



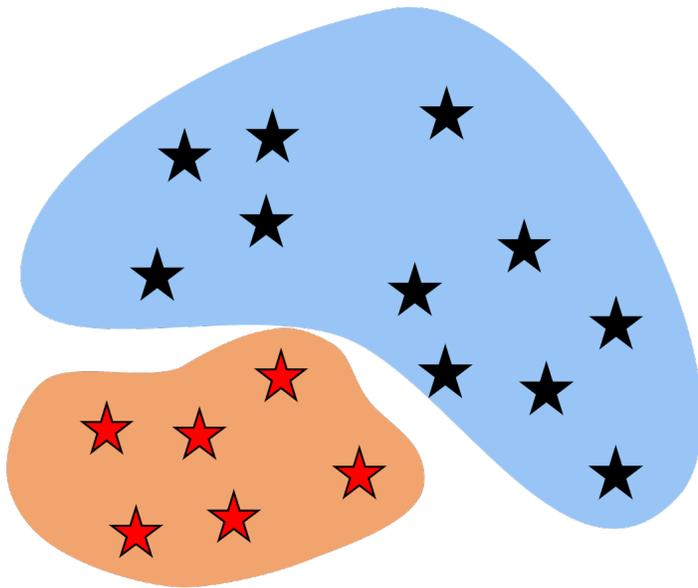
Noisy-Pair Robust Representation Alignment for Positive-Unlabeled Learning

Hengwei Zhao* Zhengzhong Tu Zhuo Zheng Wei Wang Junjue Wang

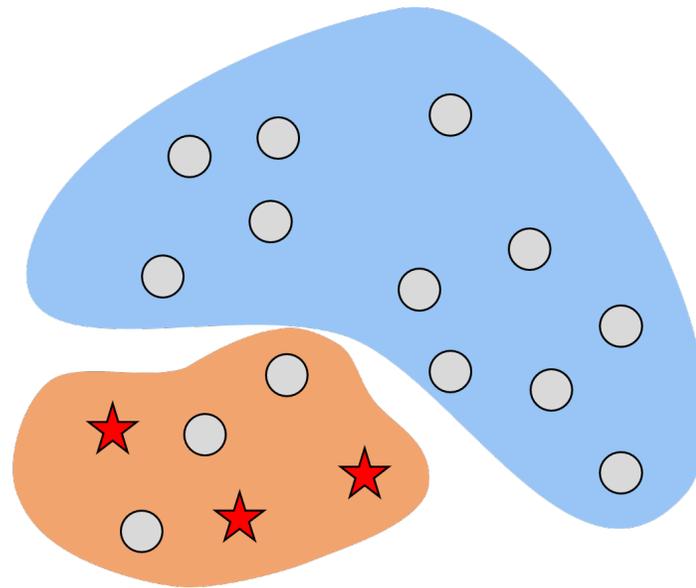
Rusty Feagin Wenzhe Jiao*

What is Positive-Unlabeled (PU) Learning?

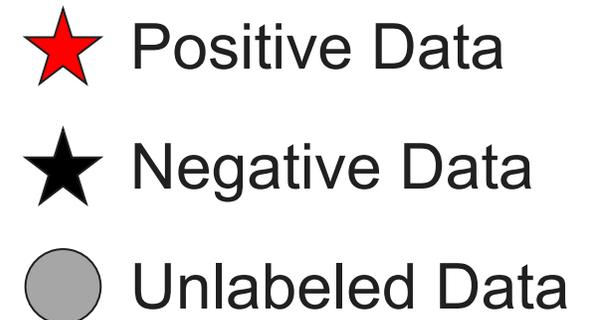
- PU learning aims to learn a binary classifier with limited positive data and a large pool of unlabeled data



Supervised Learning

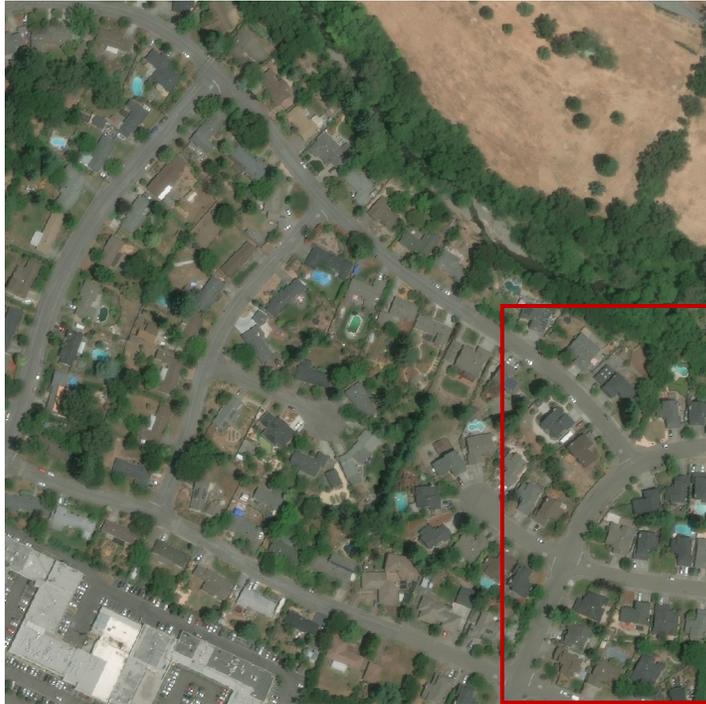


PU Learning

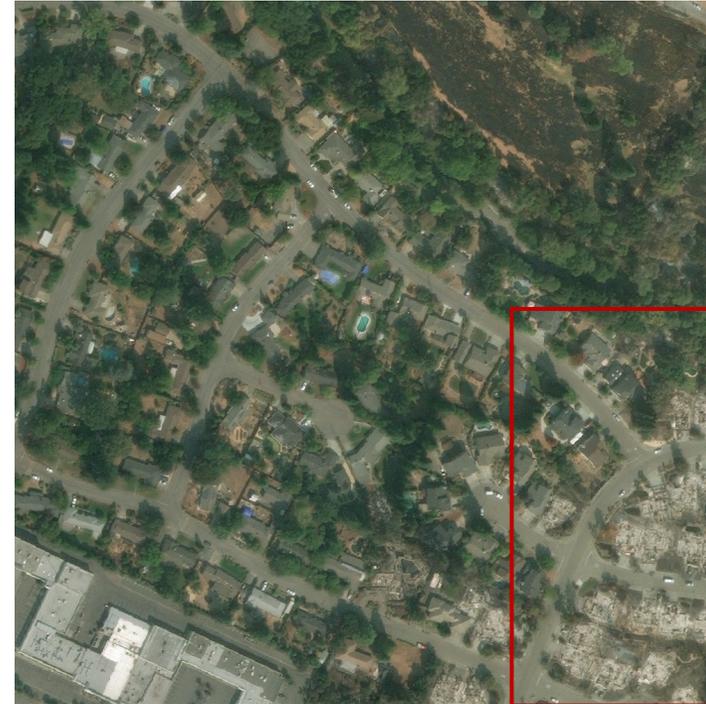


Why is PU Learning Important?

- In numerous real-world applications, positive data are readily obtainable, whereas negative labels are scarce or even infeasible to acquire



Pre-disaster image

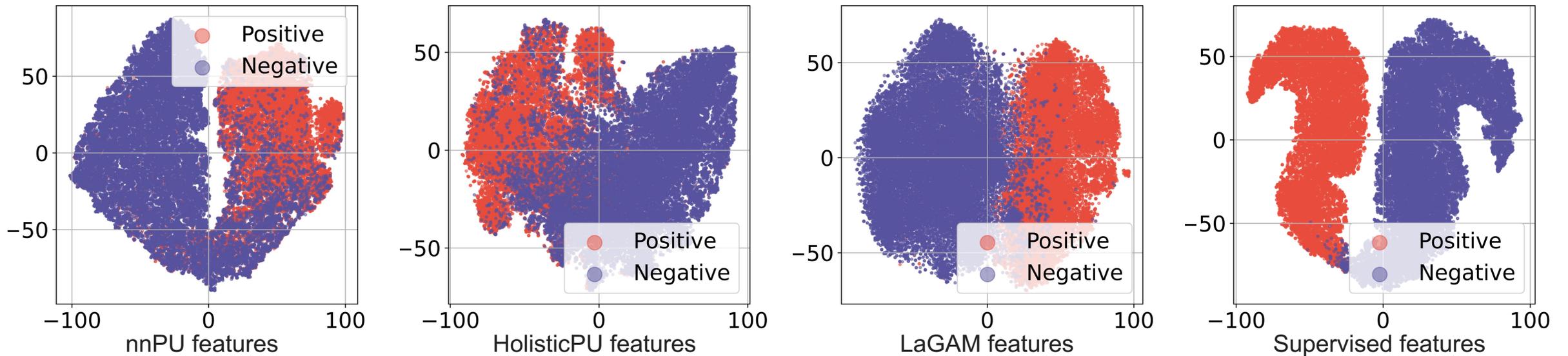


Post-disaster image

Example: After a disaster, training a classification model using only positive (damaged buildings) data can significantly shorten the time required for building damage mapping

What Limits PU Classifier vs. Supervised Methods?

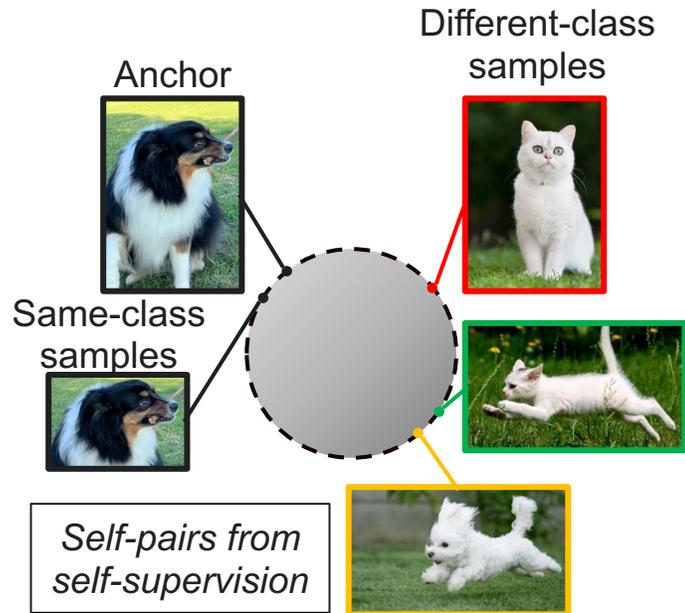
- Inaccurate annotations impede learning discriminative representations



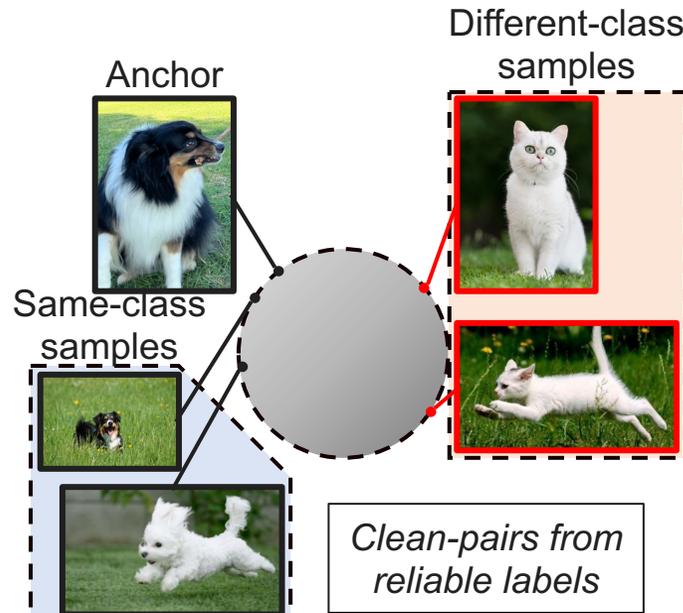
t-SNE visualizations of the representations on CIFAR-10 dataset

Challenge: How can discriminative representations be learned under unreliable supervision?

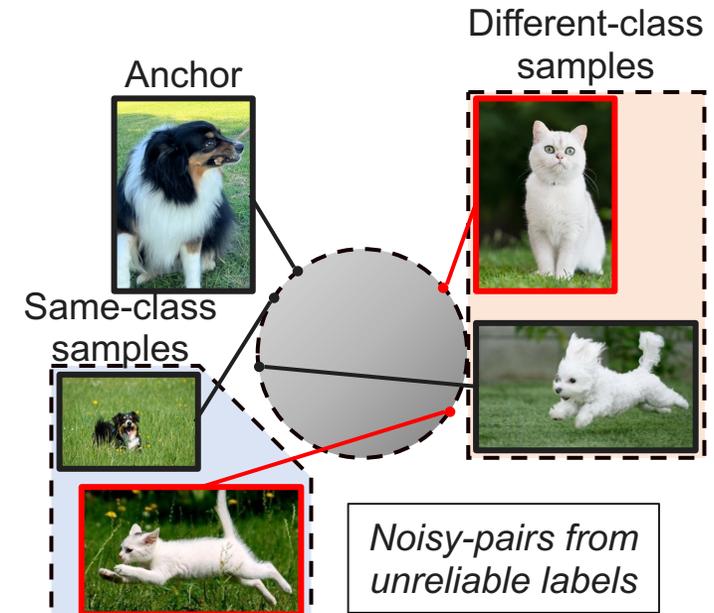
Noisy-Pair Robust Representation Learning



Self-Supervised Representation Learning



Supervised Representation Learning



Noisy-Pair Robust Representation Learning

Representation learning can acquire discriminative representations either by pulling same-class samples closer to the anchor and pushing different-class samples apart (contrastive representation learning), or by pulling same-class samples closer to the anchor (non-contrastive representation learning):

- ❖ Self-supervised representation learning: same-class pairs from augmented anchor
- ❖ Supervised representation learning: same-class pairs from reliable labels
- ❖ **Noisy-pair representation learning: same-class pairs from unreliable labels**



Analysis of Gradients

- Gradient magnitudes from noisy pairs overwhelm those from clean pairs in \mathcal{L}_r

- In noisy-pair robust representation learning, the representations may have:

$$\text{Noisy pairs: } (\mathbf{x}_i, \mathbf{x}_m) \rightarrow \tilde{\mathbf{q}}_i^T \tilde{\mathbf{q}}_m \approx 0 \quad \text{Clean pairs: } (\mathbf{x}_i, \mathbf{x}_j) \rightarrow \tilde{\mathbf{q}}_i^T \tilde{\mathbf{q}}_j \rightarrow 1$$

- **Supervised Non-Contrastive Loss (SNCL \mathcal{L}_r)**

$$\mathcal{L}_r(\mathbf{x}_i, \mathbf{x}_j) = 2(1 - \langle \tilde{\mathbf{q}}_i, \tilde{\mathbf{k}}_j \rangle) \mathbb{1}\{y_i = y_j\} \Rightarrow \left\| \frac{\partial \mathcal{L}_r(\mathbf{x}_i, \mathbf{x}_m)}{\partial \mathbf{q}_i} \right\|_2^2 = \frac{4}{\|\mathbf{q}_i\|_2^2} (1 - (\tilde{\mathbf{q}}_i^T \tilde{\mathbf{q}}_m)^2) > \frac{4}{\|\mathbf{q}_i\|_2^2} (1 - (\tilde{\mathbf{q}}_i^T \tilde{\mathbf{q}}_j)^2) = \left\| \frac{\partial \mathcal{L}_r(\mathbf{x}_i, \mathbf{x}_j)}{\partial \mathbf{q}_i} \right\|_2^2$$

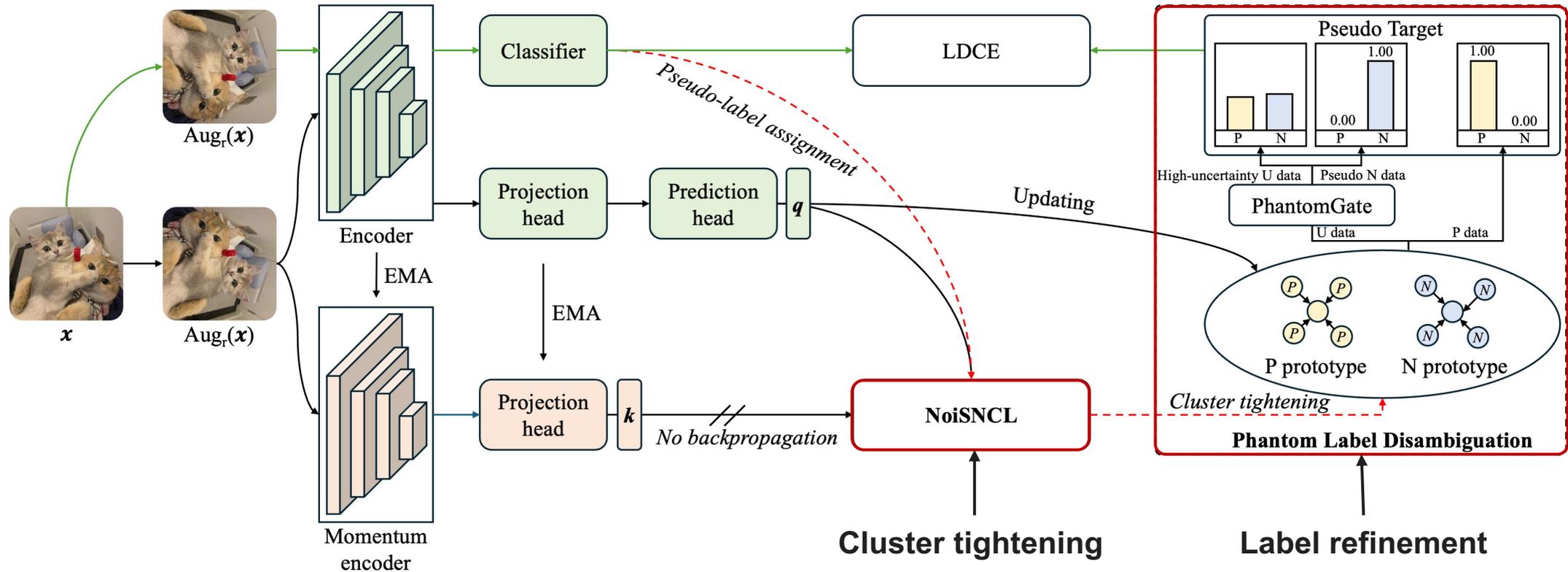
- **Noisy-Pair Robust Supervised Non-Contrastive Loss (NoiSNCL $\tilde{\mathcal{L}}_r$)**

$$\tilde{\mathcal{L}}_r(\mathbf{x}_i, \mathbf{x}_j) = 2\sqrt{1 - \langle \tilde{\mathbf{q}}_i, \tilde{\mathbf{k}}_j \rangle} \mathbb{1}\{y_i = y_j\} \Rightarrow \left\| \frac{\partial \tilde{\mathcal{L}}_r(\mathbf{x}_i, \mathbf{x}_m)}{\partial \mathbf{q}_i} \right\|_2^2 = \frac{1}{\|\mathbf{q}_i\|_2^2} (1 + (\tilde{\mathbf{q}}_i^T \tilde{\mathbf{q}}_m)) < \frac{1}{\|\mathbf{q}_i\|_2^2} (1 + (\tilde{\mathbf{q}}_i^T \tilde{\mathbf{q}}_j)) = \left\| \frac{\partial \tilde{\mathcal{L}}_r(\mathbf{x}_i, \mathbf{x}_j)}{\partial \mathbf{q}_i} \right\|_2^2$$

In \mathcal{L}_r , the gradient magnitudes contributed by noisy pairs dominate those from clean pairs; this issue is effectively mitigated by the proposed $\tilde{\mathcal{L}}_r$

Noisy-Pair Robust Non-Contrastive PU Learning

■ Noisy-Pair Robust Representation Alignment Framework (*NcPU*)



- NoiSNCL: Learning discriminative representations under unreliable supervision
- Phantom Label Disambiguation (PLD): Refines supervision by regret-based label updating



Theoretical Analysis

- The NoiSNCL and PLD modules can be theoretically justified to iteratively benefit each other from the perspective of the EM framework

- **E-Step** (Label refinement). At the E-step, each unlabeled example is assigned to one specific cluster. Given a network $g(\cdot)$ parameterized by θ , the objective is to find θ^* that maximizes the log-likelihood function:

$$\theta^* = \arg \max_{\theta} \sum_{x \in \mathcal{U}} \sum_{y \in \mathcal{Z}} \mathbb{1}(\tilde{y} = y) \log p(x, y | \theta)$$

- **M-Step** (Cluster tightening). At the M-step, minimizing $\tilde{\mathcal{L}}_r$ encourages embeddings to concentrate around their cluster centers. Under some mild assumptions, minimizing $\tilde{\mathcal{R}}_r(\mathbf{x})$ is equivalent to maximizing a lower bound of the likelihood:

$$\tilde{\mathcal{R}}_r(\mathbf{x}) = \frac{1}{n_u} \sum_{x \in \mathcal{U}} \frac{1}{|\mathcal{Q}|} \sum_{k_+ \in \mathcal{Q}} \tilde{\mathcal{L}}_r$$

- ❖ **Theorem 1** Assume the distribution of each class in the representation space follows a d-variate von Mises-Fisher (vMF) distribution, which leads to: $h(\mathbf{x} | \tilde{\mathbf{v}}_c, \kappa) = c_d(\kappa) e^{\kappa \tilde{\mathbf{v}}_c^T \tilde{\mathbf{g}}(\mathbf{x})}$, where $\tilde{\mathbf{v}}_c = \mathbf{v}_c / \|\mathbf{v}_c\|$, κ is the concentration parameter, and $c_d(\kappa)$ is the normalization constant. Under the assumption of a uniform class prior, optimizing $\tilde{\mathcal{R}}_r(\mathbf{x})$ and log-likelihood function is equivalent to maximizing L_1 and L_2 below, respectively.

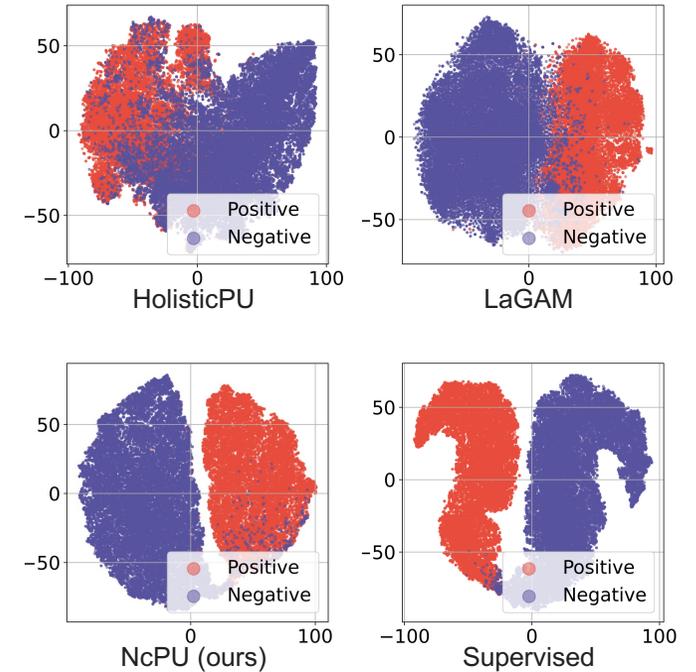
$$L_1 = \sum_{s_c \in \mathcal{U}} \frac{|S_c|}{n_u} \|\mathbf{v}_c\|^2 \leq \sum_{s_c \in \mathcal{U}} \frac{|S_c|}{n_u} \|\mathbf{v}_c\| = L_2$$



Experiments

■ Main results

Method	Additional N Data	CIFAR-10		CIFAR-100		STL-10		ABCD		xBD	
		OA	F1								
CE		60.45 \pm 0.1	2.42 \pm 0.4	50.36 \pm 0.0	1.86 \pm 0.2	50.30 \pm 0.0	1.19 \pm 0.2	55.70 \pm 0.2	20.93 \pm 0.4	84.08 \pm 0.2	25.70 \pm 2.4
uPU		65.52 \pm 0.2	26.82 \pm 0.9	61.44 \pm 0.9	43.12 \pm 2.1	57.08 \pm 0.4	25.88 \pm 1.3	83.76 \pm 2.1	81.47 \pm 2.9	86.82 \pm 0.3	55.43 \pm 2.8
nnPU		87.29 \pm 0.5	83.71 \pm 0.6	72.00 \pm 0.8	74.93 \pm 0.4	80.62 \pm 0.1	79.28 \pm 0.2	87.73 \pm 0.4	88.36 \pm 0.3	82.60 \pm 0.6	59.66 \pm 0.7
vPU		85.94 \pm 0.6	82.98 \pm 0.9	69.01 \pm 1.2	70.78 \pm 0.2	75.76 \pm 5.5	70.52 \pm 10.4	84.06 \pm 3.0	84.13 \pm 3.4	73.60 \pm 1.8	50.30 \pm 1.2
ImbPU		87.29 \pm 0.4	83.80 \pm 0.4	72.07 \pm 0.7	75.05 \pm 0.6	80.68 \pm 0.6	79.41 \pm 0.6	88.14 \pm 0.6	88.69 \pm 0.5	82.51 \pm 0.5	59.72 \pm 0.5
TEDn		86.29 \pm 2.4	80.70 \pm 4.6	69.85 \pm 0.9	61.73 \pm 1.9	66.26 \pm 4.9	49.90 \pm 10.7	88.90 \pm 0.9	89.10 \pm 0.9	85.40 \pm 0.8	52.65 \pm 4.6
PUET		78.51 \pm 0.4	73.85 \pm 0.5	62.81 \pm 0.2	71.09 \pm 0.1	75.36 \pm 0.2	73.56 \pm 0.1	78.09 \pm 2.9	66.52 \pm 24.9	74.92 \pm 0.1	38.38 \pm 0.6
HolisticPU		84.20 \pm 2.1	78.10 \pm 3.9	64.01 \pm 6.5	51.94 \pm 15.1	72.81 \pm 6.4	66.06 \pm 14.9	65.49 \pm 1.5	51.60 \pm 1.5	81.98 \pm 4.1	53.35 \pm 2.4
DistPU		85.29 \pm 2.6	83.96 \pm 2.2	67.63 \pm 0.8	73.68 \pm 0.8	85.62 \pm 1.5	85.41 \pm 0.9	86.25 \pm 1.7	87.36 \pm 1.2	82.94 \pm 0.8	57.58 \pm 0.2
PiCO		89.72 \pm 0.1	87.40 \pm 0.0	69.98 \pm 0.4	72.71 \pm 0.3	60.71 \pm 0.6	71.04 \pm 0.3	74.07 \pm 2.2	79.27 \pm 1.3	49.36 \pm 0.5	39.52 \pm 0.2
LaGAM	✓	95.78 \pm 0.5	94.90 \pm 0.6	84.82 \pm 0.1	84.42 \pm 0.2	88.64 \pm 0.0	88.50 \pm 0.1	75.90 \pm 0.4	75.38 \pm 0.6	79.14 \pm 1.5	58.78 \pm 1.7
WSC		90.55 \pm 0.3	87.92 \pm 0.8	75.39 \pm 2.1	73.76 \pm 4.0	79.06 \pm 4.5	74.16 \pm 7.0	80.10 \pm 2.8	76.12 \pm 4.3	84.89 \pm 0.8	62.17 \pm 1.3
NcPU(ours)		97.36\pm0.1	96.67\pm0.2	88.28\pm0.6	88.14\pm0.9	91.40\pm0.4	90.82\pm0.6	91.10\pm0.6	91.21\pm0.5	87.60\pm1.0	64.84\pm1.0
Supervised	✓	96.96 \pm 0.2	96.24 \pm 0.2	89.65 \pm 0.3	89.78 \pm 0.4	—	—	92.00 \pm 0.2	91.96 \pm 0.2	88.47 \pm 0.3	73.32 \pm 0.4



t-SNE visualizations of the representations on CIFAR-10 dataset

NcPU outperforms SOTA methods, achieving performance comparable to its supervised counterpart

Experiments

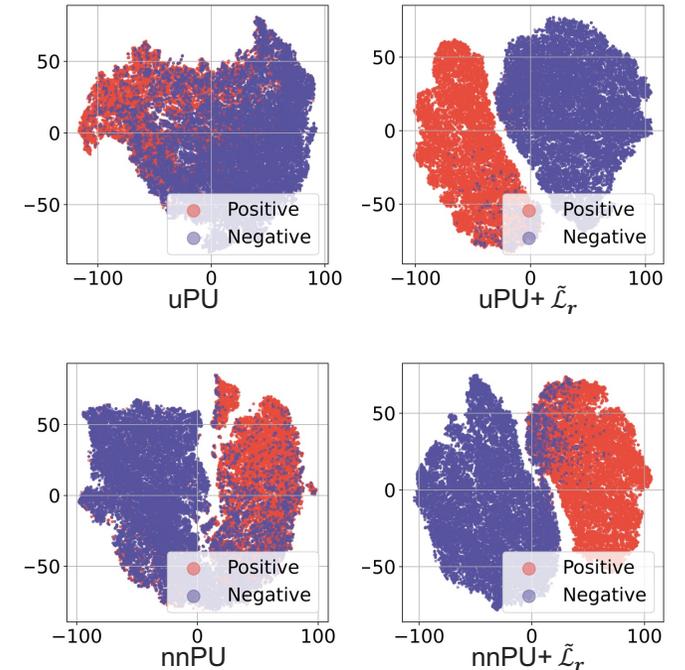
■ Ablation studies and analyses

❖ Ablation and comparative analyses on CIFAR-100 dataset

NCL	LD	OA	F1	P	R
	s	61.54 ± 7.8	40.58 ± 22.9	84.36 ± 7.1	30.69 ± 25.0
$\tilde{\mathcal{L}}_r$		50.27 ± 0.1	1.09 ± 0.4	97.97 ± 1.8	0.55 ± 0.2
$\mathcal{L}_{\text{self-r}}$	s	73.22 ± 1.7	72.75 ± 1.6	74.10 ± 2.2	71.47 ± 1.9
\mathcal{L}_r	s	84.58 ± 0.8	85.90 ± 0.6	79.12 ± 1.1	93.96 ± 0.3
$\tilde{\mathcal{L}}_r$	s'	75.14 ± 2.7	79.91 ± 1.7	67.15 ± 2.6	98.73 ± 0.5
$\tilde{\mathcal{L}}_r$	$s'+\text{SAT}$	50.25 ± 0.0	1.01 ± 0.1	97.85 ± 3.7	0.51 ± 0.1
$\tilde{\mathcal{L}}_r$	s	88.28 ± 0.6	88.14 ± 0.9	89.12 ± 1.7	87.27 ± 3.2

❖ Performance analysis of $\tilde{\mathcal{L}}_r$ under risk estimation with the real π_p and supervised methods

Method	CIFAR-10		CIFAR-100	
	OA	F1	OA	F1
uPU	69.43 ± 0.3	41.32 ± 1.0	61.68 ± 0.9	44.18 ± 2.6
uPU+ $\tilde{\mathcal{L}}_r$	97.35 ± 0.1	96.66 ± 0.2	83.71 ± 1.6	81.40 ± 2.2
nnPU	83.25 ± 0.2	76.94 ± 0.5	71.22 ± 0.5	68.12 ± 1.0
nnPU+ $\tilde{\mathcal{L}}_r$	97.03 ± 0.2	96.37 ± 0.2	87.81 ± 0.3	87.23 ± 0.4
Supervised+ \mathcal{L}_r	98.53 ± 0.0	98.17 ± 0.1	94.45 ± 0.1	94.52 ± 0.1
Supervised+ $\tilde{\mathcal{L}}_r$	98.75 ± 0.0	98.43 ± 0.1	94.56 ± 0.1	94.64 ± 0.1



t-SNE visualizations of the representations on CIFAR-10 dataset

$\tilde{\mathcal{L}}_r$ enhances the quality of learned representations, enabling simple PU methods to achieve competitive performance

Application

- Building damage mapping from satellite imagery





Thank you for your attention!



Homepage



Code