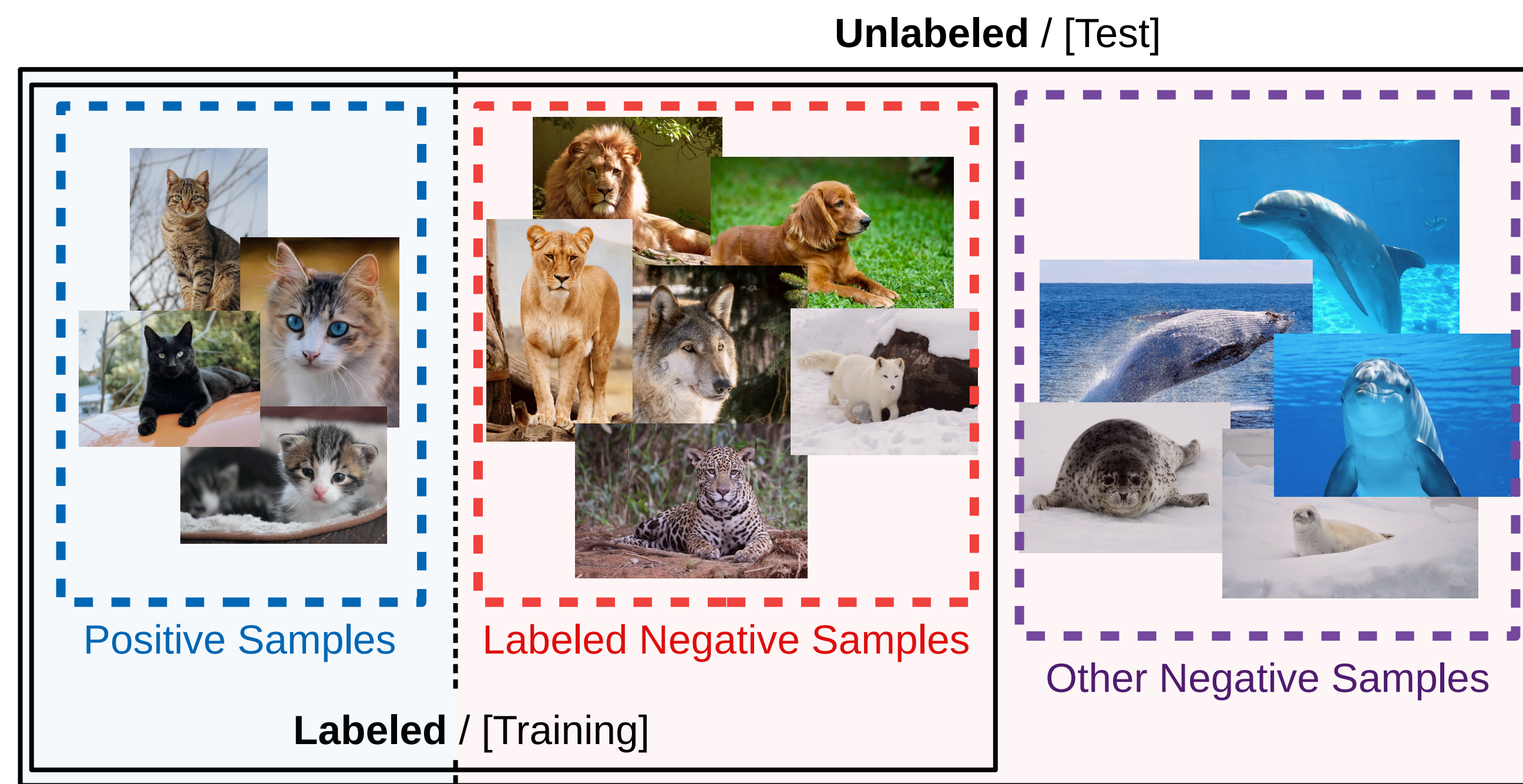


# Classification from Positive, Unlabeled and Biased Negative Data

Yu-Guan Hsieh<sup>1</sup> Gang Niu<sup>2</sup> Masashi Sugiyama<sup>2,3</sup>

<sup>1</sup> ENS Paris, France <sup>2</sup> RIKEN, Japan <sup>3</sup> The University of Tokyo, Japan

## Motivation

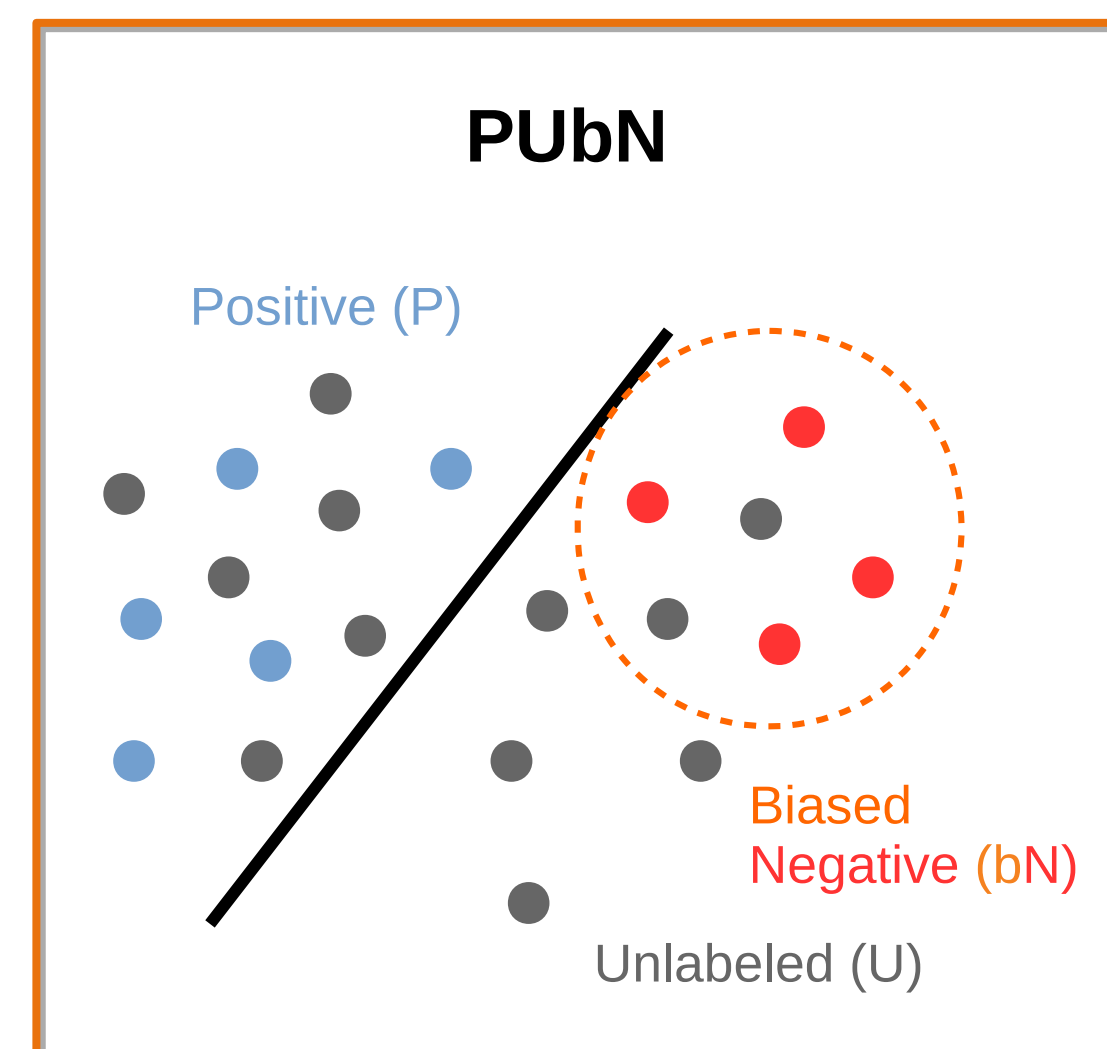
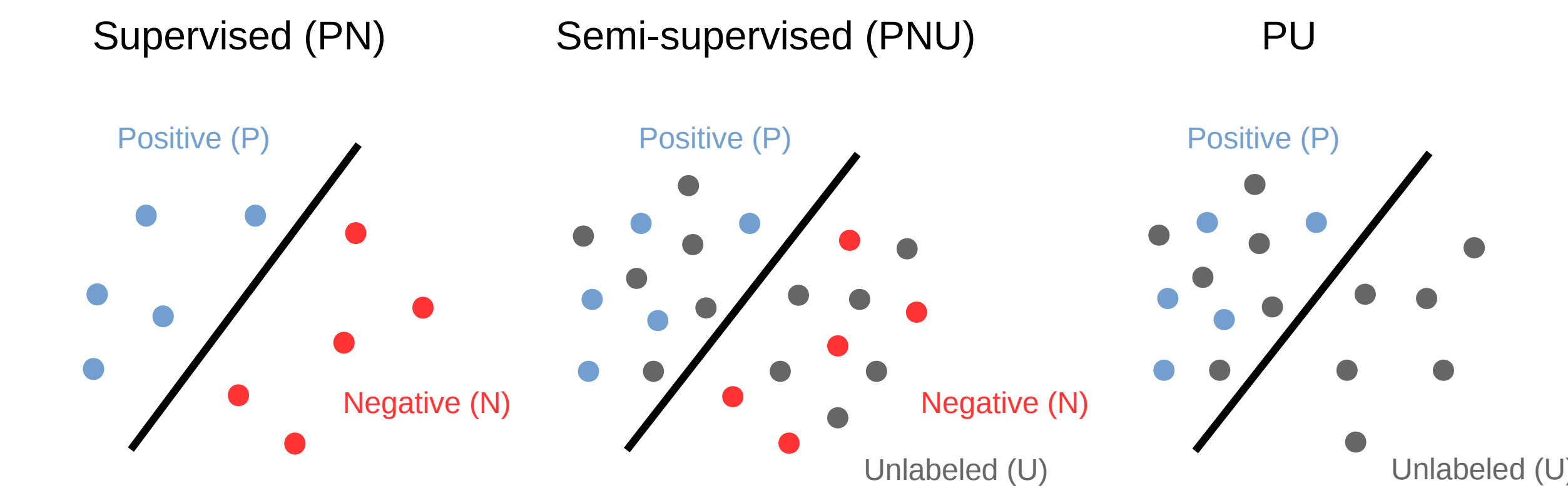


- Information retrieval, text classification
- Medical diagnosis: healthy population that goes through physical exams is biased

## Related Work

- Semi-supervised learning:**  
N data in general implicitly assumed unbiased ; U data used for regularization
- Dataset shift:** a variant that has been rarely studied  
Ex. Covariate shift ; Source component shift
- PU learning:** add bN data
- Pseudo-labeling / Importance-weighting

## Problem Setting



$x$  : feature  
 $y \in \{+1, -1\}$  : label  
 $s \in \{+1, -1\}$  : latent variable causing the bias  
**N data with selection bias**  
 $p(s = +1 | x, y = +1) = 1$

- $\bullet \sim p_P(x) := p(x | y = +1)$
- $\bullet \sim p(x)$
- $\bullet \sim p_{bN}(x) := p(x | y = -1, s = +1)$

## Method

### Empirical Risk Minimization

$$\min_{g \in \mathcal{G}} \underbrace{\mathbb{E}_{(x,y) \sim p(x,y)} [\ell(yg(x))]}_{R(g)} \quad \leftarrow \quad \underbrace{\min_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \ell(y_i g(x_i))}_{\hat{R}(g)}$$

Risk Minimization      Unbiased Estimator      #P data + #N data      Empirical Risk Minimization

Q: The N data are biased

$$R(g) = \underbrace{\pi R_P^+(g)}_{\text{\#P data}} + \underbrace{\rho R_{bN}^-(g)}_{\text{\#bN data}} + \underbrace{(1 - \pi - \rho) R_{s=-1}^-(g)}_{\text{\#U data, Partial risk for samples with low probability of being labeled}}$$

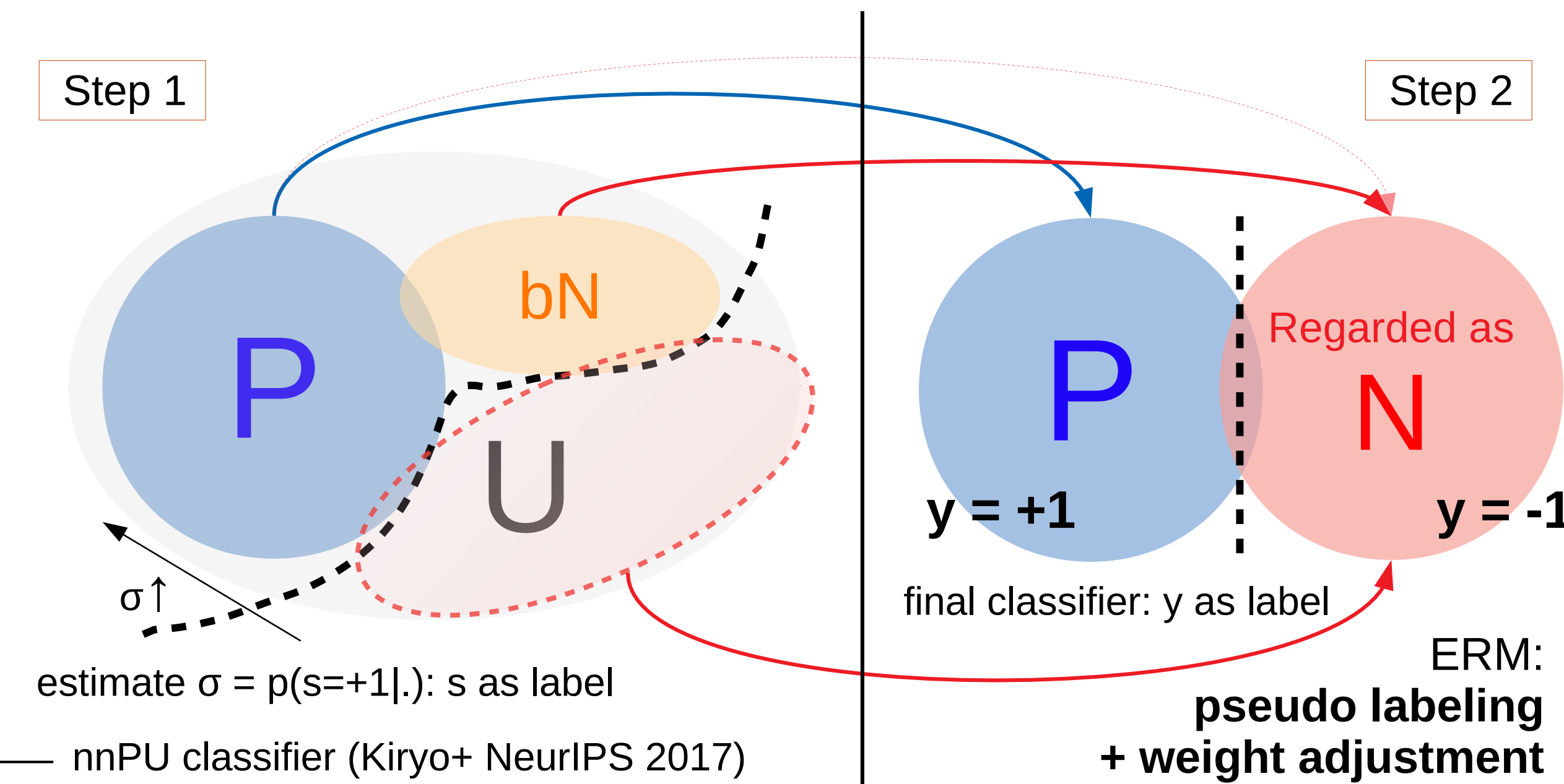
Idea:  $\sigma(x)$  probability of  $x$  being labeled  $\eta > 0$  determining how much attention paid to U data

$$\tilde{R}_{s=-1}^-(g) = \underbrace{\mathbb{E}_{x \sim p(x)} [1_{\sigma(x) \leq \eta} \ell(-g(x)) (1 - \sigma(x))]}_{\text{\#U data, Partial risk for samples with low probability of being labeled}} + \underbrace{\pi \mathbb{E}_{x \sim p_P(x)} \left[ 1_{\sigma(x) > \eta} \ell(-g(x)) \frac{1 - \sigma(x)}{\sigma(x)} \right]}_{\text{\#P data}} + \underbrace{\rho \mathbb{E}_{x \sim p_{bN}} \left[ 1_{\sigma(x) > \eta} \ell(-g(x)) \frac{1 - \sigma(x)}{\sigma(x)} \right]}_{\text{\#bN data}}$$

Partial risk for samples with high probability of being labeled

$$\begin{aligned} R_P^+(g) &:= \mathbb{E}_{x \sim p_P(x)} [\ell(g(x))] & \pi &:= p(y = +1) & \ell &: \text{loss function} \\ R_{bN}^-(g) &:= \mathbb{E}_{x \sim p_{bN}(x)} [\ell(-g(x))] & \rho &:= p(y = -1, s = +1) \\ R_{s=-1}^-(g) &:= \mathbb{E}_{x \sim p(x|s=-1)} [\ell(-g(x))] & \sigma(x) &:= p(s = +1 | x) \end{aligned}$$

### Algorithm Outline



#### PU risk estimator

$$R(g) = \underbrace{\mathbb{E}_{x \sim p(x)} [\ell(-g(x))]}_{\text{\#U data}} + \underbrace{\pi \mathbb{E}_{x \sim p_P(x)} [\ell(g(x)) - \ell(-g(x))]}_{\text{\#P data}}$$

Q: Severe overfitting      A: Avoid regarding all U as N

#### Non-negative correction

$$\tilde{R}_{pu}(g) = \underbrace{\frac{\pi}{n_P} \sum_{x \in \mathcal{X}_P} [\ell(g(x))]}_{\text{\#P data}} + \max \left\{ 0, \underbrace{\frac{1}{n_U} \sum_{x \in \mathcal{X}_U} \ell(-g(x))}_{\text{\#U data}} - \underbrace{\frac{\pi}{n_P} \sum_{x \in \mathcal{X}_P} [\ell(-g(x))]}_{\text{\#P data}} \right\}$$

N partial risk  $\geq 0$

## Estimation Error Bound

With probability at least  $1 - \delta$

$$R(\hat{g}) - R(g^*) \leq \underbrace{\frac{4\pi L_\ell \mathfrak{R}_{n_P, p_P}(\mathcal{G})}{\eta} + \frac{2\pi C_\ell}{\eta} \sqrt{\frac{\ln(6/\delta)}{2n_P}}}_{\text{\#P data}} + \underbrace{\frac{4\rho L_\ell \mathfrak{R}_{n_{bN}, p_{bN}}(\mathcal{G})}{\eta} + \frac{2\rho C_\ell}{\eta} \sqrt{\frac{\ln(6/\delta)}{2n_{bN}}}}_{\text{\#bN data}} + \underbrace{4L_\ell \mathfrak{R}_{n_U, p}(\mathcal{G}) + 2C_\ell \sqrt{\frac{\ln(6/\delta)}{2n_U}}}_{\text{\#U data}} + \underbrace{2C_\ell \sqrt{\zeta \epsilon} + \frac{2C_\ell}{\eta} \sqrt{(1 - \zeta) \epsilon}}_{\text{Bias due to inexact approximation of } \sigma}$$

$\mathfrak{R}_{n,q}(\mathcal{G})$  : Rademacher Complexity

$\hat{\sigma}$  : estimate of  $\sigma$        $\zeta := p(\hat{\sigma}(x) \leq \eta)$        $\epsilon := \mathbb{E}_{x \sim p(x)} [|\hat{\sigma}(x) - \sigma(x)|^2]$

Assumption     $\ell$  is  $L_\ell$ -Lipschitz     $\sup_{g \in \mathcal{G}} \|g\|_\infty \leq C_g$      $\sup_{|z| \leq C_g} \ell(z) \leq C_\ell$

1) If RC terms vanish asymptotically, it holds a.s.

$$\limsup R(\hat{g}) - R(g^*) \leq \underbrace{2C_\ell \sqrt{\zeta \epsilon} + 2(C_\ell/\eta) \sqrt{(1 - \zeta) \epsilon}}_{\text{Bias due to inexact approximation of } \sigma}$$

2) Classical convergence rate + bias

3) To control  $\epsilon$ : approximation error + estimation error

## Experiments

### Setting

- Models: ConvNet / ResNet / FCN + Training: Amsgrad
- VALIDATION!** equally composed of P+U+bN
- 1/10 #U ~ #P = #bN ; Same model for the two steps

### Baselines

- nnPNU (Sakai+ 2017 ICML):  
linear combination of PU and PN risk
- PU  $\rightarrow$  PN: one classifier for s and one to separate P from bN

### Results

bN data helps

Dataset	P	$\pi$	bN	$\rho$	nnPU/nnPNU	PUBN(N)	PU $\rightarrow$ PN
MNIST	2, 4, 6, 8, 10	0.49	Not given	NA	5.76 $\pm$ 1.04	<b>4.64 <math>\pm</math> 0.62</b>	NA
			1, 3, 5	0.3	5.33 $\pm$ 0.97	<b>4.05 <math>\pm</math> 0.27</b>	<b>4.00 <math>\pm</math> 0.30</b>
			9 > 5 > others	0.2	4.60 $\pm$ 0.65	<b>3.91 <math>\pm</math> 0.66</b>	<b>3.77 <math>\pm</math> 0.31</b>
CIFAR-10	Airplane, automobile, ship, truck	0.4	Not given	NA	12.02 $\pm$ 0.65	<b>10.70 <math>\pm</math> 0.57</b>	NA
			Cat, dog, horse	0.3	10.25 $\pm$ 0.38	<b>9.71 <math>\pm</math> 0.51</b>	10.37 $\pm$ 0.65
			Horse > deer = frog > others	0.25	9.98 $\pm$ 0.53	<b>9.92 <math>\pm</math> 0.42</b>	10.17 $\pm$ 0.35
CIFAR-10	Cat, deer, dog, horse	0.4	Not given	NA	23.78 $\pm$ 1.04	<b>21.13 <math>\pm</math> 0.90</b>	NA
			Bird, frog	0.2	22.00 $\pm$ 0.53	<b>18.83 <math>\pm</math> 0.71</b>	19.88 $\pm$ 0.62
			Car, truck	0.2	22.00 $\pm$ 0.74	<b>20.19 <math>\pm</math> 1.06</b>	21.83 $\pm$ 1.36
20 Newsgroups	alt., comp., misc., rec.	0.56	Not given	NA	14.67 $\pm$ 0.87	<b>13.30 <math>\pm</math> 0.53</b>	NA
			sci.	0.21	14.69 $\pm$ 0.46	<b>13.10 <math>\pm</math> 0.90</b>	13.58 $\pm$ 0.97
			talk.	0.17	14.38 $\pm$ 0.74	<b>12.61 <math>\pm</math> 0.75</b>	13.76 $\pm$ 0.66
			soc. > talk. > sci.	0.1	14.41 $\pm$ 0.76	<b>12.18 <math>\pm</math> 0.59</b>	12.92 $\pm$ 0.51