
Imbalanced Adversarial Training with Reweighting

Wentao Wang*

Michigan State University
wangw116@msu.edu

Han Xu*

Michigan State University
xuhan1@msu.edu

Xiaorui Liu

Michigan State University
xiaorui@msu.edu

Yaxin Li

Michigan State University
liyaxin1@msu.edu

Bhavani Thuraisingham

The University of Texas at Dallas
bxt043000@utdallas.edu

Jiliang Tang

Michigan State University
tangjili@msu.edu

Abstract

Adversarial training has been empirically proven to be one of the most effective and reliable defense methods against adversarial attacks. However, almost all existing studies about adversarial training are focused on balanced datasets, where each class has an equal amount of training examples. Research on adversarial training with imbalanced training datasets is rather limited. As the initial effort to investigate this problem, we reveal the facts that adversarially trained models present two distinguished behaviors from naturally trained models in imbalanced datasets: (1) Compared to natural training, adversarially trained models can suffer much worse performance on under-represented classes, when the training dataset is extremely imbalanced. (2) Traditional *reweighting* strategies may lose efficacy to deal with the imbalance issue for adversarial training. For example, upweighting the under-represented classes will drastically hurt the model’s performance on well-represented classes, and as a result, finding an optimal reweighting value can be tremendously challenging. In this paper, to further understand our observations, we theoretically show that the poor data separability is one key reason causing this strong tension between under-represented and well-represented classes. Motivated by this finding, we propose *Separable Reweighted Adversarial Training* (SRAT) to facilitate adversarial training under imbalanced scenarios, by learning more separable features for different classes. Extensive experiments on various datasets verify the effectiveness of the proposed framework.

1 Introduction

The existence of adversarial samples [32, 14] has risen huge concerns on applying deep neural network (DNN) models into security-critical applications, such as autonomous driving [8] and video surveillance systems [22]. As countermeasures against adversarial attacks, adversarial training [26, 41] has been empirically proven to be one of the most effective and reliable defense methods. In general, adversarial training can be formulated to minimize the model’s average error on adversarially perturbed input examples [26, 41, 30]. Although promising to improve the model’s robustness, most existing adversarial training methods [41, 35] assume that the number of training examples from each class is equally distributed. However, datasets collected from real-world applications typically have imbalanced distribution [13, 25, 34]. Hence, it is natural to ask: *what is the behavior of adversarial training under imbalanced scenarios? Can we directly apply existing imbalanced learning strategies in natural training to tackle the imbalance issue for adversarial training?* Recent studies find that adversarial training usually presents distinct properties from natural training [31, 38]. For example,

*Equal Contribution

compared to natural training, adversarially trained models suffer more from the overfitting issue [31]. Moreover, it is evident from a recent study [38] that the adversarially trained models tend to present strong class-wise performance disparities, even if the training examples are uniformly distributed over different classes. Imagine that if the training data distribution is highly imbalanced, these properties of adversarial training can be greatly exaggerated and make it extremely difficult to be applied in practice. Therefore, it is important but challenging to answer aforementioned questions.

As the initial effort to study the imbalanced problem in adversarial training, in this work, we first investigate the performance of existing adversarial training under imbalanced settings. As a preliminary study shown in Section 2.1, we apply both natural training and PGD adversarial training [26] on multiple imbalanced image datasets constructed from the CIFAR10 dataset [21] using the ResNet18 architecture [17] and evaluate trained models’ performance on class-balanced test datasets. From the preliminary results, we observe that, compared to naturally trained models, adversarially trained models always present very low standard accuracy and robust accuracy² on under-represented classes. For example, a naturally trained model can achieve around 40% and 60% standard accuracy on under-represented classes “frog” and “truck” separately, while an adversarially trained model gets both 0% standard & robust accuracy on these two classes. This observation suggests that adversarial training is more sensitive to imbalanced data distribution than natural training. Thus, when applying adversarial training in practice, imbalance learning strategies should always be considered for help.

As a result, we explore the potential solutions which can handle the imbalance issues for adversarial training. In this work, we focus on studying the behavior of the *reweighting* strategy [16] and leave other strategies such as resampling [12] for one future work. In Section 2.2, we apply the reweighting strategy to existing adversarial training with varied weights assigning to one under-represented class and evaluate trained models’ performance. From the results, we observe that, in adversarial training, increasing weights for an under-represented class can substantially improve the standard & robust accuracy on this class, but drastically hurt the model’s performance on the well-represented class. For example, the robust accuracy of the adversarially trained model on the under-represented class “horse” can be greatly improved when setting a relatively large weight, like 200, to its examples, but the model’s robust accuracy on the well-represented class “cat” is dropped to even lower than the class “horse” and, hence, the overall robust performance of the model is also decreased. These facts indicate that the performance of adversarially trained models is very sensitive to the reweighting manipulations and it could be very hard to figure out an eligible reweighting strategy which is optimal for all classes.

It is also worth noting that this phenomenon is absent in natural training under the same settings. In natural training, from the results in Section 2.2, we find that upweighting the under-represented class increases model’s standard accuracy on this class but only slightly hurts the accuracy on other classes, even when adopting a large weight for under-represent class. To further investigate the possible reasons leading to different behaviors of the reweighting strategy in natural and adversarial training, we visualize their learned features via t-SNE [33]. As shown in Figure 3, we observe that features learned by the adversarially trained model of different classes tend to mix together while they are well separated for the naturally trained model. This observation motivates us to theoretically show that when the given data distribution has poor data separability, upweighting under-represented classes will hurt the model’s performance on well-represented classes. Motivated by this theoretical understanding, we propose a novel algorithm *Separable Reweighted Adversarial Training (SRAT)* to facilitate the reweighting strategy in imbalanced adversarial training by enhancing the separability of learned features. Through extensive experiments, we validate the effectiveness of SRAT.

2 Preliminary Study

2.1 The Behavior of Adversarial Training

In this subsection, we conduct preliminary studies to examine the performance of PGD adversarial training [26], under an imbalanced training dataset which is resampled from CIFAR10 dataset [21]. Following previous imbalanced learning works [10, 4], we consider to construct an imbalanced

²In this work, we denote *standard accuracy* as model’s accuracy on the input samples without perturbations and *robust accuracy* as model’s accuracy on the input samples which are adversarially perturbed. Without clear clarification, we consider the perturbation is constrained by l_∞ -norm 8/255.

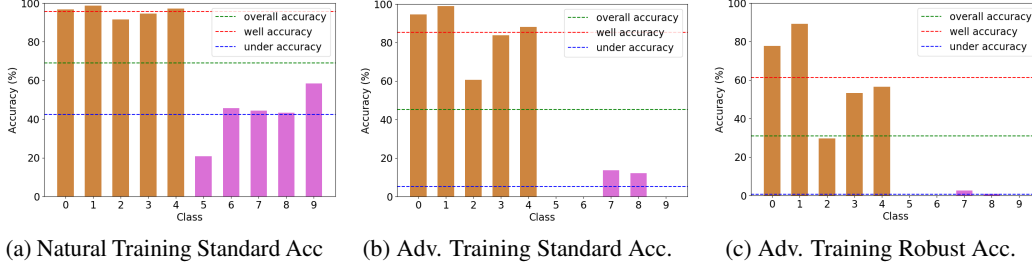


Figure 1: Class-wise performance of natural & adversarial training under an imbalanced CIFAR10.

training dataset where each of the first 5 classes (well-represented classes) has 5,000 training examples, and each of the last 5 classes (under-represented classes) has 50 training examples. Figure 1 shows the performance of naturally and adversarially trained models under a ResNet18 [17] architecture. From the figure, we can observe that, comparing with natural training, PGD adversarial training will result in a larger performance gap between well-represented classes and under-represented classes. For example, in natural training, the ratio between the average standard accuracy of well-represented classes (brown) and under-represented classes (violet) is about 2:1, while in adversarial training, this ratio expands to 16:1. Moreover, for adversarial training, although it can achieve good standard & robust accuracy on well-represented classes, it has extremely poor performance on under-represented classes. There are 3 out of the 5 under-represented classes with 0% standard & robust accuracy. As a conclusion, the performance of adversarial training is easier to be affected by imbalanced distribution than natural training and suffers more on under-represented classes. In Appendix A.1, we provide more implementation details of this experiment, as well as additional results of the same experiment under other imbalanced settings. The results in Appendix A.1 further support our findings.

2.2 The Reweighting Strategy in Natural Training v.s. in Adversarial Training

The preliminary study in Section 2.1 demonstrates that it is highly demanding to adjust the original adversarial training methods to accommodate class-imbalanced data. Therefore, in this subsection, we investigate the effectiveness of existing imbalanced learning strategies in natural training when adopted in adversarial training. In this paper, we focus on the reweighting strategy [16] as the initial effort to study this problem and leave other methods such as resampling [7] for future investigation. In this subsection, we conduct experiments under a binary classification problem, where the training dataset contains two classes that are randomly selected from CIFAR10 dataset, with each class having 5,000 and 50 training examples respectively. Under this training dataset, we arrange multiple trails of (reweighted) natural training and (reweighted) adversarial training, with the weight ratio between the under-represented class and well-represented class ranging from 1:1 to 200:1.

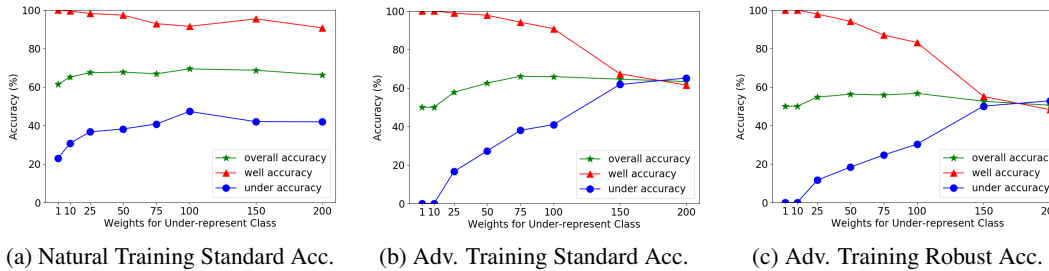


Figure 2: Class-wise performance of reweighted natural & adversarial training in binary classification.

Figure 2 shows the experimental results with data from the classes “horse” and “cat”. As demonstrated in Figure 2, increasing the weight of the under-represented class will drastically increase the model’s performance of the under-represented class, while also immensely decreasing the performance of the well-represented class. For example, when increasing the weight ratio between two classes from 1:1 to 150:1, the under-represented class’s standard accuracy can be improved from 0% to ~ 60% and its robust accuracy from 0% to ~ 50%. However, the standard & robust accuracy of the well-represented class is also drastically decreasing. For instance, the well-represented class’s standard accuracy drops

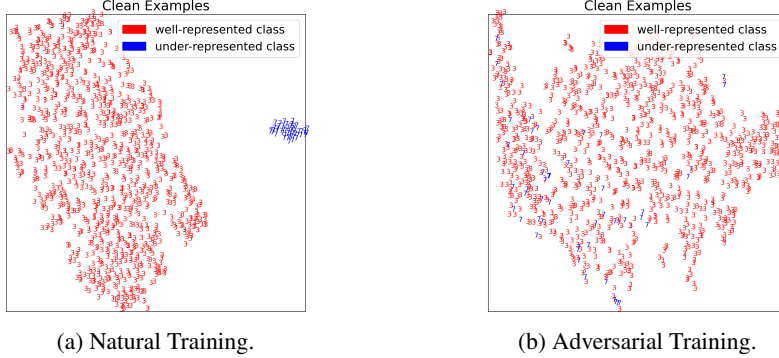


Figure 3: t-SNE visualization of penultimate layer features.

from 100% to 60%, and its robust accuracy drops from 100% to 50%. These results illustrate that adversarial training’s performance can be significantly affected by the reweighting strategy. As a result, the reweighting strategy in this setting can hardly help improve the overall performance no matter which weight ratio is chosen, because the model’s performance always presents a strong tension between these two classes. As a comparison, for the naturally trained models (Figure 2a), increasing the weights for the under-represented examples will only slightly decrease the performance on the well-represented class. More experiments using different binary imbalanced datasets are reported in Appendix A.2, where we have similar observations.

3 Theoretical Analysis

In Section 2.2, we observe that in natural training, the reweighting strategy can only make a small impact on the two classes’ performance. This phenomenon has been extensively studied by recent works [3, 37], which investigate the decision boundaries of perfectly fitted DNNs. In particular, they consider the case where the data is linearly (or nonlinearly) separable and study the behavior of linear (or nonlinear) models optimized by reweighted SGD algorithms. Interestingly, they conclude that over the process of training, these models’ decision boundaries will eventually converge to weight-agnostic solutions. For example, a linear classifier optimized by SGD on a linearly separable data will converge to the solution of the *hard-margin support vector machine* [28]. In other words, as long as the data can be well separated, reweighting will not make huge influence on the finally trained models, which is consistent with what we observed above.

Although these studies only focus on natural training, their interpretations and conclusions motivate our hypothesis in adversarial training. For adversarial training, we conjecture that it is because the models separate the data poorly, thus, their performance is highly sensitive to the reweighting strategy. As a direct validation of this hypothesis, in Figure 3, we visualize the learned (penultimate layer) features of the imbalanced training examples used in the binary classification problem in Section 2.2. We find that adversarially trained models do present obviously poorer separability on the learned features. This suggests that, compared to naturally trained models, adversarially trained models have a weaker ability to separate training data and could potentially make themselves sensitive to reweighting. Next, we will theoretically analyze the impact of reweighting on linear models which are optimized under poorly separable data. Since our empirical study shows that adversarially trained models usually poorly separate the data (see Figure 3), the analysis can hopefully shed light on the behavior of reweighting in adversarial trained models in practice.

Binary Classification Problem. To construct the theoretical study, we focus on a binary classification problem, with a Gaussian mixture distribution \mathcal{D} which is defined as:

$$y \sim \{-1, +1\}, \quad x \sim \begin{cases} \mathcal{N}(\mu, \sigma^2 I), & \text{if } y = +1 \\ \mathcal{N}(-\mu, \sigma^2 I), & \text{if } y = -1 \end{cases} \quad \text{and } \mu = \overbrace{(\eta, \dots, \eta)}^{\text{dim} = d}, \quad (1)$$

where the two classes’ centers ($\pm\mu \in \mathbb{R}^d$) with each dimension has mean value $\pm\eta$ ($\eta > 0$), variance σ^2 . Formally, we define the data *separability* as $S = \eta/\sigma^2$. Intuitively, if the separability term S is larger, it suggests that two classes have farther distance or data examples of each class are

more concentrated, so these two classes can be well separated. Previous works [3] also closely studied this term to describe data separability. Besides, we particularly define the imbalanced training dataset satisfying the condition $\Pr.(y = +1) = K \cdot \Pr.(y = -1)$ and $K > 1$ which indicates the imbalance ratio between the two classes. During test, we assume that two classes have the equal probability to appear. Under data distribution \mathcal{D} , we will discuss the performance of linear classifiers $f(x) = \text{sign}(w^T x - b)$ where w and b are the weight and bias term of model f . If a reweighting strategy is involved, we define the model will upweight the under-represented class “-1” by ρ . In the following lemma, we first derive the solution of the optimized linear classifier f training on this imbalanced dataset. Then we will extend the result of Lemma 3.1 to analyze the impact of data separability on the performance of model f .

Lemma 3.1 *Under the data distribution \mathcal{D} as defined in Eq. (1), with an imbalanced ratio K and a reweight ratio ρ , the optimal classifier which minimizes the (reweighted) empirical risk:*

$$f^* = \arg \min_f \left(\Pr.(f(x) \neq y|y = -1) \cdot \Pr.(y = -1) \cdot \rho \right. \\ \left. + \Pr.(f(x) \neq y|y = +1) \cdot \Pr.(y = +1) \right) \quad (2)$$

has the solution: $w = \mathbf{1}$ and $b = \frac{1}{2} \log(\frac{\rho}{K}) \frac{d\sigma^2}{\eta} = \frac{1}{2} \log(\frac{\rho}{K}) \frac{d}{S}$.

The proof of Lemma 3.1 can be found at Appendix A.3.1. Note that the final optimized classifier has a weight vector equal to $\mathbf{1}$ and its bias term b only depends on K , ρ and the data separability S . In the following, our first theorem is focused on one special setting when $\rho = 1$, which is the original ERM model without reweighting. Specifically, Theorem 3.1 calculates and compares the model’s performance under data distributions: \mathcal{D}_1 (with a higher separability S_1) and \mathcal{D}_2 (with a lower separability S_2). From Theorem 3.1, we aim to compare the behavior of linear models when they can poorly separate data (like adversarial trained models) or they can well separate data (like naturally trained models).

Theorem 3.1 *Under two data distributions $(x^{(1)}, y^{(1)}) \in \mathcal{D}_1$ and $(x^{(2)}, y^{(2)}) \in \mathcal{D}_2$ with the separability $S_1 > S_2$, let f_1^* and f_2^* be the optimal non-reweighted classifiers ($\rho = 1$) under \mathcal{D}_1 and \mathcal{D}_2 , respectively. Given the imbalance ratio K is large enough, we have:*

$$\Pr.(f_1^*(x^{(1)}) \neq y^{(1)}|y^{(1)} = -1) - \Pr.(f_1^*(x^{(1)}) \neq y^{(1)}|y^{(1)} = +1) \\ < \Pr.(f_2^*(x^{(2)}) \neq y^{(2)}|y^{(2)} = -1) - \Pr.(f_2^*(x^{(2)}) \neq y^{(2)}|y^{(2)} = +1). \quad (3)$$

The proof of Theorem 3.1 is provided at Appendix A.3.2. Intuitively, Theorem 3.1 suggests that when the data separability S is low (such as \mathcal{D}_2), the optimized classifier (without reweighting) can intrinsically have a larger error difference between the under-represented class “-1” and the well-represented class “+1”. Similar to the observation in Section 2.1 and Figure 3, adversarially trained models also present a weak ability to separate data, and it also presents a strong performance gap between the well-represented class and under-represented class. Conclusively, Theorem 3.1 indicates that the poor ability to separate the training data can be one important reason which leads to the strong performance gap of adversarially trained models.

Next, we consider the case when the reweighting strategy is applied. Similar to Theorem 3.1, we also calculate the models’ classwise error under \mathcal{D}_1 and \mathcal{D}_2 with different levels of separability. In particular, Theorem 3.2 focuses on the well-represented class “+1” and calculates its error increase when upweighting the under-represented class “-1” by ρ . Through the analysis in Theorem 3.2, we compare the impact of upweighting the under-represented class on the performance of well-represented class.

Theorem 3.2 *Under two data distributions $(x^{(1)}, y^{(1)}) \in \mathcal{D}_1$ and $(x^{(2)}, y^{(2)}) \in \mathcal{D}_2$ with different separability $S_1 > S_2$, let f_1^* and f_2^* be the optimal non-reweighted classifiers ($\rho = 1$) under \mathcal{D}_1 and \mathcal{D}_2 respectively, and let $f_1'^*$ and $f_2'^*$ be the optimal reweighted classifiers under \mathcal{D}_1 and \mathcal{D}_2 given the optimal reweighting ratio ($\rho = K$). Given the imbalance ratio K is large enough, we have:*

$$\Pr.(f_1'^*(x^{(1)}) \neq y^{(1)}|y^{(1)} = +1) - \Pr.(f_1^*(x^{(1)}) \neq y^{(1)}|y^{(1)} = +1) \\ < \Pr.(f_2'^*(x^{(2)}) \neq y^{(2)}|y^{(2)} = +1) - \Pr.(f_2^*(x^{(2)}) \neq y^{(2)}|y^{(2)} = +1). \quad (4)$$

The detail the proof of Theorem 3.2 at Appendix A.3.3. The theorem shows that, when the data distribution has poorer data separability, such as \mathcal{D}_2 , upweighting the under-represented class can cause greater hurt on the performance of the well-represented class. It is also consistent with our empirical findings about adversarial training models. Since the adversarially trained models poorly separate the data (Figure 3), upweighting the under-represented class always drastically decreases the performance of well-represented class (Section 2.2). Through the discussions in both Theorem 3.1 and Theorem 3.2, we can conclude that the poor separability can be one important reason which makes adversarial training and its reweighted variants extremely difficult to achieve good performance under imbalance data distribution. Therefore, in the next section, we explore potential solutions which can facilitate the reweighting strategy in adversarial training.

4 Separable Reweighted Adversarial Training (SRAT)

The observations from both preliminary studies and theoretical understandings indicate that more separable data will advance the reweighting strategy in adversarial training under imbalanced scenarios. Thus, in this section, we present a framework, Separable Reweighted Adversarial Training (SRAT), that enables the effectiveness of the reweighting strategy in adversarial training under imbalanced scenarios by increasing the separability in the learned latent feature space.

4.1 Reweighted Adversarial Training

Given an input example (\mathbf{x}, y) , adversarial training [26] aims to obtain a robust model f_θ that can make the same prediction y for an adversarial example \mathbf{x}' , generated by applying an adversarially perturbation on \mathbf{x} . The adversarially perturbations are typically bounded by a small value ϵ under L_p -norm, i.e., $\|\mathbf{x}' - \mathbf{x}\|_p \leq \epsilon$. More formally, adversarial training can be formulated as solving a min-max optimization problem, where a DNN model is trained on minimizing the prediction error on adversarial examples generated by iteratively maximizing some loss function.

As indicated in Section 2.1, adversarial training cannot be applied in imbalanced scenarios directly, as it presents very low performance on under-represented classes. To tackle this problem, a natural idea is to integrate existing imbalanced learning strategies proposed in natural training, such as reweighting, into adversarial training to improve the trained model’s performance on those under-represented classes. Hence, the reweighted adversarial training can be defined as

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \max_{\|\mathbf{x}'_i - \mathbf{x}_i\|_p \leq \epsilon} w_i \mathcal{L}(f_\theta(\mathbf{x}'_i), y_i), \quad (5)$$

where w_i is a reweighting value assigned for each input sample (\mathbf{x}_i, y_i) based on the example size of the class (\mathbf{x}_i, y_i) belongs to or some properties of (\mathbf{x}_i, y_i) . In most existing adversarial training methods [26, 41, 35], the cross entropy (CE) loss is adopted as the loss function $\mathcal{L}(\cdot, \cdot)$. However, the CE loss could be suboptimal in imbalanced settings and some new loss functions designed for imbalanced settings specifically, such as Focal loss [24] and LDAM loss [4], have been prove superiority in natural training. Hence, besides CE loss, Focal loss and LDAM loss can also be adopted as the loss function $\mathcal{L}(\cdot, \cdot)$ in Eq. (5).

4.2 Increasing Feature Separability

Our preliminary study indicates that only reweighted adversarial training cannot work well under imbalanced scenarios. Moreover, the reweighting strategy behaves very differently between natural training and adversarial training. Meanwhile, our theoretical analysis suggests that the poor separability of the feature space produced by the adversarially trained model can be one reason to understand these observations. Hence, in order to facilitate the reweighting strategy in adversarial training under imbalanced scenarios, we equip a feature separation loss with our SRAT method. We aim to enforce the learned feature space as separable as possible. More specifically, the goal of the feature separation loss is to make (1) the learned features of examples from the same class well clustered, and (2) the features of examples from different classes well separated. By achieving this goal, the model is able to learn more discriminative features for each class. Correspondingly adjusting the decision boundary via the reweighting strategy to fit under-represented classes’ examples more will not hurt

well-represented classes drastically. The feature separation loss is formally defined as:

$$\mathcal{L}_{sep}(\mathbf{x}'_i) = -\frac{1}{|P(i)|} \sum_{p \in P(i)} \log \frac{\exp(\mathbf{z}'_i \cdot \mathbf{z}'_p / \tau)}{\sum_{a \in A(i)} \exp(\mathbf{z}'_i \cdot \mathbf{z}'_a / \tau)}, \quad (6)$$

where \mathbf{z}'_i is the feature representation of the adversarial example \mathbf{x}'_i of \mathbf{x}_i , $\tau \in \mathcal{R}^+$ is a scalar temperature parameter, $P(i)$ denotes the set of input examples belonging to the same class with \mathbf{x}_i and $A(i)$ indicates the set of all input examples excepts \mathbf{x}'_i . When minimizing the feature separation loss during training, the learned features of examples from the same class will tend to aggregate together in the latent feature space, and, hence, result in a more separable latent feature space. Our proposed feature separation loss $\mathcal{L}_{sep}(\cdot)$ is inspired by the supervised contrastive loss proposed in [20]. The main difference is, instead of applying data augmentation techniques to generate two different views of each data example and feeding the model with augmented data examples, our feature separation loss directly takes the adversarial example \mathbf{x}'_i of each data example \mathbf{x}_i as input.

4.3 Training Schedule

By combining the feature separation loss with the reweighted adversarial training, the final object function for Separable reweighted Adversarial Training (SRAT) can be defined as:

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \max_{\|\mathbf{x}'_i - \mathbf{x}_i\|_p \leq \epsilon} w_i \mathcal{L}(f_{\theta}(\mathbf{x}'_i), y_i) + \lambda \mathcal{L}_{sep}(\mathbf{x}'_i), \quad (7)$$

where we use a hyper-parameter λ to balance the contributions from the reweighted adversarial training and the feature separation loss.

In practice, in order to better take advantage of the reweighting strategy in our SRAT method, we adopt a deferred reweighting training schedule [4]. Specifically, before annealing the learning rate, our SRAT method first trains a model guided by Eq. (7) without introducing the reweighting strategy, i.e., setting $w_i = 1$ for every input example \mathbf{x}'_i , and then applies reweighting into model training process with a smaller learning rate. Our SRAT method enables to learn more separable feature space, thus comparing with applying the reweighting strategy from the beginning of training, this deferred re-balancing training schedule enables the reweighting strategy to obtain more benefits from our SRAT method, and as a result, it can boost the performance of our SRAT method with the help of the reweighting strategy. The detailed training algorithm for SRAT is shown in Appendix A.4.

5 Experiment

In this section, we perform comprehensive experiments to validate the effectiveness of our proposed SRAT method. We first compare our method with several representative imbalanced learning methods in adversarial training under various imbalanced scenarios and then conduct ablation study to understand our method more deeply.

5.1 Experimental Settings

Datasets. We conduct experiments on multiple imbalanced training datasets artificially created from two benchmark image datasets CIFAR10 [21] and SVHN [27] with diverse imbalanced distributions. Specifically, we consider two types of imbalance types: Exponential (Exp) imbalance [10] and Step imbalance [2]. For Exp imbalance, the number of training examples of each class will be reduced according to an exponential function $n = n_i \tau^i$, where i is the class index, n_i is the number of training examples in the original CIFAR10/SVHN training dataset for class i and $\tau \in (0, 1)$. We categorize five most frequent classes in the constructed imbalanced training dataset as well-represented classes and the remaining five classes as under-represented classes. For Step imbalance, we follow the same process adopted in Section 2.1 to construct imbalanced training datasets based on CIFAR10 and SVHN, separately. Moreover, in both imbalanced types, we denote *imbalance ratio* K as the ratio between training example sizes of the most frequent and least frequent class. In our experiments, we construct four different imbalanced datasets, named as ‘‘Step-100’’, ‘‘Step-10’’, ‘‘Exp-100’’ and ‘‘Exp-10’’, by adopting different imbalanced types (Step or Exp) with different imbalanced ratios ($K = 100$ or $K = 10$) to train models, and evaluate model’s performance on the original uniformly

distributed test datasets of CIFAR10 and SVHN correspondingly. More detailed information about imbalanced training sets used in our experiments can be found in Appendix A.5.

Baseline methods. We implement several representative and state-of-the-art imbalanced learning methods (or their combinations) into adversarial training as baseline methods. These methods include: (1) Focal loss (Focal); (2) LDAM loss (LDAM); (3) Class-balanced reweighting (CB-Reweight) [10], where each example is reweighted proportionally by the inverse of the effective number³ of its class; (4) Class-balanced Focal loss (CB-Focal) [10], a combination of Class-balanced method [10] and Focal loss [24], where well-classified examples will be down-weighted while hard-classified examples will be up-weighted controlled by their corresponding effective numbers; (5) deferred reweighted CE loss (DRCB-CE), where a deferred reweighting training schedule is applied based on the CE loss; (6) deferred reweighted Class-balanced Focal loss (DRCB-Focal), where a deferred reweighting training schedule is applied based on the CB-Focal loss; (7) deferred reweighted Class-balanced LDAM loss (DRCB-LDAM) [4], where a deferred reweighting training schedule is applied based on the CB-LDAM loss. In addition, we also include the original PGD adversarial training method using cross entropy loss (CE) in our experiments.

Our proposed methods. We evaluate three variants of our proposed SRAT method with different implementations of the prediction loss $\mathcal{L}(\cdot, \cdot)$ in Eq. (5), i.e., CE loss, Focal loss and LDAM loss. The variant utilizing CE loss is denoted as SRAT-CE, and, similarly, other two variants are denoted as SRAT-Focal and SRAT-LDAM, respectively. For all these three variants, Class-balanced method [10] is adopted to set reweighting values within the deferred reweighting training schedule.

Implementation details. We implement all methods used in our experiments based on a Pytorch library DeepRobust [23]. For CIFAR10 based imbalanced datasets, the adversarial examples used in training are calculated by PGD-10, with a perturbation budget $\epsilon = 8/255$, and step size $\gamma = 2/255$. For robustness evaluation, we report robust accuracy under l_∞ -norm $8/255$ attacks generated by PGD-20 on Resnet-18 [17] models. For SVHN based imbalanced datasets, the setting is similar with CIFAR10 based datasets, excepts we set step size γ to $1/255$ in both training and test phases, as suggested in [36]. For the deferred reweighting training schedule used in our methods and some baseline methods, we set the number of the training epochs to 200 and the initial learning rate to 0.1, and then decay the learning rate at epoch 160 and 180 with the ratio 0.01. The reweighting strategy will be applied starting from epoch 160.

5.2 Performance Comparison

Table 1 and 2 show the performance comparison on various imbalanced CIFAR10 datasets with different imbalance types and imbalance ratios. In these two tables, we use bold values to denote the highest accuracy among all methods and use the underline values to indicate our SRAT variants which achieve the highest accuracy among their corresponding baseline methods utilizing the same loss function for making predictions. Due to the limited space, we report the performance comparison on SVHN based imbalanced datasets in Appendix A.6.

From Table 1 and Table 2, we can make the following observations. First, compared to baseline methods, our SRAT methods can obtain improved performance in terms of both overall standard & robust accuracy under almost all imbalanced settings. More importantly, our SRT methods make significantly improvement on those under-represented classes, especially under the extremely imbalanced setting. For example, on the Step imbalanced dataset with imbalance ratio $K = 100$, our SRAT-Focal method improves the standard accuracy on under-represented classes from 21.81% achieved by the best baseline method utilizing Focal loss to 51.83% and robust accuracy from 3.24% to 15.89%. These results demonstrate that our proposed SRAT method is able to obtain more robustness under imbalanced settings. Second, the performance gap among three variants SRAT-CE, SRAT-Focal and SRAT-LDAM are mainly caused by the gap between the loss functions in these methods. As shown in Table 1 and 2, DRCB-LDAM typically performs better than DRCE-CE and DRCB-Focal, and similarly, SRAT-LDAM outperforms SRAT-CE and SRAT-Focal under corresponding imbalanced settings.

³The effective number is defined as the volume of examples and can be calculated by $(1 - \beta^{n_i}) / (1 - \beta)$, where $\beta \in [0, 1)$ is a hyperparameter and n_i denotes the number of examples of class i .

Table 1: Performance comparison on imbalanced CIFAR10 datasets (Imbalanced Type: Step)

Imbalance Ratio	10				100			
	Standard Accuracy		Robust Accuracy		Standard Accuracy		Robust Accuracy	
Method	Overall	Under	Overall	Under	Overall	Under	Overall	Under
CE	63.26	40.62	36.96	14.23	47.29	9.03	30.39	1.62
Focal	63.57	41.17	36.89	14.25	47.36	9.03	30.12	1.45
LDAM	57.08	31.09	37.18	12.44	42.49	0.85	30.80	0.05
CB-Reweight	73.30	74.80	41.34	42.15	37.68	19.64	25.58	10.33
CB-Focal	73.47	73.69	41.19	41.02	15.44	0.00	14.46	0.00
DRCB-CE	75.89	70.55	39.93	33.33	53.40	22.86	28.31	3.35
DRCB-Focal	74.61	67.06	37.91	29.50	52.75	21.81	27.78	3.24
DRCB-LDAM	72.95	75.42	45.23	44.98	61.60	50.69	31.37	16.25
SRAT-CE	76.32	73.20	41.71	37.86	59.10	40.24	30.02	11.72
SRAT-Focal	75.41	74.91	42.05	41.28	62.93	51.83	28.38	15.89
SRAT-LDAM	73.99	76.63	45.60	45.96	63.13	52.73	33.51	18.89

Table 2: Performance comparison on imbalanced CIFAR10 datasets (Imbalanced Type: Exp).

Imbalance Ratio	10				100			
	Standard Accuracy		Robust Accuracy		Standard Accuracy		Robust Accuracy	
Method	Overall	Under	Overall	Under	Overall	Under	Overall	Under
CE	71.95	64.09	37.94	26.79	48.40	23.04	26.94	6.17
Focal	72.06	63.99	37.62	26.27	49.16	23.69	26.84	5.88
LDAM	67.39	58.01	41.35	28.65	48.39	25.69	29.51	8.95
CB-Reweight	75.17	76.87	41.02	41.67	57.49	56.47	29.01	26.53
CB-Focal	74.73	76.67	38.86	42.41	50.35	60.05	27.15	33.56
DRCB-CE	76.25	75.83	40.02	37.93	57.30	37.90	26.97	10.57
DRCB-Focal	75.36	72.72	37.76	33.83	54.76	31.79	25.24	7.81
DRCB-LDAM	73.92	78.53	46.29	48.81	62.65	57.19	31.66	22.11
SRAT-CE	76.94	79.50	41.50	43.08	64.93	64.34	29.68	25.42
SRAT-Focal	75.26	80.52	42.37	47.22	62.57	64.88	30.34	28.66
SRAT-LDAM	74.63	79.82	46.72	50.38	63.11	65.60	34.22	32.55

5.3 Ablation Study

In this subsection, we provide ablation study to understand our SRAT method more comprehensively.

Feature space visualization. In order to facilitate the reweighting strategy in adversarial training under the imbalanced setting, we present a feature separation loss in our SRAT method. The main goal of the feature separation loss is to enforce the learned feature space as much separated as possible. For checking whether the feature separation loss can work as expected, we apply t-SNE [33] to visualize the latent feature space learned by our SRAT-LDAM method in Figure 4. As a comparison, we also provide the visualization of feature space learned by the original PGD adversarial training method (CE) and DRCB-LDAM method.

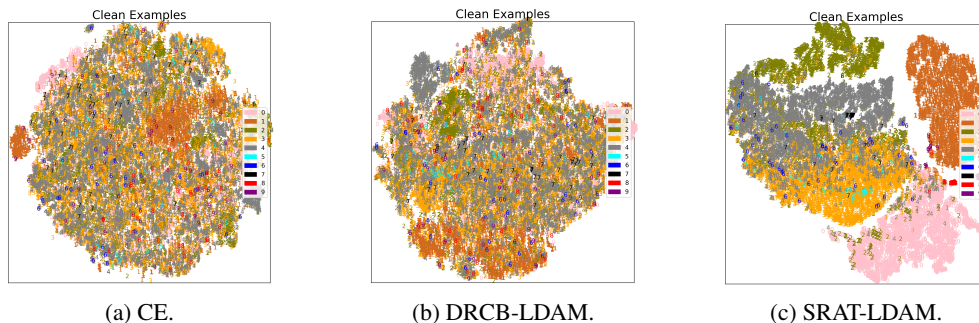


Figure 4: t-SNE feature visualization of training examples learned by SRAT and two baseline methods using imbalanced training datasets “Step-100”.

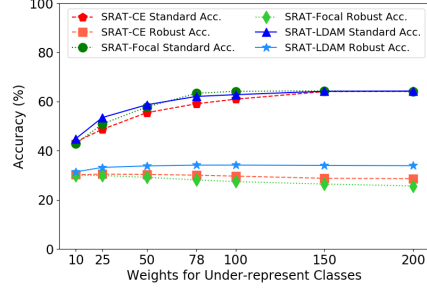


Figure 5: The impact of reweighting values using an imbalanced training dataset “Step-100”.

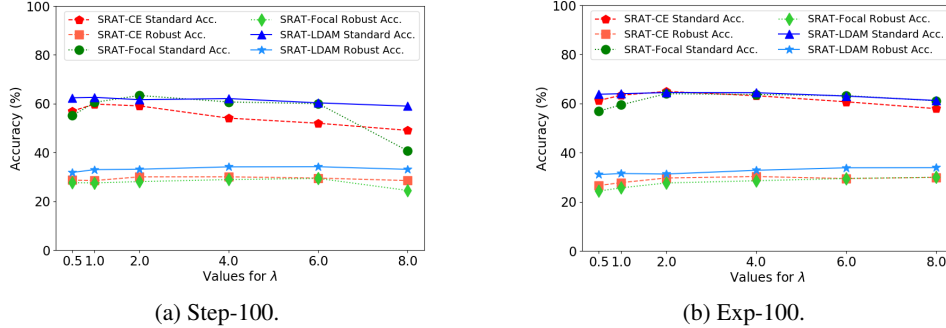


Figure 6: The impact of the hyper-parameter λ using imbalanced training datasets “Step-100” and “Exp-100”.

As shown in Figure 4, the feature space learned by our SRAT-LDAM method is more separable than two baseline methods. This observation demonstrates that, with our proposed feature separation loss, the adversarially trained model is able to learn much better features and thus our SRAT method can achieve superiority performance.

Impact of reweighting values. As in all SRAT variants, we adopt the Class-balanced method [10] to assign different weights to different classes based on their effective number. To explore how the assigned weights impact the performance of our proposed SRAT method, we conduct experiments on a Step-imbalanced CIFAR10 dataset with imbalance ratio $K = 100$ to see the change of model’s performance using different reweighting values. In our experiments, we assign five well-represented classes with weight 1 and change the weight for remaining five under-represented classes from 10 to 200. The experimental results are shown in Figure 5. Here, we use an approximation integer 78 to denote the weight calculated by the Class-balanced method when the imbalance ratio equals 100.

From Figure 5, we can observe that, for all SRAT variants, the model’s standard accuracy is increased with the increase of the weights assigning to under-represented classes. However, the robust accuracy for these three methods do not synchronize with the change of their standard accuracy. When increasing the weights for under-represented classes, robust accuracy of SRAT-LDAM is almost unchanged and robust accuracy of SRAT-CE and SRAT-Focal even has slight decrease. As a trade-off, using a relative large weights, such as 78 or 100, in our SRAT method can obtain satisfactory performance on both standard & robust accuracy, where the former is calculated by the Class-balanced method and the latter equals the imbalance ratio K .

Impact of hyper-parameter λ . In our proposed SRAT method, the contributions of feature separation loss and prediction loss are controlled by a hyper-parameter γ . In this part, we study how this hyper-parameter affects the performance of our SRAT method. In our experiments, we evaluate the models’ performance of all SRAT variants with different values of λ used in training process on both Step-imbalanced CIFAR10 dataset and Exp-imbalanced CIFAR10 dataset with imbalance ratio $K = 100$.

As shown in Figure 6, the performance of all SRAT variants are not very sensitive with the choice of λ . However, a large value of λ , such as 8, may hurt the model’s performance.

6 Related Work

Adversarial Robustness. The vulnerability of DNN models to adversarial examples has been verified by many existing successful attack methods [14, 5, 26]. To improve model robustness against adversarial attacks, various defense methods have been proposed [14, 26, 29, 9]. Among them, adversarial training has been proven to be one of the most effective defense methods [1]. Adversarial training can be formulated as solving a min-max optimization problem where the outer minimization process enforces the model to be robust to adversarial examples, generated by the inner maximization process via some existing attacking methods like PGD [26]. Based on adversarial training, several variants, such as TRADES [41], MART [35] and FAT [42], have been presented to improve the model’s performance further. More details about adversarial robustness can be found in recent surveys [6, 39]. Since almost all studies of adversarial training are focused on balanced datasets, it’s worthwhile to investigate the performance of adversarial training methods on imbalanced training datasets.

Imbalanced Learning. Most existing works of imbalanced training can be roughly classified into two categories, i.e., re-sampling and reweighting. *Re-sampling* methods aim to reduce the level of imbalance through either over-sampling data examples from under-represented classes [2, 3] or under-sampling data examples from well-represented classes [18, 11, 15, 40]. *reweighting* methods allocate different weights for different classes or even different data examples. For example, Focal loss [24] enlarges the weights of wrongly-classified examples while reducing the weights of well-classified examples in the standard cross entropy loss; and LDAM loss [4] regularizes the under-represented classes more strongly than the over-represented classes to attain good generalization performance on under-represented classes. More information about imbalanced learning can be found in recent surveys [16, 19]. The majority of existing methods focused on the nature training scenario and their trained models will be crashed when facing adversarial attacks [32, 14]. Hence, in this paper, we develop a novel method that can defend adversarial attacks and achieve well-pleasing performance under the imbalance setting.

7 Conclusion

In this work, we first empirically investigate the behavior of adversarial training under imbalanced settings and explore the potential solutions to assist adversarial training in tackling the imbalanced issues. As neither adversarial training method itself nor adversarial training with reweighting strategy can work well under imbalanced scenarios, we further theoretically verify that the poor data separability is one key reason causing the failure of adversarial training based methods under imbalanced scenarios. Based on our findings, we propose the Separable Reweighted Adversarial Training (SRAT) framework to facilitate the reweighting strategy in imbalanced adversarial training by enhancing the separability of learned features. Through extensive experiments, we validate the effectiveness of SRAT. In the future, we plan to examine how other types of defense methods perform under imbalanced scenarios and how other types of balanced learning strategies in natural training behavior under adversarial training.

References

- [1] Anish Athalye, Nicholas Carlini, and David Wagner. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In *International Conference on Machine Learning*, pages 274–283. PMLR, 2018.
- [2] Mateusz Buda, Atsuto Maki, and Maciej A Mazurowski. A systematic study of the class imbalance problem in convolutional neural networks. *Neural Networks*, 106:249–259, 2018.
- [3] Jonathon Byrd and Zachary Lipton. What is the effect of importance weighting in deep learning? In *International Conference on Machine Learning*, pages 872–881. PMLR, 2019.
- [4] Kaidi Cao, Colin Wei, Adrien Gaidon, Nikos Archiga, and Tengyu Ma. Learning imbalanced datasets with label-distribution-aware margin loss. *arXiv preprint arXiv:1906.07413*, 2019.
- [5] Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy (SP)*, pages 39–57. IEEE, 2017.

- [6] Anirban Chakraborty, Manaar Alam, Vishal Dey, Anupam Chattopadhyay, and Debdeep Mukhopadhyay. Adversarial attacks and defences: A survey. *arXiv preprint arXiv:1810.00069*, 2018.
- [7] Nitesh V Chawla, Kevin W Bowyer, Lawrence O Hall, and W Philip Kegelmeyer. Smote: synthetic minority over-sampling technique. *Journal of artificial intelligence research*, 16:321–357, 2002.
- [8] Chenyi Chen, Ari Seff, Alain Kornhauser, and Jianxiong Xiao. Deepdriving: Learning affordance for direct perception in autonomous driving. In *Proceedings of the IEEE international conference on computer vision*, pages 2722–2730, 2015.
- [9] Jeremy Cohen, Elan Rosenfeld, and Zico Kolter. Certified adversarial robustness via randomized smoothing. In *International Conference on Machine Learning*, pages 1310–1320. PMLR, 2019.
- [10] Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie. Class-balanced loss based on effective number of samples. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9268–9277, 2019.
- [11] Chris Drummond, Robert C Holte, et al. C4. 5, class imbalance, and cost sensitivity: why under-sampling beats over-sampling. In *Workshop on learning from imbalanced datasets II*, volume 11, pages 1–8. Citeseer, 2003.
- [12] Andrew Estabrooks, Taeho Jo, and Nathalie Japkowicz. A multiple resampling method for learning from imbalanced data sets. *Computational intelligence*, 20(1):18–36, 2004.
- [13] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88(2):303–338, 2010.
- [14] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- [15] Haibo He and Eduardo A Garcia. Learning from imbalanced data. *IEEE Transactions on knowledge and data engineering*, 21(9):1263–1284, 2009.
- [16] Haibo He and Yunqian Ma. Imbalanced learning: foundations, algorithms, and applications. 2013.
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [18] Nathalie Japkowicz and Shaju Stephen. The class imbalance problem: A systematic study. *Intelligent data analysis*, 6(5):429–449, 2002.
- [19] Justin M Johnson and Taghi M Khoshgoftaar. Survey on deep learning with class imbalance. *Journal of Big Data*, 6(1):1–54, 2019.
- [20] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschiot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *arXiv preprint arXiv:2004.11362*, 2020.
- [21] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [22] Alexey Kurakin, Ian J Goodfellow, and Samy Bengio. Adversarial examples in the physical world. *arXiv preprint arXiv:1607.02533*, 2016.
- [23] Yaxin Li, Wei Jin, Han Xu, and Jiliang Tang. Deeprobust: A pytorch library for adversarial attacks and defenses. *arXiv preprint arXiv:2005.06149*, 2020.
- [24] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017.
- [25] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.
- [26] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.

- [27] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011.
- [28] William S Noble. What is a support vector machine? *Nature biotechnology*, 24(12):1565–1567, 2006.
- [29] Aditi Raghunathan, Jacob Steinhardt, and Percy Liang. Certified defenses against adversarial examples. *arXiv preprint arXiv:1801.09344*, 2018.
- [30] Leslie Rice, Eric Wong, and Zico Kolter. Overfitting in adversarially robust deep learning. In *International Conference on Machine Learning*, pages 8093–8104. PMLR, 2020.
- [31] Ludwig Schmidt, Shibani Santurkar, Dimitris Tsipras, Kunal Talwar, and Aleksander Mądry. Adversarially robust generalization requires more data. *arXiv preprint arXiv:1804.11285*, 2018.
- [32] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- [33] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008.
- [34] Grant Van Horn and Pietro Perona. The devil is in the tails: Fine-grained classification in the wild. *arXiv preprint arXiv:1709.01450*, 2017.
- [35] Yisen Wang, Difan Zou, Jinfeng Yi, James Bailey, Xingjun Ma, and Quanquan Gu. Improving adversarial robustness requires revisiting misclassified examples. In *International Conference on Learning Representations*, 2019.
- [36] Dongxian Wu, Shu-Tao Xia, and Yisen Wang. Adversarial weight perturbation helps robust generalization. *Advances in Neural Information Processing Systems*, 33, 2020.
- [37] Da Xu, Yuting Ye, and Chuanwei Ruan. Understanding the role of importance weighting for deep learning. *arXiv preprint arXiv:2103.15209*, 2021.
- [38] Han Xu, Xiaorui Liu, Yaxin Li, and Jiliang Tang. To be robust or to be fair: Towards fairness in adversarial training. *arXiv preprint arXiv:2010.06121*, 2020.
- [39] Han Xu, Yao Ma, Hao-Chen Liu, Debayan Deb, Hui Liu, Ji-Liang Tang, and Anil K Jain. Adversarial attacks and defenses in images, graphs and text: A review. *International Journal of Automation and Computing*, 17(2):151–178, 2020.
- [40] Show-Jane Yen and Yue-Shi Lee. Cluster-based under-sampling approaches for imbalanced data distributions. *Expert Systems with Applications*, 36(3):5718–5727, 2009.
- [41] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric Xing, Laurent El Ghaoui, and Michael Jordan. Theoretically principled trade-off between robustness and accuracy. In *International Conference on Machine Learning*, pages 7472–7482. PMLR, 2019.
- [42] Jingfeng Zhang, Xilie Xu, Bo Han, Gang Niu, Lizhen Cui, Masashi Sugiyama, and Mohan Kankanhalli. Attacks which do not kill training make adversarial learning stronger. In *International Conference on Machine Learning*, pages 11278–11287. PMLR, 2020.

A Appendix

A.1 The Behavior of Adversarial Training

In order to examine the performance of PGD adversarial training under imbalanced scenarios, we adversarially train ResNet18 [17] models on multiple imbalanced training datasets based on CIFAR10 dataset [21]. Similar with observations we discussed in Section 2.1, as shown in Figure 7, Figure 8 and Figure 9, adversarial training produces larger performance gap between well-represented classes and under-represented classes than natural training. Especially, in all imbalanced scenarios, adversarially trained models obtain very low robust accuracy on under-represented classes, which proves again that adversarial training cannot be applied in practical imbalanced scenarios directly.

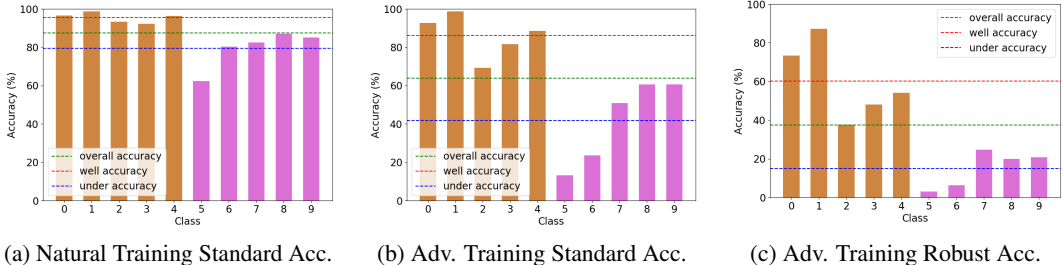


Figure 7: Class-wise performance of natural & adversarial training under an imbalanced CIFAR10 dataset “Step-10”.

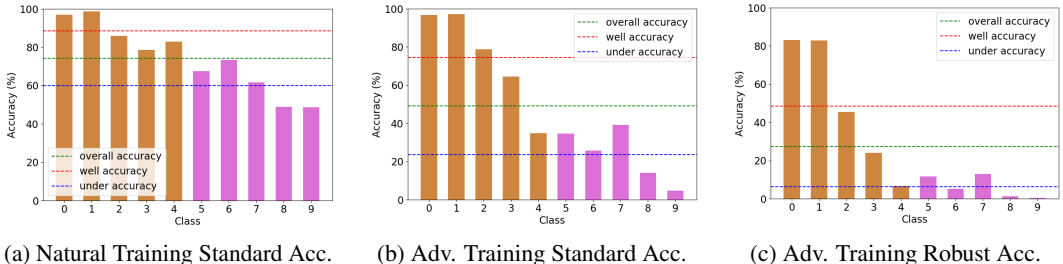


Figure 8: Class-wise performance of natural & adversarial training under an imbalanced CIFAR10 dataset “Exp-100”.

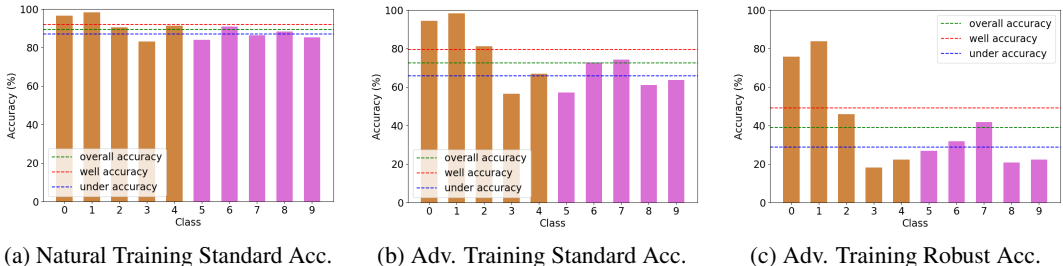


Figure 9: Class-wise performance of natural & adversarial training under an imbalanced CIFAR10 dataset “Exp-10”.

A.2 Reweighting Strategy in Natural Training v.s. in Adversarial Training

For exploring whether the reweighting strategy can help adversarial training deal with imbalanced issues, we evaluate performance of adversarial trained models using diverse binary imbalanced training datasets with different weights assigning to under-represented class. As shown in Figure 10, Figure 11, Figure 12, for adversarially trained models, increasing the weights assigning to under-represented class will improve models’ performance on under-represented class. However, as

the same time, the models' performance on well-represented class will be drastically decreased. As a comparison, adopting larger weights in naturally trained models will also improve models' performance on under-represented class but only result in slight drop in performance on well-represented class. In other words, the reweighting strategy proposed in natural training to handle imbalanced problem may only provide limited help in adversarial training, and, hence, new techniques are needed for adversarial training under imbalanced scenarios.

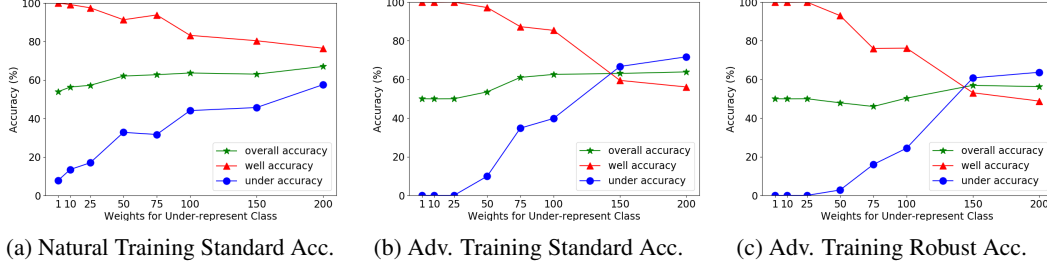


Figure 10: Class-wise performance of reweighted natural & adversarial training in binary classification. (“auto” as well-represented class and “truck” as under-represented class).

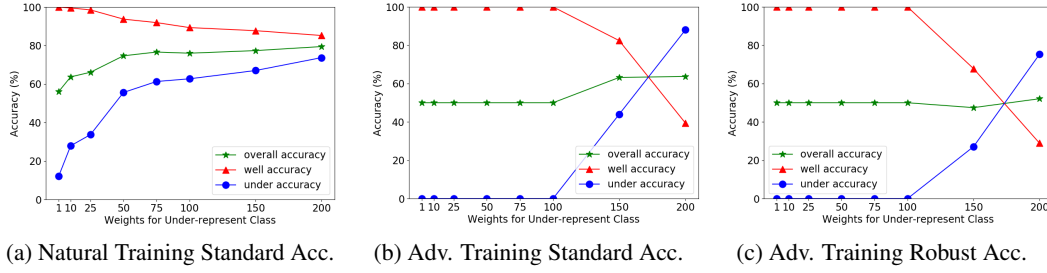


Figure 11: Class-wise performance of reweighted natural & adversarial training in binary classification. (“bird” as well-represented class and “frog” as under-represented class).

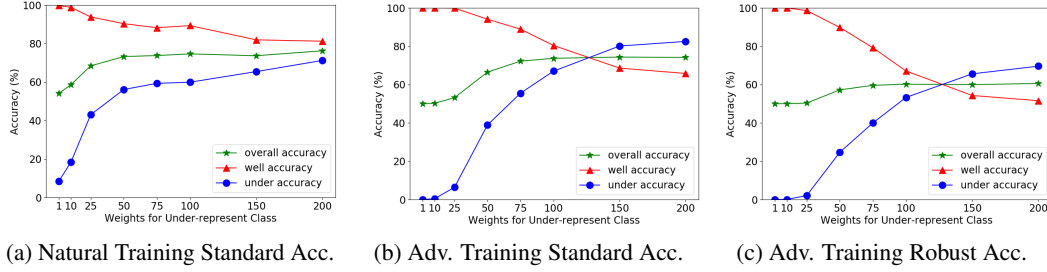


Figure 12: Class-wise performance of reweighted natural & adversarial training in binary classification. (“dog” as well-represented class and “deer” as under-represented class).

A.3 Proofs of the Theorems in Section 3

A.3.1 Proof of Lemma 3.1

Lemma 3.1 Under the data distribution \mathcal{D} as defined in Eq. (1), with an imbalanced ratio K and a reweight ratio ρ , the optimal classifier which minimizes the (reweighted) empirical risk:

$$f^* = \arg \min_f \left(\Pr(f(x) \neq y | y = -1) \cdot \Pr(y = -1) \cdot \rho + \Pr(f(x) \neq y | y = +1) \cdot \Pr(y = +1) \right) \quad (2)$$

has the solution: $w = \mathbf{1}$ and $b = \frac{1}{2} \log\left(\frac{\rho}{K}\right) \frac{d\sigma^2}{\eta} = \frac{1}{2} \log\left(\frac{\rho}{K}\right) \frac{d}{S}$.

Proof 1 (Proof of Lemma 3.1) We will first prove that the optimal model f^* has parameters $w_1 = w_2 = \dots = w_d$ (or $w = \mathbf{1}$) by contradiction. We define $G = \{1, 2, \dots, d\}$ and make the following assumption: for the optimal w and b , we assume if there exist $w_i < w_j$ for $i \neq j$ and $i, j \in G$. Then we obtain the following standard errors for two classes of this classifier f with weight w :

$$\begin{aligned} \Pr(f^*(x) \neq y|y = -1) &= \Pr\left\{ \sum_{k \neq i, k \neq j} w_k \mathcal{N}(-\eta, \sigma^2) - b + w_i \mathcal{N}(-\eta, \sigma^2) + w_j \mathcal{N}(-\eta, \sigma^2) > 0 \right\}, \\ \Pr(f^*(x) \neq y|y = +1) &= \Pr\left\{ \sum_{k \neq i, k \neq j} w_k \mathcal{N}(+\eta, \sigma^2) - b + w_i \mathcal{N}(+\eta, \sigma^2) + w_j \mathcal{N}(+\eta, \sigma^2) < 0 \right\}. \end{aligned} \quad (8)$$

However, if we define a new classifier \tilde{f} whose weight \tilde{w} uses w_j to replace w_i , we obtain the errors for the new classifier:

$$\begin{aligned} \Pr(\tilde{f}(x) \neq y|y = -1) &= \Pr\left\{ \sum_{k \neq i, k \neq j} w_k \mathcal{N}(-\eta, \sigma^2) - b + w_j \mathcal{N}(-\eta, \sigma^2) + w_j \mathcal{N}(-\eta, \sigma^2) > 0 \right\}, \\ \Pr(\tilde{f}(x) \neq y|y = +1) &= \Pr\left\{ \sum_{k \neq i, k \neq j} w_k \mathcal{N}(+\eta, \sigma^2) - b + w_j \mathcal{N}(+\eta, \sigma^2) + w_j \mathcal{N}(+\eta, \sigma^2) < 0 \right\}. \end{aligned} \quad (9)$$

By comparing the errors in Eq. (8) and Eq. (9), it can imply the classifier \tilde{f} has smaller error in each class. Therefore, it contradicts with the assumption that f is the optimal classifier with smallest error. Thus, we conclude for an optimal linear classifier in natural training, it must satisfies $w_1 = w_2 = \dots = w_d$ (or $w = \mathbf{1}$) if we do not consider the scale of w .

Next, we calculate the optimal bias term b given $w = \mathbf{1}$, where we find an optimal b can minimize the (reweighted) empirical risk:

$$\begin{aligned} & \text{Error}_{\text{train}}(f^*) \\ &= \Pr(f^*(x) \neq y|y = -1) \cdot \Pr(y = -1) \cdot \rho + \Pr(f^*(x) \neq y|y = +1) \cdot \Pr(y = +1) \\ &\propto \Pr(f^*(x) \neq y|y = -1) \cdot \rho + \Pr(f^*(x) \neq y|y = +1) \cdot K \\ &= \rho \cdot \Pr\left(\sum_{i=1}^d \mathcal{N}(-\eta, \sigma^2) - b > 0\right) + K \cdot \Pr\left(\sum_{i=1}^d \mathcal{N}(\eta, \sigma^2) - b < 0\right) \\ &= \rho \cdot \Pr(\mathcal{N}(0, 1) < -\frac{b + d\eta}{d\sigma}) + K \cdot \Pr(\mathcal{N}(0, 1) < \frac{b - d\eta}{d\sigma}), \end{aligned}$$

and we take the derivative with respect to b :

$$\frac{\partial \text{Error}_{\text{train}}}{\partial b} = \frac{\rho}{\sqrt{2\pi}} \cdot \left(-\frac{1}{d\sigma}\right) \exp\left(-\frac{1}{2}\left(-\frac{b + d\eta}{d\sigma}\right)^2\right) + \frac{K}{\sqrt{2\pi}} \cdot \left(\frac{1}{d\sigma}\right) \exp\left(-\frac{1}{2}\left(\frac{b - d\eta}{d\sigma}\right)^2\right).$$

When $\partial \text{Error}_{\text{train}} / \partial b = 0$, we can calculate the optimal b which gives the minimum value of the empirical error, and we have:

$$b = \frac{1}{2} \log\left(\frac{\rho}{K}\right) \frac{d\sigma^2}{\eta} = \frac{1}{2} \log\left(\frac{\rho}{K}\right) \frac{d}{S}.$$

A.3.2 Proof of Theorem 3.1

Theorem 3.1 Under two data distributions $(x^{(1)}, y^{(1)}) \in \mathcal{D}_1$ and $(x^{(2)}, y^{(2)}) \in \mathcal{D}_2$ with the separability $S_1 > S_2$, let f_1^* and f_2^* be the optimal non-reweighted classifiers ($\rho = 1$) under \mathcal{D}_1 and \mathcal{D}_2 , respectively. Given the imbalance ratio K is large enough, we have:

$$\begin{aligned} & \Pr(f_1^*(x^{(1)}) \neq y^{(1)}|y^{(1)} = -1) - \Pr(f_1^*(x^{(1)}) \neq y^{(1)}|y^{(1)} = +1) \\ &< \Pr(f_2^*(x^{(2)}) \neq y^{(2)}|y^{(2)} = -1) - \Pr(f_2^*(x^{(2)}) \neq y^{(2)}|y^{(2)} = +1). \end{aligned} \quad (3)$$

Proof 2 (Proof of Theorem 3.1) Without loss of generality, for distribution $\mathcal{D}_1, \mathcal{D}_2$ with different mean-variance pairs $(\pm\eta_1, \sigma_1^2)$ and $(\pm\eta_2, \sigma_2^2)$, we can only consider the case $\eta_1 = \eta_2$ and $\sigma_1^2 < \sigma_2^2$. Otherwise, we can simply rescale one of them to match the mean vector of the other and will not impact the results. Under this definition, the optimal classifier f_1^* and f_2^* has weight vector $w_1 = w_2 = \mathbf{1}$ and bias term b_1, b_2 , with the value as demonstrated in Lemma 3.1. Next, we will prove the Theorem 3.1 by 2 steps.

Step 1. For the error of class “-1”, we have:

$$\begin{aligned}
\Pr(f_1^*(x^{(1)}) \neq y^{(1)} | y^{(1)} = -1) &= \Pr\left(\sum_{i=1}^d \mathcal{N}(-\eta, \sigma_1^2) - b_1 > 0\right) \\
&< \Pr\left(\sum_{i=1}^d \mathcal{N}(-\eta, \sigma_1^2) - b_2 > 0\right) \quad (\text{because } S_1 > S_2) \\
&< \Pr\left(\sum_{i=1}^d \mathcal{N}(-\eta, \sigma_2^2) - b_2 > 0\right) \quad (\text{because } \sigma_1^2 < \sigma_2^2) \\
&= \Pr(f_2^*(x^{(2)}) \neq y^{(2)} | y^{(2)} = -1).
\end{aligned}$$

Step 2. For the error of class “+1”, we have:

$$\begin{aligned}
\Pr(f_1^*(x^{(1)}) \neq y^{(1)} | y^{(1)} = +1) &= \Pr\left(\sum_{i=1}^d \mathcal{N}(\eta, \sigma_1^2) - b_1 < 0\right) \\
&= \Pr\left(\mathcal{N}(0, 1) < \frac{b_1 - d\eta}{d\sigma_1}\right) \tag{10} \\
&= \Pr\left(\mathcal{N}(0, 1) < \frac{-\log(K) \cdot \sigma_1}{2\eta} - \frac{\eta}{\sigma_1}\right),
\end{aligned}$$

and similarly,

$$\Pr(f_2^*(x^{(2)}) \neq y^{(2)} | y^{(2)} = +1) = \Pr\left(\mathcal{N}(0, 1) < \frac{-\log(K) \cdot \sigma_2}{2\eta} - \frac{\eta}{\sigma_2}\right). \tag{11}$$

Note that when K is large enough, i.e., $\log(K) > \frac{2 \cdot \eta^2}{\sigma_1 \cdot \sigma_2}$, we can get the Z-score in Eq. (10) is larger than Eq. (11). As a result, we have:

$$\Pr(f_1^*(x^{(1)}) \neq y^{(1)} | y^{(1)} = +1) > \Pr(f_2^*(x^{(2)}) \neq y^{(2)} | y^{(2)} = +1). \tag{12}$$

By combining Step 1 and Step 2, we can get the inequality in Theorem 3.1.

A.3.3 Proof of Theorem 3.2

Theorem 3.2 Under two data distributions $(x^{(1)}, y^{(1)}) \in \mathcal{D}_1$ and $(x^{(2)}, y^{(2)}) \in \mathcal{D}_2$ with different separability $S_1 > S_2$, let f_1^* and f_2^* be the optimal non-reweighted classifiers ($\rho = 1$) under \mathcal{D}_1 and \mathcal{D}_2 respectively, and let $f_1'^*$ and $f_2'^*$ be the optimal reweighted classifiers under \mathcal{D}_1 and \mathcal{D}_2 given the optimal reweighting ratio ($\rho = K$). Given the imbalance ratio K is large enough, we have:

$$\begin{aligned}
&\Pr(f_1'^*(x^{(1)}) \neq y^{(1)} | y^{(1)} = +1) - \Pr(f_1^*(x^{(1)}) \neq y^{(1)} | y^{(1)} = +1) \\
&< \Pr(f_2'^*(x^{(2)}) \neq y^{(2)} | y^{(2)} = +1) - \Pr(f_2^*(x^{(2)}) \neq y^{(2)} | y^{(2)} = +1).
\end{aligned} \tag{4}$$

Proof 3 (Proof of Theorem 3.2) We first show that under both distribution \mathcal{D}_1 and \mathcal{D}_2 , the optimal reweighting ratio ρ is equal to the imbalance ratio K . Based on the results in Eq. (8) and calculated model parameters w and b , we have the test error (given the model trained by reweight value ρ):

$$\begin{aligned}
&\text{Error}_{\text{test}}(f^*) \\
&= \Pr(f^*(x) \neq y | y = -1) \cdot \Pr(y = -1) + \Pr(f^*(x) \neq y | y = +1) \cdot \Pr(y = +1) \\
&\propto \Pr(\mathcal{N}(0, 1) < -\frac{b + d\eta}{d\sigma}) + \Pr(\mathcal{N}(0, 1) < \frac{b - d\eta}{d\sigma}) \\
&= \Pr(\mathcal{N}(0, 1) < -\frac{1}{2} \log\left(\frac{\rho}{K}\right) - \frac{\sigma}{\eta}) + \Pr(\mathcal{N}(0, 1) < \frac{1}{2} \log\left(\frac{\rho}{K}\right) - \frac{\sigma}{\eta}).
\end{aligned}$$

The value of taking the minimum when its derivative with respect to ρ is equal to 0, where we can get $\rho = K$ and the bias term $b = 0$. Note that the variance values have the relation: $\sigma_1^2 < \sigma_2^2$. Therefore, it is easy to get that:

$$\begin{aligned}
\Pr(f_1'^*(x^{(1)}) \neq y^{(1)} | y^{(1)} = +1) &= \Pr\left(\sum_{i=1}^d \mathcal{N}(\eta, \sigma_1^2) < 0\right) \\
&< \Pr\left(\sum_{i=1}^d \mathcal{N}(\eta, \sigma_2^2) < 0\right) = \Pr(f_2'^*(x^{(2)}) \neq y^{(2)} | y^{(2)} = +1).
\end{aligned} \tag{13}$$

Combining the results in Eq. (12) and (13), we have proved the inequality in Theorem 3.2.

A.4 Algorithm of SRAT

The algorithm of our proposed SRAT framework is shown in Algorithm 1. Specifically, in each training iteration, we first generate adversarial examples using PGD for examples in the current batch (Line 5). If the current training iteration does not reach a predefined starting reweighting epoch T_d , we will assign same weights, i.e., $w_i = 1$ for all adversarial examples \mathbf{x}_i in the current batch (Line 6). Otherwise, the reweighting strategy will be adopted in the final loss function (Line 15), where a specific weight w_i will be assigned for each adversarial example \mathbf{x}_i if its corresponding clean example \mathbf{x}_i comes from an under-represented class.

Algorithm 1 Separable Reweighted Adversarial Training (SRAT).

Input: imbalanced training dataset $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$, number of total training epochs T , starting reweighting epoch T_d , batch size N , number of batches M , learning rate γ

Output: An adversarially robust model f_θ

- 1: Initialize the model parameters θ randomly;
 - 2: **for** epoch = $1, \dots, T_d - 1$ **do**
 - 3: **for** mini-batch = $1, \dots, M$ **do**
 - 4: Sample a mini-batch $\mathcal{B} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ from D ;
 - 5: Generate adversarial example \mathbf{x}'_i for each $\mathbf{x}_i \in \mathcal{B}$;
 - 6: $\mathcal{L}(f_\theta) = \frac{1}{N} \sum_{i=1}^N \max_{\|\mathbf{x}'_i - \mathbf{x}_i\|_p \leq \epsilon} \mathcal{L}(f_\theta(\mathbf{x}'_i), y_i) + \lambda \mathcal{L}_{sep}(\mathbf{x}'_i)$
 - 7: $\theta \leftarrow \theta - \gamma \nabla_{\theta} \mathcal{L}(f_\theta)$
 - 8: **end for**
 - 9: Optional: $\gamma \leftarrow \gamma / \kappa$
 - 10: **end for**
 - 11: **for** epoch = T_d, \dots, T **do**
 - 12: **for** mini-batch = $1, \dots, M$ **do**
 - 13: Sample a mini-batch $\mathcal{B} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ from D ;
 - 14: Generate adversarial example \mathbf{x}'_i for each $\mathbf{x}_i \in \mathcal{B}$;
 - 15: $\mathcal{L}(f_\theta) = \frac{1}{N} \sum_{i=1}^N \max_{\|\mathbf{x}'_i - \mathbf{x}_i\|_p \leq \epsilon} w_i \mathcal{L}(f_\theta(\mathbf{x}'_i), y_i) + \lambda \mathcal{L}_{sep}(\mathbf{x}'_i)$
 - 16: $\theta \leftarrow \theta - \gamma \nabla_{\theta} \mathcal{L}(f_\theta)$
 - 17: **end for**
 - 18: Optional: $\gamma \leftarrow \gamma / \kappa$
 - 19: **end for**
-

A.5 Data Distribution of Imbalanced Training Datasets

In our experiments, we construct multiple imbalanced training datasets to simulate various kinds of imbalanced scenarios by combining different imbalance types (i.e., Exp and Step) with different imbalanced ratios (i.e., $K = 10$ and $K = 100$). Figure 13 and Figure 14 show the data distribution of all ten-classes imbalanced training datasets used in our preliminary studies and experiments based on CIFAR10 [21] and SVHN [27] datasets, respectively.

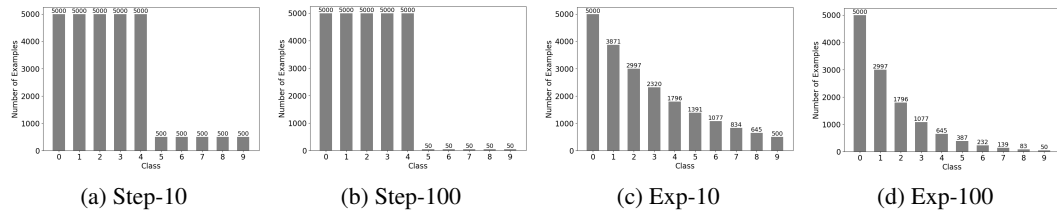


Figure 13: Data distribution of imbalanced training datasets constructed from CIFAR10 dataset.

A.6 Performance Comparison on Imbalanced SVHN Datasets

Table 3 and Table 4 show the performance comparison on various imbalanced SVHN datasets with different imbalance types and imbalance ratios. We use bold values to denote the highest accuracy among all methods and use the underline values to indicate our SRAT variants which achieve the

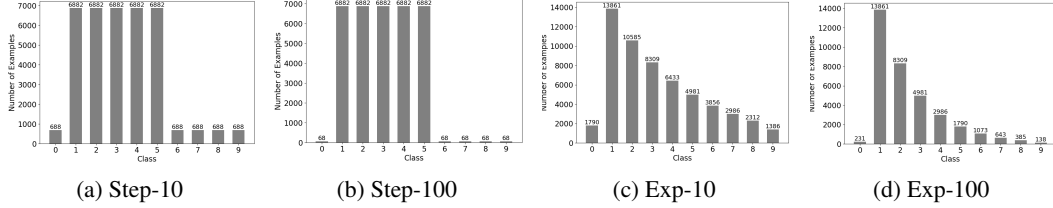


Figure 14: Data distribution of imbalanced training datasets constructed from SVHN dataset.

highest accuracy among their corresponding baseline methods utilizing the same loss function for making predictions.

From Table 3 and Table 4, we get similar observation that, comparing with baseline methods, our proposed SRAT method can produce a robust model which can achieve improved overall performance when the training dataset is imbalanced. In addition, based on the experimental results in Table 1 to Table 4, we find that, compared with the performance improvement between DRCB-LDAM and SRAT-LDAM, the improvement between DRCB-CE and SRAT-CE and the improvement between DRCB-Focal and SRAT-Focal are more obviously. The possible reason behind this phenomenon is, the LDAM loss can also implicitly produce a more separable feature space [4] while CE loss and Focal loss do not conduct any specific operations on the latent feature space. Hence, the feature separation loss contained in SRAT-CE and SRAT-Focal could be more effective on learning separable feature space and facilitate the Focal loss on prediction. However, in SRAT-LDAM, the feature separation loss and LDAM loss may affect each other on learning feature representations and, hence, the effectiveness of the feature separation loss may be counteracted or weakened.

In conclusion, experiments conducted on multiple imbalanced datasets verify the effectiveness of our proposed SRAT method under various imbalanced scenarios.

Table 3: Performance Comparison on Imbalanced SVHN Datasets (Imbalanced Type: Step)

Imbalance Ratio	10				100			
	Standard Accuracy		Robust Accuracy		Standard Accuracy		Robust Accuracy	
Method	Overall	Under	Overall	Under	Overall	Under	Overall	Under
CE	79.88	67.04	37.62	22.08	59.61	26.19	29.57	5.03
Focal	79.96	67.03	37.83	22.47	60.58	28.17	30.27	5.83
LDAM	84.55	74.96	45.80	31.23	65.61	37.13	33.34	8.36
CB-Reweight	79.48	66.07	37.38	21.66	60.23	27.68	29.54	5.32
CB-Focal	80.29	67.56	38.10	23.00	60.73	28.37	30.09	5.75
DRCB-CE	80.62	68.74	37.25	22.79	60.67	28.36	30.02	5.59
DRCB-Focal	79.11	65.72	37.01	22.02	61.65	30.29	30.78	7.06
DRCB-LDAM	87.83	82.63	46.45	35.15	63.78	33.99	33.60	7.28
SRAT-CE	<u>82.89</u>	<u>72.79</u>	<u>38.23</u>	<u>24.70</u>	<u>63.39</u>	<u>33.85</u>	29.64	<u>6.11</u>
SRAT-Focal	<u>85.05</u>	<u>77.10</u>	<u>39.51</u>	<u>28.06</u>	<u>70.12</u>	<u>47.44</u>	<u>32.18</u>	<u>11.08</u>
SRAT-LDAM	87.65	82.62	46.03	34.75	<u>71.56</u>	<u>50.33</u>	<u>33.54</u>	<u>11.63</u>

Table 4: Performance Comparison on Imbalanced SVHN Datasets (Imbalanced Type: Exp)

Imbalance Ratio	10				100			
	Standard Accuracy		Robust Accuracy		Standard Accuracy		Robust Accuracy	
Method	Overall	Under	Overall	Under	Overall	Under	Overall	Under
CE	87.54	82.67	44.12	35.33	72.51	56.30	33.34	16.93
Focal	87.82	83.01	44.88	35.97	72.61	56.48	34.09	17.62
LDAM	90.06	86.69	51.84	43.73	79.11	66.86	40.42	25.18
CB-Reweight	87.66	82.79	44.39	35.53	72.25	55.97	33.36	17.16
CB-Focal	87.86	82.96	44.61	35.55	73.23	57.34	34.25	17.90
DRCB-CE	88.49	84.51	43.82	36.28	73.74	58.03	33.52	17.68
DRCB-Focal	87.47	82.78	42.52	34.31	71.95	55.11	33.43	17.63
DRCB-LDAM	91.24	89.65	52.39	46.71	80.29	69.23	40.16	24.64
SRAT-CE	<u>88.70</u>	<u>84.94</u>	<u>44.54</u>	<u>36.59</u>	<u>77.11</u>	<u>64.47</u>	<u>34.48</u>	<u>19.91</u>
SRAT-Focal	<u>89.51</u>	<u>85.42</u>	<u>45.37</u>	<u>37.20</u>	<u>80.04</u>	<u>69.54</u>	<u>35.25</u>	<u>23.04</u>
SRAT-LDAM	91.27	89.55	52.10	46.13	80.71	70.49	40.33	25.11