

## **ProMix: Combating Label Noise via Maximizing Clean Sample Utility**

Ruixuan Xiao<sup>1</sup>, Yiwen Dong<sup>1</sup>, Haobo Wang<sup>1\*</sup>, Lei Feng<sup>2</sup>, Runze Wu<sup>3</sup>, Gang Chen<sup>1</sup> and Junbo Zhao<sup>1</sup> <sup>1</sup>Zhejiang University, Hangzhou, China <sup>2</sup>Nanyang Technological University, Singapore <sup>3</sup>NetEase Fuxi AI Lab, Hangzhou, China {xiaoruixuan, dyw424, wanghaobo, cg, j.zhao}@zju.edu.cn, lfengqaq@gmail.com, wurunze1@corp.netease.com

**IJCAI 2023** 



## **1. Learning with Noisy Labels**:

The vanilla CE is proved to easily overfit corrupted data. Existing methods can be roughly categorized into two types:

- the estimated noise transition matrix to explicitly correct the loss function, but hard in the presence of heavy noise and a large number of classes.
- filter out a clean subset for robust training, e.g. DivideMix(to leverage the unselected examples): existing a quality-quantity trade-off.

Motivation: our work targets to maximize the utility of clean samples.

### 2. Debiased Learning:

In long-tailed learning, models can be easily biased towards dominant classes. As for SSL, confirmation bias(unreliable pseudo-labels) may hurt generalization performance.

Motivation: we attempt to alleviate the biases in our selection and pseudo-labeling procedures.

## Method





(a) Overall Framework of ProMix.

(b) Progressive Selection.





(b) Progressive Selection.

Motivation:

- In the training, the easy classes are usually fitted.
- The loss value of example with different observed labels may not be comparable.

Datasets D to C sets according to the noisy labels at each epoch:

 $\mathcal{S}_j = \{(\boldsymbol{x}_i, \tilde{y}_i) \in \mathcal{D} | \tilde{y}_i = j\}$ 

Produce a roughly balanced base set:

 $k = \min(\left\lceil \frac{n}{C} \times R \right\rceil, |\mathcal{S}_j|)$ 

$$\mathcal{D}_{CSS} = \cup_{j=1}^{C} \mathcal{C}_{j}$$



Motivation: Find out the potentially clean samples missed by CSS

Inspiration:

the DNNs can assign an arbitrary label with high confidence and select them as clean, which results in a vicious cycle.

we choose those samples to have high confidence in their given labels, which tend to be originally clean.

$$e_i = \max_j p_i^j \qquad \qquad y_i' = \arg\max_j p_i^j$$

 $(\max_{j} f_{1}^{j}(\boldsymbol{x}) \geq \tau) \wedge (\max_{j} f_{2}^{j}(\boldsymbol{x}) \geq \tau) \wedge$   $(\arg\max_{j} f_{1}^{j}(\boldsymbol{x}) = \arg\max_{j} f_{2}^{j}(\boldsymbol{x}))$ (8) Label Guessing by Agreement (LGA)

$$\mathcal{D}_{MHCS} = \{ (\boldsymbol{x}_i, \tilde{y}_i) \in \mathcal{D} | e_i \ge \tau, y'_i = \tilde{y}_i \}$$
(3)

 $\mathcal{D}_l = \mathcal{D}_{CSS} \cup \mathcal{D}_{MHCS}$ 



(b) Progressive Selection.

# Method/Debiased Semi-Supervised Training







#### Motivation:

- 1. to utilizes the remaining noisy samples to boost performance.
- 2. Distribution Bias: The selected clean samples may imbalance, since some labels are typically more ambiguous than others. Bias towards domain classes.
- 3. Confirmation Bias: Pseudo-labels bring confirmation bias.







Motivation:

- 1. skewed label distribution.
- 2. pseudo-labels can naturally bias towards some easy classes even if training data are balanced.

a Debiased Margin-based Loss (DML) to encourage a larger margin between the sample-rich classes and the sample-deficient classes.

$$l_{\text{DML}} = -\sum_{j=1}^{C} \tilde{y_i}^j \log \frac{e^{f^j(\boldsymbol{x}_i) + \alpha \cdot \log \pi_j}}{\sum_{k=1}^{C} e^{f^k(\boldsymbol{x}_i) + \alpha \cdot \log \pi_k}} \quad (5)$$

we suppress the logits on those easy classes and ensure other classes are fairly learned:

$$\tilde{f}_i = f(\boldsymbol{x}_i) - \alpha \log \boldsymbol{\pi} \tag{6}$$

$$\boldsymbol{\pi} = m\boldsymbol{\pi} + (1-m)\frac{1}{|\boldsymbol{\mathcal{B}}|} \sum_{\boldsymbol{x}_i \in \boldsymbol{\mathcal{B}}} \boldsymbol{p}_i \tag{7}$$

$$\mathcal{L}_{total} = \mathcal{L}_{cls} + \gamma (\mathcal{L}_{cr} + \mathcal{L}_{mix}) \tag{9}$$



$$\mathcal{L}_{\text{LDAM}}((x,y);f) = -\log \frac{e^{z_y - \Delta_y}}{e^{z_y - \Delta_y} + \sum_{j \neq y} e^{z_j}}$$
(12)  
where  $\Delta_j = \frac{C}{n_j^{1/4}} \text{ for } j \in \{1, \dots, k\}$ (13)

[1] Learning imbalanced datasets with label-distribution-aware margin loss. NeurlPS2019

## **Experiments**



Dataset	CIFAR-10					CIFAR-100			
	Sym.				Asym.	Sym.			
Methods\Noise Ratio	20%	50%	80%	90%	40%	20%	50%	80%	90%
CE	86.8	79.4	62.9	42.7	85.0	62.0	46.7	19.9	10.1
Co-Teaching+	89.5	85.7	67.4	47.9	-	65.6	51.8	27.9	13.7
JoCoR	85.7	79.4	27.8	-	76.4	53.0	43.5	15.5	-
M-correction	94.0	92.0	86.8	69.1	87.4	73.9	66.1	48.2	24.3
PENCIL	92.4	89.1	77.5	58.9	88.5	69.4	57.5	31.1	15.3
DivideMix	96.1	94.6	93.2	76.0	93.4	77.3	74.6	60.2	31.5
ELR+	95.8	94.8	93.3	78.7	93.0	77.6	73.6	60.8	33.4
LongReMix	96.2	95.0	93.9	82.0	94.7	77.8	75.6	62.9	33.8
MOIT	94.1	91.1	75.8	70.1	93.2	75.9	70.1	51.4	24.5
SOP+	96.3	95.5	94.0	-	93.8	78.8	75.9	63.3	-
PES(semi)	95.9	95.1	93.1	-	-	77.4	74.3	61.6	-
ULC	96.1	95.2	94.0	86.4	94.6	77.3	74.9	61.2	34.5
ProMix(last) ProMix(best)	97.59 97.69	97.30 97.40	95.05 95.49	91.13 93.36	96.51 96.59	82.39 82.64	79.72 80.06	68.95 69.37	42.74 42.93

Table 1: Accuracy comparisons on CIFAR-10/100 with symmetric (20%-90%) and asymmetric noise (40%). We report both the averaged test accuracy over last 10 epochs and the best accuracy of ProMix. Results of previous methods are the best test accuracies cited from their original papers, where the blank ones indicate that the corresponding results are not provided. **Bold entries** indicate superior results.



40.21%	40.20%
--------	--------

- 1 LGA could boost performance, especially under severe noise
- ② Because of inadequate labeled samples.
- (3) Employs clean and disregard unselected.
- (4) Generating and utilizing the pseudo-labels on the same classification head.
- (5) Sticks to the vanilla pseudo-labeling and adopts the standard cross entropy loss.

Ablation	LGA	Selection Strategy	CIFAR-10 Sym.80%	CIFAR-100 Sym.80%	CIFAR-10N Worst	CIFAR-100N Noisy Fine
ProMix	<ul> <li>✓</li> </ul>	CSS+MHCS	95.05	68.95	96.34	73.79
w/o LGA	×	CSS+MHCS	94.32	61.19	95.57	73.10
w/o MHCS	×	CSS	93.17	60.35	95.11	70.40
w/o Base Selection	×	MHCS	39.59	21.01	74.09	40.75
w/o CBR	×	CSS+MHCS	93.96	60.72	95.26	72.61
w/o DBR	×	CSS+MHCS	94.07	60.91	95.07	72.51
with Only Clean	×	CSS+MHCS	93.82	60.35	95.18	72.06

Table 5: Ablation study of ProMix on CIFAR-10-Symmetric 80%, CIFAR-100-Symmetric 80%, CIFAR-10N-Worst and CIFAR-100N-Noisy. ProMix with Only Clean denotes training with merely selected samples. CBR/DBR indicates Confirmation/Distribution Bias Removal.

## **Experiments**/Ablation





(a) CIFAR-10-Symmetric 80%. (b) CIFAR-100-Symmetric 80%.

Figure 3: Comparison of clean sample selection on CIFAR-10/100 dataset with 80% symmetric noise. The threshold of DivideMix and forget rate of JoCoR have been re-tuned to get the best F1 score.



#### (a) CIFAR-10-Symmetric 80%. (b) CIFAR-100-Symmetric 80%.

Figure 4: Accuracy of pseudo-labels for the unlabeled data with and without debiasing on CIFAR-10/100 with 80% Symmetric noise. LGA is disabled in order to avoid unexpected interference.





NANJING UNIVERSITY OF AERONAUTICS AND ASTRONAUTICS

# Thank you