

Accepted to ICML2017

# Semi-Supervised Classification based on Classification from Positive and Unlabeled Data

Tomoya Sakai,<sup>1,2</sup> Marthinus Christoffel du Plessis,  
Gang Niu,<sup>1,2</sup> and Masashi Sugiyama<sup>2,1</sup>

<sup>1</sup>The University of Tokyo, Japan

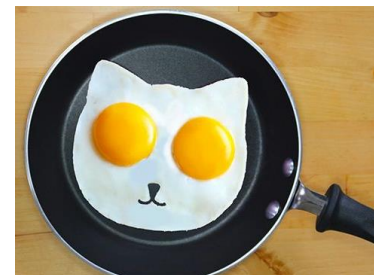
<sup>2</sup>RIKEN, Japan

# Classification Problem

Identify class or category of data points

Examples

- Image is a cat (**positive** class) or not (**negative** class)



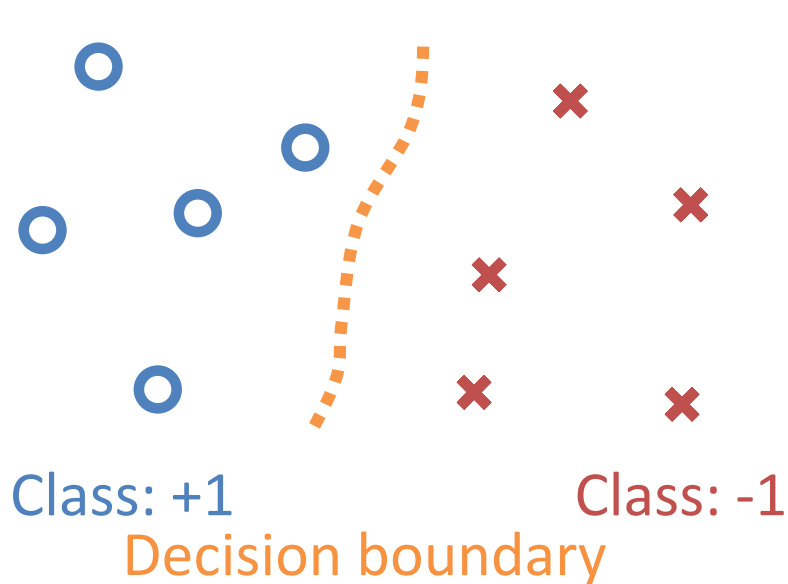
Class: +1

Decision boundary

Class: -1

# Supervised Learning (PN Learning) <sup>3</sup>

Learn from **labeled** (positive and negative) data



Pattern:  $\mathbf{x} \in \mathbb{R}^d$   
Label:  $y \in \{\pm 1\}$

○ : positive (P) data

$$\{\mathbf{x}_i^P\}_{i=1}^{n_P} \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x} \mid y = +1)$$

✕ : negative (N) data

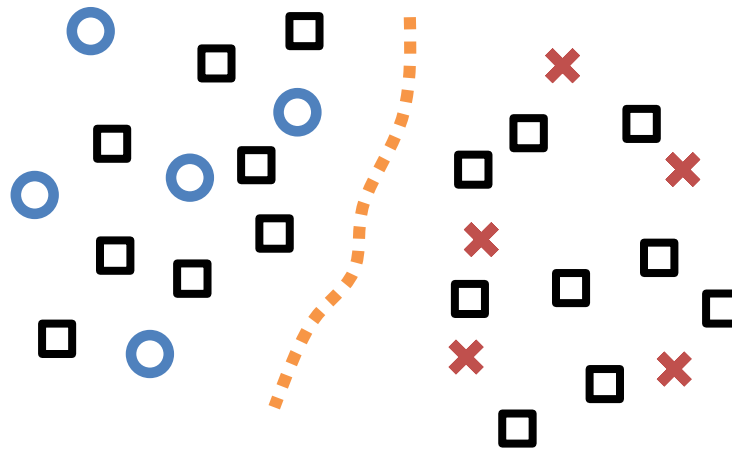
$$\{\mathbf{x}_i^N\}_{i=1}^{n_N} \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x} \mid y = -1)$$

😊 Better performance with many labeled data

😞 Collecting labeled data is **costly**

# Semi-Supervised Learning (SSL)

Learn from a **small** amount of **labeled** data  
and a **large** amount of **unlabeled** data



○ : positive data

× : negative data

□ : unlabeled data

$$\{\mathbf{x}_i^U\}_{i=1}^{n_U} \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x})$$

😊 Labeling cost becomes cheap

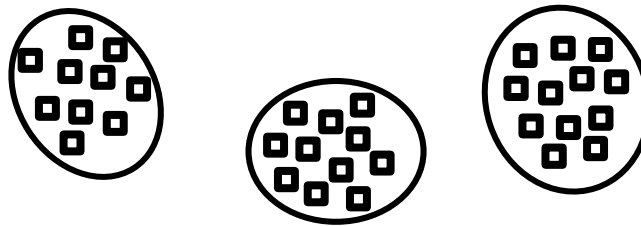
# This Work

Existing methods:

- Require **strong distributional assumptions** for utilizing unlabeled data

- ◆ Ex. the cluster assumption requires the samples in the same cluster be likely to share the same label

(Chapelle et al., NIPS, 2002)



- ◆ If the distributional assumptions are **not** satisfied, the performance of the existing methods **decreases**

Propose method:

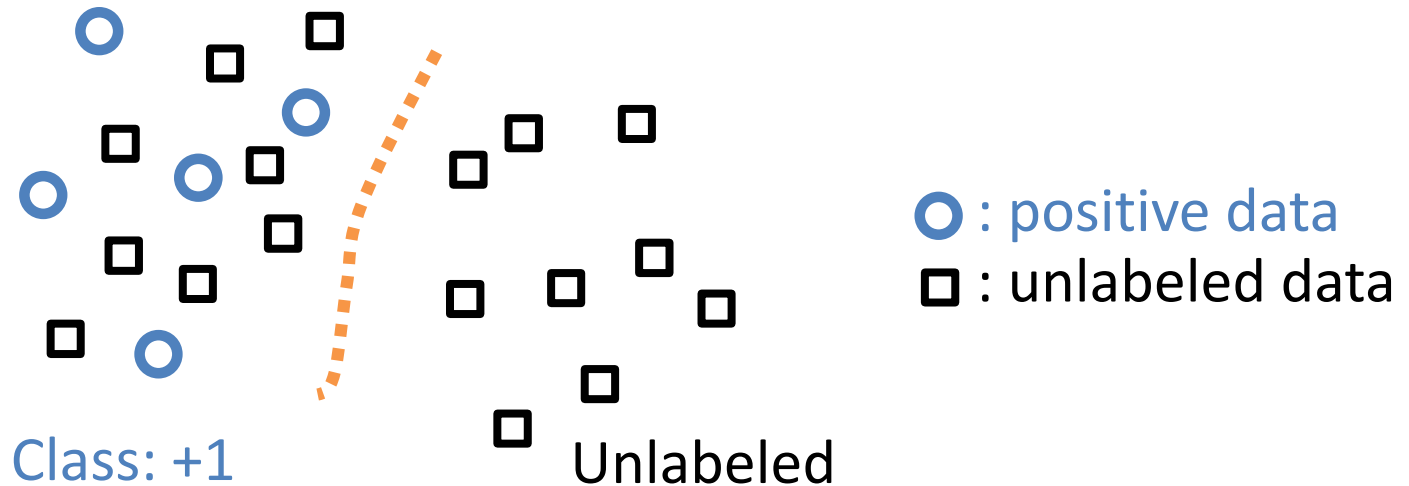
- **NOT** require the strong distributional assumptions

# Outline

1. Introduction
2. **Proposed Method**
3. Experiment
4. Conclusions

# Our Idea: Use of PU Learning

Learning from **positive (P)** and **unlabeled (U)** data

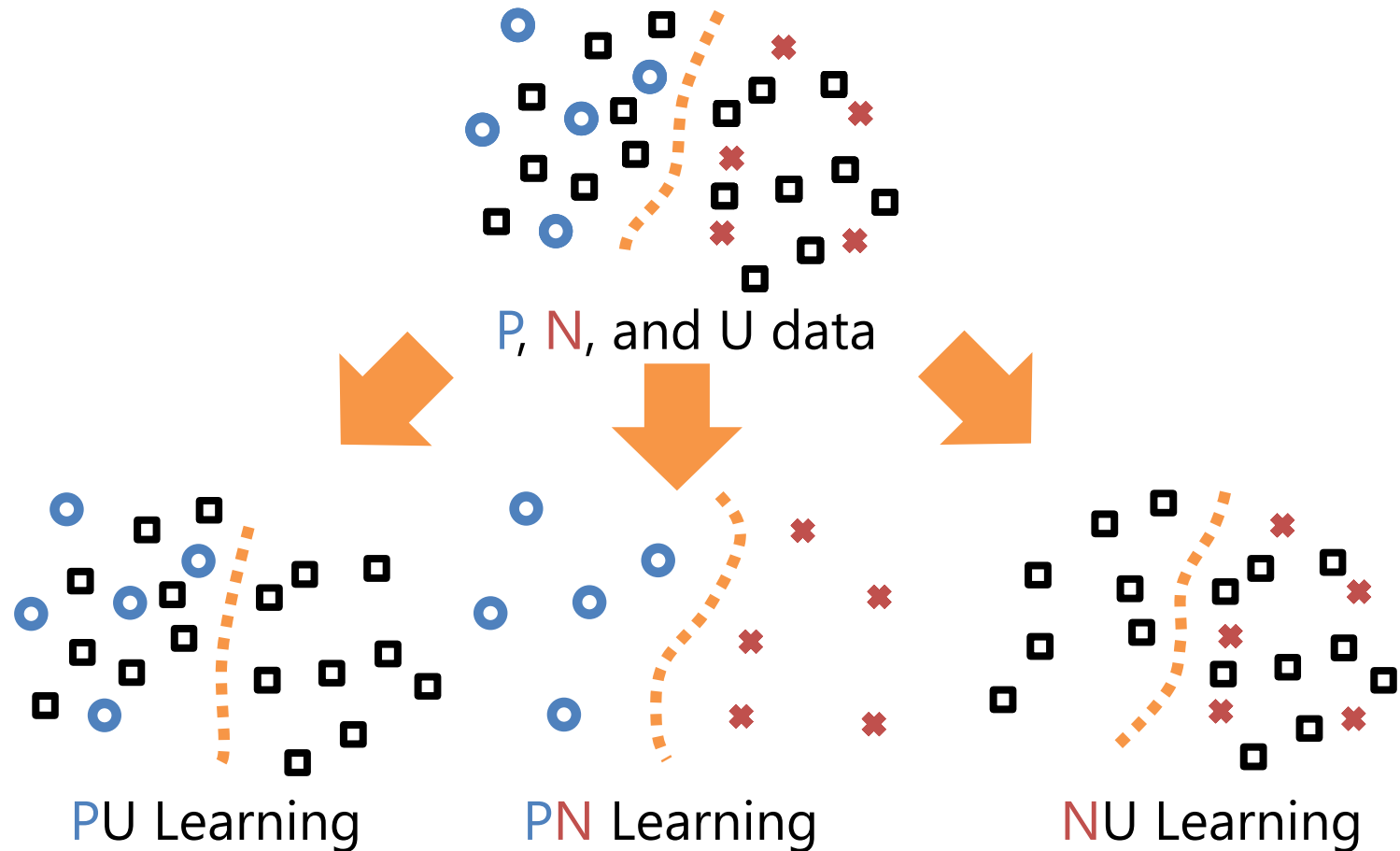


😊 Utilize unlabeled data **without** the distributional assumption (du Plessis et al., NIPS, 2014)

How to use PU learning for semi-supervised learning?

# Our Idea (cont.): SSL Decomposition <sup>8</sup>

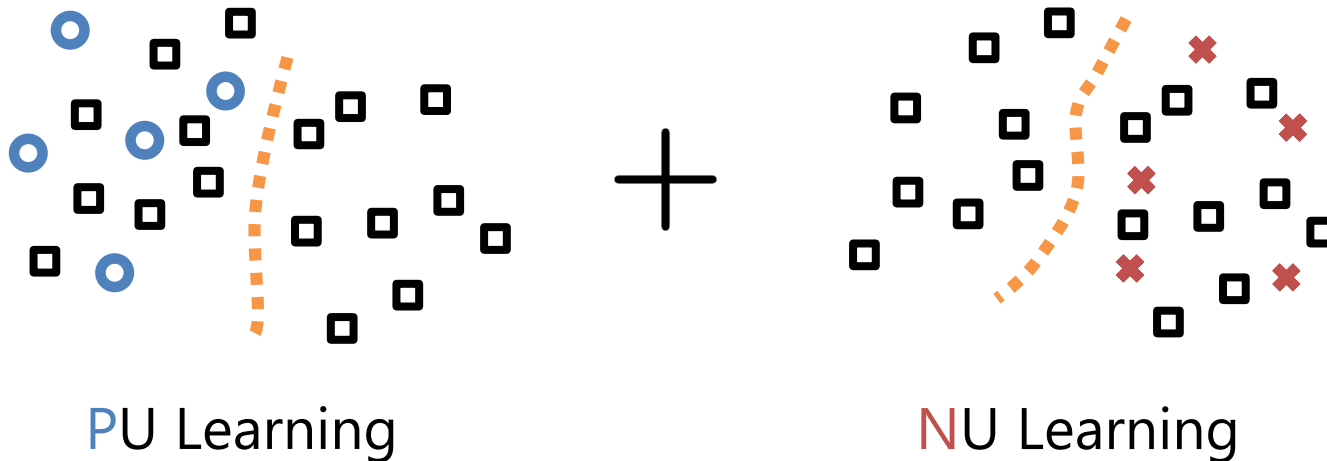
Decompose Semi-Supervised Learning (SSL)  
into PU, PN, and NU learning



➤ Reconstruct SSL from PU, PN, and NU learning

# PU+NU: PUNU Learning

Since PU and NU are symmetric,  
it might be natural to combine **PU** and **NU** learning



Is this really good?

# Is PUNU Learning Better?

From Niu et al. NIPS, (2016), we have

(Case I) The size of unlabeled data is sufficiently large

$$\left\{ \begin{array}{l} \text{PU} > \text{PN} > \text{NU} \text{ (classification accuracy)} \\ \text{or} \\ \text{NU} > \text{PN} > \text{PU} \end{array} \right. \quad \text{➤ PU or NU is the best}$$

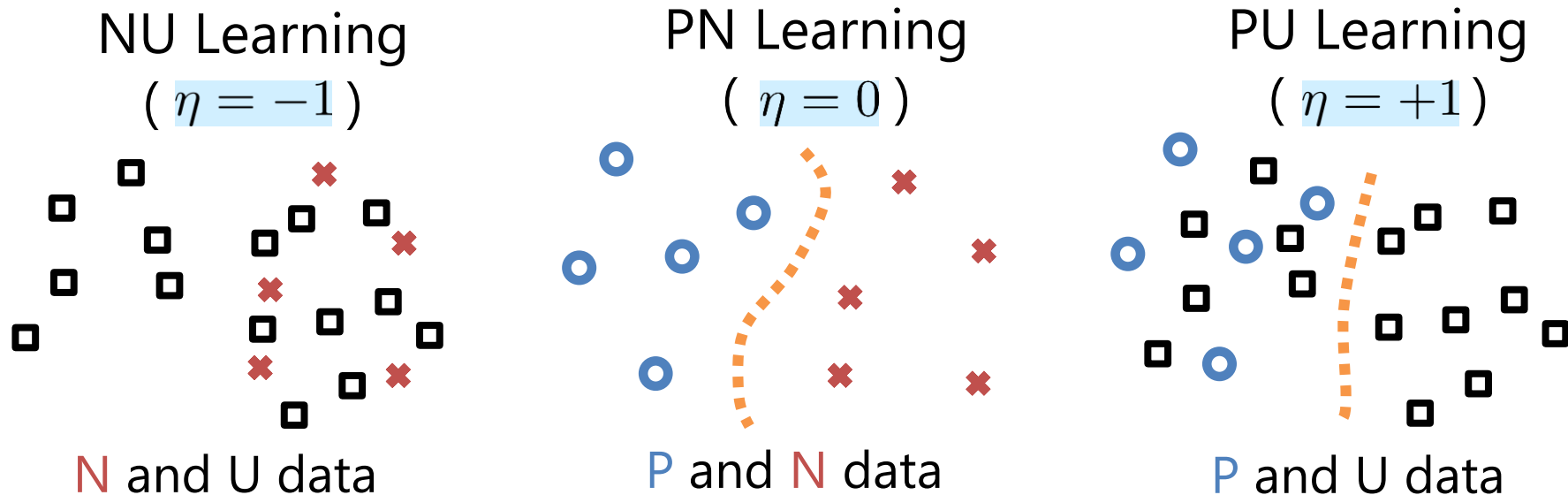
(Case II) The size of unlabeled data is small

$$\left\{ \begin{array}{l} \text{PN} > \text{PU} > \text{NU} \\ \text{or} \\ \text{PN} > \text{NU} > \text{PU} \end{array} \right. \quad \text{➤ PN is always the best}$$

- PUNU (=PU+NU) is **not** a good idea  
since it contains the worst one in the combination
- Combine PN with PU/NU would be promising  
(PN+PU or PN+NU is better)

# Proposed Method: PNU Learning

Combine PN with PU/NU learning (PN+PU or PN+NU)



➤ The PNU risk:

$$R_{\text{PNU}}^{\eta}(g) := \begin{cases} (1 - \eta)R_{\text{PN}}(g) + \eta R_{\text{PU}}(g) & (\eta \geq 0) \\ (1 + \eta)R_{\text{PN}}(g) - \eta R_{\text{NU}}(g) & (\eta < 0) \end{cases}$$

◆ Obtain the decision rule by minimizing the PNU risk

$\eta \in [-1, 1]$   $g$ : classifier

$R_{\text{PN}}, R_{\text{PU}}, R_{\text{NU}}$ : risks in PN, PU, and NU learning

# Comparison with Existing Methods 12

Existing approach:

- Design regularizer based on the distributional assumption
- Use unlabeled data for **regularization**

$$\hat{g} = \operatorname{argmin}_g \hat{R}_{\text{PN}}(g) + \lambda \hat{W}(g)$$

Entropy regularizer

Manifold regularizer

⋮

Unlabeled data are used

Our approach:

$R_{\text{PN}}$ : risk in supervised learning

- **Not** require the strong distributional assumptions
- Utilize unlabeled data for **risk evaluation**

$$\hat{g} = \operatorname{argmin}_g \hat{R}_{\text{PNU}}^{\eta}(g)$$

Labeled and unlabeled data are used

# Theoretical Analyses

Without the distributional assumptions, we prove the follows:

## ◆ Generalization error bound

$$\mathbb{E}_{p(\mathbf{x}, y)}[\ell_{0-1}(yg(\mathbf{x}))] \leq 2\hat{R}_{\text{PNU}}^\eta(g) + \mathcal{O}_p\left(\frac{1}{\sqrt{n_{\text{P}}}} + \frac{1}{\sqrt{n_{\text{N}}}} + \frac{1}{\sqrt{n_{\text{U}}}}\right) \forall g \in \mathcal{G}$$

➤ Unlabeled data help reduce the bound 

➤ Optimal parametric convergence rate

## ◆ Variance reduction

$\text{Var}[\hat{R}_{\text{PNU}}^\eta(g)] < \text{Var}[\hat{R}_{\text{PN}}(g)]$  for some  $\eta$  if  $n_{\text{U}}$  is sufficiently large

➤ The PNU risk is **stable** in terms of the variance

➤ More **stable cross-validation**

$$\mathcal{G} = \{g(\mathbf{x}) = \langle \mathbf{w}, \phi(\mathbf{x}) \rangle \mid \|\mathbf{w}\| \leq C_w, \|\phi(\mathbf{x})\| \leq C_\phi\}$$

$R_{\text{PN}}$ : risk in supervised learning     $\ell_{0-1}(m)$ : 0-1 loss

# Outline

1. Introduction
2. Proposed Method
- 3. Experiment**
4. Conclusions

# Image Classification

(smaller is better)

Average with the standard error of misclassification rate  $n_L = 100$

Data set	$n_U$	$\theta_P$	$\hat{\theta}_P$	PNU	ER	LapSVM	SMIR	WellSVM
Arts	1000	0.50	0.49 (0.01)	<b>27.4 (1.3)</b>	<b>26.6 (0.5)</b>	<b>26.1 (0.7)</b>	40.1 (3.9)	<b>27.5 (0.5)</b>
	5000	0.50	0.50 (0.01)	<b>24.8 (0.6)</b>	26.1 (0.5)	26.1 (0.4)	30.1 (1.6)	N/A
	10000	0.50	0.52 (0.01)	<b>25.6 (0.7)</b>	<b>25.4 (0.5)</b>	<b>25.5 (0.6)</b>	N/A	N/A
Deserts	1000	0.73	0.67 (0.01)	<b>13.0 (0.5)</b>	15.3 (0.6)	16.7 (0.8)	17.2 (0.8)	18.2 (0.7)
	5000	0.73	0.67 (0.01)	<b>13.4 (0.4)</b>	<b>13.3 (0.5)</b>	16.6 (0.6)	24.4 (0.6)	N/A
	10000	0.73	0.68 (0.01)	<b>13.3 (0.5)</b>	<b>13.7 (0.6)</b>	16.8 (0.8)	N/A	N/A
Fields	1000	0.65	0.57 (0.01)	<b>22.4 (1.0)</b>	26.2 (1.0)	26.6 (1.3)	28.2 (1.1)	26.6 (0.8)
	5000	0.65	0.57 (0.01)	<b>20.6 (0.5)</b>	22.6 (0.6)	24.7 (0.8)	29.6 (1.2)	N/A
	10000	0.65	0.57 (0.01)	<b>21.6 (0.6)</b>	<b>22.5 (0.6)</b>	25.0 (0.9)	N/A	N/A
Stadiums	1000	0.50	0.50 (0.01)	<b>11.4 (0.4)</b>	<b>11.5 (0.5)</b>	12.5 (0.5)	<b>17.4 (3.6)</b>	<b>11.7 (0.4)</b>
	5000	0.50	0.50 (0.01)	<b>11.0 (0.5)</b>	<b>10.9 (0.3)</b>	<b>11.1 (0.3)</b>	13.4 (0.7)	N/A
	10000	0.50	0.51 (0.00)	<b>10.7 (0.3)</b>	<b>10.9 (0.3)</b>	<b>11.2 (0.2)</b>	N/A	N/A
Platforms	1000	0.27	0.33 (0.01)	<b>21.8 (0.5)</b>	23.9 (0.6)	24.1 (0.5)	30.1 (2.3)	26.2 (0.8)
	5000	0.27	0.34 (0.01)	<b>23.3 (0.8)</b>	<b>24.4 (0.7)</b>	<b>24.9 (0.7)</b>	26.6 (0.3)	N/A
	10000	0.27	0.34 (0.01)	<b>21.4 (0.5)</b>	24.3 (0.6)	24.8 (0.5)	N/A	N/A

➤ PNU learning outperforms the existing methods

Class-prior is estimated by

the energy distance minimization method (Kawakubo et al., IEICE-ED, 2016)

\* Colored cells indicate the best and comparable method in terms of t-test (sig. lev. 5%)

\* The methods taking 2 hours were omitted and indicated as "N/A"

# Outline

1. Introduction
2. Proposed Method
3. Experiment
4. Conclusions

# World of SSL based on PU Learning <sup>17</sup>

---

	PN	PU
MR Minimization	Vapnik, Springer, 2000.	du Plessis et al., NIPS, 2014.

---

MR: misclassification rate

- Obtain classification rules based on MR minimization

# World of SSL based on PU Learning<sup>18</sup>

	PN	PU	Semi-Supervised
MR Minimization	Vapnik, Springer, 2000.	du Plessis et al., NIPS, 2014.	This talk!

MR: misclassification rate

- Obtain classification rules based on MR minimization

SSL: Semi-Supervised Learning

# World of SSL based on PU Learning <sup>19</sup>

	PN	PU	Semi-Supervised
MR Minimization	Vapnik, Springer, 2000.	du Plessis et al., NIPS, 2014.	This talk!
AUC Maximization	Herschtal and Raskutti, ICML, 2004.	Tomorrow poster session! Sakai et al., MLJ, 2017.	

AUC: area under the receiver operating characteristic curve

➤ Useful for imbalanced classification when  $n_P \ll n_N$

$$\{\mathbf{x}_i^P\}_{i=1}^{n_P} \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x} \mid y = +1)$$

$$\{\mathbf{x}_i^N\}_{i=1}^{n_N} \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x} \mid y = -1)$$

# World of SSL based on PU Learning<sup>20</sup>

	PN	PU	Semi-Supervised
MR Minimization	Vapnik, Springer, 2000.	du Plessis et al., NIPS, 2014.	This talk!
AUC Maximization	Herschtal and Raskutti, ICML, 2004.	Tomorrow poster session! Sakai et al., MLJ, 2017.	
SMI Estimation	Suzuki et al., BMC bioinform, 2009.	Sakai et al., arXiv, 2017.	To be explored

SMI: squared-loss mutual information

- Statistical dependency measure
- Useful for dimension reduction, feature selection, independence test, and object matching

# World of SSL based on PU Learning <sup>21</sup>

	PN	PU	Semi-Supervised
MR Minimization	Vapnik, Springer, 2000.	du Plessis et al., NIPS, 2014.	This talk!
AUC Maximization	Herschtal and Raskutti, ICML, 2004.	Tomorrow poster session! Sakai et al., MLJ, 2017.	
SMI Estimation	Suzuki et al., BMC bioinform, 2009.	Sakai et al., arXiv, 2017.	To be explored

Code available at

<https://github.com/t-sakai-kure/PNU>

