



Nanjing University of Aeronautics and Astronautics



### Contrastive Label Disambiguation for Partial Label Learning

Submit to ICLR2022 Average score: 8 (Top1)



### **Partial Label Learning**

• Partial label learning vs. ordinary supervised learning

#### candidate label set





#### Ordinary supervised learning (only one ground-truth label)

#### **Partial label learning**

(a candidate label set, among which only one label is valid)

### **Partial Label Learning: Applications**



• Multi-modal large-scale image annotation



顿扣篮得分,勇士将分差缩小至1分。紧接着沙梅特命中三分,佩顿连得4分,佩恩命中三

分,库里回<mark>敬三分球,麦基补扣得手,三节结束太阳80-78领先。</mark>

#### Candidate set (accurate one in black)



### **Partial Label Learning: Applications**



Crowdsourcing labeling



Candidate set (accurate one in black) Horse Donkey Mule

Annotator 1: Horse Annotator 2: Donkey Annotator 3: Mule

## **Partial Label Learning Framework**



Learning from partial labels



- ✓ Each instance is associated
  with multiple candidate labels
- ✓ Only one of the candidate
  label is the unknown groundtruth label
- Learning a classifier from the candidate set  $Y_i$

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

- Solution: disambiguation:
  - Recover the ground-truth confidence s from  $Y_i$

#### 6

### **Existing Disambiguating Strategies**

- Regularization term
  - Entropy regularizer:  $\mathcal{L}_d(\hat{\mathbf{Y}}) = -$
  - Smooth assumption:  $\mu \sum_{i,j}^{m} s_{ij} \left\| \frac{\mathbf{p}_i}{\sqrt{d_{ii}}} \frac{\mathbf{p}_j}{\sqrt{d_{jj}}} \right\|_2^2$
  - Shioothassamption
  - Self training: 🛛 🖤

0

Noisy label identification



$$-\frac{1}{n}\sum_{i=1}^{n}\hat{\mathbf{y}}_{i}^{\top}\log\hat{\mathbf{y}}_{i} \qquad [AAAI'20]$$

g: 
$$w_{ij} = \begin{cases} g_j(\boldsymbol{x}_i) / \sum_{k \in s_i} g_k(\boldsymbol{x}_i) & \text{if } j \in s_i, \\ 0 & \text{otherwise.} \end{cases}$$

[AAAI'18]

#### 7

## The PiCO Framework

• Contrastive label disambiguation





## The PiCO Framework (cont.)



#### Supervised contrastive learning



Self-Supervised contrastive loss:

Supervised contrastive loss:

$$-\sum_{i\in I} \log \frac{\exp\left(\boldsymbol{z}_{i} \cdot \boldsymbol{z}_{j(i)} / \tau\right)}{\sum_{a\in A(i)} \exp\left(\boldsymbol{z}_{i} \cdot \boldsymbol{z}_{a} / \tau\right)}$$

$$\sum_{i \in I} \frac{-1}{|P(i)|} \sum_{p \in P(i)} \log \frac{\exp\left(\mathbf{z}_i \cdot \mathbf{z}_p / \tau\right)}{\sum_{a \in A(i)} \exp\left(\mathbf{z}_i \cdot \mathbf{z}_a / \tau\right)}$$

### The PiCO Framework (cont.)



- Contrastive representation learning
  - Positive set selection:  $P(x) = \{k' | k' \in A(x), \tilde{y}' = \tilde{y}\}$
  - The contrastive loss:

$$\mathcal{L}_{\text{cont}}(g; \boldsymbol{x}, \tau, A) = -\frac{1}{|P(\boldsymbol{x})|} \sum_{\boldsymbol{k}_{+} \in P(\boldsymbol{x})} \log \frac{\exp(\boldsymbol{q}^{\top} \boldsymbol{k}_{+}/\tau)}{\sum_{\boldsymbol{k}' \in A(\boldsymbol{x})} \exp(\boldsymbol{q}^{\top} \boldsymbol{k}'/\tau)}$$

• The overall loss:  $\mathcal{L} = \mathcal{L}_{cls} + \lambda \mathcal{L}_{cont}$   $\mathcal{L}_{cls}(f; \boldsymbol{x}_i, Y_i) = \sum_{j=1}^{C} -s_{i,j} \log(f^j(\boldsymbol{x}_i))$ 



## The PiCO Framework (cont.)

- Prototype-based label disambiguation
  - Pseudo target updating:

$$s = \phi s + (1 - \phi) z$$
,  $z_c = \begin{cases} 1 & \text{if } c = \arg \max_{j \in Y} q^\top \mu_j \\ 0 & \text{else} \end{cases}$ 

Prototype updating

 $\boldsymbol{\mu}_c = \operatorname{Normalize}(\gamma \boldsymbol{\mu}_c + (1 - \gamma) \boldsymbol{q}), \quad \text{if } c = \arg \max_{j \in Y} f^j(\operatorname{Aug}_q(\boldsymbol{x}))),$ 





## Experiments



#### • PiCO achieves SOTA results

Table 1: Accuracy comparisons on benchmark datasets. Bold indicates superior results. Notably PiCO achieves comparable results to the fully supervised learning (less than 1% in accuracy with 1 false candidate).

Dataset	Method	q = 0.1	q = 0.3	q = 0.5
CIFAR-10	PiCO (ours)	<b>94.39</b> ± 0.18%	$\textbf{94.18} \pm 0.12\%$	$\textbf{93.58} \pm 0.06\%$
	LWS	$90.30 \pm 0.60\%$	$88.99 \pm 1.43\%$	$86.16 \pm 0.85\%$
	PRODEN	$90.24 \pm 0.32\%$	$89.38 \pm 0.31\%$	$87.78 \pm 0.07\%$
	CC	$82.30 \pm 0.21\%$	$79.08 \pm 0.07\%$	$74.05 \pm 0.35\%$
	MSE	$79.97 \pm 0.45\%$	$75.64\pm0.28\%$	$67.09 \pm 0.66\%$
	EXP	$79.23 \pm 0.10\%$	$75.79 \pm 0.21\%$	$70.34 \pm 1.32\%$
	Fully Supervised	$94.91 \pm 0.07\%$		
Dataset	Method	q = 0.01	q = 0.05	q = 0.1
CIFAR-100	PiCO (ours)	<b>73.09</b> ± 0.34%	$\textbf{72.74} \pm 0.30\%$	$\textbf{69.91} \pm 0.24\%$
	LWS	$65.78 \pm 0.02\%$	$59.56 \pm 0.33\%$	$53.53 \pm 0.08\%$
	PRODEN	$62.60 \pm 0.02\%$	$60.73 \pm 0.03\%$	$56.80 \pm 0.29\%$
	CC	$49.76 \pm 0.45\%$	$47.62 \pm 0.08\%$	$35.72 \pm 0.47\%$
	MSE	$49.17 \pm 0.05\%$	$46.02 \pm 1.82\%$	$43.81 \pm 0.49\%$
	EXP	$44.45 \pm 1.50\%$	$41.05 \pm 1.40\%$	$29.27\pm2.81\%$

# Experiments



• PiCO learns more distinguishable representations



Figure 3: T-SNE visualization of the image representation on CIFAR-10 with q = 0.5. Different colors represent the corresponding classes.

$$s_j = 1/|Y| \ (j \in Y)$$

# Experiments



• Effective of contrastive loss and label disambiguation

Table 2: Ablation study on CIFAR-10 with $q = 0.5$ and CIFAR-100 with $q = 0.05$ .								
Ablation	$\mathcal{L}_{cont}$	Label Disambiguation	$\begin{array}{c} \text{CIFAR-10} \\ (q = 0.5) \end{array}$	$\begin{array}{l} \text{CIFAR-100} \\ (q=0.05) \end{array}$				
PiCO	$\checkmark$	Ours	93.58	72.74				
PiCO w/o Disambiguation	$\checkmark$	Uniform Pseudo Target	84.50	64.11				
PiCO w/o $\mathcal{L}_{cont}$	×	Uniform Pseudo Target	76.46	56.87				
PiCO with $\phi = 0$	$\checkmark$	Soft Prototype Probs	91.60	71.07				
PiCO with $\phi = 0$	$\checkmark$	One-hot Prototype	91.41	70.10				
PiCO	$\checkmark$	MA Soft Prototype Probs	81.67	63.75				

$$\boldsymbol{s} = \phi \boldsymbol{s} + (1 - \phi) \boldsymbol{z}, \quad z_c = \begin{cases} 1 & \text{if } c = \arg \max_{j \in Y} \boldsymbol{q}^\top \boldsymbol{\mu}_j \\ 0 & \text{else} \end{cases}$$

s = z

$$s_i = \frac{\exp(\boldsymbol{q}^{\top} \boldsymbol{\mu}_i / \tau)}{\sum_{j \in Y} \exp(\boldsymbol{q}^{\top} \boldsymbol{\mu}_j / \tau)}$$

# Discussion



- Smooth assumption
  - Single-label



| |

• Multi-Label



# Discussion



• Global Average Pooling (GAP)









Nanjing University of Aeronautics and Astronautics



#### ΤΗΑΝΚS