



**ICLR**



吉林大学



重慶大學  
CHONGQING UNIVERSITY

# Who Is Your Right Mixup Partner in Positive and Unlabeled Learning?

**Changchun Li<sup>1,2</sup>, Ximing Li<sup>1,2</sup>, Lei Feng<sup>3,4</sup>, Jihong Ouyang<sup>1,2</sup>**

<sup>1</sup>College of Computer Science and Technology, Jilin University, China

<sup>2</sup>Key Laboratory of Symbolic Computation and Knowledge Engineering of Ministry of Education, China

<sup>3</sup>College of Computer Science, Chongqing University, China

<sup>4</sup>Imperfect Information Learning Team, RIKEN Center for Advanced Intelligence Project, Japan

# Background

- Positive and Unlabeled (PU) Learning

PU learning aims to induce a **binary classifier** from weak training datasets of **positive and unlabeled instances**.

- Mixup Technique

The mixup is approximately equivalent to applying adversarial training, enabling to improve **robustness** with even **scarce and noisy** supervision.

$$\hat{\mathbf{x}} = \lambda \mathbf{x}_i + (1 - \lambda) \mathbf{x}_j, \quad \hat{y} = \lambda y_i + (1 - \lambda) y_j, \quad \lambda \sim \text{Beta}(\alpha, \alpha), \quad \alpha \in (0, \infty)$$

# Motivation

- Disambiguation-free objective of PU learning

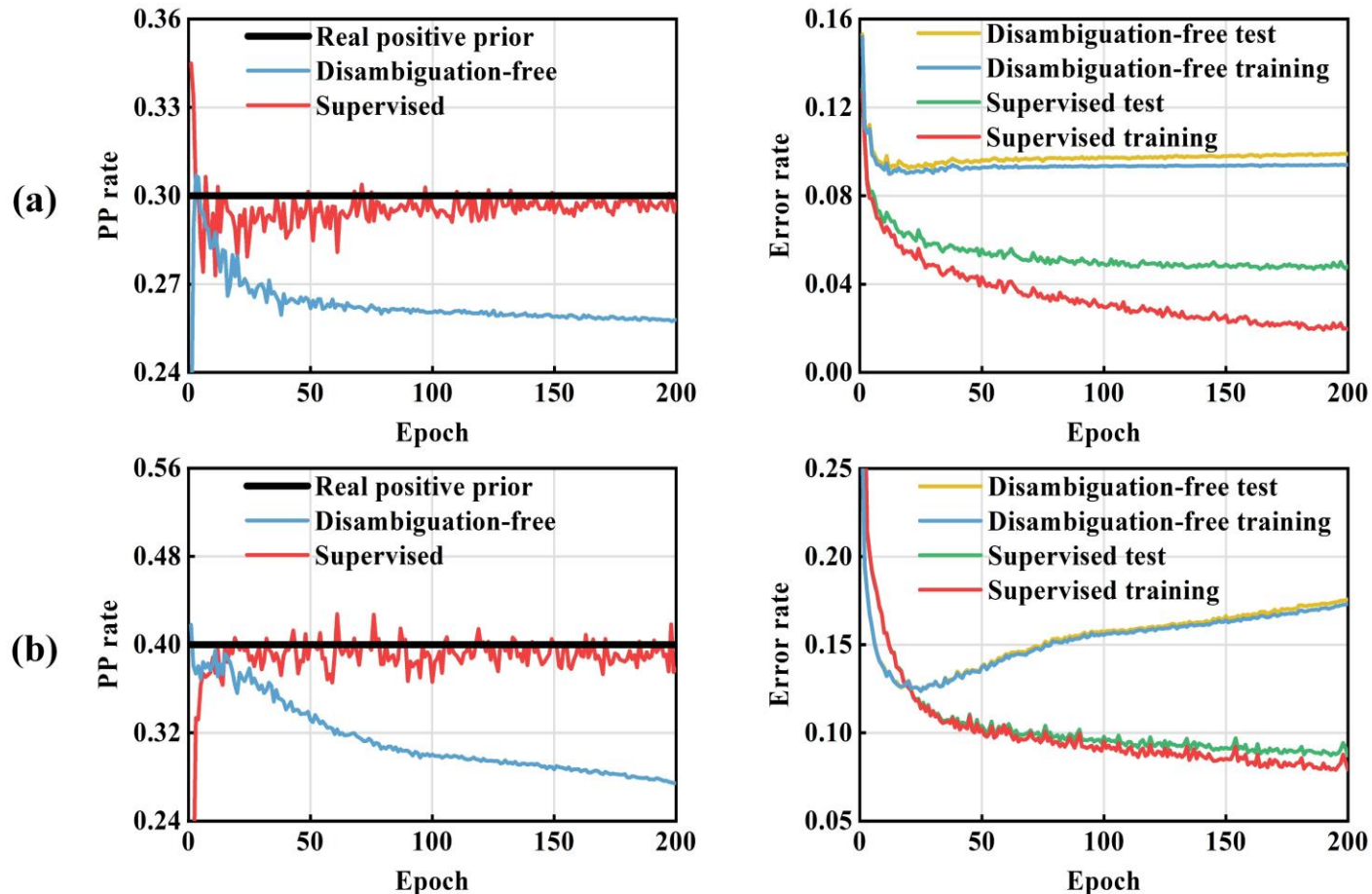
All unlabeled instances are treated as **pseudo-negative** instances, and the binary classifier is trained based on **positive and pseudo-negative** instances.

$$\mathcal{L}(\mathcal{X}_p, \mathcal{X}_u; \Theta) = \frac{1}{|\mathcal{X}_p|} \sum_{(\mathbf{x}, y) \in \mathcal{X}_p} \ell(f(\mathbf{x}; \Theta), y) + \frac{\beta}{|\mathcal{X}_u|} \sum_{(\mathbf{x}, y) \in \mathcal{X}_u} \ell(f(\mathbf{x}; \Theta), y)$$

# Motivation

- Decision boundary deviation phenomenon in PU learning

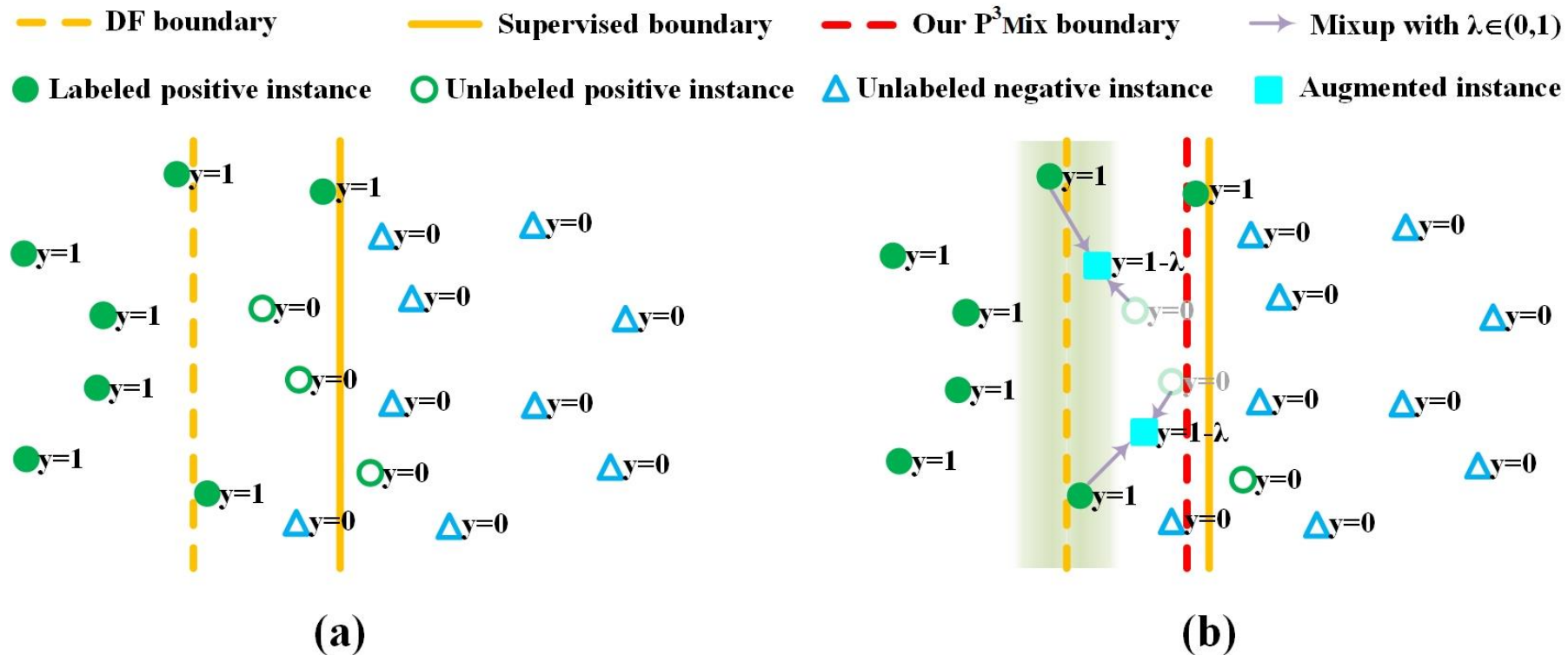
The number of training instances predicted as positive (PP rate) by the disambiguation-free classifier tends to be **smaller** than usual.



# Motivation

- Decision boundary deviation phenomenon in PU learning

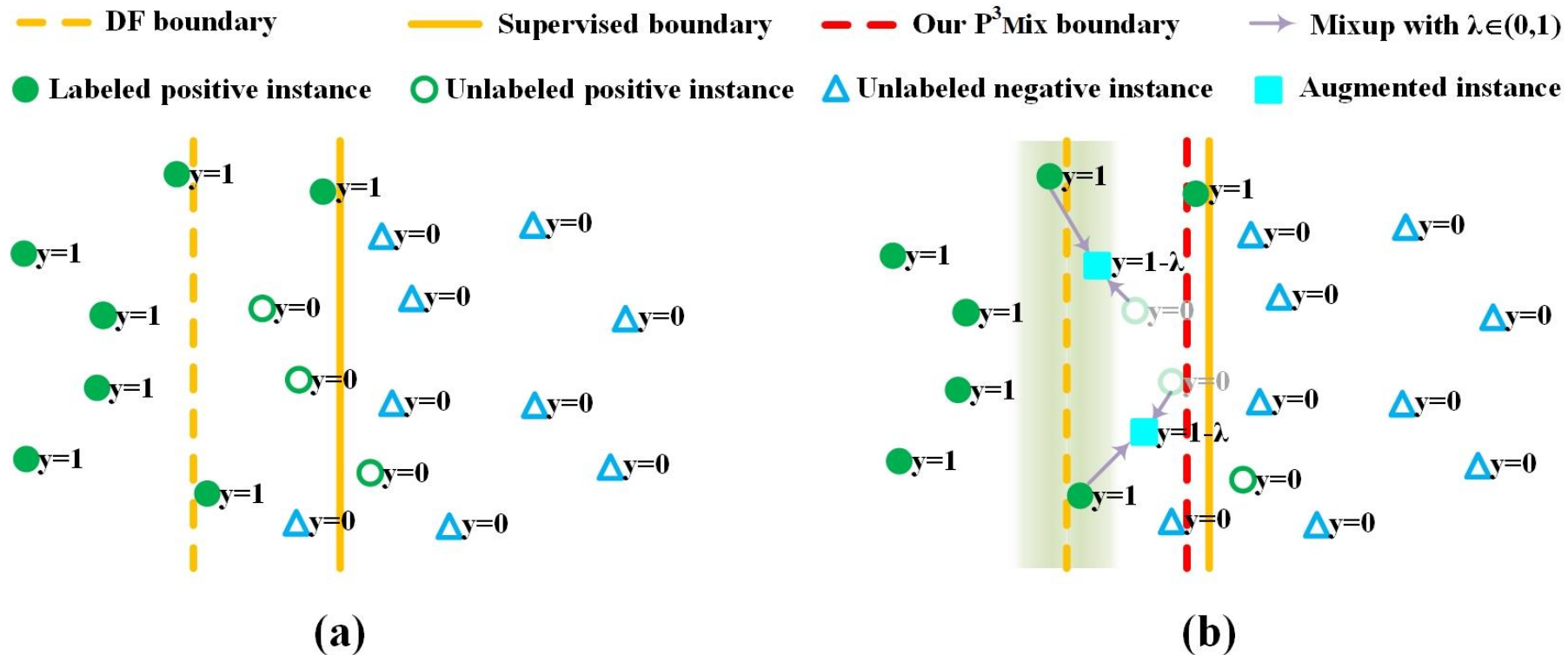
The disambiguation-free boundary tends to deviate from the fully supervised boundary towards the **positive** side.



# Motivation

- Decision boundary deviation phenomenon in PU learning

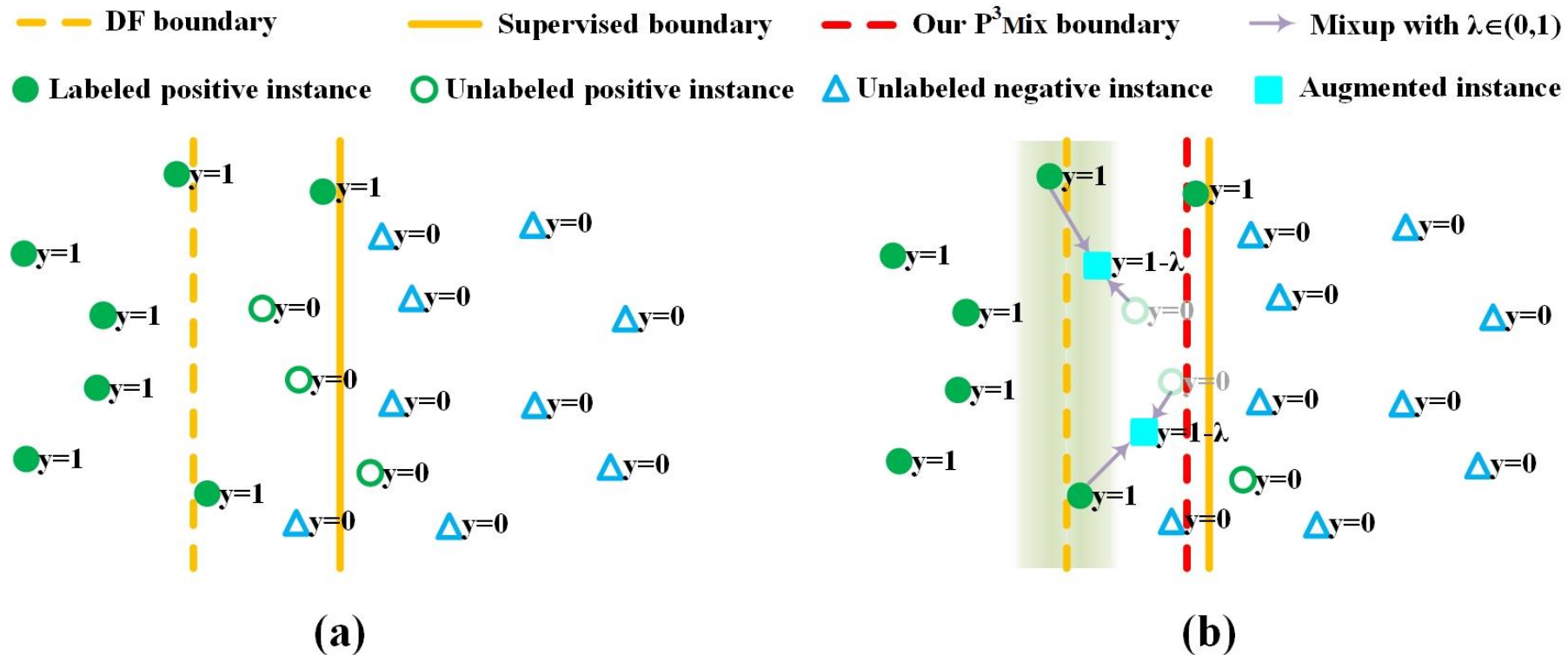
This is mainly caused by the **marginal pseudo-negative instances**, which lie between the two boundaries, and are more likely to be positive but actually annotated by negative.



# Motivation

- Heuristic mixup for PU learning

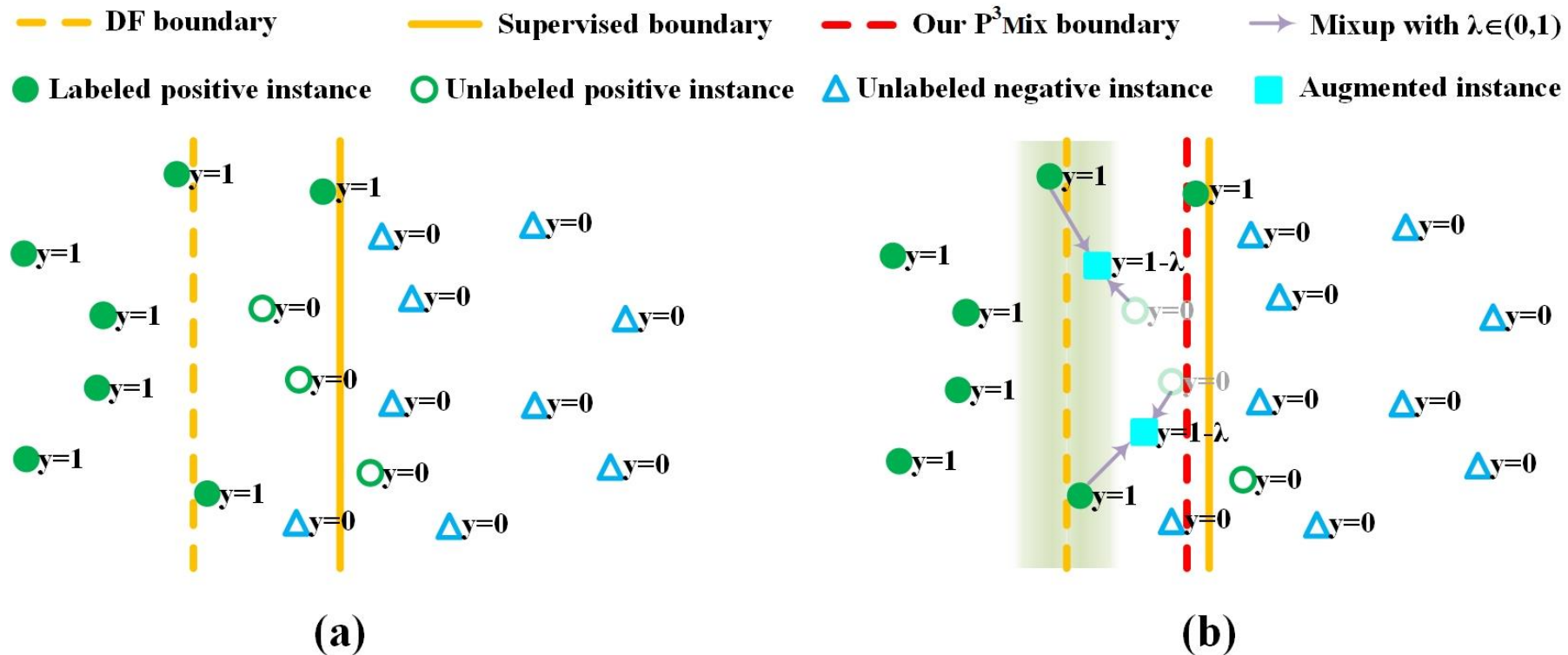
We extend mixup to a specific **heuristic** version for PU learning, enabling to achieve both data augmentation and supervision correction.



# Motivation

- Heuristic mixup for PU learning

It transforms the marginal pseudo-negative instances into **augmented instances** which are partially positive and yet also lie between the two boundaries, so as to push the learned boundary towards the fully supervised one.

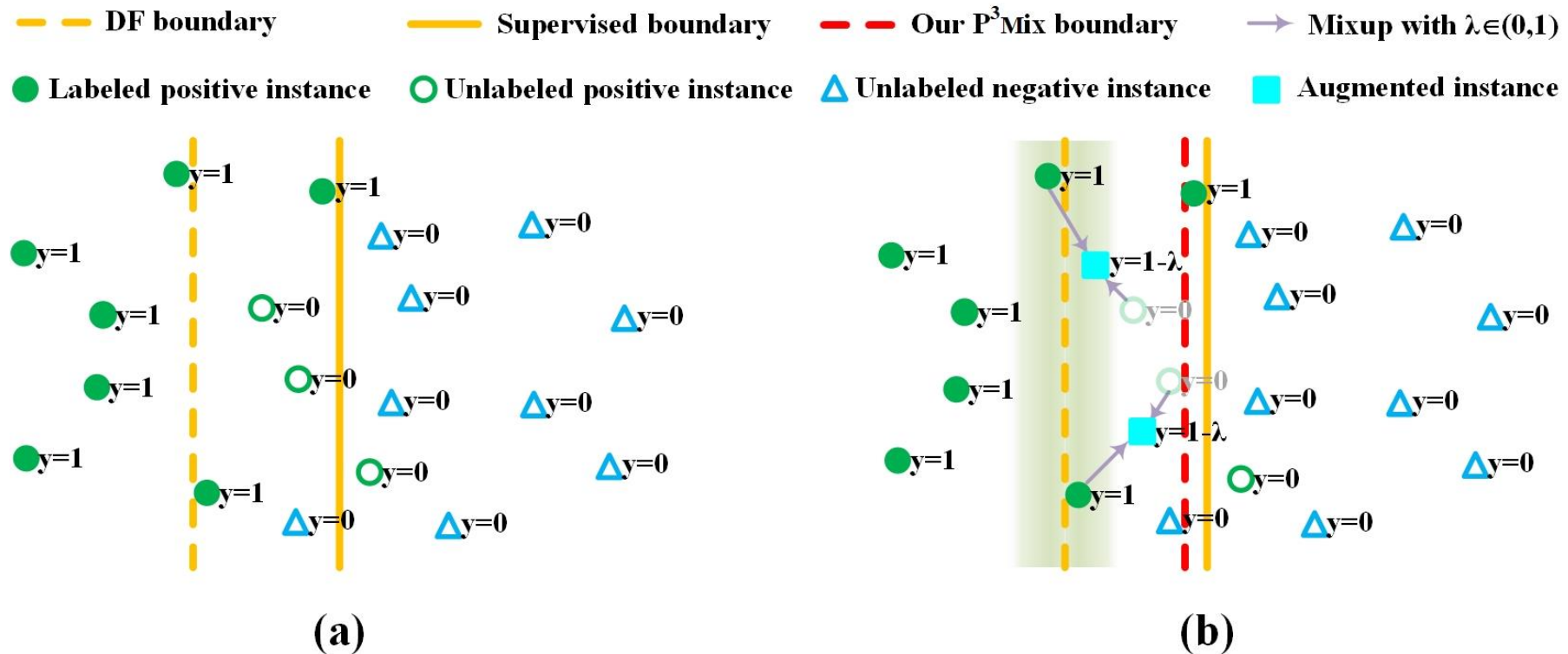




# Motivation

- Heuristic mixup for PU learning

This can be achieved by selecting the mixup partners for marginal pseudo-negative instances from the **positive** instances that are **around** the learned boundary.



# P<sup>3</sup>Mix

- Objective of P<sup>3</sup>Mix

P<sup>3</sup>Mix transforms batches of positive instances  $\mathcal{X}_p \subset \mathcal{P}$  and pseudo-negative ones  $\mathcal{X}_u \subset \mathcal{U}$  into the batches of augmented instances  $\hat{\mathcal{X}}_p$  and  $\hat{\mathcal{X}}_u$  using the proposed heuristic mixup, and its objective is given by:

$$\mathcal{L}(\hat{\mathcal{X}}_p, \hat{\mathcal{X}}_u; \Theta) = \frac{1}{|\hat{\mathcal{X}}_p|} \sum_{(\hat{\mathbf{x}}, \hat{y}) \in \hat{\mathcal{X}}_p} \ell(f(\hat{\mathbf{x}}; \Theta), \hat{y}) + \frac{\beta}{|\hat{\mathcal{X}}_u|} \sum_{(\hat{\mathbf{x}}, \hat{y}) \in \hat{\mathcal{X}}_u} \ell(f(\hat{\mathbf{x}}; \Theta), \hat{y}),$$
$$\hat{\mathcal{X}}_p, \hat{\mathcal{X}}_u = \text{HeuristicMixup}(\mathcal{X}_p, \mathcal{X}_u, \alpha),$$

# P<sup>3</sup>Mix

- Heuristic mixup

1. For each instance  $(\mathbf{x}_i, y_i) \in \mathcal{X}_p \cup \mathcal{X}_u$  we select a mixup partner  $(\mathbf{x}_j, y_j)$  to generate an augmented instance  $(\hat{\mathbf{x}}_i, \hat{y}_i)$  by using the modified mixup operator:

$$\hat{\mathbf{x}}_i = \lambda' \mathbf{x}_i + (1 - \lambda') \mathbf{x}_j, \quad \hat{y}_i = \lambda' y_i + (1 - \lambda') y_j, \quad \lambda' = \max(\lambda, 1 - \lambda), \\ \lambda \sim \text{Beta}(\alpha, \alpha), \quad \alpha \in (0, \infty),$$

2. Select the mixup partner for each of **marginal pseudo-negative instances**  $\mathcal{X}_{mpn} \subset \mathcal{X}_u$  from the **candidate mixup pool**  $\mathcal{X}_{cnd} \subset \mathcal{P}$  of positive instances that are around the current learned boundary; Select the mixup partners for positive instances  $\mathcal{X}_p$  and other pseudo-negative instances  $\mathcal{X}_u \setminus \mathcal{X}_{mpn}$  from both positive and pseudo-negative ones  $\mathcal{X}_p \cup \mathcal{X}_u$

$$(\mathbf{x}_j, y_j) \sim \begin{cases} \text{Uniform}(\mathcal{X}_{cnd}) & \text{if } (\mathbf{x}_i, y_i) \in \mathcal{X}_{mpn}, \\ \text{Uniform}(\mathcal{X}_p \cup \mathcal{X}_u) & \text{if } (\mathbf{x}_i, y_i) \in \mathcal{X}_p \cup \mathcal{X}_u \setminus \mathcal{X}_{mpn}. \end{cases}$$

# P<sup>3</sup>Mix

- Marginal pseudo-negative instance estimation

We define the marginal pseudo-negative instances as the “unreliable” pseudo-negative instances measured by the predictive scores with thresholding  $\gamma \in [0.5, 1]$

$$\mathcal{X}_{mpn} = \{(\mathbf{x}, y = 0) | (\mathbf{x}, y = 0) \in \mathcal{X}_u, 1 - \gamma \leq f(\mathbf{x}; \Theta) \leq \gamma\},$$

- Candidate mixup pool

For each positive instance we compute its entropy value of the predictive score, and update the candidate mixup pool with the positive instances with the top- $k$  maximum entropy values as follows:

$$\mathcal{X}_{cnd} = \{(\mathbf{x}, y = 1) | (\mathbf{x}, y = 1) \in \mathcal{P}, \mathcal{H}(f(\mathbf{x}; \Theta)) \in \text{Rank}(\{\mathcal{H}(f(\mathbf{x}_i; \Theta))\}_{i=1}^{n_p})\},$$

# Robustness of P<sup>3</sup>Mix

- Early-learning regularization and P<sup>3</sup>Mix-E

We employ the early-learning regularization to **prevent the memorization of imprecise supervision**.

$$\mathcal{R}_{elr}(\{(\hat{\mathbf{x}}_i, \tilde{\mathbf{y}}_i)\}_{i=1}^{|\hat{\mathcal{X}}_p|+|\hat{\mathcal{X}}_u|}; \Theta) = \frac{1}{|\hat{\mathcal{X}}_p|+|\hat{\mathcal{X}}_u|} \sum_{i=1}^{|\hat{\mathcal{X}}_p|+|\hat{\mathcal{X}}_u|} \log(1 - \langle f(\hat{\mathbf{x}}_i; \Theta), \tilde{\mathbf{y}}_i \rangle)$$

- Pseudo-negative instance correction and P<sup>3</sup>Mix-C

We focus on the pseudo-negative instances with **high confidence to be positive**:

$$\{(\mathbf{x}, y = 0) | (\mathbf{x}, y = 0) \in \mathcal{X}_u, f(\mathbf{x}; \Theta) > \gamma\}$$

and directly **revise** their labels to **positive** before their corresponding mixup operators.

# Experiment 1: Classification Accuracy

P<sup>3</sup>Mix-E and P<sup>3</sup>Mix-C consistently **outperform** all PU learning baselines on all benchmark datasets, indicating their superior performance.

Table 2: Results of classification accuracy (mean $\pm$ std). The highest scores among PU learning methods are indicated in **bold**.

Dataset	F-MNIST-1	F-MNIST-2	CIFAR-10-1	CIFAR-10-2	STL-10-1	STL-10-2
uPU	71.3 $\pm$ 1.4	84.0 $\pm$ 4.0	76.5 $\pm$ 2.5	71.6 $\pm$ 1.4	76.7 $\pm$ 3.8	78.2 $\pm$ 4.1
nnPU	89.7 $\pm$ 0.8	88.8 $\pm$ 0.9	84.7 $\pm$ 2.4	83.7 $\pm$ 0.6	77.1 $\pm$ 4.5	80.4 $\pm$ 2.7
nnPU+mixup	91.4 $\pm$ 0.3	88.2 $\pm$ 0.7	87.2 $\pm$ 0.6	85.8 $\pm$ 1.2	79.8 $\pm$ 0.8	82.2 $\pm$ 0.9
Self-PU	90.8 $\pm$ 0.4	89.1 $\pm$ 0.7	85.1 $\pm$ 0.8	83.9 $\pm$ 2.6	78.5 $\pm$ 1.1	80.8 $\pm$ 2.1
PAN	88.7 $\pm$ 1.2	83.6 $\pm$ 2.5	87.0 $\pm$ 0.3	82.8 $\pm$ 1.0	77.7 $\pm$ 2.5	79.8 $\pm$ 1.4
VPU	90.6 $\pm$ 1.2	86.8 $\pm$ 0.8	86.8 $\pm$ 1.2	82.5 $\pm$ 1.1	78.4 $\pm$ 1.1	82.9 $\pm$ 0.7
MixPUL	87.5 $\pm$ 1.5	89.0 $\pm$ 0.5	87.0 $\pm$ 1.9	87.0 $\pm$ 1.1	77.8 $\pm$ 0.7	78.9 $\pm$ 1.9
PULNS	90.7 $\pm$ 0.5	87.9 $\pm$ 0.5	87.2 $\pm$ 0.6	83.7 $\pm$ 2.9	80.2 $\pm$ 0.8	83.6 $\pm$ 0.7
P <sup>3</sup> Mix-E	<b>91.9<math>\pm</math>0.3</b>	<b>89.5<math>\pm</math>0.5</b>	88.2 $\pm$ 0.4	84.7 $\pm$ 0.5	80.2 $\pm$ 0.9	<b>83.7<math>\pm</math>0.7</b>
P <sup>3</sup> Mix-C	<b>92.0<math>\pm</math>0.4</b>	<b>89.4<math>\pm</math>0.3</b>	<b>88.7<math>\pm</math>0.4</b>	<b>87.9<math>\pm</math>0.5</b>	<b>80.7<math>\pm</math>0.7</b>	<b>84.1<math>\pm</math>0.3</b>
Supervised	95.2 $\pm$ 0.2	95.2 $\pm$ 0.2	91.3 $\pm$ 0.3	91.3 $\pm$ 0.3	85.6 $\pm$ 0.6	85.6 $\pm$ 0.6



# Experiment 2: Ablation Study

1. The proposed heuristic mixup can **significantly improve** the classification performance.
2. Both early-learning regularization and pseudo-negative instance correction contribute to the improvement of the classification performance in all cases.

Table 3: Results of ablative study (mean $\pm$ std). The highest scores are indicated in **bold**.

Dataset	F-MNIST-1	F-MNIST-2	CIFAR-10-1	CIFAR-10-2	STL-10-1	STL-10-2
DF	75.2 $\pm$ 1.2	62.7 $\pm$ 2.8	72.0 $\pm$ 3.2	57.4 $\pm$ 3.7	78.1 $\pm$ 0.6	80.6 $\pm$ 2.4
DF+mixup	78.4 $\pm$ 1.7	72.4 $\pm$ 1.4	79.2 $\pm$ 3.0	67.4 $\pm$ 2.5	78.9 $\pm$ 0.3	80.7 $\pm$ 1.9
P <sup>3</sup> Mix	87.0 $\pm$ 1.1	79.0 $\pm$ 1.6	87.0 $\pm$ 1.1	84.3 $\pm$ 0.6	79.8 $\pm$ 0.7	83.4 $\pm$ 0.7
DF-E	90.1 $\pm$ 0.7	74.2 $\pm$ 5.5	82.4 $\pm$ 1.6	69.4 $\pm$ 3.0	67.3 $\pm$ 2.0	75.0 $\pm$ 3.7
DF-E+mixup	90.6 $\pm$ 0.7	86.1 $\pm$ 2.5	85.7 $\pm$ 0.7	76.4 $\pm$ 0.9	78.3 $\pm$ 1.1	79.3 $\pm$ 2.3
P <sup>3</sup> Mix-E	<b>91.9<math>\pm</math>0.3</b>	<b>89.5<math>\pm</math>0.5</b>	<b>88.2<math>\pm</math>0.4</b>	<b>84.7<math>\pm</math>0.5</b>	<b>80.2<math>\pm</math>0.9</b>	<b>83.7<math>\pm</math>0.7</b>
DF-C	89.6 $\pm$ 1.8	87.4 $\pm$ 2.4	87.2 $\pm$ 0.8	84.7 $\pm$ 1.1	80.2 $\pm$ 3.0	82.7 $\pm$ 2.6
DF-C+mixup	91.6 $\pm$ 0.3	88.3 $\pm$ 1.2	87.7 $\pm$ 1.1	81.3 $\pm$ 3.6	79.9 $\pm$ 3.1	81.6 $\pm$ 2.8
P <sup>3</sup> Mix-C	<b>92.0<math>\pm</math>0.4</b>	<b>89.4<math>\pm</math>0.3</b>	<b>88.7<math>\pm</math>0.4</b>	<b>87.9<math>\pm</math>0.5</b>	<b>80.7<math>\pm</math>0.7</b>	<b>84.1<math>\pm</math>0.3</b>

# **Who Is Your Right Mixup Partner in Positive and Unlabeled Learning?**

**Thank you!**

**Changchun Li**  
**changchunli93@gmail.com**