

1 We sincerely thank the three reviewers for their constructive comments and supports.

2 **Response to Reviewer #1:**

3 **Q1:** Compressing on the user part. **A:** GPUs are essential to doing effective deep learning. Compared with setting up  
4 their own servers, many users tend to spin up cloud instances with GPUs by balancing the flexibility and the investment,  
5 especially when the GPUs are only needed for several hours. In addition, not every user is a deep learning expert, and  
6 thus a cloud service would be expected to produce efficient deep neural networks according to users' needs.

7 **Q2:** Cloud platform is used, no compression is required. **A:** The compressed networks are often deployed on low-end  
8 computing devices, e.g. digital cameras and mobile phones. By doing this, the cloud service latency can be avoided and  
9 the application can be well executed even there is no Internet connection, which will improve the user experience.

10 **Q3:** Novelty. **A:** The novelty of this paper is twofold. Firstly, compressing network with little labeled data is a  
11 challenging task that has rarely been investigated by existing works in the field. To compensate the lack of labeled data,  
12 we introduce PU learning to select the most related data from a large pool of unlabeled data for the compression task.  
13 Secondly, we enhance the robustness of knowledge distillation to deal with data imbalance problem and noise. An  
14 attention based multi-scaled feature extractor is developed to cope with PU data better.

15 **Q4:** Comparison on using the attention-based feature extraction or not. **A:** The results are shown in Figure 2. The blue  
16 line is the result on using attention-based feature extractor. The red line is the result without attention-based feature  
17 extractor. The superiority of the proposed architecture leads to the accuracy improvement from 91.39% to 91.57%.

18 **Q5:** Comparison on part of, or all the unlabeled data set with KD method. Using robust KD or not. **A:** We include new  
19 experiments on CIFAR-10. Randomly choosing 50,000 samples leads to a 87.02% accuracy (91.56% for PU method  
20 choosing 50,000 samples). The result is bad since many negative data are selected in this way and is used to train the  
21 student network. Using all 1.2M unlabeled data leads to a 94.01% accuracy since all the positive data are guaranteed to  
22 be selected, but is about 12 times slower than using PU data (93.75% for PU method choosing 0.1M samples). Robust  
23 KD leads to 0.6% increase.

24 **Q6:** Minor problems. **A:** Thanks for this nice concern. All these typos will be corrected in the final version.

25 **Response to Reviewer #2:**

26 **Q1:** Figure 4. **A:** The first row represents the data in the original dataset. The second is the positive data selected by PU  
27 classifier. The third is the negative data selected by PU classifier. Spider is in both selected data and negative data, since  
28 the PU classifier does not classify unlabeled data with 100% accuracy, and spider represents the noise in positive data.

29 **Response to Reviewer #3:**

30 **Q1:** Advantages over related works. **A:** Theoretically, we utilize the strength of PU learning in augmenting data, and  
31 the strength of KD method in compressing neural networks, which is suitable for solving compression problem with  
32 little labeled data. Besides, we propose multi-feature network with attention and robust KD method to better solve  
33 the problem. Experimentally, we compare with the state-of-the-art methods in Table 3, and the accuracy results show  
34 the priority of the proposed method. Generally speaking, the proposed method is robust on the number of positive  
35 samples, and performs better. In other methods, the student network is compressed by pruning [1] or re-normalization  
36 [2] the giant teacher network with unlabeled data, which means that the detail architecture of the teacher network is  
37 required. In the proposed method, we only need the input and output interface of the teacher network instead of the  
38 whole architecture, which is more flexible for users to protect their own privacy.

39 **Q2:** The motivation to tackle imbalanced data problem is unclear. **A:** PU classifier treats all samples in the training set  
40 as 'positive', the specific category is undistinguishable. PU may select lots of data for some positive categories, while  
41 little for others. This causes data imbalanced problem. Eq.(7) and (8) are used to tackle the imbalanced data problem.  
42 We utilize the output of teacher network, and assign larger weights to categories with fewer samples. The weighted  
43 surrogate KD loss (Eq.(8)) can alleviate the imbalanced data problem.

44 **Q3:** Class prior  $\pi_p$ . **A:** In CIFAR-10 experiments, it is set equal to the ratio of manually selected data. In MNIST  
45 experiments, it is set to the real ratio of numbers in EMNIST. In ImageNet experiments, which is much like the reality  
46 settings,  $\pi_p$  is estimated by the prior estimation methods, such as [3]. In Figure 2, we analysis the relationship between  
47 classification accuracy and class prior. And it shows that our method is robust to the choice of class prior when it is not  
48 far from the true class prior.

49 [1] Tang Y, You S, Xu C, et al. Bringing Giant Neural Networks Down to Earth with Unlabeled Data. arXiv, 2019.

50 [2] He X, Cheng J. Learning Compression from Limited Unlabeled Data. ECCV, 2018.

51 [3] M. C. du Plessis and M. Sugiyama. Class prior estimation from positive and unlabeled data. NIPS, 2014.