



ParNeC

模式识别与神经计算研究组

PAttern Recognition and NEural Computing

---

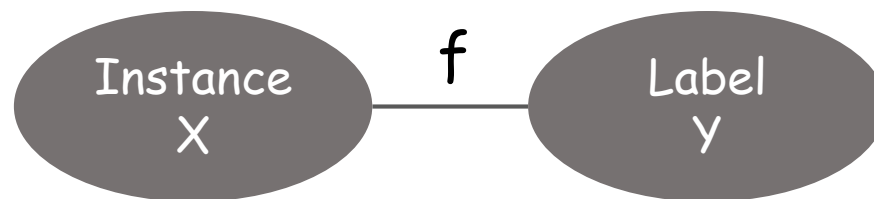
# Revisiting Consistency Regularization for Deep Partial Label Learning

---

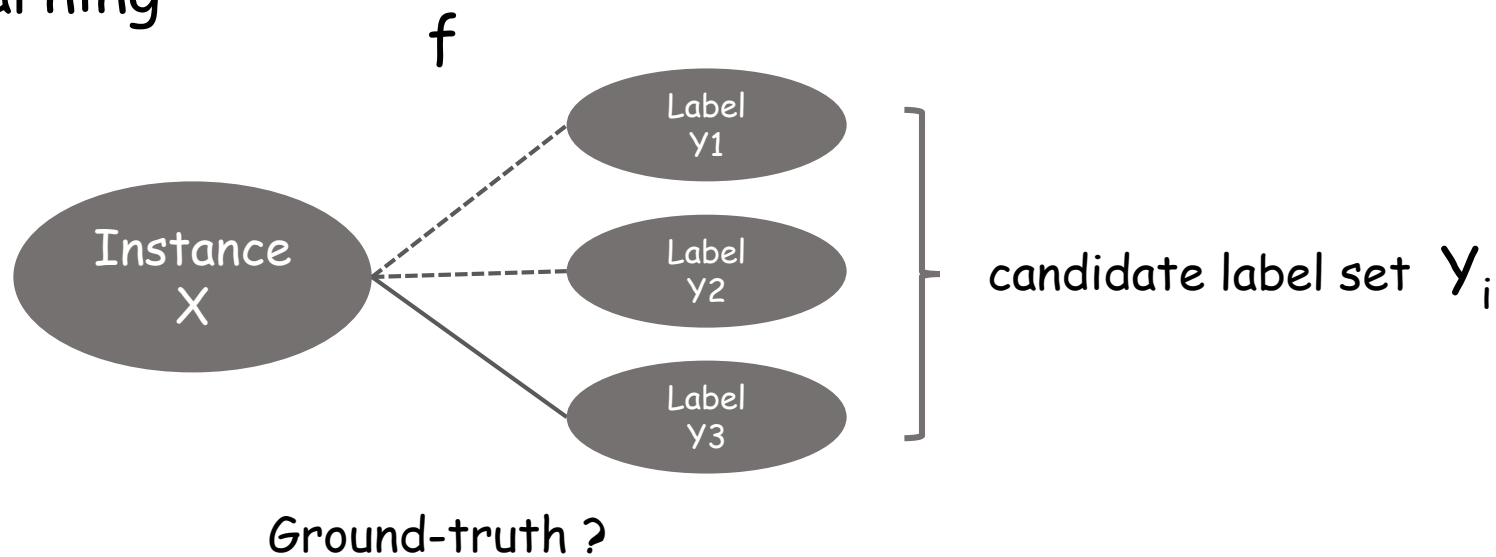
*ICML 2022*

# Partial Label Learning

- Supervised Learning



- Partial Label Learning



# Partial Label Learning



---

A dog image  $x_i$  with  
 $Y_i = \{\text{Husky}, \underline{\text{Malamute}}, \text{Samoyed}\}$

---

An input image with three candidate labels,  
where the ground-truth is **Malamute**

# Partial Label Learning

- The key challenge of PLL is to identify the ground-truth label from the candidate label set.
- Learning a classifier from the candidate set  $Y_i$

$$\mathcal{L}_{\text{cls}}(f; \mathbf{x}_i, Y_i) = \sum_{j=1}^C -s_{i,j} \log(f^j(\mathbf{x}_i)) \quad \text{s.t.} \quad \sum_{j \in Y_i} s_{i,j} = 1 \text{ and } s_{i,j} = 0, \forall j \notin Y_i, \quad (1)$$

- **Disambiguation:**

Recover the ground-truth confidence  $s$  from  $Y_i$

# Method

- initialization

$$p_k = \begin{cases} \frac{1}{|\mathcal{S}|} & \text{if } k \in \mathcal{S}, \\ 0 & \text{otherwise.} \end{cases} \quad (6)$$

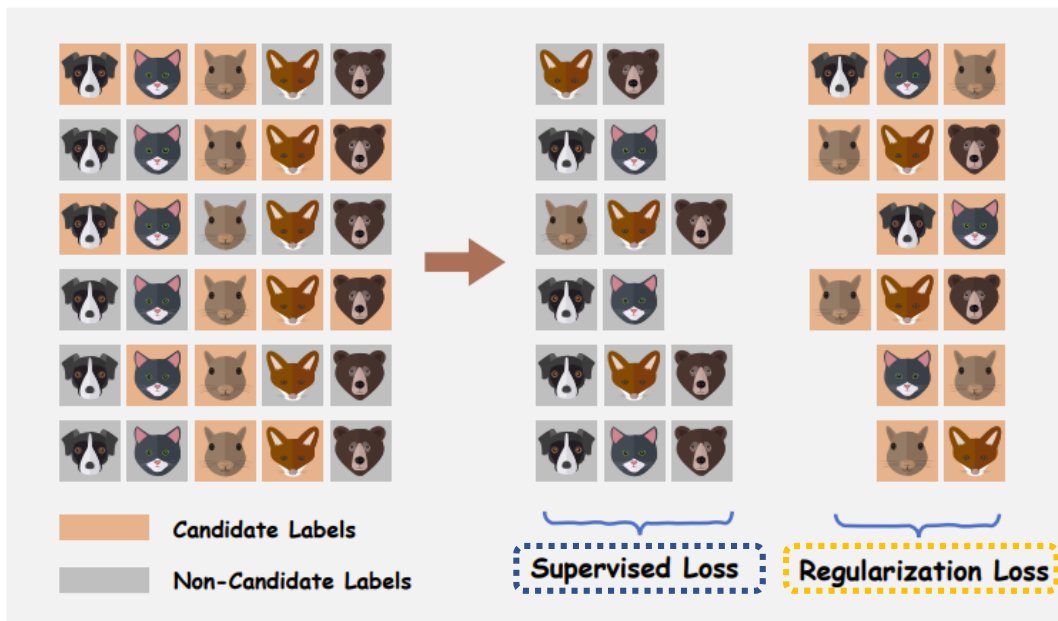


Figure 1. The label matrix is decoupled into two complementary parts including candidate labels and non-candidate labels, on which the supervised loss and regularization loss could be calculated respectively.

$$\mathcal{L}_{\text{Sup}}(\mathbf{x}, \mathcal{S}) = - \sum_{k \notin \mathcal{S}} \log(1 - g_k(\mathbf{x})). \quad (2)$$

- The complements of the candidate labels would never be the ground-truth labels
- Objective function

$$\mathcal{L}(\mathbf{x}, \mathcal{S}) = \mathcal{L}_{\text{Sup}}(\mathbf{x}, \mathcal{S}) + \lambda \Psi(\mathbf{x}, \mathcal{S}), \quad (3)$$

# Consistency Regularization

- Minimizing the divergence of each output pair given a set of **random augmented** instances in

$$\mathcal{A}(\mathbf{x}) = \{\text{Aug}_i(\mathbf{x}) | 1 \leq i \leq K\}$$

- Drawbacks
  - Optimizing the output divergence of each instance pair is **inefficient**, especially with large K
  - Some random augmentations which could cause significant **semantic shift**

- Conformal label distribution **p**

$$\sum_{k \in \mathcal{S}} p_k = 1 \text{ and } p_k = 0, \forall k \notin \mathcal{S}$$

$$\Psi(\mathbf{x}, \mathcal{S}) = \sum_{z \in \mathcal{A}(\mathbf{x})} \text{KL}(\mathbf{p} || g(z)). \quad (4)$$

# Obtain the conformal label distribution

- Given this regularization term, we need to obtain **the conformal label distribution** of each instance before optimization the network's parameters
- $p$  could be treated as **latent variable** and optimized simultaneously with the parameters  $\theta$  from the objective function.
- Bi-level optimization.

$$\arg \min_{\theta} \mathcal{L}(\theta, \mathbf{p}^*)$$

$$\text{subject to } \mathbf{p}^* = \arg \min_{\mathbf{p}} \mathcal{L}(\theta, \mathbf{p});$$

$$\sum_{k \in \mathcal{S}} p_k = 1; p_k = 0, \forall k \notin \mathcal{S}.$$

(5)

$$p_k^* = \frac{\left( \prod_{z \in \mathcal{A}(\mathbf{x})} g_k(z) \right)^{\frac{1}{|\mathcal{A}(\mathbf{x})|}}}{\sum_{j \in \mathcal{S}} \left( \prod_{z \in \mathcal{A}(\mathbf{x})} g_j(z) \right)^{\frac{1}{|\mathcal{A}(\mathbf{x})|}}}, \quad (7)$$

- Loss

$$\mathcal{L}(\boldsymbol{x}, \mathcal{S}) = \mathcal{L}_{\text{Sup}}(\boldsymbol{x}, \mathcal{S}) + \gamma(t) \cdot \Psi(\boldsymbol{x}, \mathcal{S}), \quad (8)$$

- Dynamic balancing

$$\gamma(t) = \min\left\{\frac{t}{T}, \lambda\right\}$$

---

## Algorithm 1 Our Regularized Training Method

---

**Input:** Training dataset  $\mathcal{D} = \{x_i, S_i\}_{i=1}^n$ ;  
The classifier  $g$  and its initial parameters  $\theta$ ;  
Epochs  $T$  and iterations  $I$ ;  
The number of augmentations  $K$ ;  
Maximum balancing factor  $\lambda$ ;

**Procedure:**

- 1: Initialize  $p$  for each instance by Eq.(6).
- 2: **for**  $t = 1$  **to**  $T$  **do**
- 3:   **for**  $i = 1$  **to**  $I$  **do**
- 4:     Fetch a random batch  $\mathcal{B}$  from  $\mathcal{D}$ ;
- 5:     Obtain the conformal label distributions by Eq.(7);
- 6:     Calculate the loss on the current batch by Eq. (8);
- 7:     Update network parameter  $\theta$  via gradient descent;
- 8:   **end for**
- 9: **end for**

**Output:** Learned multi-class classifier  $g$ .

---

# Experiments

Dataset	Method	$q = 0.1$	$q = 0.3$	$q = 0.5$	$q = 0.7$
Fashion-MNIST	Ours	<b>93.79 ± 0.05%</b>	<b>93.72 ± 0.20%</b>	<b>93.38 ± 0.08%</b>	<b>92.19 ± 0.03%</b>
	PiCO	93.36 ± 0.09%	93.41 ± 0.08%	92.88 ± 0.03%	91.73 ± 0.07%
	PRODEN	89.15 ± 0.58%	89.10 ± 0.26%	88.22 ± 0.35%	85.87 ± 0.28%
	VALEN	86.52 ± 0.16%	85.80 ± 0.17%	85.01 ± 0.30%	81.92 ± 0.30%
	LWS	91.44 ± 0.13%	91.85 ± 0.14%	90.59 ± 0.18%	89.46 ± 0.16%
	RC	92.64 ± 0.14%	92.08 ± 0.03%	92.01 ± 0.04%	90.83 ± 0.35%
	CC	92.26 ± 0.12%	91.75 ± 0.04%	90.92 ± 0.06%	89.73 ± 0.23%
	Fully Supervised	93.92 ± 0.07%			
Kuzushiji-MNIST	Ours	<b>98.27 ± 0.07%</b>	<b>98.08 ± 0.03%</b>	<b>97.44 ± 0.04%</b>	<b>95.93 ± 0.11%</b>
	PiCO	97.68 ± 0.06%	97.34 ± 0.07%	97.15 ± 0.03%	91.90 ± 0.04%
	PRODEN	94.61 ± 0.39%	93.08 ± 0.46%	90.15 ± 0.51%	81.10 ± 0.78%
	VALEN	85.59 ± 0.42%	82.96 ± 0.26%	78.13 ± 0.89%	68.09 ± 0.80%
	LWS	96.22 ± 0.10%	96.15 ± 0.24%	95.43 ± 0.02%	93.63 ± 0.04%
	RC	96.84 ± 0.09%	96.31 ± 0.15%	96.17 ± 0.06%	95.84 ± 0.12%
	CC	96.45 ± 0.04%	96.16 ± 0.02%	95.62 ± 0.10%	95.33 ± 0.14%
	Fully Supervised	98.31 ± 0.05%			

Table 1. Accuracy (mean±std) comparisons on Fashion-MNIST, Kuzushiji-MNIST, SVHN and CIFAR-10 with uniform partial labels on different ambiguity levels. The best result among each column is highlighted in bold.

# Experiments



SVHN	Ours	<b>97.56 ± 0.02%</b>	<b>97.19 ± 0.06%</b>	<b>97.58 ± 0.05%</b>	<b>95.46 ± 0.23%</b>
	PiCO	96.01 ± 0.02%	96.24 ± 0.09%	95.89 ± 0.06%	94.71 ± 0.17%
	PRODEN	96.35 ± 0.25%	96.07 ± 0.23%	95.61 ± 0.21%	94.35 ± 0.29%
	VALEN	90.80 ± 0.16%	90.92 ± 0.30%	89.19 ± 0.64%	32.72 ± 0.17%
	LWS	96.12 ± 0.05%	95.72 ± 0.18%	33.75 ± 0.09%	19.59 ± 0.11%
	RC	96.20 ± 0.03%	95.54 ± 0.03%	96.03 ± 0.05%	95.73 ± 0.08%
	CC	96.15 ± 0.02%	95.79 ± 0.05%	95.35 ± 0.03%	94.55 ± 0.15%
	Fully Supervised	97.58 ± 0.03%			
CIFAR-10	Ours	<b>97.45 ± 0.04%</b>	<b>97.28 ± 0.02%</b>	<b>97.05 ± 0.05%</b>	<b>95.77 ± 0.08%</b>
	PiCO	95.78 ± 0.05%	95.25 ± 0.06%	94.73 ± 0.11%	92.73 ± 0.08%
	VALEN	68.97 ± 0.23%	63.61 ± 0.17%	52.57 ± 0.32%	38.02 ± 2.16%
	PRODEN	91.94 ± 0.32%	91.10 ± 0.50%	89.82 ± 0.47%	86.48 ± 0.47%
	LWS	86.47 ± 0.20%	84.31 ± 0.14%	54.73 ± 0.19%	38.49 ± 0.24%
	RC	88.96 ± 0.06%	87.49 ± 0.17%	83.48 ± 0.19%	75.01 ± 0.19%
	CC	88.78 ± 0.05%	86.69 ± 0.40%	83.75 ± 0.28%	77.60 ± 0.22%
	Fully Supervised	97.57 ± 0.03%			

Table 1. Accuracy (mean±std) comparisons on Fashion-MNIST, Kuzushiji-MNIST, SVHN and CIFAR-10 with uniform partial labels on different ambiguity levels. The best result among each column is highlighted in bold.

# Experiments



Table 2. Accuracy (mean $\pm$ std) comparisons on CIFAR-100 with uniform partial labels on different ambiguity levels.

Dataset	Method	$q = 0.01$	$q = 0.05$	$q = 0.1$	$q = 0.2$
CIFAR-100	Ours	<b>83.12 <math>\pm</math> 0.20%</b>	<b>82.77 <math>\pm</math> 0.10%</b>	<b>82.24 <math>\pm</math> 0.07%</b>	<b>80.97 <math>\pm</math> 0.29%</b>
	PiCO	74.39 $\pm$ 0.15%	73.97 $\pm$ 0.09%	51.94 $\pm$ 0.11%	20.29 $\pm$ 0.04%
	PRODEN	72.55 $\pm$ 0.77%	71.55 $\pm$ 0.94%	70.84 $\pm$ 0.87%	58.86 $\pm$ 0.85%
	VALEN	43.68 $\pm$ 0.50%	43.19 $\pm$ 0.56%	35.00 $\pm$ 0.57%	4.51 $\pm$ 0.61%
	LWS	58.54 $\pm$ 0.12%	55.19 $\pm$ 0.23%	40.12 $\pm$ 0.34%	23.90 $\pm$ 0.18%
	RC	64.95 $\pm$ 0.23%	62.48 $\pm$ 0.14%	57.48 $\pm$ 0.04%	44.13 $\pm$ 0.23%
	CC	63.74 $\pm$ 0.17%	61.22 $\pm$ 0.21%	58.65 $\pm$ 0.06%	51.65 $\pm$ 0.49%
	Fully Supervised			83.16 $\pm$ 0.19%	

# Experiments

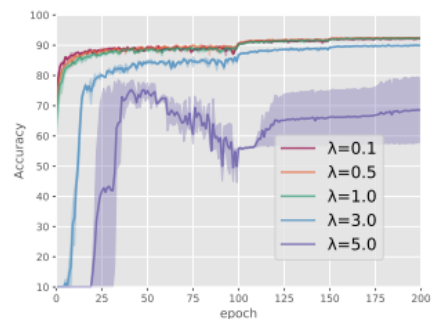
- Instance-dependent partial labels

- The flipping probability of each incorrect label is computed by  $q_j(x) = \frac{\hat{g}_j(x)}{\max_{k \in S} \hat{g}_k(x)}$
- $\hat{g}$  denotes a pre-trained neural network.

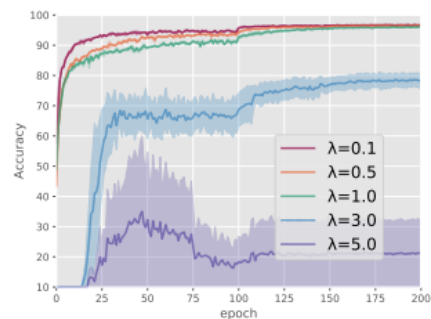
*Table 3. Accuracy (mean±std) comparisons on Kuzushiji-MNIST, Fashion-MNIST, CIFAR-10 with instance-dependent partial labels.*

Method	Kuzushiji-MNIST	Fashion-MNIST	CIFAR10
Ours	<b>95.07 ± 0.06%</b>	<b>89.21 ± 0.21%</b>	87.80 ± 0.11%
PiCO	92.87 ± 0.08%	86.93 ± 0.20%	<b>92.64 ± 0.07%</b>
PRODEN	87.71 ± 0.62%	84.25 ± 0.61%	76.51 ± 0.69%
VALEN	83.16 ± 0.33%	85.40 ± 0.10%	62.37 ± 0.35%
LWS	91.17 ± 0.18%	86.14 ± 0.34%	44.08 ± 0.15%
RC	94.00 ± 0.12%	89.10 ± 0.24%	75.90 ± 0.34%
CC	94.01 ± 0.05%	88.33 ± 0.17%	79.58 ± 0.22%

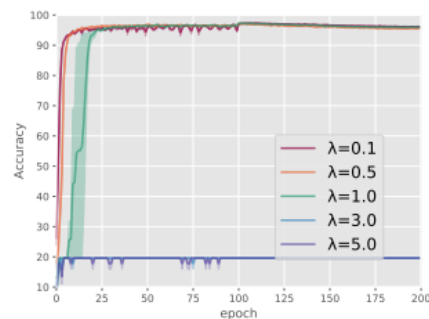
# Experiments(Ablation Study)



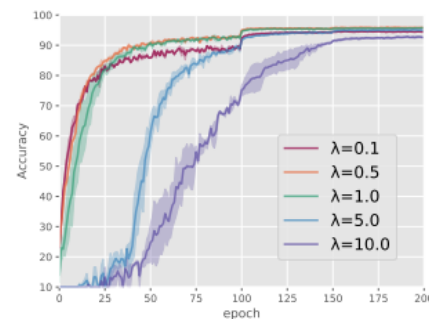
(a) Fixed  $\lambda$  Fashion-MNIST.



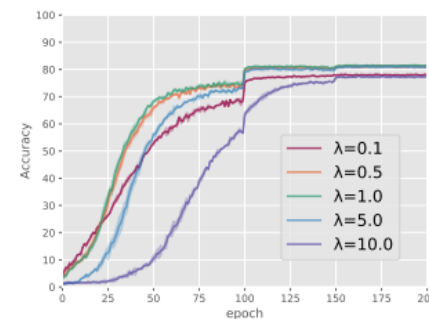
(b) Fixed  $\lambda$  on Kuzushiji-MNIST.



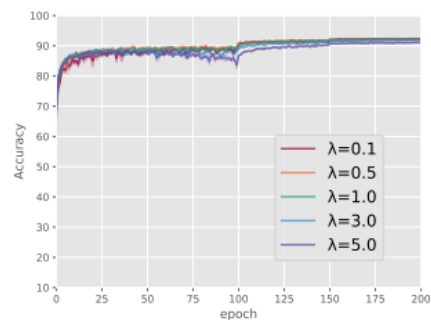
(c) Fixed  $\lambda$  on SVHN.



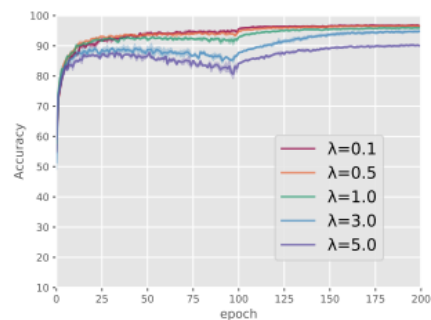
(d) Fixed  $\lambda$  on CIFAR-10.



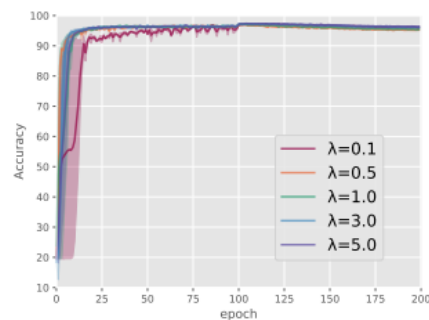
(e) Fixed  $\lambda$  on CIFAR-100.



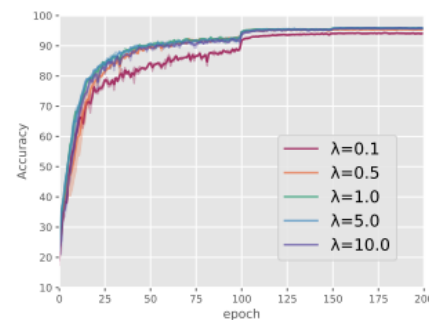
(f) Dynamic  $\lambda$  Fashion-MNIST.



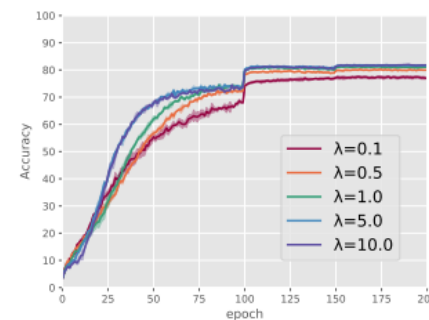
(g) Dynamic  $\lambda$  on Kuzushiji-MNIST.



(h) Dynamic  $\lambda$  on SVHN.



(i) Dynamic  $\lambda$  on CIFAR-10.



(j) Dynamic  $\lambda$  on CIFAR-100.

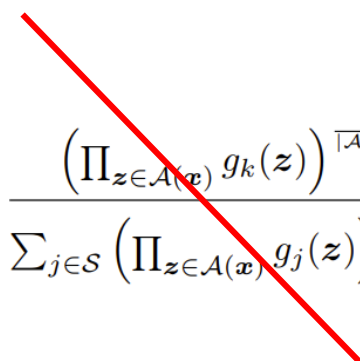
Figure 2. Accuracy curves with different balancing strategies. The top row shows results with fixed  $\lambda$  and the bottom row shows the results with dynamic  $\lambda$ . Dark colors show the mean accuracy of 5 trials and light colors show standard deviation.

# Experiments(Ablation Study)

- The importance of regularization term.

Table 4. Accuracy (mean $\pm$ std) of the degenerated method with uniform partial labels.

	$q = 0.5$	$q = 0.7$
Fashion-MNIST	91.49 $\pm$ 0.71% $\downarrow$ (1.89%)	90.61 $\pm$ 0.66% $\downarrow$ (1.58%)
Kuzushiji-MNIST	95.99 $\pm$ 0.06% $\downarrow$ (1.45%)	92.28 $\pm$ 0.20% $\downarrow$ (3.65%)
SVHN	96.60 $\pm$ 0.71% $\downarrow$ (0.98%)	95.02 $\pm$ 0.09% $\downarrow$ (0.44%)
CIFAR-10	95.39 $\pm$ 1.47% $\downarrow$ (1.66%)	93.29 $\pm$ 1.23% $\downarrow$ (2.48%)
	$q = 0.05$	$q = 0.1$
CIFAR-100	76.07 $\pm$ 0.15% $\downarrow$ (6.70%)	74.93 $\pm$ 0.19% $\downarrow$ (7.31%)


$$p_k^* = \frac{\left(\prod_{z \in \mathcal{A}(x)} g_k(z)\right)^{\frac{1}{|\mathcal{A}(x)|}}}{\sum_{j \in \mathcal{S}} \left(\prod_{z \in \mathcal{A}(x)} g_j(z)\right)^{\frac{1}{|\mathcal{A}(x)|}}}, \quad (7)$$

- Specifically, we replace the conformal label distribution inferred by Eq. (7) with simply re-normalized model outputs, and remaining the augmentation techniques.

# Experiments(Ablation Study)

- The influence of backbone network.

Table 5. Accuracy (mean $\pm$ std) comparison between PiCO and Our method with different backbones. WRN-34-10 is short for Wide-ResNet-34-10.

		PiCO	Ours
ResNet-18	Fully Supervised	73.56 $\pm$ 0.10%	79.88 $\pm$ 0.03%
	$p = 0.01$	73.09 $\pm$ 0.34% $\downarrow$ ( <b>0.47%</b> )	79.54 $\pm$ 0.12% $\downarrow$ ( <b>0.34%</b> )
	$p = 0.05$	72.74 $\pm$ 0.30% $\downarrow$ ( <b>0.82%</b> )	78.96 $\pm$ 0.06% $\downarrow$ ( <b>0.92%</b> )
	$p = 0.10$	69.91 $\pm$ 0.24% $\downarrow$ ( <b>3.65%</b> )	77.72 $\pm$ 0.08% $\downarrow$ ( <b>2.15%</b> )
WRN-34-10	Fully Supervised	74.47 $\pm$ 0.14%	83.16 $\pm$ 0.19%
	$p = 0.01$	74.39 $\pm$ 0.15% $\downarrow$ ( <b>0.08%</b> )	83.12 $\pm$ 0.20% $\downarrow$ ( <b>0.04%</b> )
	$p = 0.05$	73.97 $\pm$ 0.09% $\downarrow$ ( <b>0.50%</b> )	82.77 $\pm$ 0.10% $\downarrow$ ( <b>0.39%</b> )
	$p = 0.10$	51.94 $\pm$ 0.11% $\downarrow$ ( <b>22.53%</b> )	82.24 $\pm$ 0.07% $\downarrow$ ( <b>0.92%</b> )



Thanks