

Note: Self-Distillation from the Last Mini-Batch for Consistency Regularization.[1]

sgc

May 1, 2022

1 Introduction

Conventional KD requires a complex teacher resulting its low-efficiency. And recent self-KD need extra network or hard to parallelize. This paper proposed Self-Distillation from Last Mini-Batch (DLB). Rather than storing network configuration, this method only stores the soft labels.

Contributions:

1, With no network architecture modifications, DLB method requires very few additional computation costs as well as the run-time memory to implement.

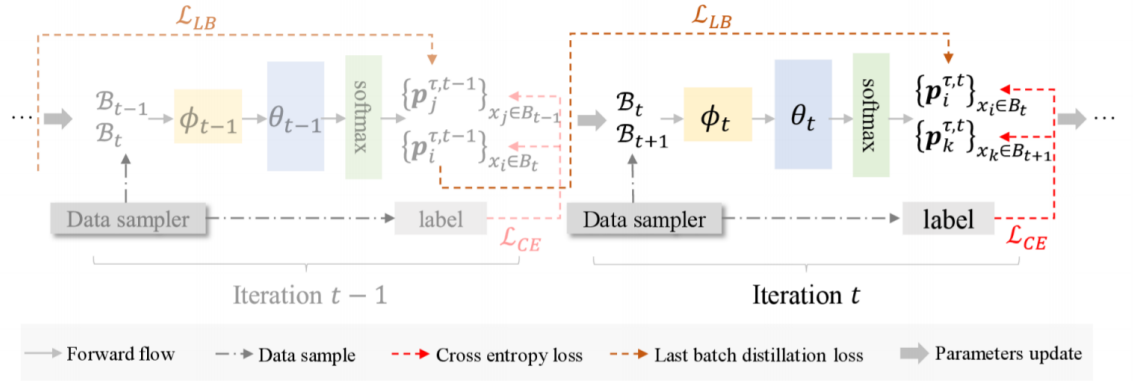


Figure 1. The overall architecture of our DLB. We write \mathcal{B}_t , ϕ_t , θ_t for a mini-batch of data samples, random augmentation and the trainable parameters indexed in the t^{th} iteration.

2 Method

2.1 Self-Distillation from Last Batch

The total loss consists of two parts:

$$L = L_{CE} + \alpha L_{LB}$$

Last-batch consistency regularization loss:

$$\mathcal{L}_{LB} = \frac{1}{n} \sum_{i=1}^n \tau^2 \cdot D_{KL}(\mathbf{p}_i^{\tau, t-1} \| \mathbf{p}_i^{\tau, t}).$$

where $p_i^{\tau, t-1}$ is the prediction of $t-1$ iteration, and $p_i^{\tau, t}$ is t 's.

Both the predictions from B_{t-1} and B_t in $t-1^{th}$ iteration update the L_{CE} .

Algorithm 1 Pseudo code for DLB.

Input: balancing coefficient α
Input: distillation temperature τ
Require: data_loader samples batches as in Figure 1

```
1: last_logits = None # initialization
2: for (x, gt_labels) in data_loader do
3:    $\hat{n} = \text{gt\_labels.size}(0)$  # batch size
4:   logits = model.forward(x)
5:   loss = CELoss(logits, gt_labels)
6:   if last_logits != None then
7:     soft_targets = Softmax(last_logits/ $\tau$ )
8:     pred = Softmax(logits[: $\hat{n}/2$ ]/ $\tau$ )
9:     loss +=  $\alpha * \text{LBLoss}(\text{pred}, \text{soft\_targets}) * \tau^2$ 
10:  end if
11:  loss.backward() # update parameters
12:  last_logits = logits[: $\hat{n}/2$ ].detach() # no gradient
13: end for
```

3 Results

DLB consistently improved the performance on various backbones, outperformed the SOTA. DLB and augmentation regularization work orthogonally.

4 Summary

This paper proposed DLB which does not require a complex teacher or ensemble peer students. It boosts the robustness to label noise. It uses half samples from the last iteration to keep consistency of samples and thus make it robust to noise.

References

- [1] Y. Shen, L. Xu, Y. Yang, Y. Li, and Y. Guo, “Self-distillation from the last mini-batch for consistency regularization,” *arXiv preprint arXiv:2203.16172*, 2022. (document)