

1 We thank the reviewers for their insightful comments and constructive feedback. As the reviewers mentioned, our work  
2 shows the following strengths. (1) The theoretical analysis is “profound and novel” [R3,R4,R5]. (2) Experiments are  
3 designed “thoroughly” and “carefully” which “verifies the feasibility” [R3,R4]. (3) The paper is “well-written and  
4 organized” [R3,R4]. We will answer the major points below and address all remaining ones in the final version.

5 **[R3]:** “For eq. (3) (4) (5), the first item on the right-hand side,  $\sqrt{\frac{C(F)}{n_j}}$  or  $\frac{C(F)}{\sqrt{n_j}}\gamma$ ”

6 • This depends on how  $C(F)$  is defined. If it is defined to be the Rademacher complexity, then the former is correct.

7 **[R3]:** I suggest the authors to polish up the Figure 1.

8 • Thanks for the suggestion! We’ll update with a better one for the final version.

9 **[R3]:** ““Hinge loss (HG) does not work well with 100 classes”, what you mean by not work well?”

10 • When trained on CIFAR-100, Hinge loss seems to suffer from optimization issues — the training accuracy is at most  
11 about 80%. Thus we didn’t report the test accuracy because the failure here is of a different nature.

12 **[R4,R6]:** “It is unclear to me why the loss function (10) enforces the desired margin in (9).”; “Provide a strong  
13 justification for the equation (10)”

14 • The Hinge loss in (10) achieves its minimum value zero only if the margin is at least  $\Delta_y$ . Recall that the margin is  
15 defined to be  $\gamma = z_y - \max_{j \neq y} z_j$ . Therefore, Hinge loss =  $\max\{\Delta_y - \gamma, 0\} = 0$  if and only if  $\gamma \geq \Delta_j$ . Hinge loss is  
16 a standard loss that encourages margins in the context of SVM.<sup>1</sup> We extend it to allow label-dependent margins.

17 **[R4]:** “I wonder what exactly is showing in Figure 2.”

18 • We visualize the distributions of the last-but-one layer of the neural network, which are referred to as the features.  
19 Please refer to the details in L230-L235. We will clarify more in the final version.

20 **[R4,R6]:** “If the second stage does not move the weight by much, shouldn’t the ERM with LDAM loss work well  
21 enough?”; “Provide a better why DRW is important?”

22 • We believe that the second stage with smaller learning rate serves as a fine-tuning-like process to capture sophisticated  
23 details in each class. Thus in the second stage, emphasizing rare examples are important, because without it, the training  
24 accuracies for all the classes can not be approximately 100%. (Relatively smaller movements in the second stage could  
25 also change the performance by more than a few percents.) With the initial large learning rate in the first stage, by  
26 contrast, the network learns the shared patterns/features shared across all tasks, and therefore it would be better to train  
27 with all the examples with uniform weights. Such phenomenon/intuitions were also observed<sup>2</sup> and justified in recent  
28 works<sup>3</sup>. We realized this from the ablation study in Fig. 6 in Appendix, which shows that the features learned in the  
29 first stage with ERM are better than those with re-weighting.

30 **[R5]:** “How to decide the hyperparameter C? How is the LDAM-HG-DRS in Table.1 implemented?”

31 • We tune  $C$  as a hyper-parameter for each dataset. In particular, we use  $C = 0.5$  for all CIFAR-10 and CIFAR-100  
32 experiments, and  $C = 0.3$  for all iNaturalist experiments. Regarding the LDAM-HG-DRS implementation, we follow  
33 Eq. (10) to implement Hinge loss. Here DRS means the delayed re-sampling strategy.

34 **[R5]:** “CB+Softmax and LDAM seem to be quite similar”; “it seems that the main boost of performance is stemmed  
35 from the DRW (deferred re-weighting)”; additional baseline CB+DRW.

36 • We’d like first to clarify that CB only re-weights the losses, and therefore is a re-weighting scheme more similar to  
37 vanilla re-weighting than to LDAM (which is a new loss). DRW, a deferred re-weighting scheme that we proposed, is  
38 an improved version of CB or vanilla re-weighting, and is orthogonal to LDAM. In Tab. 2, we see that either using  
39 LDAM alone (4th row), or DRW alone (3rd row), on top of the ERM baseline, can outperform prior work. LDAM alone  
40 (3.5% improvement) is slightly more useful than DRW alone (2.6%), and together, they give 6.8% improvement. Thus  
41 we don’t agree that the main boost stems from DRW. We found CB+DRW does not outperform DRW alone, which also  
42 suggests that DRW is a better re-weighting scheme.

43 **[R6]:** Test the proposed method for more general machine learning tasks.

44 • Thank you for your suggestion. We selected these datasets (1) to compare with related works, (2) because they are  
45 challenging, (3) because they are representative of ubiquitous real-world dataset imbalance issues. Nonetheless, we add  
46 one additional sentiment analysis experiment on the Large Movie Review (IMDB) Dataset, a popular and standard  
47 task in NLP. We manually created an imbalanced training set by removing 90% of negative reviews. We train a 2-layer  
48 bidirectional LSTM with Adam optimizer. Test accuracy of different methods are listed as follows: ERM: 63.18,  
49 Re-weight: 76.34, Re-sample: 73.50, LDAM: 82.16. Thus our conclusions hold on other tasks. We will add this result  
50 to the final version of the paper.

---

<sup>1</sup>Wikipedia contributors. “Hinge loss.” Wikipedia, The Free Encyclopedia.

<sup>2</sup>Nakkiran, Preetum, et al. “SGD on Neural Networks Learns Functions of Increasing Complexity.”

<sup>3</sup>Li, Yuanzhi, et al. “Towards Explaining the Regularization Effect of Initial Large Learning Rate in Training Neural Networks.”