



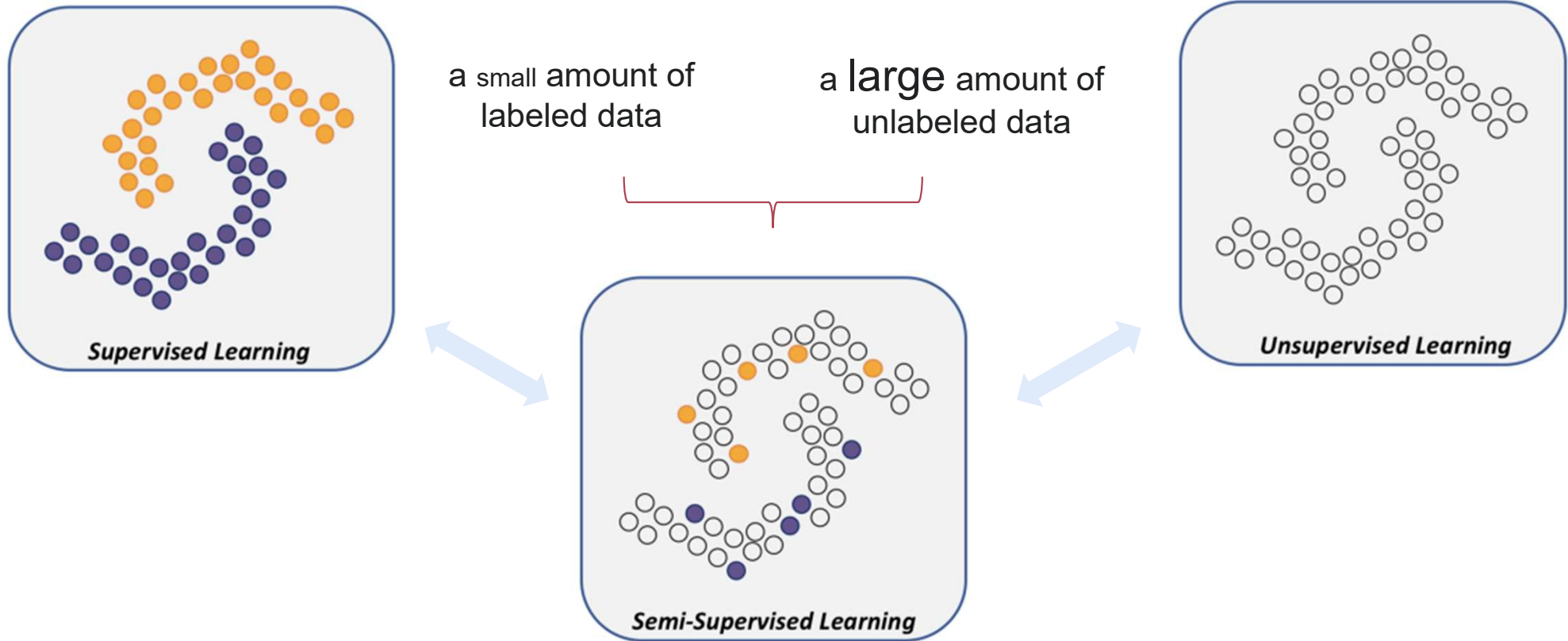
模式分析与机器智能
工业和信息化部重点实验室
MIT Key Laboratory of
Pattern Analysis & Machine Intelligence

ParNeC | 模式识别与神经计算研究组
Pattern Recognition and Neural Computing

IMBALANCED SEMI-SUPERVISED LEARNING

Background

Semi-supervised Learning (SSL)



Imbalanced Semi-supervised Learning (SSL)

assumption: The distribution of labeled and/or unlabeled data are class-balanced.



assumption: The labeled and unlabeled data are in the same class-imbalance distribution.

The test set is class-balanced

ABC: Auxiliary Balanced Classifier for Class-imbalanced Semi-supervised Learning

Hyuck Lee Seungjae Shin Heeyoung Kim
Department of Industrial and Systems Engineering, KAIST
{dlgur0921, tmdwo0910, heeyoungkim}@kaist.ac.kr

NIPS 2021

Related work

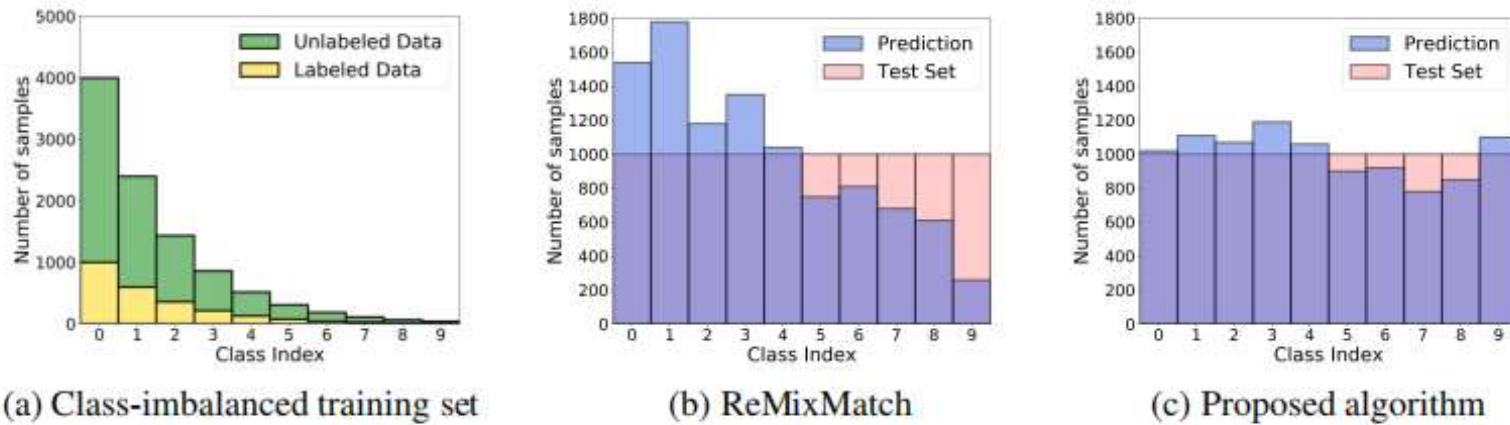


Figure 1: Predictions on a class-balanced test set using ReMixMatch (b) and the proposed algorithm (c) trained on a class-imbalanced training set (a).

Experiments on CIFAR-10-LT. Training set is class-imbalanced with **imbalance ratio $\gamma = 100$** .

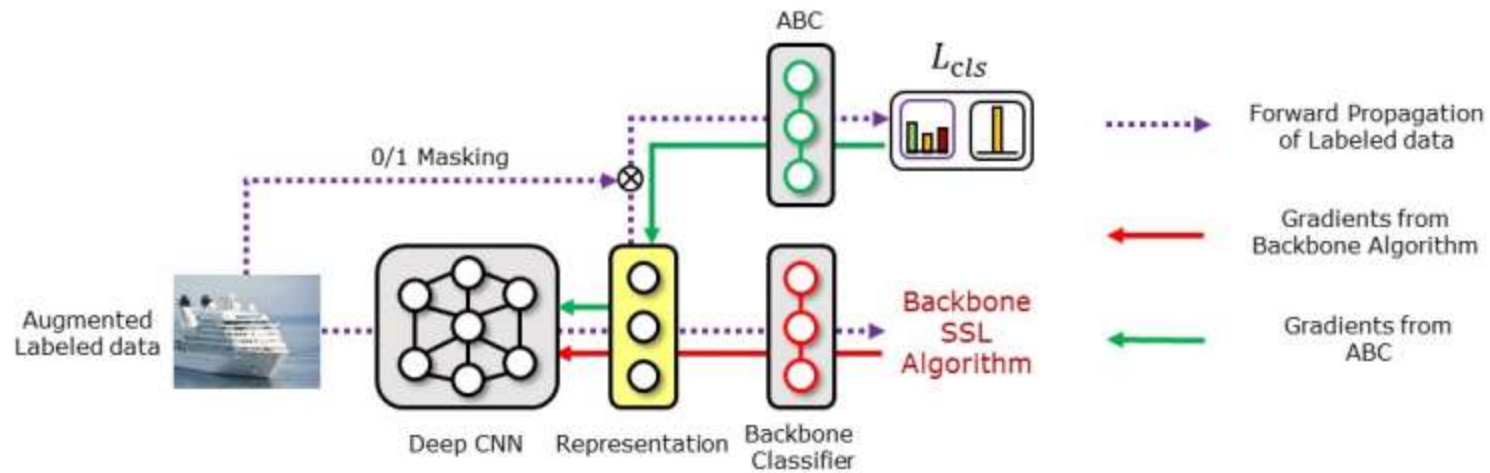


Figure 2: Overall procedure for balanced training of the ABC with a 0/1 mask

$$L_{cls} = \frac{1}{B} \sum_{b=1}^B M(x_b) \mathbf{H}(p_s(y|\alpha(x_b)), p_b),$$

$$M(x_b) = \mathcal{B}\left(\frac{N_L}{N_{y_b}}\right),$$

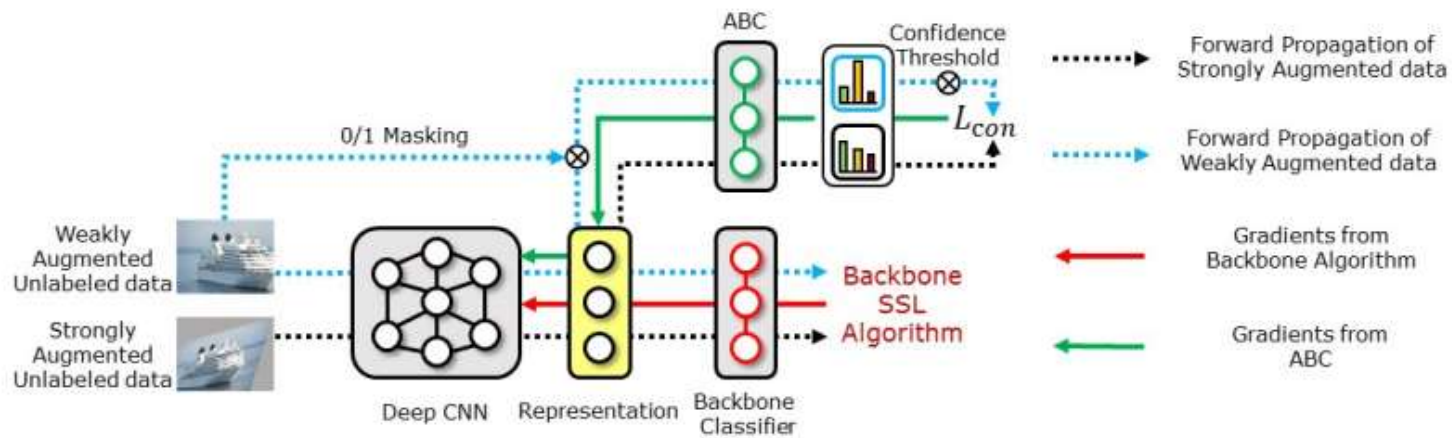


Figure 5: Overall procedure of consistency regularization for the ABC

$$L_{con} = \frac{1}{B} \sum_{b=1}^B \sum_{k=1}^2 M(u_b) \mathbf{I}(\max(q_b) \geq \tau) \mathbf{H}(p_s(y|\mathcal{A}_k(u_b)), q_b),$$

$$M(u_b) = \mathcal{B}\left(\frac{N_L}{N_{\hat{q}_b}}\right),$$

Method

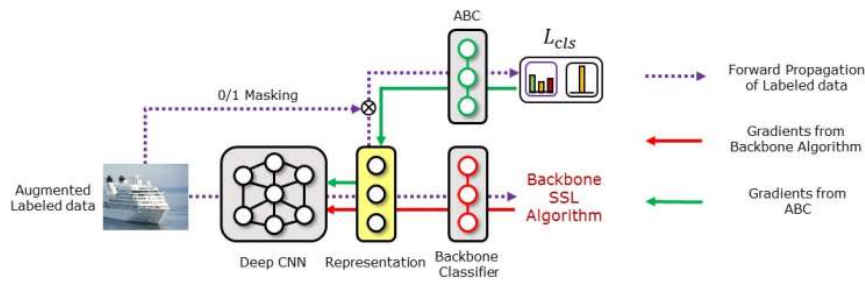


Figure 2: Overall procedure for balanced training of the ABC with a 0/1 mask

$$L_{cls} = \frac{1}{B} \sum_{b=1}^B M(x_b) \mathbf{H}(p_s(y|\alpha(x_b)), p_b),$$

$$M(x_b) = \mathcal{B}\left(\frac{N_L}{N_{y_b}}\right),$$

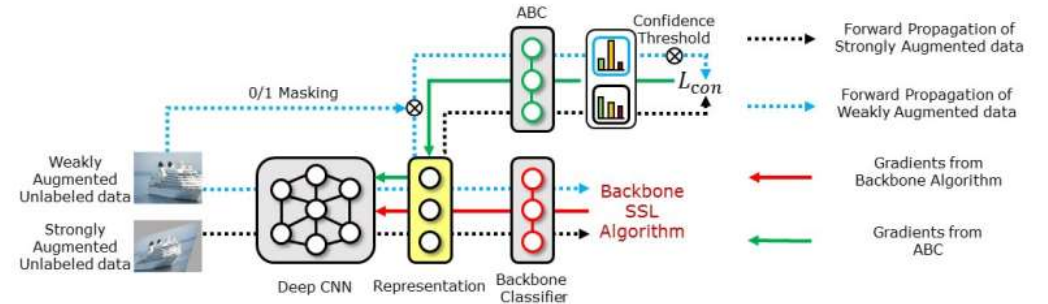


Figure 5: Overall procedure of consistency regularization for the ABC

$$L_{con} = \frac{1}{B} \sum_{b=1}^B \sum_{k=1}^2 M(u_b) \mathbf{I}(\max(q_b) \geq \tau) \mathbf{H}(p_s(y|\mathcal{A}_k(u_b)), q_b),$$

$$M(u_b) = \mathcal{B}\left(\frac{N_L}{N_{q_b}}\right),$$

$$L_{total} = L_{cls} + L_{con} + L_{back}.$$

Algorithm 1 Pseudo code of the proposed algorithm

Input: $\mathcal{MB}_{\mathcal{X}} = \{(x_b, y_b) : b \in (1, \dots, B)\} \subset \mathcal{X}$, $\mathcal{MB}_{\mathcal{U}} = \{(u_b) : b \in (1, \dots, B)\} \subset \mathcal{U}$
Output: Classification model $f : \mathbb{R}^d \rightarrow \{1, \dots, L\}$
Parameters : θ (Parameters of Wide ResNet-28-2 and ABC)

- 1: **while** Training **do**
- 2: **for** $b = 1$ to B **do**
- 3: $\alpha(x_b) = \text{Augment}(x_b)$
- 4: $\alpha(u_b) = \text{WeakAugment}(u_b)$
- 5: $\mathcal{A}_k(u_b) = \text{StrongAugment}_k(x_b)$, $k = 1, 2$
- 6: Predicted class distribution for $\alpha(x_b) = p_s(y|\alpha(x_b))$
- 7: Generate 0/1 mask $M(x_b)$.
- 8: Soft pseudo label $q_b = p_s(y|\alpha(u_b))$
- 9: **if** $\max(q_b) \geq 0.95$ **then**
- 10: Predicted class distribution for $\mathcal{A}_k(u_b) = p_s(y|\mathcal{A}_k(u_b))$, $k = 1, 2$
- 11: Generate 0/1 mask $M(u_b)$.
- 12: **end if**
- 13: Loss from the backbone $L_{back} += \text{backbone}(\alpha(x_b), \alpha(u_b), \mathcal{A}_k(u_b))$
- 14: **end for**
- 15: Calculate the classification loss L_{cls} and consistency regularization loss L_{con} .
- 16: Total Loss $L_{total} = L_{cls} + L_{con} + L_{back}$
- 17: $\Delta\theta \propto \nabla_{\theta} L_{total}$, $\theta \leftarrow \theta + \Delta\theta$
- 18: **end while**

Table 1: Overall accuracy/minority-class-accuracy under the main setting

Algorithm	CIFAR-10-LT	SVHN-LT	CIFAR-100-LT
	$\gamma = 100, \beta = 20\%$	$\gamma = 100, \beta = 20\%$	$\gamma = 20, \beta = 40\%$
Vanilla	55.3 \pm 1.30 / 33.9 \pm 1.88	77.0 \pm 0.67 / 63.3 \pm 1.25	40.1 \pm 1.15 / 25.2 \pm 0.95
VAT [24]	55.3 \pm 0.88 / 28.2 \pm 1.55	81.3 \pm 0.47 / 68.2 \pm 0.88	40.4 \pm 0.34 / 24.8 \pm 0.38
BALMS [27]	70.7 \pm 0.59 / 69.8 \pm 1.03	87.6 \pm 0.53 / 85.0 \pm 0.67	50.2 \pm 0.54 / 42.9 \pm 1.03
FixMatch [29]	72.3 \pm 0.33 / 53.8 \pm 0.63	88.0 \pm 0.30 / 79.4 \pm 0.54	51.0 \pm 0.20 / 32.8 \pm 0.41
w/ CReST+PDA [34]	76.6 \pm 0.46 / 61.4 \pm 0.85	89.1 \pm 0.69 / 81.7 \pm 1.18	51.6 \pm 0.29 / 36.4 \pm 0.46
w/ DARP [18]	73.7 \pm 0.98 / 57.0 \pm 2.12	88.6 \pm 0.19 / 80.5 \pm 0.54	51.4 \pm 0.37 / 33.9 \pm 0.77
w/ DARP+cRT [18]	78.1 \pm 0.89 / 66.6 \pm 1.55	89.9 \pm 0.44 / 83.5 \pm 0.61	54.7 \pm 0.46 / 41.2 \pm 0.42
w/ ABC	81.1\pm0.82 / 72.0\pm1.77	92.0\pm0.38 / 87.9\pm0.73	56.3\pm0.19 / 43.4\pm0.42
ReMixMatch [3]	73.7 \pm 0.39 / 55.9 \pm 0.87	89.8 \pm 0.42 / 82.8 \pm 0.68	54.0 \pm 0.29 / 37.1 \pm 0.37
w/ CReST+PDA [34]	75.7 \pm 0.34 / 59.6 \pm 0.76	90.9 \pm 0.20 / 85.2 \pm 0.39	54.6 \pm 0.48 / 38.1 \pm 0.69
w/ DARP [18]	74.4 \pm 0.41 / 56.9 \pm 0.67	90.2 \pm 0.22 / 83.5 \pm 0.40	54.5 \pm 0.33 / 37.7 \pm 0.58
w/ DARP+cRT [18]	78.5 \pm 0.61 / 66.4 \pm 1.68	92.1 \pm 0.48 / 87.6 \pm 0.75	55.1 \pm 0.45 / 43.6 \pm 0.58
w/ ABC	82.4\pm0.45 / 75.7\pm1.18	93.9\pm0.16 / 92.5\pm0.4	57.6\pm0.26 / 46.7\pm0.50

Table 2: Overall accuracy/minority-class accuracy for CIFAR-10 under various settings

CIFAR-10-LT				
Algorithm	$\gamma = 100, \beta = 10\%$	$\gamma = 100, \beta = 30\%$	$\gamma = 50, \beta = 20\%$	$\gamma = 150, \beta = 20\%$
FixMatch [29]	70.0 \pm 0.59 / 48.9 \pm 1.04	74.9 \pm 0.63 / 58.2 \pm 1.28	81.2 \pm 0.07 / 70.7 \pm 0.36	68.5 \pm 0.60 / 45.8 \pm 1.15
w/ CReST+PDA [34]	73.9 \pm 0.40 / 58.9 \pm 1.14	77.6 \pm 0.73 / 64.0 \pm 1.39	83.3 \pm 0.10 / 75.7 \pm 0.39	70.0 \pm 0.82 / 49.4 \pm 1.52
w/ DARP+cRT [18]	74.6 \pm 0.98 / 59.2 \pm 2.12	79.0 \pm 0.25 / 67.7 \pm 0.95	83.6 \pm 0.42 / 77.1 \pm 1.19	73.2 \pm 0.85 / 57.1 \pm 1.13
w/ ABC	77.2\pm1.60 / 65.7\pm2.85	81.5\pm0.29 / 72.9\pm0.96	85.2\pm0.51 / 80.2\pm0.64	77.1\pm0.46 / 64.4\pm0.92
ReMixMatch [3]	71.5 \pm 0.51 / 52.2 \pm 1.08	75.8 \pm 0.10 / 59.4 \pm 0.17	81.5 \pm 0.17 / 70.7 \pm 0.32	69.9 \pm 0.23 / 48.4 \pm 0.60
w/ CReST+PDA [34]	73.8 \pm 0.32 / 56.6 \pm 0.43	78.6 \pm 0.73 / 64.8 \pm 1.49	83.9 \pm 0.26 / 75.4 \pm 0.52	71.3 \pm 0.77 / 50.8 \pm 1.59
w/ DARP+cRT [18]	75.9 \pm 1.20 / 62.1 \pm 3.10	81.0 \pm 0.16 / 70.7 \pm 0.72	84.5 \pm 0.80 / 77.8 \pm 1.67	73.9 \pm 0.59 / 57.4 \pm 1.45
w/ ABC	79.8\pm0.36 / 70.8\pm0.92	84.3\pm1.03 / 80.6\pm0.97	87.5\pm0.31 / 84.6\pm1.19	80.6\pm0.66 / 72.1\pm1.51

Table 4: Overall accuracy/minority-class accuracy for the large-scale LSUN dataset

LSUN, $\gamma = 100, \beta = 20\%$					
Algorithm	w/ -	w/ cRT [17]	w/ DARP [18]	w/ DARP+cRT [18]	w/ ABC
FixMatch [29]	73.1 / 55.3	77.0 / 71.5	71.0 / 51.8	75.8 / 69.5	78.9 / 75.5
ReMixMatch [3]	69.4 / 49.1	75.4 / 69.5	65.6 / 44.1	72.1 / 67.5	76.9 / 69.5

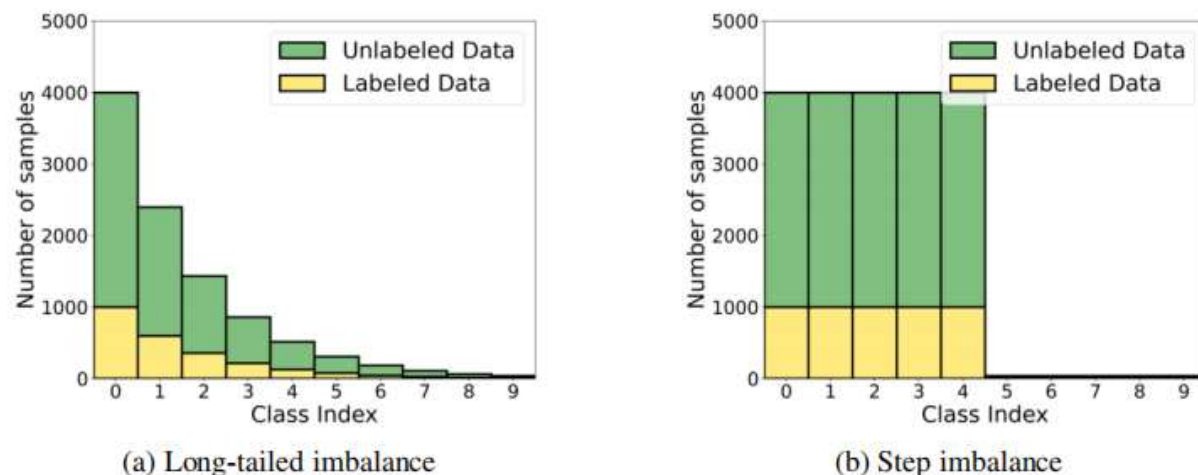


Figure 6: Long-tailed imbalance and step imbalance

Table 3: Overall accuracy/minority-class accuracy on CIFAR-10 under a step imbalance setting

CIFAR-10-Step, $\gamma = 100$, $\beta = 20\%$				
Algorithm	w/ -	w/ CReST+PDA [34]	w/ DARP+cRT [18]	w/ ABC
FixMatch [29]	54.0 \pm 0.84 / 11.8 \pm 1.71	71.1 \pm 0.78 / 48.2 \pm 2.26	69.8 \pm 1.51 / 45.1 \pm 2.70	75.9\pm0.49 / 57.0\pm1.07
ReMixMatch [3]	60.8 \pm 0.10 / 25.1 \pm 1.28	64.6 \pm 0.97 / 33.5 \pm 2.05	72.3 \pm 1.77 / 50.6 \pm 3.53	76.4\pm1.70 / 65.7\pm1.30

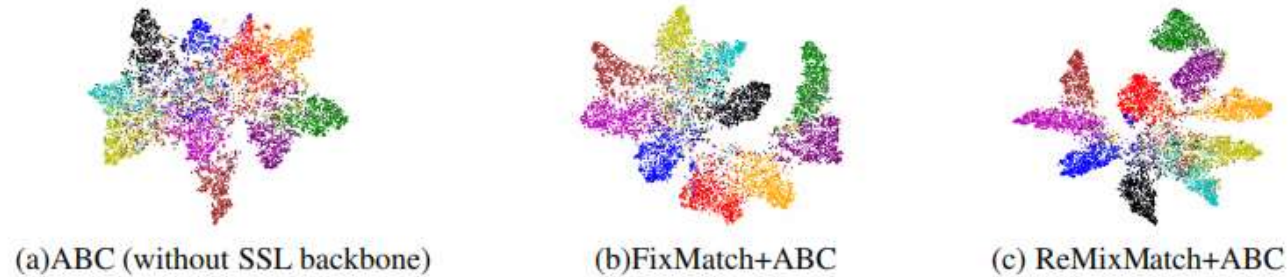


Figure 3: t-SNE of the proposed algorithm and the ABC (without SSL backbone)

Table 5: Ablation study for ReMixMatch+ABC on CIFAR-10-LT, $\gamma = 100$, $\beta = 20\%$

Ablation study	Overall	Minority
ReMixMatch+ABC (proposed algorithm)	82.4	75.7
Without gradually decreasing the parameter of $\mathcal{B}(\cdot)$ for consistency regularization	81.8	74.6
Without consistency regularization for the ABC	79.4	66.9
Without using the 0/1 mask for the consistency regularization loss L_{con}	79.0	69.2
Without using the 0/1 mask for the classification loss L_{cls}	74.4	57.8
Without using the confidence threshold τ for consistency regularization	74.3	75.4
Using hard pseudo labels for consistency regularization	70.2	75.1
Without training backbone (ABC without SSL backbone)	68.7	56.2
Training the ABC with a re-weighting technique	81.2	74.1
Decoupled training of the backbone and ABC	79.5	72.3

Imbalanced Semi-supervised Learning (SSL)

assumption: The distribution of labeled and/or unlabeled data are class-balanced.



assumption: The labeled and unlabeled data are in the different class-imbalance distribution.

The test set is class-balanced



模式分析与机器智能
工业和信息化部重点实验室
MIT Key Laboratory of
Pattern Analysis & Machine Intelligence

ParNeC | 模式识别与神经计算研究组
Pattern Recognition and Neural Computing

IMBALANCED SEMI-SUPERVISED LEARNING WITH BIAS ADAPTIVE CLASSIFIER

Renzhen Wang¹, Xixi Jia², Quanziang Wang¹, Yichen Wu³, Deyu Meng^{1,4,5*}

¹Xi'an Jiaotong University, ²Xidian University, ³City University of Hong Kong

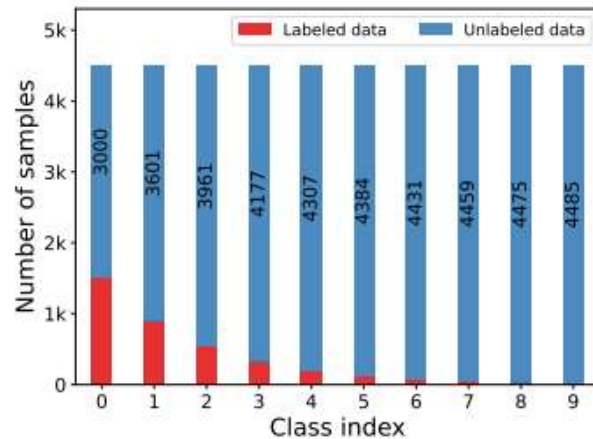
⁴Macau University of Science and Technology, ⁵Peng Cheng Laboratory

{rzwang, dymeng}@mail.xjtu.edu.cn

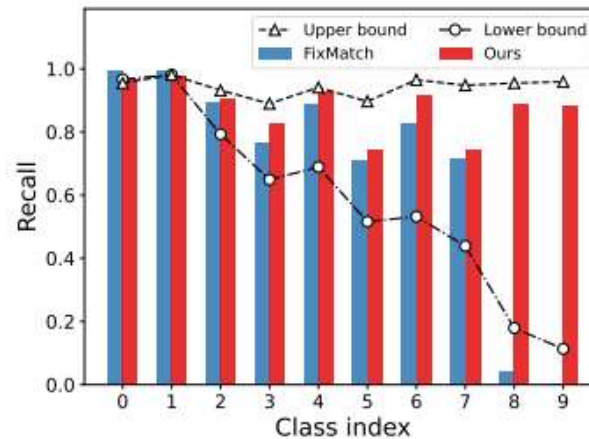
ICLR 2023

Related work

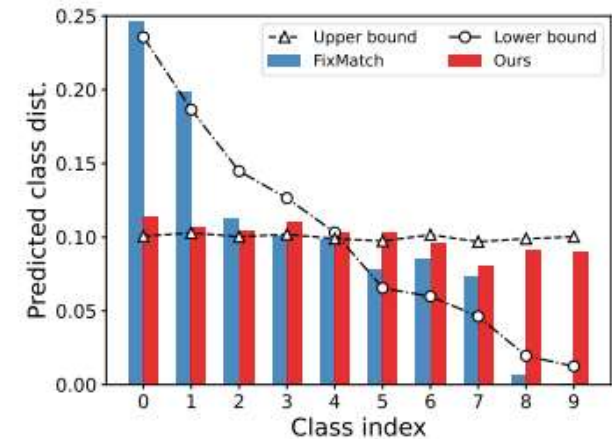
Imbalanced Semi-supervised Learning



(a) Class distribution



(b) Per-class recall

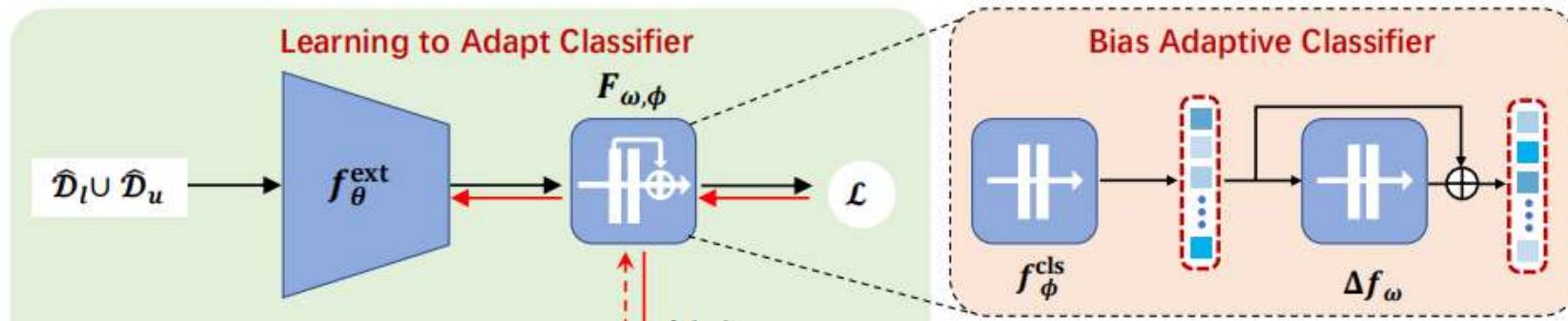


(c) Predicted class distribution

Experiments on CIFAR-10-LT. Labeled set is class-imbalanced with **imbalance ratio $\gamma = 100$** , while the **whole training data remains balanced**.

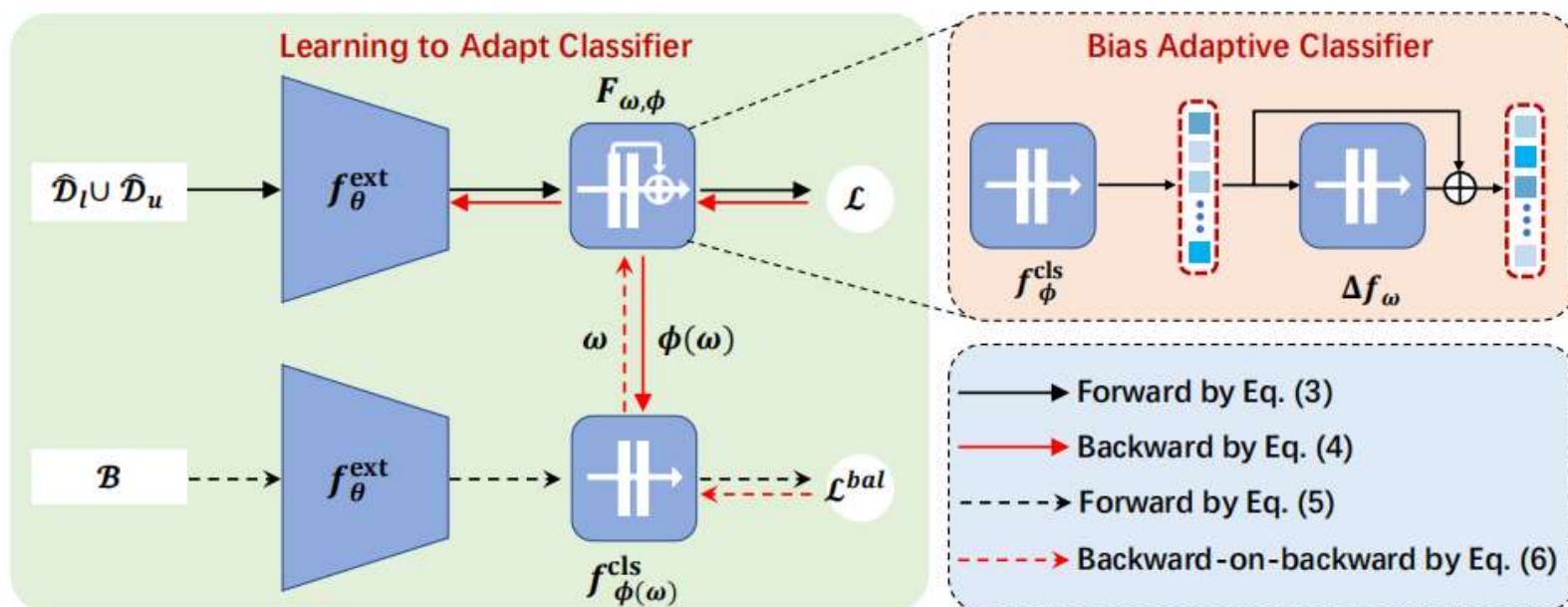
Bias Adaptive Classifier

$$F_{\omega, \phi}(\mathbf{z}) = (\mathbf{I} + \Delta f_{\omega}) \circ f_{\phi}^{\text{cls}}(\mathbf{z}),$$



$$\mathcal{L} = \frac{1}{|\hat{\mathcal{D}}_l|} \sum_{\mathbf{x}_i \in \hat{\mathcal{D}}_l} H(\hat{\Psi}(\mathbf{x}_i), \mathbf{y}_i) + \frac{1}{|\hat{\mathcal{D}}_u|} \sum_{\mathbf{x}_i \in \hat{\mathcal{D}}_u} \lambda_i H(\hat{\Psi}(\mathbf{x}_i), \hat{\mathbf{y}}_i).$$

Bias Adaptive Classifier



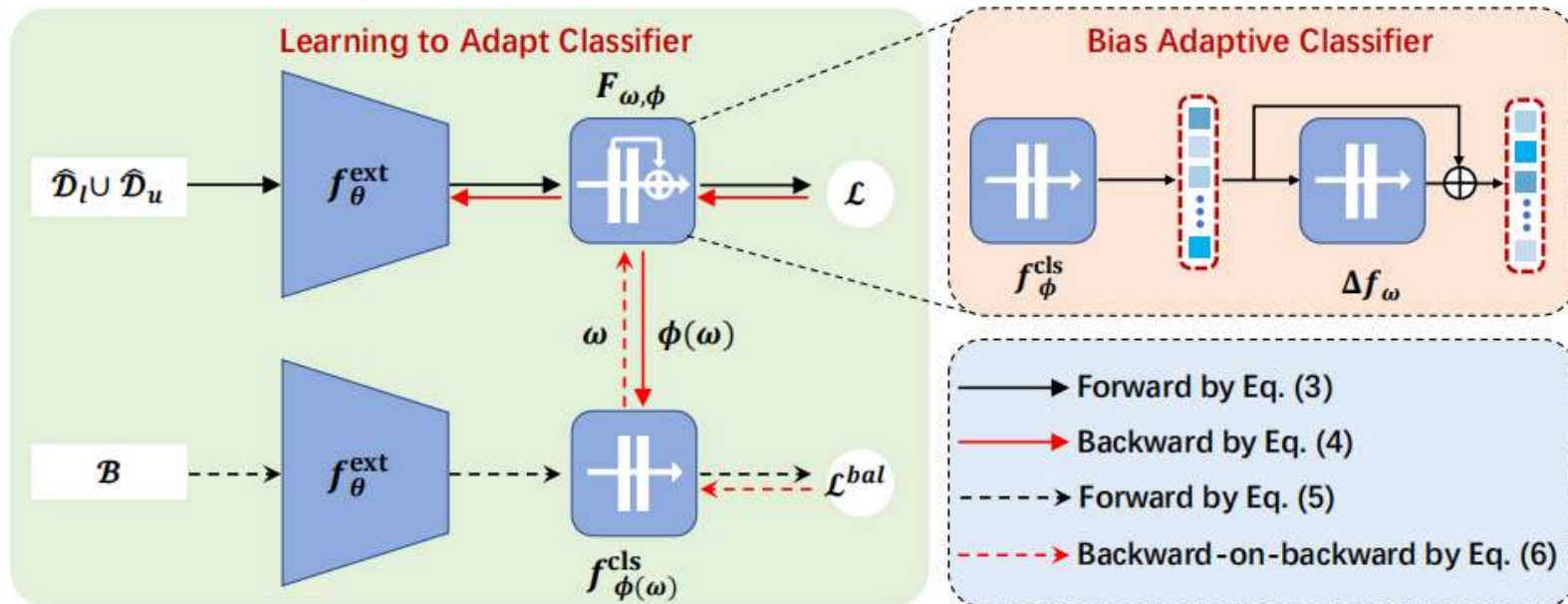
$$\mathcal{L} = \frac{1}{|\hat{\mathcal{D}}_l|} \sum_{\mathbf{x}_i \in \hat{\mathcal{D}}_l} H(\hat{\Psi}(\mathbf{x}_i), \mathbf{y}_i) + \frac{1}{|\hat{\mathcal{D}}_u|} \sum_{\mathbf{x}_i \in \hat{\mathcal{D}}_u} \lambda_i H(\hat{\Psi}(\mathbf{x}_i), \hat{\mathbf{y}}_i). \quad (3)$$

$$(\theta^{t+1}, \phi^{t+1}(\omega)) = (\theta^t, \phi^t) - \alpha \nabla_{\theta, \phi} \mathcal{L}, \quad (4)$$

$$\mathcal{L}^{\text{bal}} = \frac{1}{|\mathcal{B}|} \sum_{\mathbf{x}_i \in \mathcal{B}} H(f_{\phi^{t+1}}^{\text{cls}} \circ f_{\theta^{t+1}}^{\text{ext}}(\mathbf{x}_i), \mathbf{y}_i), \quad (5)$$

$$\omega^{t+1} = \omega^t - \eta \nabla_{\omega} \mathcal{L}^{\text{bal}}, \quad (6)$$

Complexity Analysis



However, in the backward pass, the second-order gradient of ω only requires to unroll the gradient graph of the linear classifier f_{ϕ}^{cls} .

$$\frac{\#\text{Params}(f_{\phi}^{\text{cls}})}{\#\text{Params}(\Psi)} \times \text{training time of one full backward pass.}$$

Theoretical Analysis

$$\mathcal{L}^{bal} = \frac{1}{|\mathcal{B}|} \sum_{\mathbf{x}_i \in \mathcal{B}} \text{H}(f_{\phi^{t+1}(\omega)}^{\text{cls}} \circ f_{\theta^{t+1}}^{\text{ext}}(\mathbf{x}_i), \mathbf{y}_i),$$

As the training loss is

$$\mathcal{L}_i(\theta, \phi) = \log \sum_{k=1}^d e^{z_{i,k} + \Delta f_{i,k}} - z_{i,c_i} - \Delta f_{i,c_i},$$

where c_i is the class label of the i -th sample. Therefore

$$\Xi_{i,k} = \begin{cases} \frac{e^{z_{i,k} + \Delta f_{i,k}}}{\sum_{s=1}^d e^{z_{i,s} + \Delta f_{i,s}}}, & k \neq c_i \\ \frac{e^{z_{i,k} + \Delta f_{i,k}}}{\sum_{s=1}^d e^{z_{i,s} + \Delta f_{i,s}}} - 1, & k = c_i \end{cases} = \mathbf{p}_i - \mathbf{y}_i$$

Meanwhile, the upper level loss is defined as

$$\mathcal{L}^{bal}(\theta, \phi) = -\frac{1}{m} \sum_{j=1}^m \mathcal{L}_j^{bal}(\theta, \phi),$$



$$G_i = \frac{\partial(\mathbf{p}_i - \mathbf{y}_i)}{\partial \phi} \Big|_{\phi^t}^T \left(\frac{1}{m} \sum_{j=1}^m \frac{\partial \mathcal{L}_j^{bal}(\theta, \phi)}{\partial \phi} \Big|_{\phi^{t+1}} \right)$$

In consequence, we have

$$\begin{aligned} \omega^{t+1} &= \omega^t - \eta \nabla \mathcal{L}^{bal}(\theta, \phi^{t+1}) \Big|_{\omega^t} \\ &= \omega^t + \eta \alpha \frac{1}{n} \sum_{i=1}^n \frac{\partial^2 \mathcal{L}_i(\theta, \phi)}{\partial \phi \partial \omega} \Big|_{\phi^t, \omega^t} \frac{\partial \mathcal{L}^{bal}(\theta, \phi)}{\partial \phi} \Big|_{\phi^{t+1}} \\ &= \omega^t + \eta \alpha \frac{1}{n} \sum_{i=1}^n \frac{\partial \Delta f_i}{\partial \omega} \Big|_{\omega^t} \left(\frac{\partial^2 \mathcal{L}_i(\theta, \phi)}{\partial \phi \partial \Delta f_i} \Big|_{\phi^t} \frac{\partial \mathcal{L}^{bal}(\theta, \phi)}{\partial \phi} \Big|_{\phi^{t+1}} \right), \end{aligned}$$

Denote by $\Xi_i = \frac{\partial \mathcal{L}_i(\theta, \phi)}{\partial \Delta f_i}$, then Eq. (9) becomes

$$\omega^{t+1} = \omega^t + \eta \alpha \frac{1}{n} \sum_{i=1}^n G_i \frac{\partial \Delta f_i}{\partial \omega} \Big|_{\omega^t},$$

where

$$G_i = \left\langle \frac{\partial \Xi_i}{\partial \phi} \Big|_{\phi^t}, \frac{\partial \mathcal{L}^{bal}(\theta, \phi)}{\partial \phi} \Big|_{\phi^{t+1}} \right\rangle.$$

Algorithm 1 learning to adapt classifier during training

Input: labeled / unlabeled training data $\mathcal{D}_l / \mathcal{D}_u$, labeled / unlabeled batch size n / m , max iterations T

Output: classification network parameters $\{\theta, \phi\}$

- 1: Initialize $\{\theta^0, \phi^0\} \leftarrow \{\theta, \phi\}$ and $\omega^0 \leftarrow \omega$.
 - 2: **for** $t = 0$ **to** T **do**
 - 3: $\hat{\mathcal{D}}_l = \{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^n \leftarrow \text{SampleMiniBatch}(\mathcal{D}_l, n)$.
 - 4: $\hat{\mathcal{D}}_u = \{\mathbf{x}_i\}_{i=1}^m \leftarrow \text{SampleMiniBatch}(\mathcal{D}_u, m)$.
 - 5: $\mathcal{B} = \{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^n \leftarrow \text{SampleMiniBatch}(\mathcal{D}_l, n)$.
 - 6: Estimate pseudo-label \hat{y}_i for $\mathbf{x}_i \in \hat{\mathcal{D}}_u$.
 - 7: Compute lower-level loss \mathcal{L} by Eq. (3).
 - 8: Update network parameters $\{\theta^{t+1}, \phi^{t+1}\}$ by Eq. (4).
 - 9: Compute upper-level loss \mathcal{L}^{bal} by Eq. (5).
 - 10: Update bias attractor parameters ω^{t+1} by Eq. (6).
 - 11: **end for**
-

Experiment

We adopt balanced accuracy (bACC) and geometric mean scores (GM) as the evaluation metrics.

Methods	CIFAR-10 ($\gamma_l = \gamma_u$)		CIFAR-10 ($\gamma_l \neq \gamma_u$)	
	$\gamma = 100$	$\gamma = 150$	$\gamma_u = 1$ (uniform)	$\gamma_u = 100$ (reversed)
Vanilla	58.8 \pm 0.13 / 51.0 \pm 0.11	55.6 \pm 0.43 / 44.0 \pm 0.98	58.8 \pm 0.13 / 51.0 \pm 0.11	58.8 \pm 0.13 / 51.0 \pm 0.11
w/ Re-sampling	55.8 \pm 0.47 / 45.1 \pm 0.30	52.2 \pm 0.05 / 38.2 \pm 1.49	55.8 \pm 0.47 / 45.1 \pm 0.30	55.8 \pm 0.47 / 45.1 \pm 0.30
w/ LDAM-DRW	62.8 \pm 0.17 / 58.9 \pm 0.60	57.9 \pm 0.20 / 50.4 \pm 0.30	62.8 \pm 0.17 / 58.9 \pm 0.60	62.8 \pm 0.17 / 58.9 \pm 0.60
w/ cRT	63.2 \pm 0.45 / 59.9 \pm 0.40	59.3 \pm 0.10 / 54.6 \pm 0.72	63.2 \pm 0.45 / 59.9 \pm 0.40	63.2 \pm 0.45 / 59.9 \pm 0.40
MixMatch	64.8 \pm 0.28 / 49.0 \pm 2.05	62.5 \pm 0.31 / 42.5 \pm 1.68	41.5 \pm 0.76 / 12.0 \pm 1.34	47.9 \pm 0.09 / 20.5 \pm 0.85
w/ DARP	67.9 \pm 0.14 / 61.2 \pm 0.15	65.8 \pm 0.52 / 56.5 \pm 2.08	86.7 \pm 0.80 / 86.2 \pm 0.82	72.9 \pm 0.24 / 71.0 \pm 0.32
w/ SaR	66.8 \pm 0.92 / 59.9 \pm 1.32	64.4 \pm 2.21 / 57.3 \pm 1.95	68.4 \pm 3.20 / 62.0 \pm 2.17	65.5 \pm 1.01 / 64.2 \pm 0.95
w/ DASO	69.8 \pm 1.10 / 69.3 \pm 1.07	66.5 \pm 1.99 / 65.4 \pm 2.25	75.5 \pm 0.48 / 74.6 \pm 0.67	65.7 \pm 1.01 / 62.0 \pm 1.23
w/ ABC	75.7 \pm 0.76 / 74.7 \pm 0.47	68.5 \pm 0.40 / 56.4 \pm 1.50	72.1 \pm 0.53 / 41.2 \pm 4.40	62.9 \pm 0.36 / 59.9 \pm 0.60
w/ L2AC (ours)	76.6\pm0.73 / 75.7\pm1.08	72.1\pm0.62 / 70.3\pm0.93	87.2\pm0.09 / 86.7\pm0.08	74.0\pm0.82 / 72.9\pm1.01
FixMatch	71.5 \pm 0.72 / 66.8 \pm 1.51	68.4 \pm 0.15 / 59.9 \pm 0.43	68.9 \pm 1.95 / 42.8 \pm 8.11	65.5 \pm 0.05 / 26.0 \pm 0.44
w/ DARP	75.5 \pm 0.05 / 73.0 \pm 0.09	70.4 \pm 0.25 / 64.9 \pm 0.17	85.4 \pm 0.55 / 85.0 \pm 0.65	74.9 \pm 0.51 / 72.3 \pm 1.13
w/ CReST+	77.5 \pm 0.15 / 76.1 \pm 0.15	72.1 \pm 0.74 / 68.9 \pm 1.29	N/A	N/A
w/ SaR	77.6 \pm 0.42 / 75.9 \pm 0.76	71.5 \pm 0.23 / 66.9 \pm 0.25	85.9 \pm 0.68 / 85.3 \pm 0.53	78.3 \pm 0.34 / 76.1 \pm 0.21
w/ DASO	78.3 \pm 0.55 / 76.5 \pm 0.57	74.6 \pm 0.74 / 71.7 \pm 0.52	87.9 \pm 0.41 / 87.7 \pm 0.43	79.5 \pm 0.91 / 78.9 \pm 0.96
w/ ABC	80.2 \pm 0.42 / 78.9 \pm 1.29	74.7 \pm 1.04 / 72.2 \pm 1.45	81.3 \pm 0.34 / 80.2 \pm 0.36	70.3 \pm 0.50 / 67.9 \pm 0.70
w/ L2AC (ours)	82.1\pm0.57 / 81.5\pm0.64	77.6\pm0.53 / 75.8\pm0.71	89.5\pm0.18 / 89.2\pm0.19	82.2\pm1.23 / 81.7\pm1.36

It can be observed that L2AC significantly improves MixMatch and FixMatch at least **9%** absolute gain on bACC and at least **14%** on GM for all settings.

Experiment

Methods	CIFAR-100 ($\gamma_l = \gamma_u$)		STL-10 ($\gamma_u = \text{N/A}$)	
	$\gamma_l = 10$	$\gamma_l = 20$	$\gamma_l = 10$	$\gamma_l = 20$
FixMatch	55.1 \pm 0.09 / 46.7 \pm 0.53	49.5 \pm 0.38 / 34.2 \pm 1.01	69.6 \pm 0.60 / 62.6 \pm 1.11	65.5 \pm 0.05 / 26.0 \pm 0.44
w/ DARP	56.3 \pm 0.25 / 48.2 \pm 0.73	50.2 \pm 0.18 / 36.0 \pm 0.60	72.9 \pm 0.24 / 69.5 \pm 0.18	74.9 \pm 0.51 / 72.3 \pm 1.13
w/ ABC	58.2 \pm 0.08 / 51.8 \pm 0.25	53.1 \pm 0.19 / 42.2 \pm 0.82	78.2 \pm 0.35 / 77.3 \pm 0.30	72.7 \pm 0.08 / 70.6 \pm 0.22
w/ DASO	58.3 \pm 0.39 / 51.4 \pm 0.80	53.0 \pm 0.27 / 39.5 \pm 1.45	78.2 \pm 0.63 / 77.4 \pm 0.53	75.4 \pm 0.81 / 74.4 \pm 1.00
w/ L2AC (ours)	57.8 \pm 0.19 / 52.1 \pm 0.31	52.6 \pm 0.13 / 43.0 \pm 0.45	79.9 \pm 0.52 / 79.1 \pm 0.49	77.0 \pm 0.65 / 75.8 \pm 0.68

It is worth noting that the unlabeled set of STL-10 is noisy as it contains samples that do not belong to any of the classes in the labeled set.

L2AC significantly outperforms ABC and DASO on both bACC and GM, which demonstrates that it has greater potential to be applied in the practical SSL scenarios.

Experiment

Methods	bACC/GM	Methods	bACC/GM
Vanilla	$38.3_{\pm 0.05} / 29.9_{\pm 0.08}$	DARP (Kim et al., 2020a)	$45.5_{\pm 0.32} / 37.5_{\pm 0.04}$
cRT (Kang et al., 2020)	$39.3_{\pm 0.21} / 33.7_{\pm 0.37}$	ABC (Lee et al., 2021)	$47.0_{\pm 0.26} / 39.2_{\pm 0.34}$
FixMatch (Sohn et al., 2020)	$44.9_{\pm 0.11} / 35.7_{\pm 0.66}$	L2AC (ours)	$48.8_{\pm 0.19} / 40.6_{\pm 0.17}$

This further verifies the efficacy of our proposed method toward the real-world imbalanced SSL applications.

Experiment

CIFAR-10 ($\gamma_l = \gamma_u = 100$)	$\beta = 1$	$\beta = 5$	$\beta = 10$	$\beta = 20$	$\beta = 30$
FixMatch	54.9 / 16.5	65.1 / 35.5	69.0 / 53.9	72.0 / 62.2	76.5 / 74.3
w/ L2AC (ours)	62.8 / 55.8	75.9 / 74.1	79.3 / 78.4	80.8 / 79.9	83.6 / 83.2
STL-10 ($\gamma_l = 10, \gamma_u = \text{N/A}$)	$\beta = 5$	$\beta = 10$	$\beta = 20$	$\beta = 40$	$\beta = 60$
FixMatch	46.5 / 19.9	48.8 / 27.0	58.2 / 39.8	67.2 / 60.7	69.2 / 67.6
w/ L2AC (ours)	62.8 / 57.1	66.5 / 62.9	72.6 / 70.6	77.0 / 75.7	78.8 / 77.9

In such an extremely biased scenario, our L2AC significantly improves FixMatch by around **16%** over bACC and **37%** over GM.

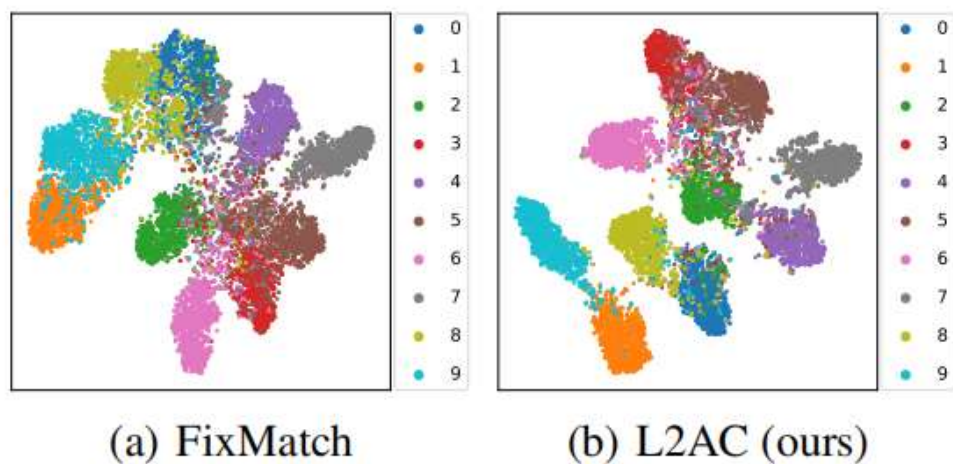


Figure 5: t-SNE visualization of training data for (a) FixMatch and (b) L2AC. L2AC helps to discriminate tail classes from majority ones.

Table 5: Ablation study.

Methods	CIFAR-10 ($\gamma_l = 100$)	
	$\gamma_u = 100$	1/100 (reversed)
FixMatch	71.5 / 68.8	65.5 / 26.0
FixMatch w/ bias attractor	73.9 / 70.7	66.6 / 44.8
L2AC w/o bi-level training	78.4 / 76.6	79.3 / 78.0
L2AC (ours)	82.1 / 81.5	82.2 / 81.7

Experiment

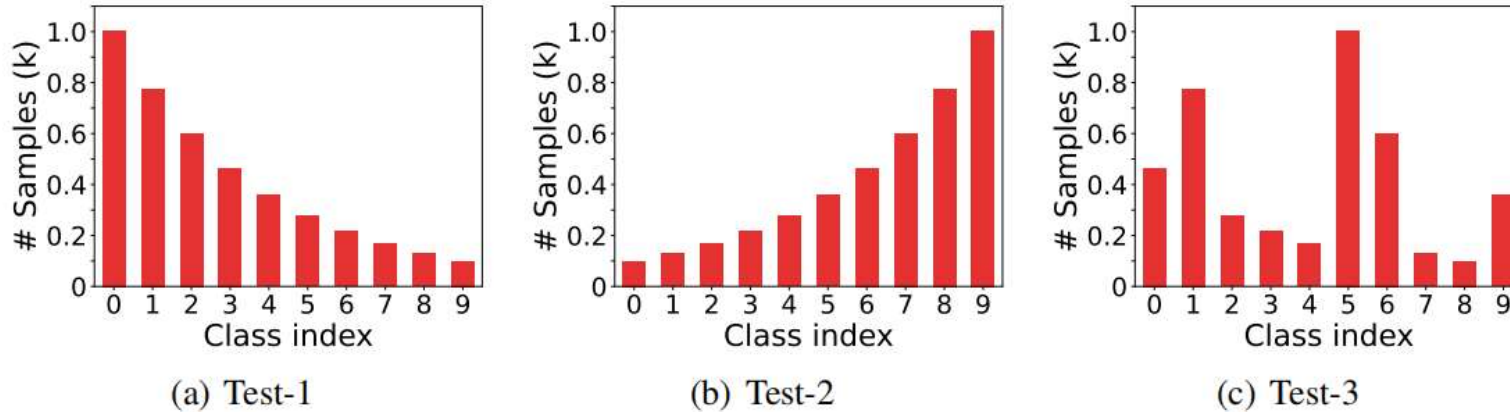


Table 6: Imbalanced test set results. ACC: accuracy for all samples.

Methods	Test-1			Test-2			Test-3		
	bACC	GM	ACC	bACC	GM	ACC	bACC	GM	ACC
FixMatch (Sohn et al., 2020)	72.4	66.3	86.1	72.7	66.9	56.6	72.3	65.7	75.6
w/ DARP (Kim et al., 2020a)	74.8	72.5	86.3	75.5	73.2	63.9	75.6	73.3	77.4
w/ CReST (Wei et al., 2021)	77.8	76.5	86.3	77.2	74.8	68.9	77.5	76.3	80.3
w/ ABC (Lee et al., 2021)	80.2	79.2	88.1	80.2	79.0	71.7	80.1	79.0	82.8
w/ L2AC (ours)	82.6	82.0	87.2	82.4	81.8	78.6	82.7	82.1	83.9

All these models are trained on CIFAR-10 with imbalanced ratio $\gamma = \gamma_l = \gamma_u = 100$.

Thanks