from the Chinese Ministry of Education, the Key Research and Development Program in Shaanxi Province of China(2019ZDLGY03-06), the National Science Basic Research Plan in Shaanxi Province of China(2022JQ-607), the China Postdoctoral fund(2022T150506), the Scientific Research Project of Education Department In Shaanxi Province of China (No.20JY023), the National Natural Science Foundation of China (No. 61977052).

REFERENCES

- [1] Sumyeong Ahn, Jongwoo Ko, and Se-Young Yun. 2023. CUDA: Curriculum of Data Augmentation for Long-tailed Recognition. In *The Eleventh International Conference on Learning Representations*. <https://openreview.net/forum?id=RgUPdudkWIN>
- [2] Shaden Alshammari, Yu-Xiong Wang, Deva Ramanan, and Shu Kong. 2022. Long-tailed recognition via weight balancing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6897–6907.
- [3] Kaidi Cao, Colin Wei, Adrien Gaidon, Nikos Arachiga, and Tengyu Ma. 2019. Learning imbalanced datasets with label-distribution-aware margin loss. *Advances in neural information processing systems* 32 (2019).
- [4] Peng Chu, Xiao Bian, Shaopeng Liu, and Haibin Ling. 2020. Feature space augmentation for long-tailed data. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIX* 16. Springer, 694–710.
- [5] Jiequan Cui, Shu Liu, Zhuotao Tian, Zhisheng Zhong, and Jiaya Jia. 2022. ResL: Residual learning for long-tailed recognition. *IEEE transactions on pattern analysis and machine intelligence* 45, 3 (2022), 3695–3706.
- [6] Jiequan Cui, Zhisheng Zhong, Shu Liu, Bei Yu, and Jiaya Jia. 2021. Parametric contrastive learning. In *Proceedings of the IEEE/CVF international conference on computer vision*. 715–724.
- [7] Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie. 2019. Class-balanced loss based on effective number of samples. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 9268–9277.
- [8] Yin Cui, Yang Song, Chen Sun, Andrew Howard, and Serge Belongie. 2018. Large scale fine-grained categorization and domain-specific transfer learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 4109–4118.
- [9] Andrew Estabrooks, Taeho Jo, and Nathalie Japkowicz. 2004. A multiple resampling method for learning from imbalanced data sets. *Computational intelligence* 20, 1 (2004), 18–36.
- [10] Luca Gremontieri and Rita Fioresi. 2022. Model-centric data manifold: the data through the eyes of the model. *SIAM Journal on Imaging Sciences* 15, 3 (2022), 1140–1156.
- [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [12] Yin-Yin He, Jianxin Wu, and Xiu-Shen Wei. 2021. Distilling virtual examples for long-tailed recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 235–244.
- [13] Feng Hong, Jiangchao Yao, Zhihan Zhou, Ya Zhang, and Yanfeng Wang. 2023. Long-tailed partial label learning via dynamic rebalancing. *arXiv preprint arXiv:2302.05080* (2023).
- [14] Youngkyu Hong, Seungju Han, Kwanghee Choi, Seokjun Seo, Beomsu Kim, and Buru Chang. 2021. Disentangling label distribution for long-tailed visual recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 6626–6636.
- [15] Jaehyung Kim, Jongheon Jeong, and Jinwoo Shin. 2020. M2m: Imbalanced classification via major-to-minor translation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 13896–13905.
- [16] Na Lei, Dongsheng An, Yang Guo, Kehua Su, Shixia Liu, Zhongxuan Luo, Shing-Tung Yau, and Xianfeng Gu. 2020. A geometric understanding of deep learning. *Engineering* 6, 3 (2020), 361–374.
- [17] Mengke Li, Yiu-Ming Cheung, and Zhikai Hu. 2022. Key point sensitive loss for long-tailed visual recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 45, 4 (2022), 4812–4825.
- [18] Mengke Li, Yiu-ming Cheung, and Yang Lu. 2022. Long-tailed visual recognition via gaussian clouded logit adjustment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6929–6938.
- [19] Shuang Li, Kaixiong Gong, Chi Harold Liu, Yulin Wang, Feng Qiao, and Xinjing Cheng. 2021. Metasaug: Meta semantic augmentation for long-tailed visual recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 5212–5221.
- [20] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. 2017. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*. 2980–2988.
- [21] Bo Liu, Haoxiang Li, Hao Kang, Gang Hua, and Nuno Vasconcelos. 2021. Gistnet: a geometric structure transfer network for long-tailed recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 8209–8218.
- [22] Jialun Liu, Yifan Sun, Chuchu Han, Zhaopeng Dou, and Wenhui Li. 2020. Deep representation learning on long-tailed data: A learnable embedding augmentation perspective. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2970–2979.
- [23] Yanbiao Ma, Licheng Jiao, Fang Liu, Yuxin Li, Shuyuan Yang, and Xu Liu. 2023. Delving into Semantic Scale Imbalance. In *The Eleventh International Conference on Learning Representations*. <https://openreview.net/forum?id=07tc5kKRlo>
- [24] Sarah Parisot, Pedro M Esperança, Steven McDonagh, Tamas J Madarasz, Yongxin Yang, and Zhenguo Li. 2022. Long-tail recognition via compositional knowledge transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6939–6948.
- [25] Seulki Park, Youngkyu Hong, Byeongho Heo, Sangdoon Yun, and Jin Young Choi. 2022. The majority can help the minority: Context-rich minority oversampling for long-tailed classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6887–6896.
- [26] Jiawei Ren, Cunjun Yu, Xiao Ma, Haiyu Zhao, Shuai Yi, et al. 2020. Balanced meta-softmax for long-tailed visual recognition. *Advances in neural information processing systems* 33 (2020), 4175–4186.
- [27] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. 2015. Imagenet large scale visual recognition challenge. *International journal of computer vision* 115 (2015), 211–252.
- [28] Saptarshi Sinha, Hiroki Ohashi, and Katsuyuki Nakamura. 2022. Class-difficulty based methods for long-tailed visual recognition. *International Journal of Computer Vision* 130, 10 (2022), 2517–2531.
- [29] Kaihua Tang, Jianqiang Huang, and Hanwang Zhang. 2020. Long-tailed classification by keeping the good and removing the bad momentum causal effect. *Advances in Neural Information Processing Systems* 33 (2020), 1513–1524.
- [30] Grant Van Horn, Oisín Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam, Pietro Perona, and Serge Belongie. 2018. The inaturalist species classification and detection dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 8769–8778.
- [31] Chaozheng Wang, Shuzheng Gao, Pengyun Wang, Cuiyun Gao, Wenjie Pei, Lujia Pan, and Zenglin Xu. 2022. Label-aware distribution calibration for long-tailed classification. *IEEE Transactions on Neural Networks and Learning Systems* 99 (2022), 1–13.
- [32] Jianfeng Wang, Thomas Lukasiewicz, Xiaolin Hu, Jianfei Cai, and Zhenghua Xu. 2021. Rsg: A simple but effective module for learning imbalanced datasets. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3784–3793.
- [33] Tao Wang, Yu Li, Bingyi Kang, Junnan Li, Junhao Liew, Sheng Tang, Steven Hoi, and Jiashi Feng. 2020. The devil is in classification: A simple framework for long-tail instance segmentation. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIV* 16. Springer, 728–744.
- [34] Xudong Wang, Long Lian, Zhongqi Miao, Ziwei Liu, and Stella X Yu. 2020. Long-tailed recognition by routing diverse distribution-aware experts. *arXiv preprint arXiv:2010.01809* (2020).
- [35] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. 2017. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1492–1500.
- [36] Zhengzhou Xu, Zenghao Chai, and Chun Yuan. 2021. Towards calibrated model for long-tailed visual recognition from prior perspective. *Advances in Neural Information Processing Systems* 34 (2021), 7139–7152.
- [37] Jing-Hao Xue and Peter Hall. 2014. Why does rebalancing class-unbalanced data improve AUC for linear discriminant analysis? *IEEE transactions on pattern analysis and machine intelligence* 37, 5 (2014), 1109–1112.
- [38] Shiran Zada, Itay Benou, and Michal Irani. 2022. Pure noise to the rescue of insufficient data: Improving imbalanced classification by training on random noise images. In *International Conference on Machine Learning*. PMLR, 25817–25833.
- [39] Yuhang Zang, Chen Huang, and Chen Change Loy. 2021. Fasa: Feature augmentation and sampling adaptation for long-tailed instance segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 3457–3466.
- [40] Yifan Zhang, Bingyi Kang, Bryan Hooi, Shuicheng Yan, and Jiashi Feng. 2023. Deep long-tailed learning: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2023).
- [41] Zizhao Zhang and Tomas Pfister. 2021. Learning fast sample re-weighting without reward data. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 725–734.
- [42] Peilin Zhao, Yifan Zhang, Min Wu, Steven CH Hoi, Minghui Tan, and Junzhou Huang. 2018. Adaptive cost-sensitive online classification. *IEEE Transactions on Knowledge and Data Engineering* 31, 2 (2018), 214–228.
- [43] Zhisheng Zhong, Jiequan Cui, Shu Liu, and Jiaya Jia. 2021. Improving calibration for long-tailed recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 16489–16498.
- [44] Boyan Zhou, Quan Cui, Xiu-Shen Wei, and Zhao-Min Chen. 2020. Bbn: Bilateral-branch network with cumulative learning for long-tailed visual recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 9719–9728.
- [45] Zhi-Hua Zhou and Xu-Ying Liu. 2005. Training cost-sensitive neural networks with methods addressing the class imbalance problem. *IEEE Transactions on knowledge and data engineering* 18, 1 (2005), 63–77.

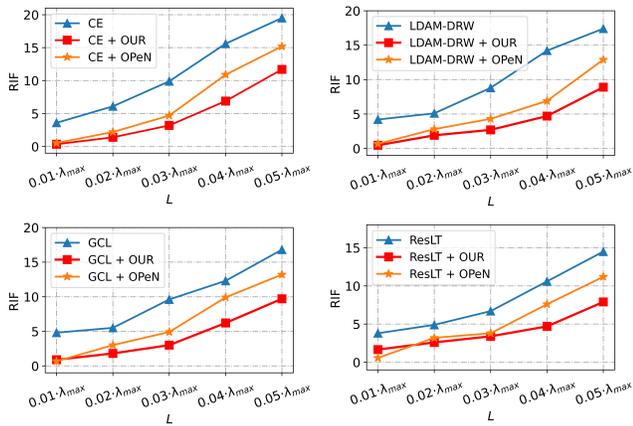


Figure 6: OUR alleviates the imbalance of model robustness. L denotes the distance between the noisy data manifold and the data manifold.

A OUR MITIGATES THE LONG TAIL OF MODEL ROBUSTNESS

To visualize the improvement effect of *OUR* on the long-tailed phenomenon of robustness, we constructed multiple noisy data manifolds corresponding to the training data of CIFAR-10-LT (IF = 100), and then calculate the values of *RIF* (Definition 2) before and after using *OUR* for CE, LDAM-DRW, GCL, and ResLT. Experimental results are illustrated in Fig.6. The imbalance degree *RIF* of model robustness is significantly reduced by adopting *OUR* to improve multiple methods, and the ability of OPeN to mitigate the long tailed phenomenon of robustness is between the original method and *OUR*. When L is small, the improvement effect of OPeN on *RIF* is close to *OUR*, but as L increases, the effect of OPeN becomes weak. This may be due to the fact that OPeN is based on pure noise and its randomness of direction leads to the inability to generate noisy data manifolds stably at long distances, which limits the performance.