
Supplementary Material for How Re-sampling Helps for Long-Tail Learning?

Anonymous Author(s)

Affiliation

Address

email

1 Contents

2	A Training Procedure	1
3	B Implementation Details for <i>context-shift augmentation</i>	1
4	C Additional Experimental Results	2
5	C.1 Effects of the context bank	2
6	C.2 Influence of augmentation variants	3
7	C.3 Comparison between different modules	3
8	C.4 The influence of Grad-CAM	4
9	C.5 Combination with self-supervised learning	4
10	C.6 Combination with the logit adjustment	4

11 A Training Procedure

12 The traing procedure of *context-shift augmentation* is summarized in Algorithm 1.

13 B Implementation Details for *context-shift augmentation*

14 For experiments on CIFAR10-LT and CIFAR100-LT, we use ResNet-32 as the backbone network
15 and train it using standard SGD with a momentum of 0.9, a weight decay of 2×10^{-4} , a batch size
16 of 128. The model is trained for 200 epochs. The initial learning rate is set to 0.2 and is annealed by
17 a factor of 10 at 160 and 180 epochs. We train each model with 1 NVIDIA GeForce RTX 3090.

18 For experiments on ImageNet-LT, we implement the proposed method on ResNet-10 and ResNet-50.
19 We use standard SGD with a momentum of 0.9, a weight decay of 5×10^{-4} , and a batch size of
20 256 to train the whole model for a total of 90 epochs. We use the cosine learning rate decay with an
21 initial learning rate of 0.2. We train each model with 2 NVIDIA Tesla V100 GPUs.

22 In all experiments, we first warm-up the uniform module for 10 epochs and then train the uniform
23 module and the balanced re-sampling module simultaneously for the rest epochs. For the uniform
24 module, we follow the simple data augmentation used in [1] with only random crop and horizontal
25 flips. For the re-sampling module, we use the proposed context-shift augmentation method. We
26 apply the trick proposed by [2] to disable the augmentation in the balanced re-sampling module at the

Algorithm 1 Training procedure of context-shift augmentation

Input: training data $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$; context memory bank Q , maximum volume size V ; model parameters ϕ, f^u, f^b ; loss functions ℓ^u, ℓ^b ;

Procedure:

```
1: Initialize model parameters  $\phi, f^u, f^b$ ;  
2: Re-sampling a class-balanced dataset  $\tilde{\mathcal{D}} = \{(\tilde{\mathbf{x}}_i, \tilde{y}_i)\}_{i=1}^N$ ;  
3: Empty memory bank  $Q$ ;  
4: for epoch = 1, ...,  $T$  do  
5:   repeat  
6:     Draw a mini-batch  $(\mathbf{x}_i, y_i)_{i=1}^B$  from  $\mathcal{D}$ ;  
7:     Draw a mini-batch  $(\tilde{\mathbf{x}}_i, \tilde{y}_i)_{i=1}^B$  from  $\tilde{\mathcal{D}}$ ;  
8:     // uniform module  
9:     for  $i = 1, \dots, B$  do  
10:      Calculate  $\mathbf{z}_i^u = f^u(\phi(\mathbf{x}_i))$  and  $\mathcal{L}_i^u = \ell^u(\mathbf{z}_i^u, y_i)$ ;  
11:      if  $p(y = y_i | \mathbf{x}_i, \phi, f^u) \geq \delta$  then  
12:        Calculate background mask  $M_i$  of  $\mathbf{x}_i$ ;  
13:        Push  $(\mathbf{x}_i, M_i)$  into  $Q$ ;  
14:      end if  
15:    end for  
16:    Calculate  $\mathcal{L}^u = \frac{1}{B} \sum_{i=1}^B \mathcal{L}_i^u$ ;  
17:    // balanced re-sampling module  
18:    if Size of  $Q$  reaches  $V$  then  
19:      Get backgrounds  $(\tilde{\mathbf{x}}_i, M_i)_{i=1}^B$  from  $Q$ ;  
20:       $\lambda \sim \text{Uniform}(0, 1)$ ;  
21:      for  $i = 1, \dots, B$  do  
22:         $\tilde{\mathbf{x}}_i = \lambda M_i \odot \tilde{\mathbf{x}}_i + (1 - \lambda M_i) \odot \tilde{\mathbf{x}}_i$ ;  
23:        Calculate  $\mathbf{z}_i^b = f^b(\phi(\tilde{\mathbf{x}}_i))$  and  $\mathcal{L}_i^b = \ell^b(\mathbf{z}_i^b, \tilde{y}_i)$ ;  
24:      end for  
25:      Calculate  $\mathcal{L}^b = \frac{1}{B} \sum_{i=1}^B \mathcal{L}_i^b$ ;  
26:    else  
27:      Assign  $\mathcal{L}^b = 0$ ;  
28:    end if  
29:    // total objective function  
30:    Calculate  $\mathcal{L} = \mathcal{L}^u + \mathcal{L}^b$ ;  
31:    Update model parameters  $\phi, f^u, f^b$  with  $\mathcal{L}$ ;  
32:  until all training data are traversed.  
33: end for
```

27 last 3 epochs to obtain further improvements, which is also applied in other baseline methods [3, 4].
28 We set the threshold δ to 0.8.

29 C Additional Experimental Results

30 C.1 Effects of the context bank

31 The context bank Q is a novel component of *context-shift augmentation* which receives diverse
32 contexts from the uniform module, and provides them to augment the data in the re-sampling module.
33 To verify the effectiveness of the context bank, we remove it from the framework and train the model
34 on CIFAR100-LT with imbalance ratio of 100. The results are reported in Table 1.

35 The results show that without context bank Q , the performance decrease by a large margin. The per-
36 formance degradation mostly comes from the medium-shot classes and the few-shot classes, which
37 indicates that the context bank can significantly improve the generalization of tail classes.

38 Moreover, we study the effect of different variants in the context bank. First, we study the threshold
39 δ for sample selecting and report the results in Table 2. On the one hand, if δ is too high, the selected
40 samples will be very few. On the other hand, if δ is too small, the selected samples might be not well

Table 1: Ablation study on the context bank Q .

	All	Many	Med.	Few
Ours	45.8	64.3	49.7	18.2
Ours w/o Q	41.2 (-4.6)	65.1 (+0.8)	41.9 (-7.8)	10.7 (-7.5)

learned. Nevertheless, as the training process progresses, most samples will fit well, so our method is not sensitive to δ . We set $\delta = 0.8$ considering its best performance.

Table 2: Influence of the threshold δ .

δ	0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9
Accuracy	45.47	45.37	45.55	44.83	45.52	45.59	45.42	45.08	45.83	44.93

Second, we study the influence of the volume size V of the context bank Q and report the results in Table 3. Since the bank Q is a first-in-first-out queue, the latest incoming backgrounds are more convincing. When the volume size is too large, the bank might contain more past samples. Besides, a larger size of V would bring more memory overhead. So we set the volume size V equal to the mini-batch size B in our method.

Table 3: Influence of the bank volume size (compared with the mini-batch size B).

Volume	$\times 1$	$\times 2$	$\times 4$	$\times 8$	$\times 16$	$\times 32$	$\times 64$
Accuracy	45.83	45.60	45.57	45.32	45.55	45.37	45.26

C.2 Influence of augmentation variants

We use a variant $\lambda \sim \text{Uniform}(0, 1)$ for generating novel samples. The value of λ can result in different proportions of foreground and background in the novel sample. Also, the size of the sampling space affects the diversity of the novel images. To explore the effect of λ , we try $\lambda \sim \text{Uniform}(a, b)$ and $\lambda \sim \text{Beta}(a, b)$ to train *context-shift augmentation* on CIFAR100-LT with imbalance ratio 100 and report the results in Figure 1.

First, when λ is close to 0, the background merely takes effect, and the performance decreases a lot. Second, when $\lambda = 1$, the background image might cover the important content in the foreground image. Also, the diversity of new samples is limited. Although the performance is better than that of $\lambda = 0$, it is still unsatisfactory. Overall, choosing $\lambda \sim \text{Uniform}(0, 1)$ or $\lambda \sim \text{Beta}(1, 1)$ lead to the best performance.

C.3 Comparison between different modules

In our framework, the uniform sampling module is only enabled in the training phase. While in the inference phase, we use the balanced re-sampling module to predict unseen instances. To verify the superiority of the re-sampling module, we compare the performance of these two modules as well as their ensemble. We report the results in Table 4. The results show that the re-sampling module is superior to the uniform module, and even achieves higher accuracy than the ensembled results. This also indicates the superiority of the proposed context-shift re-sampling method.

Moreover, we study the influence of different balance ratios on our re-sampling module and compare it with the vanilla re-sampling method. We report the results in Figure 2. For vanilla re-sampling, adopting a more balanced re-sampling would yield more severe performance degradation. In contrast, our method achieves higher performance through class-balanced re-sampling.

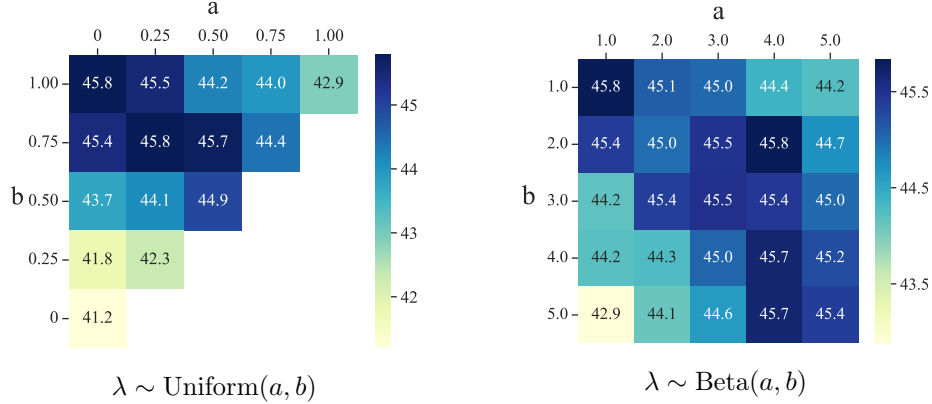


Figure 1: Influence of variants λ .

Table 4: Comparison between the uniform module and the re-sampling module in *context-shift augmentation*.

	All	Many	Med.	Few
Uniform module	39.4	68.3	37.3	6.1
Re-sampling module	45.8	64.3	49.7	18.2
Ensemble results	43.0	67.5	44.1	11.4

70 C.4 The influence of Grad-CAM

71 The Grad-CAM [5] is utilized to extract background in previous works such as open-set learning and
 72 adversarial learning [6, 7]. Also, we use Grad-CAM in the *context-shift augmentation* to generate
 73 diverse backgrounds for tail-class data. We compare Grad-CAM with CAM [8]. The results shown
 74 in Table 5 demonstrate that Grad-CAM is superior to CAM when applied to our method.

Table 5: Comparison between CAM and Grad-CAM in *context-shift augmentation*

	All	Many	Med.	Few
Ours w/ CAM	45.1	63.8	48.1	18.0
Ours w/ Grad-CAM	45.8	64.3	49.7	18.2

75 C.5 Combination with self-supervised learning

76 It is interesting to combine self-supervised methods with *context-shift augmentation*. Inspired by
 77 this, we follow the self-supervised + fine-tune method, i.e., SimSiam+rwSAM in [9] and conduct
 78 extensive experiments on CIFAR10-LT dataset. The results are shown in Table 6.

79 Note that in [9], the models are pre-trained with the long-tailed dataset, while fine-tuned with the
 80 balanced in-domain dataset. However, it is hard to achieve a balanced in-domain version of a long-
 81 tailed dataset in real-world scenarios. So we use the long-tailed dataset with balancing method
 82 including class-balanced re-sampling (CB-RS), and re-sampling with *context-shift augmentation*.
 83 The results show that our method is superior to CB-RS.

84 C.6 Combination with the logit adjustment

85 Since our work aims to study the effectiveness of re-sampling in long-tail learning, we use the
 86 Class-Balanced Re-Sampling (CB-RS) in our method. We consider combining our method with
 87 other re-balancing methods, such as Logit Adjustment (LA) [10]. Specifically, we change the class-
 88 balanced re-sampling to the uniform sampling while adopting *context-shift augmentation*. Moreover,
 89 we consider combining class-balanced re-sampling and LA simultaneously. The comparison results

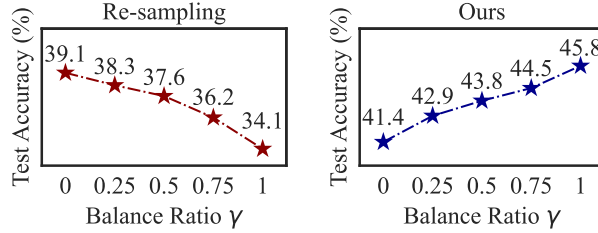


Figure 2: Comparison of re-sampling and our method under different balance ratios γ .

Table 6: Combing self-supervised method SimSiam+rwSAM with different re-balancing methods

Pre-train model	Fine-tune method	Accuracy (%)
SimSiam+rwSAM	CE	69.5
SimSiam+rwSAM	CB-RS	72.8
SimSiam+rwSAM	Ours	75.5

are shown in Table 7. The results show that our method can be combined with logit adjustment to yield a more higher performance. However, by applying the balanced loss and the balanced sampling, the model puts much focus on tail classes and results in a deterioration of overall accuracy.

Table 7: Combination with Logit Adjustment (LA)

	All	Many	Med.	Few
w/ CB-RS	45.8	64.3	49.7	18.2
w/ LA	47.2	62.4	48.8	26.1
w/ both	45.0	48.7	51.0	31.4

References

- [1] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 770–778. IEEE Computer Society, 2016.
- [2] Zhuoxun He, Lingxi Xie, Xin Chen, Ya Zhang, Yanfeng Wang, and Qi Tian. Data augmentation revisited: Rethinking the distribution gap between clean and augmented data. *CoRR*, abs/1909.09148, 2019.
- [3] Yongshun Zhang, Xiu-Shen Wei, Boyan Zhou, and Jianxin Wu. Bag of tricks for long-tailed visual recognition with deep convolutional neural networks. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Virtual Event, February 2-9, 2021*, pages 3447–3455, 2021.
- [4] Seulki Park, Youngkyu Hong, Byeongho Heo, Sangdoo Yun, and Jin Young Choi. The majority can help the minority: Context-rich minority oversampling for long-tailed classification. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 6877–6886, 2022.
- [5] Bolei Zhou, Aditya Khosla, Àgata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *CVPR*, pages 2921–2929. IEEE Computer Society, 2016.
- [6] Jie-Jing Shao, Xiao-Wen Yang, and Lan-Zhe Guo. Open-set learning under covariate shift. *Machine Learning*, pages 1–17, 2022.

- 113 [7] Qibing Ren, Yiting Chen, Yichuan Mo, Qitian Wu, and Junchi Yan. DICE: domain-attack
114 invariant causal learning for improved data privacy protection and adversarial robustness. In
115 *KDD*, pages 1483–1492. ACM, 2022.
- 116 [8] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi
117 Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-
118 based localization. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice,*
119 *Italy, October 22-29, 2017*, pages 618–626, 2017.
- 120 [9] Hong Liu, Jeff Z. HaoChen, Adrien Gaidon, and Tengyu Ma. Self-supervised learning is more
121 robust to dataset imbalance. In *ICLR*. OpenReview.net, 2022.
- 122 [10] Aditya Krishna Menon, Sadeep Jayasumana, Ankit Singh Rawat, Himanshu Jain, Andreas Veit,
123 and Sanjiv Kumar. Long-tail learning via logit adjustment. In *9th International Conference on*
124 *Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*, 2021.