

# Machine Learning from Weak Supervision:

Towards Accurate Classification with Low Labeling Costs

Slides:

<http://goo.gl/meiTwY>

Masashi Sugiyama

RIKEN Center for Advanced Intelligence Project/  
The University of Tokyo



東京大学  
THE UNIVERSITY OF TOKYO



# About Myself

2

## ■ Affiliations:

- Director: RIKEN AIP
- Professor: University of Tokyo
- Consultant: several local startups

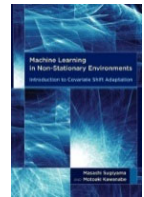
## ■ Research interests:

- Theory and algorithms of ML
- Real-world applications with partners (signal, image, language, brain, cars, robots, optics, ads, medicine, biology...)

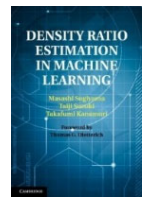
## ■ Goal:

- Develop practically useful algorithms that have theoretical support

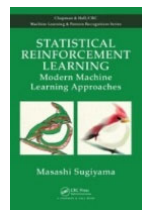
Sugiyama & Kawanabe, **Machine Learning in Non-Stationary Environments**, MIT Press, 2012



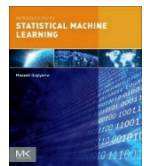
Sugiyama, Suzuki & Kanamori, **Density Ratio Estimation in Machine Learning**, Cambridge University Press, 2012



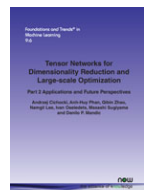
Sugiyama, **Statistical Reinforcement Learning**, Chapman and Hall/CRC, 2015



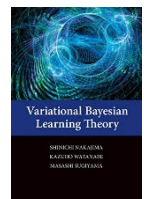
Sugiyama, **Introduction to Statistical Machine Learning**, Morgan Kaufmann, 2015



Cichocki, Phan, Zhao, Lee, Oseledets, Sugiyama & Mandic, **Tensor Networks for Dimensionality Reduction and Large-Scale Optimizations**, Now, 2017



Nakajima, Watanabe & Sugiyama, **Variational Bayesian Learning Theory**, Cambridge University Press, 2019



# What Is This Tutorial about?

3

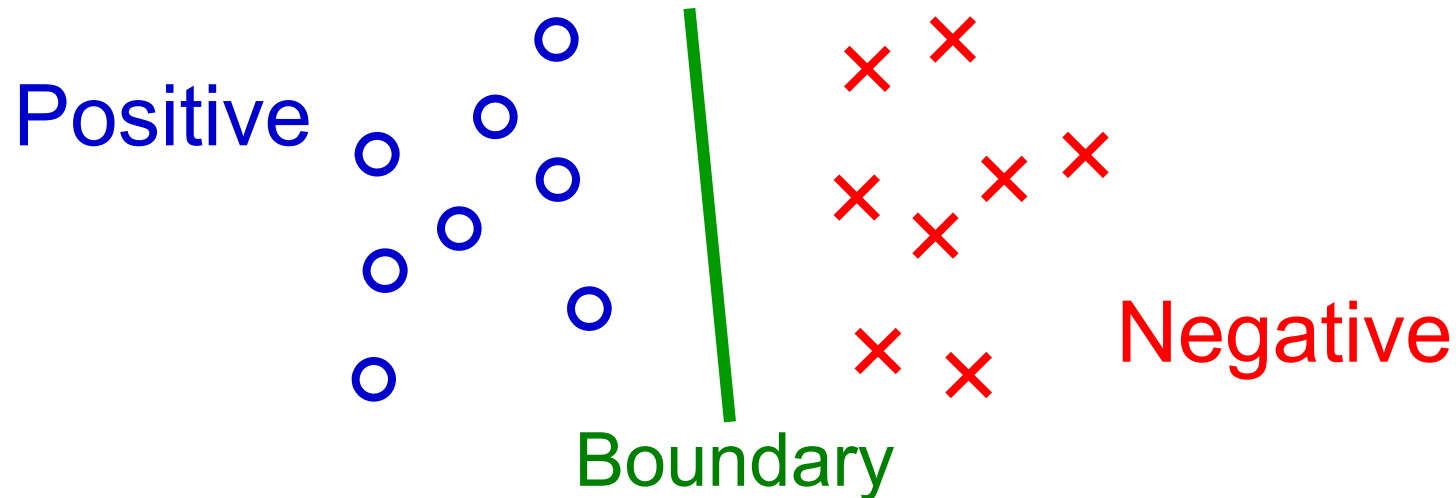
- **Machine learning from big labeled data** is highly successful.
  - Speech recognition, image understanding, natural language translation, recommendation...
- However, there are various applications where **massive labeled data is not available**.
  - Medicine, disaster, robots, brain, ...

# What Is This Tutorial about?

4

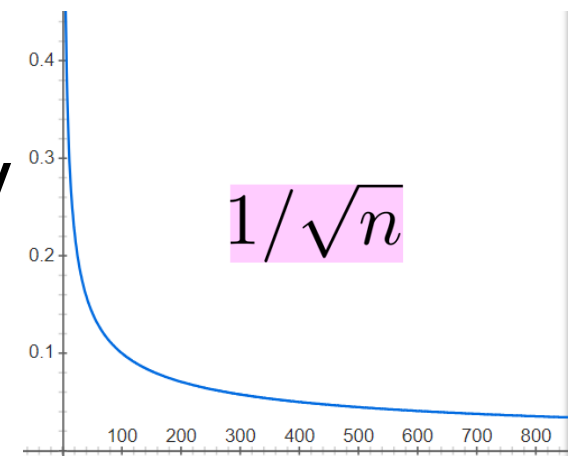
- There are many approaches to coping with the label-cost problem:
  - Improve data collection (e.g., crowdsourcing)
  - Use a simulator to generate pseudo data
  - Use domain knowledge (i.e., engineering)
  - Use cheap but weak data (e.g., unlabeled)
  
- I introduce our recent advances in **classification from weak supervision.**

# Our Target Problem: Binary Supervised Classification



- Larger amount of labeled data yields better classification accuracy.
- Estimation error of the boundary decreases in order  $1/\sqrt{n}$ .

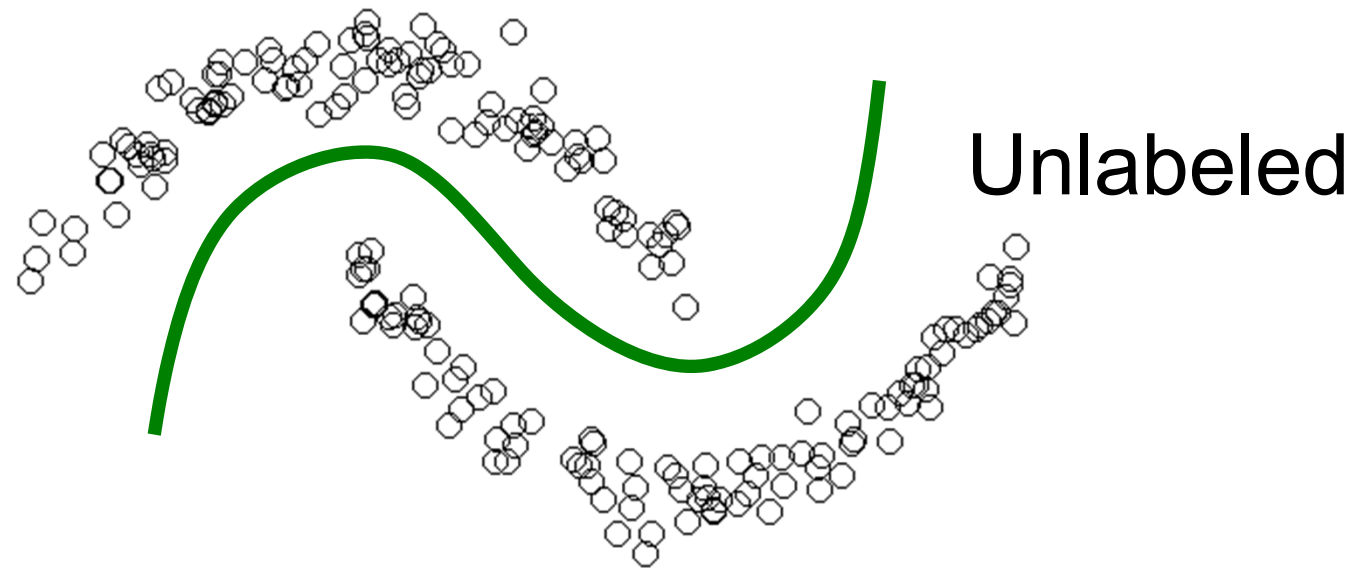
$n$  : Number of labeled samples



# Unsupervised Classification

6

- Gathering labeled data is costly. Let's use **unlabeled data** that are often cheap to collect:

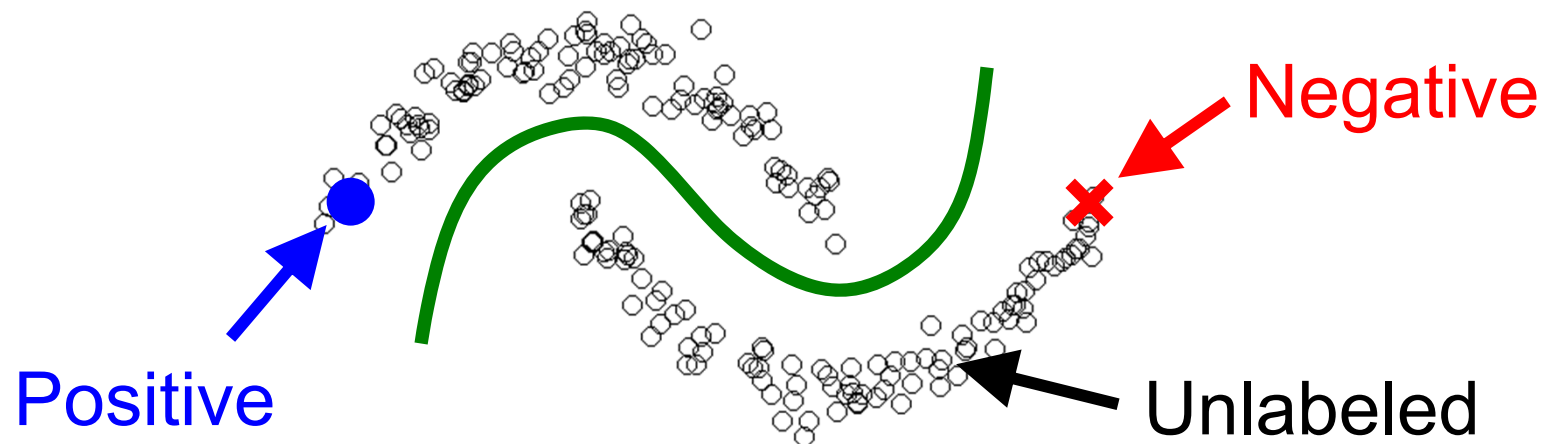


- Unsupervised classification is typically **clustering**.
- This works well only when **each cluster corresponds to a class**.

# Semi-Supervised Classification <sup>7</sup>

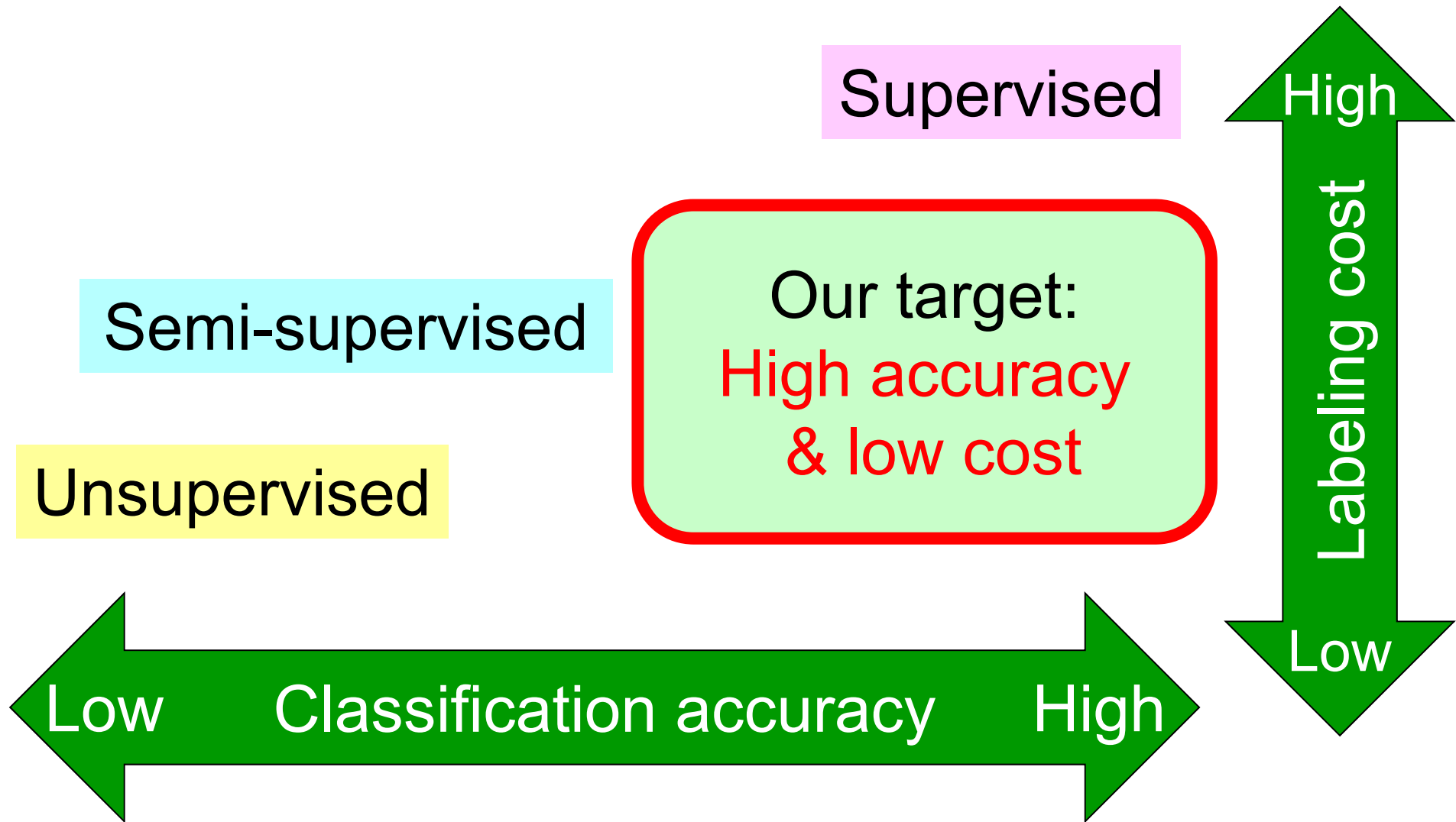
Chapelle, Schölkopf & Zien (MIT Press 2006) and many

- Use a large number of **unlabeled** samples and a small number of **labeled** samples.
- Find a boundary **along the cluster structure** induced by unlabeled samples:
  - Sometimes very useful.
  - But not that different from unsupervised classification.



# Classification of Classification

8







# This Tutorial in a Nutshell

9

1. Background
2. PN Classification
3. PU Classification
4. PNU Classification
5. Pconf Classification
6. UU Classification
7. SU Classification
8. Comp Classification
9. Summary

- **P**: Positive
- **N**: Negative
- **U**: Unlabeled
- **Conf**: Confidence
- **S**: Similar
- **Comp**: Complementary

Slides:

<http://goo.gl/meiTwY>

# Method 1: PU Classification

10

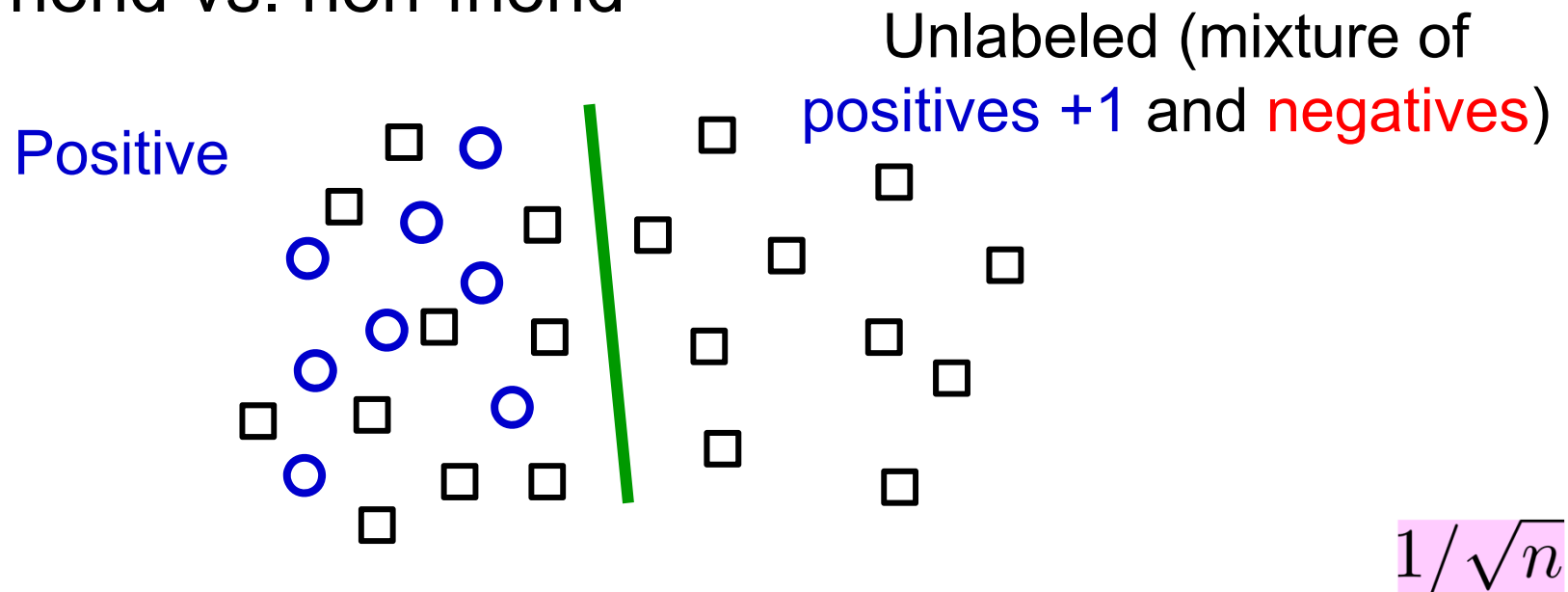
du Plessis, Niu & Sugiyama (NIPS2014, ICML2015)

Niu, du Plessis, Sakai, Ma & Sugiyama (NIPS2016)

Kiryo, Niu, du Plessis & Sugiyama (NIPS2017)

■ Only PU data is available; N data is missing:

- Click vs. non-click
- Friend vs. non-friend



■ From PU data, PN classifiers are trainable!

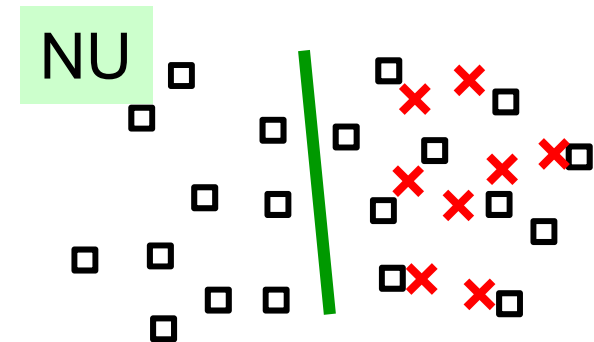
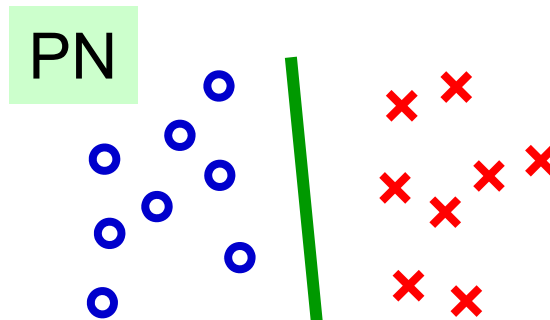
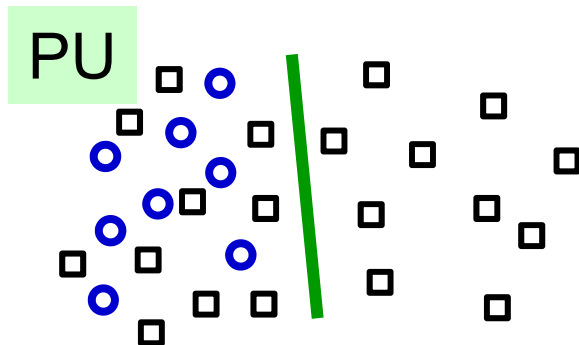
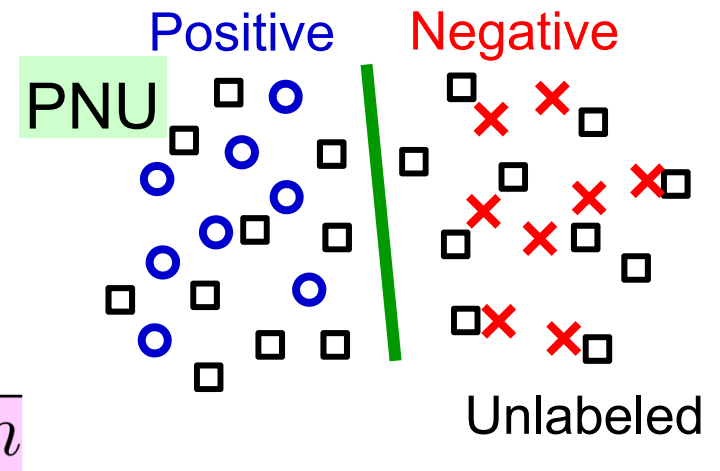
# Method 2: PNU Classification <sup>11</sup> (Semi-Supervised Classification)

Sakai, du Plessis, Niu & Sugiyama (ICML2017)

■ Let's decompose PNU into PU, PN, and NU:

- Each is solvable.
- Let's combine them!

■ Without cluster assumptions,  
PN classifiers are trainable!

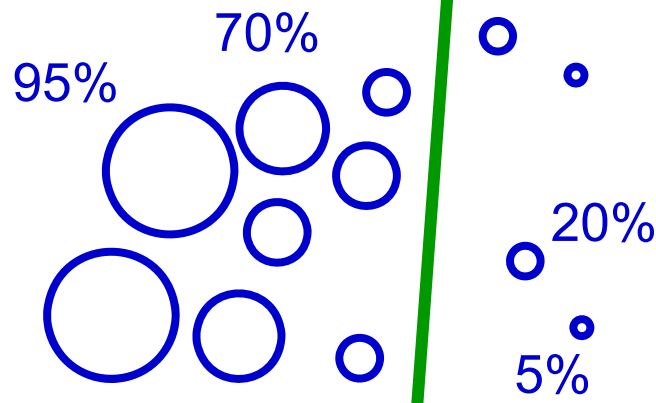


# Method 3: Pconf Classification <sup>12</sup>

Ishida, Niu & Sugiyama (NeurIPS2018)

- Only P data is available, not U data:
  - Data from rival companies cannot be obtained.
  - Only positive results are reported (publication bias).
- “Only-P learning” is unsupervised.
- From Pconf data, PN classifiers are trainable!

Positive confidence



$$1/\sqrt{n}$$

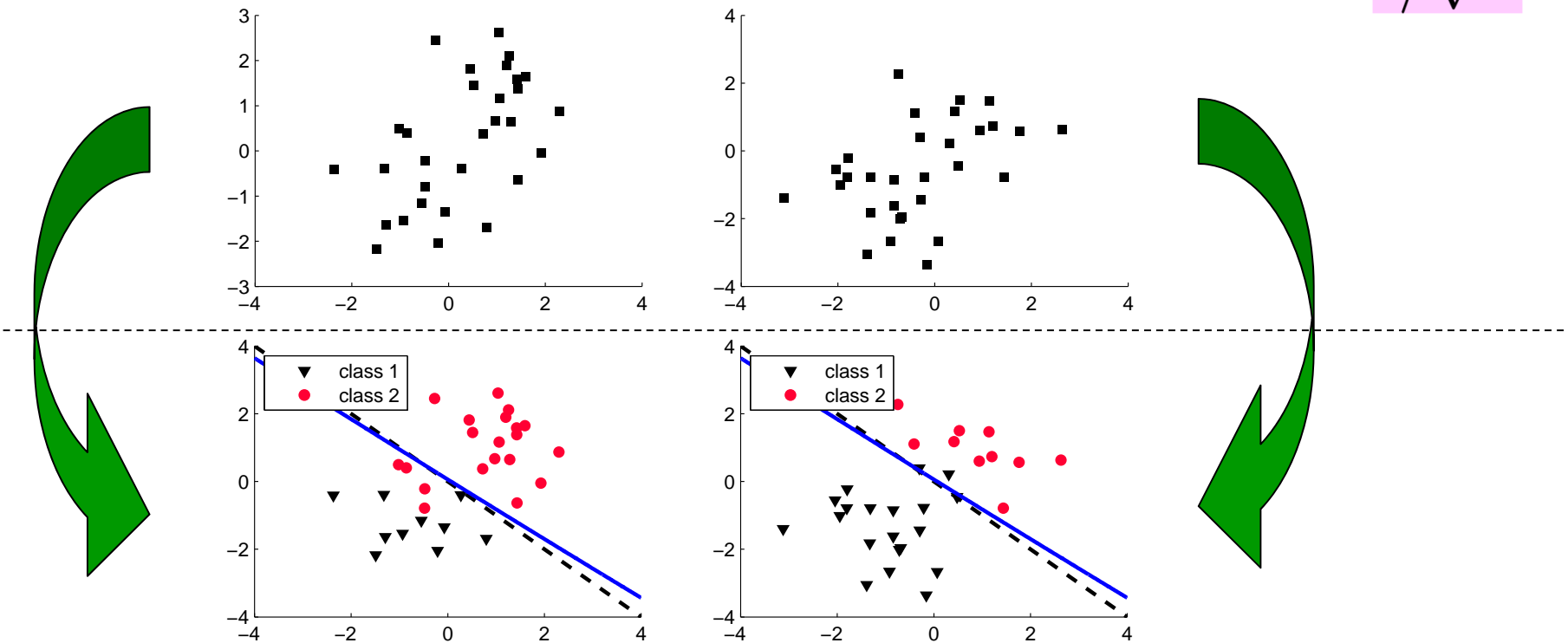
# Method 4: UU Classification

13

du Plessis, Niu & Sugiyama (TAAI2013)  
Nan, Niu, Menon & Sugiyama (ICLR2019)

- From two sets of unlabeled data with different class priors, PN classifiers are trainable!

$$1/\sqrt{n}$$



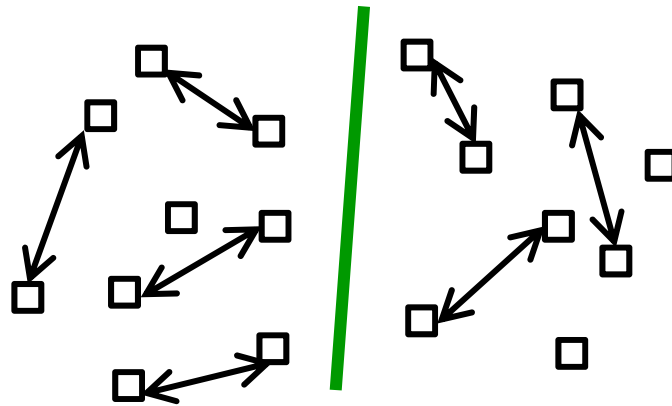
# Method 5: SU Classification

14

Bao, Niu & Sugiyama (ICML2018)

- **Delicate classification** (money, religion...):
  - Highly hesitant to directly answer questions.
  - Less reluctant to just say “**same as him/her**”.
- **From SU data, PN classifiers are trainable!**

$$1/\sqrt{n}$$



# Method 6: Comp Classification<sup>15</sup>

Ishida, Niu, Hu & Sugiyama (NIPS2017)  
Ishida, Niu, Menon & Sugiyama (ICML2019)

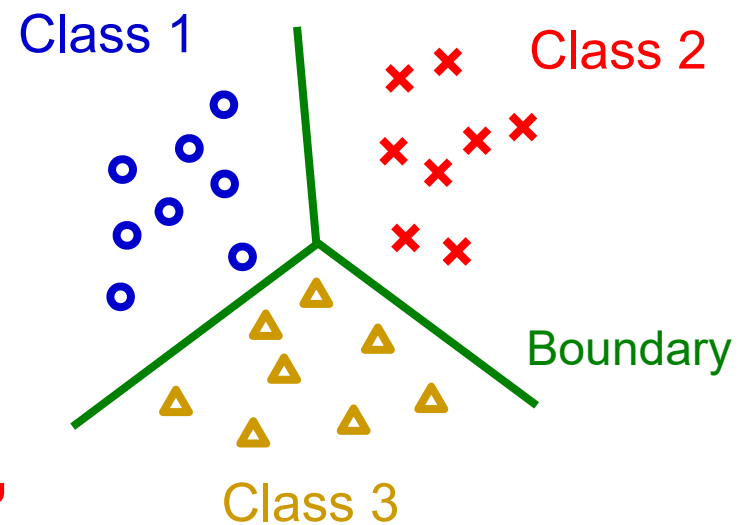
## ■ Labeling patterns in **multi-class** problems:

- Selecting a correct class from a long list of candidate classes is extremely painful.

## ■ **Complementary labels**:

- Specify a class that a pattern does **not** belong to.
- This is much easier and faster to perform!

## ■ **From complementary labels, classifiers are trainable!**



$$1/\sqrt{n}$$



# Contents

16

1. Background
2. **PN Classification**
3. PU Classification
4. PNU Classification
5. Pconf Classification
6. UU Classification
7. SU Classification
8. Comp Classification
9. Summary

- **P**: Positive
- **N**: Negative
- **U**: Unlabeled
- **Conf**: Confidence
- **S**: Similar
- **Comp**: Complementary

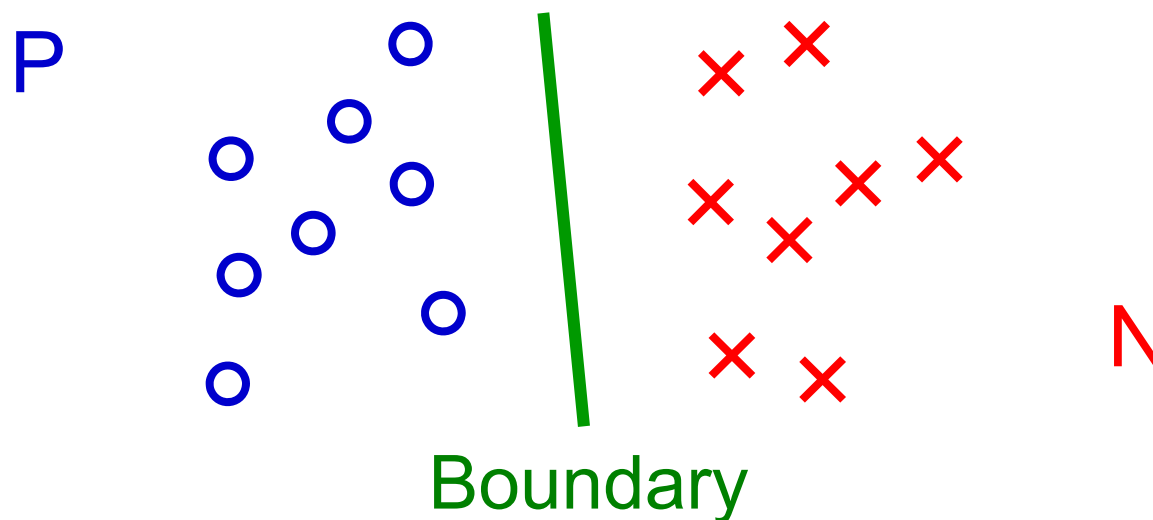
**Slides:**

<http://goo.gl/meiTwY>



# PN Classification (Ordinary Supervised Classification)

- **Labeled data:**  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}, y)$ 
  - Input  $\mathbf{x} \in \mathbb{R}^d$ :  $d$ -dimensional real vector
  - Output  $y \in \{+1, -1\}$ : Binary class label




# Some Definitions

18

■ **Classifier:**  $f : \mathbb{R}^d \rightarrow \mathbb{R}$

- Label prediction by  $\hat{y} = \text{sign}(f(\mathbf{x}))$   
(e.g., linear, additive, kernel, deep models).

■ **Margin:**  $m = yf(\mathbf{x})$   $y \in \{+1, -1\}$

- $m > 0 \implies \text{sign}(f(\mathbf{x})) = y$   
 **Classification is correct.**
- $m < 0 \implies \text{sign}(f(\mathbf{x})) \neq y$   
 **Classification is wrong.**

■ **Zero-one loss:**  $\ell_{0/1}(m) = \frac{1}{2} (1 - \text{sign}(m))$

- 1 for correct prediction.
- 0 for wrong prediction.

# Classification Error and Empirical Approximation

- **Classification error** (expected zero-one loss over all test data):

$\mathbb{E}$ : Expectation

$$R_{0/1}(f) = \mathbb{E}_{p(\mathbf{x}, y)} \left[ \ell_{0/1}(yf(\mathbf{x})) \right]$$

$$\ell_{0/1}(m) = \frac{1}{2} (1 - \text{sign}(m))$$

- **Our goal**: Find a minimizer of  $R_{0/1}(f)$ .
- But this is impossible since  $p(\mathbf{x}, y)$  is unknown:

- Let's use **samples**:  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}, y)$

i.i.d.: Independent and identically distributed

- **Empirical approximation**:

$$\hat{R}_{0/1}(f) = \frac{1}{n} \sum_{i=1}^n \ell_{0/1}(y_i f(\mathbf{x}_i)) = R_{0/1}(f) + O_p\left(\frac{1}{\sqrt{n}}\right)$$

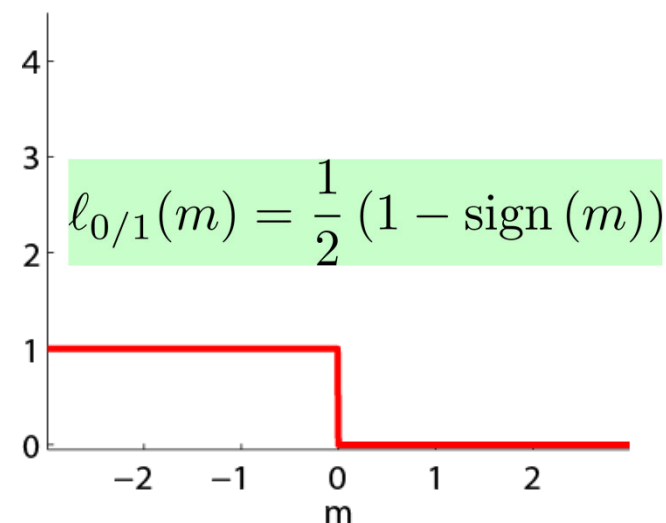
# Minimization of Empirical Classification Error

$$\hat{R}_{0/1}(f) = \frac{1}{n} \sum_{i=1}^n \ell_{0/1}(y_i f(\mathbf{x}_i))$$

- However, minimization of  $\hat{R}_{0/1}(f)$  is **NP-hard**, due to **discrete nature** of  $\ell_{0/1}$ :

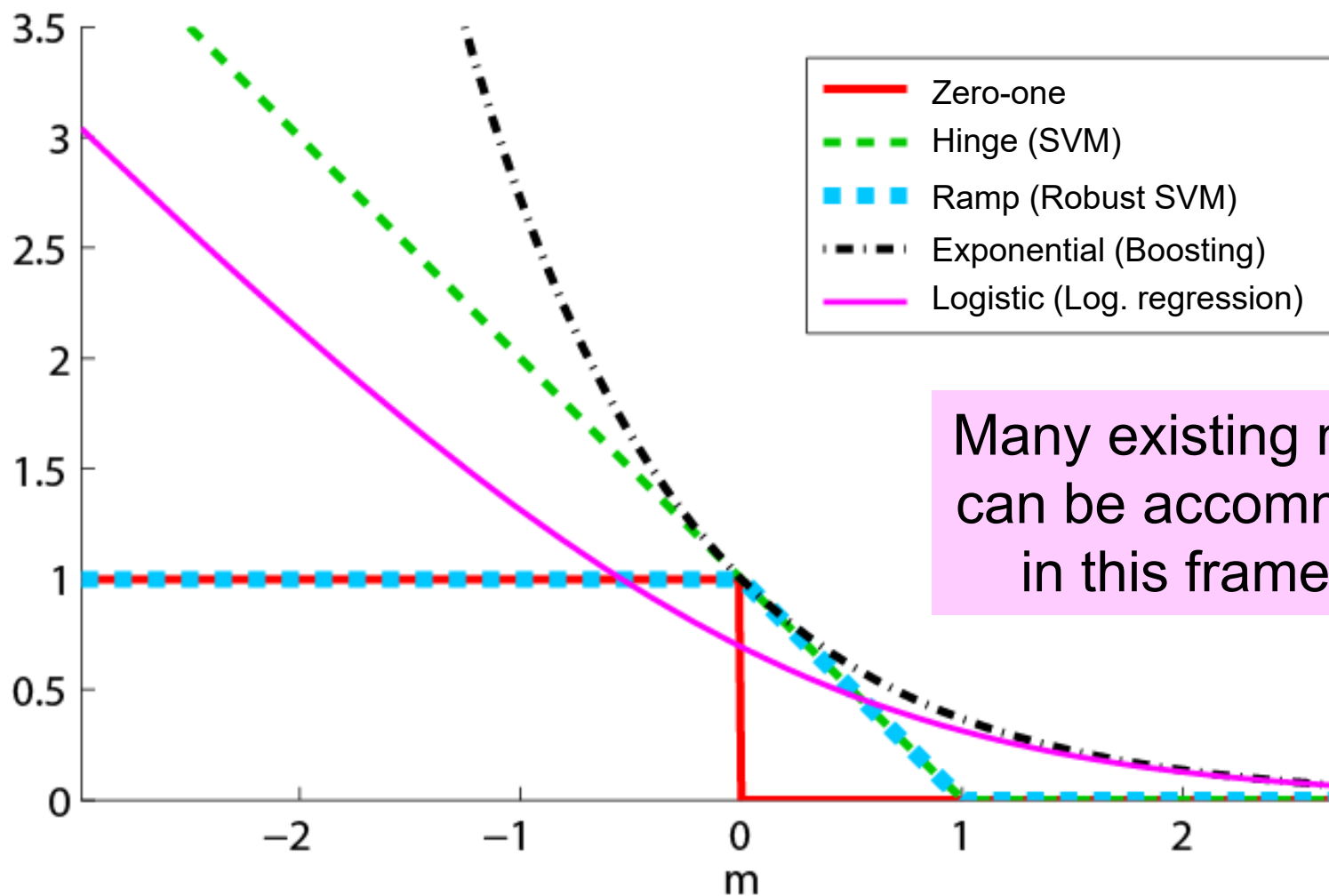
- We may not be able to obtain a global minimizer in practice.

- Let's use a **smoother loss**!



# Surrogate Loss

■ Let's use a **smoother loss** as a surrogate:



Many existing methods can be accommodated in this framework!

# PN Empirical Risk Minimization<sup>22</sup>

- **Classification risk** for loss  $\ell$ :

$$R(f) = \mathbb{E}_{p(\mathbf{x}, y)} \left[ \ell \left( y f(\mathbf{x}) \right) \right]$$

- **Empirical risk:**

- Expectation is approximated by sample average:

$$\hat{R}_{\text{PN}}(f) = \frac{1}{n} \sum_{i=1}^n \ell \left( y_i f(\mathbf{x}_i) \right) = R(f) + O_p \left( \frac{1}{\sqrt{n}} \right)$$

$$\{(\mathbf{x}_i, y_i)\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}, y)$$

- Minimize it within a certain model class (e.g., linear, additive, kernel, deep,...):

$$\hat{f}_{\text{PN}} = \operatorname{argmin}_f \hat{R}_{\text{PN}}(f)$$



# Contents

23

1. Background
2. PN Classification
3. **PU Classification**
4. PNU Classification
5. Pconf Classification
6. UU Classification
7. SU Classification
8. Comp Classification
9. Summary

- **P**: Positive
- **N**: Negative
- **U**: Unlabeled
- **Conf**: Confidence
- **S**: Similar
- **Comp**: Complementary

**Slides:**

<http://goo.gl/meiTwY>

# PU Classification: Setup

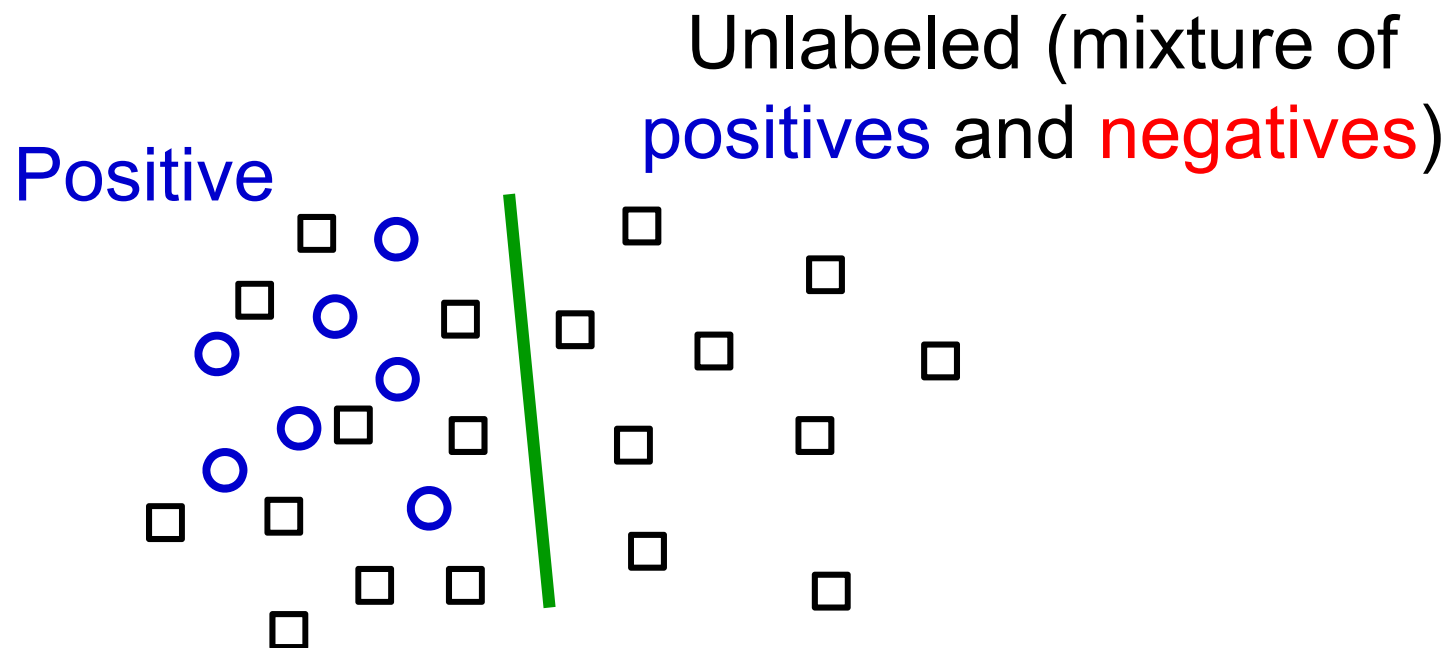
24

- **Given:** Positive and unlabeled samples

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

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

- **Goal:** Obtain a PN classifier





# PN Risk Decomposition

25

- Risk of classifier  $f$  :

$$\begin{aligned} R(f) &= \mathbb{E}_{p(\mathbf{x}, y)} \left[ \ell \left( y f(\mathbf{x}) \right) \right] \\ &= \underbrace{\pi \mathbb{E}_{p(\mathbf{x} | y = +1)} \left[ \ell \left( f(\mathbf{x}) \right) \right]}_{\text{Risk for P data}} + \underbrace{(1 - \pi) \mathbb{E}_{p(\mathbf{x} | y = -1)} \left[ \ell \left( - f(\mathbf{x}) \right) \right]}_{\text{Risk for N data}} \end{aligned}$$

$\pi = p(y = +1)$  : Class-prior probability  
(assumed known; **can be estimated**)

Scott & Blanchard (AISTATS2009)

Blanchard, Lee & Scott (JMLR2010)

du Plessis, Niu & Sugiyama (IEICE2014, MLJ2017)

Ramaswamy, Scott & Tewari (ICML2016)

- Since we do not have N data in the PU setting, the risk cannot be directly estimated.

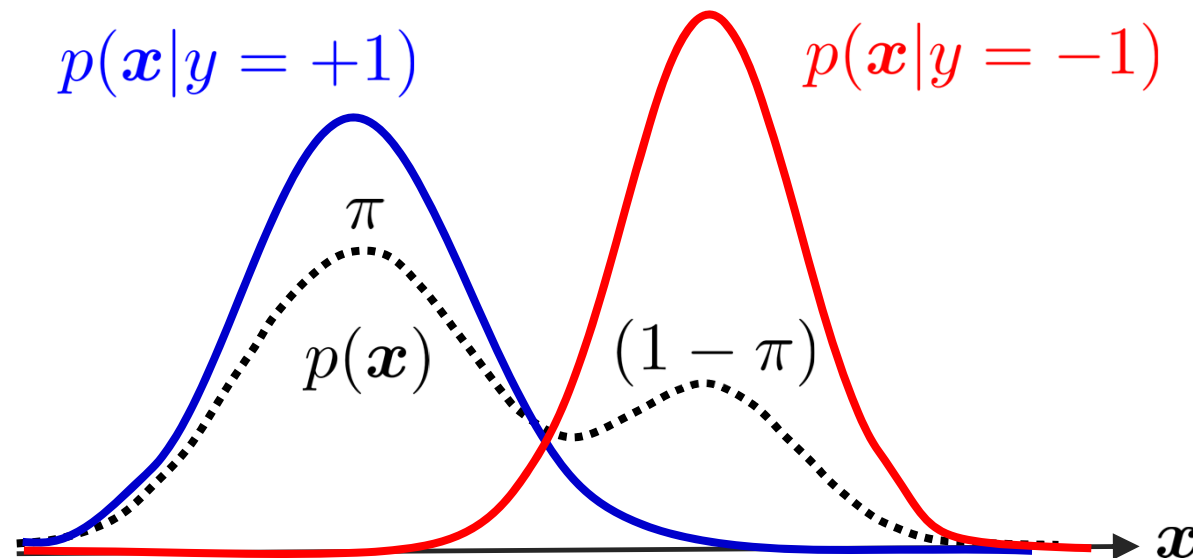
# PU Risk Estimation

26

$$R(f) = \pi \mathbb{E}_{p(\mathbf{x}|y=+1)} \left[ \ell(f(\mathbf{x})) \right] + (1 - \pi) \mathbb{E}_{p(\mathbf{x}|y=-1)} \left[ \ell(-f(\mathbf{x})) \right]$$

- U-density is a mixture of P- and N-densities:

$$p(\mathbf{x}) = \pi p(\mathbf{x}|y = +1) + (1 - \pi) p(\mathbf{x}|y = -1)$$



# PU Risk Estimation (cont.)

27

du Plessis, Niu & Sugiyama (ICML2015)

$$R(f) = \pi \mathbb{E}_{p(\mathbf{x}|y=+1)} \left[ \ell(f(\mathbf{x})) \right] + (1 - \pi) \mathbb{E}_{p(\mathbf{x}|y=-1)} \left[ \ell(-f(\mathbf{x})) \right]$$

$$p(\mathbf{x}) = \pi p(\mathbf{x}|y = +1) + (1 - \pi) p(\mathbf{x}|y = -1)$$

- This allow us to eliminate the N-density:

$$(1 - \pi) p(\mathbf{x}|y = -1) = p(\mathbf{x}) - \pi p(\mathbf{x}|y = +1)$$

$$R(f) = \pi \mathbb{E}_{p(\mathbf{x}|y=+1)} \left[ \ell(f(\mathbf{x})) \right] + \mathbb{E}_{p(\mathbf{x})} \left[ \ell(-f(\mathbf{x})) \right] - \pi \mathbb{E}_{p(\mathbf{x}|y=+1)} \left[ \ell(-f(\mathbf{x})) \right]$$

- Unbiased risk estimation is possible from PU data, just by replacing expectations by sample averages!

# PU Empirical Risk Minimization<sup>28</sup>

$$R(f) = \pi \mathbb{E}_{p(\mathbf{x}|y=+1)} [\ell(f(\mathbf{x}))] + \mathbb{E}_{p(\mathbf{x})} [\ell(-f(\mathbf{x}))] - \pi \mathbb{E}_{p(\mathbf{x}|y=+1)} [\ell(-f(\mathbf{x}))]$$

- Replacing expectations by sample averages gives an empirical risk:

$$\hat{R}_{\text{PU}}(f) = \frac{\pi}{n_{\text{P}}} \sum_{i=1}^{n_{\text{P}}} \ell(f(\mathbf{x}_i^{\text{P}})) + \frac{1}{n_{\text{U}}} \sum_{i=1}^{n_{\text{U}}} \ell(-f(\mathbf{x}_i^{\text{U}})) - \frac{\pi}{n_{\text{P}}} \sum_{i=1}^{n_{\text{P}}} \ell(-f(\mathbf{x}_i^{\text{P}}))$$

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

- Optimal convergence rate is attained:

Niu, du Plessis, Sakai, Ma & Sugiyama (NIPS2016)

$$R(\hat{f}_{\text{PU}}) - R(f^*) \leq C(\delta) \left( \frac{2\pi}{\sqrt{n_{\text{P}}}} + \frac{1}{\sqrt{n_{\text{U}}}} \right)$$

with probability  $1 - \delta$

$$\hat{f}_{\text{PU}} = \operatorname{argmin}_f \hat{R}_{\text{PU}}(f)$$

$$f^* = \operatorname{argmin}_f R(f)$$

$n_{\text{P}}, n_{\text{U}}$  : # of P, U samples

# Theoretical Comparison with PN<sup>29</sup>

Niu, du Plessis, Sakai, Ma & Sugiyama (NIPS2016)

## ■ Estimation error bounds for PU and PN:

$$R(\hat{f}_{\text{PU}}) - R(f^*) \leq C(\delta) \left( \frac{2\pi}{\sqrt{n_{\text{P}}}} + \frac{1}{\sqrt{n_{\text{U}}}} \right)$$

$$R(\hat{f}_{\text{PN}}) - R(f^*) \leq C(\delta) \left( \frac{\pi}{\sqrt{n_{\text{P}}}} + \frac{1 - \pi}{\sqrt{n_{\text{N}}}} \right)$$

$$\hat{f}_{\text{PN}} = \operatorname{argmin}_f \hat{R}_{\text{PN}}(f)$$

with probability  $1 - \delta$

$$\hat{R}_{\text{PN}}(f) = \frac{1}{n} \sum_{i=1}^n \ell(y_i f(\mathbf{x}_i))$$

$n_{\text{P}}, n_{\text{N}}, n_{\text{U}}$  : # of P, N, U samples

## ■ Comparison: PU bound is smaller than PN if

$$\frac{\pi}{\sqrt{n_{\text{P}}}} + \frac{1}{\sqrt{n_{\text{U}}}} < \frac{1 - \pi}{\sqrt{n_{\text{N}}}}$$

- PU can be better than PN, provided many PU data!

# Further Correction

30

$$R(f) = \underbrace{\pi \mathbb{E}_{p(\mathbf{x}|y=+1)} \left[ \ell \left( f(\mathbf{x}) \right) \right]}_{\text{Risk for P data}} + \underbrace{(1 - \pi) \mathbb{E}_{p(\mathbf{x}|y=-1)} \left[ \ell \left( -f(\mathbf{x}) \right) \right]}_{\text{Risk for N data } R^-(f)}$$

■ PU formulation:  $p(\mathbf{x}) = \pi p(\mathbf{x}|y = +1) + (1 - \pi)p(\mathbf{x}|y = -1)$

$$R^-(f) = \mathbb{E}_{p(\mathbf{x})} \left[ \ell \left( -f(\mathbf{x}) \right) \right] - \pi \mathbb{E}_{p(\mathbf{x}|y=+1)} \left[ \ell \left( -f(\mathbf{x}) \right) \right]$$

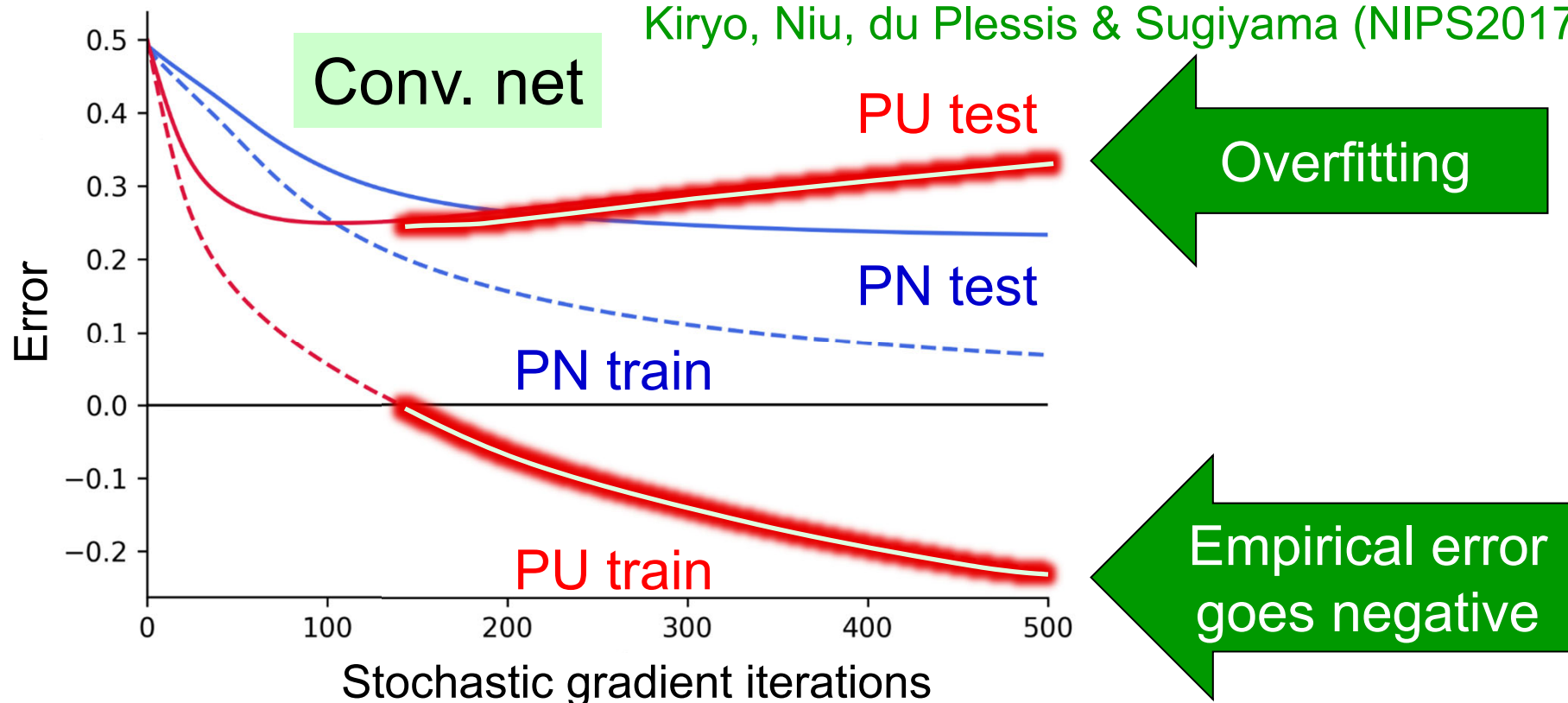
- If  $\ell(m) \geq 0, \forall m$ ,  $R^-(f) \geq 0$ .
- However, its PU empirical approximation can be **negative** due to “difference of approximations”.

$$\hat{R}_{\text{PU}}^-(f) = \frac{1}{n_U} \sum_{i=1}^{n_U} \ell \left( -f(\mathbf{x}_i^U) \right) - \frac{\pi}{n_P} \sum_{i=1}^{n_P} \ell \left( -f(\mathbf{x}_i^P) \right) \not\geq 0$$

- This problem is more critical for flexible models such as **deep nets**.

# Non-Negative PU Classification<sup>31</sup>

Kiryu, Niu, du Plessis & Sugiyama (NIPS2017)



- We constrain the sample approximation term **to be non-negative** through back-prop training:

$$\tilde{R}_{\text{PU}}(f) = \frac{\pi}{n_{\text{P}}} \sum_{i=1}^{n_{\text{P}}} \ell(f(\mathbf{x}_i^{\text{P}})) + \max \left\{ 0, \frac{1}{n_{\text{U}}} \sum_{i=1}^{n_{\text{U}}} \ell(-f(\mathbf{x}_i^{\text{U}})) - \frac{\pi}{n_{\text{P}}} \sum_{i=1}^{n_{\text{P}}} \ell(-f(\mathbf{x}_i^{\text{P}})) \right\}$$

- Now the risk estimator is biased. Is it really good?

# Theoretical Analysis

32

Kiryu, Niu, du Plessis & Sugiyama (NIPS2017)

$$\tilde{R}_{\text{PU}}(f) = \frac{\pi}{n_{\text{P}}} \sum_{i=1}^{n_{\text{P}}} \ell(f(\mathbf{x}_i^{\text{P}})) + \max \left\{ 0, \frac{1}{n_{\text{U}}} \sum_{i=1}^{n_{\text{U}}} \ell(-f(\mathbf{x}_i^{\text{U}})) - \frac{\pi}{n_{\text{P}}} \sum_{i=1}^{n_{\text{P}}} \ell(-f(\mathbf{x}_i^{\text{P}})) \right\}$$

- $\tilde{R}_{\text{PU}}(f)$  is still **consistent** and its bias decreases **exponentially**:  $\mathcal{O}(e^{-n_{\text{P}}-n_{\text{U}}})$   $n_{\text{P}}, n_{\text{U}}$ : # of P, U samples
  - In practice, we can ignore the bias of  $\tilde{R}_{\text{PU}}(f)$  !
- Mean-squared error of  $\tilde{R}_{\text{PU}}(f)$  is not more than the original one.
  - In practice,  $\tilde{R}_{\text{PU}}(f)$  is more reliable!
- Risk of  $\operatorname{argmin}_f \tilde{R}_{\text{PU}}(f)$  for linear models attains **optimal convergence rate**:  $\mathcal{O}_p \left( \frac{1}{\sqrt{n_{\text{P}}}} + \frac{1}{\sqrt{n_{\text{U}}}} \right)$ 
  - Learned function is optimal.



# Experiments

33

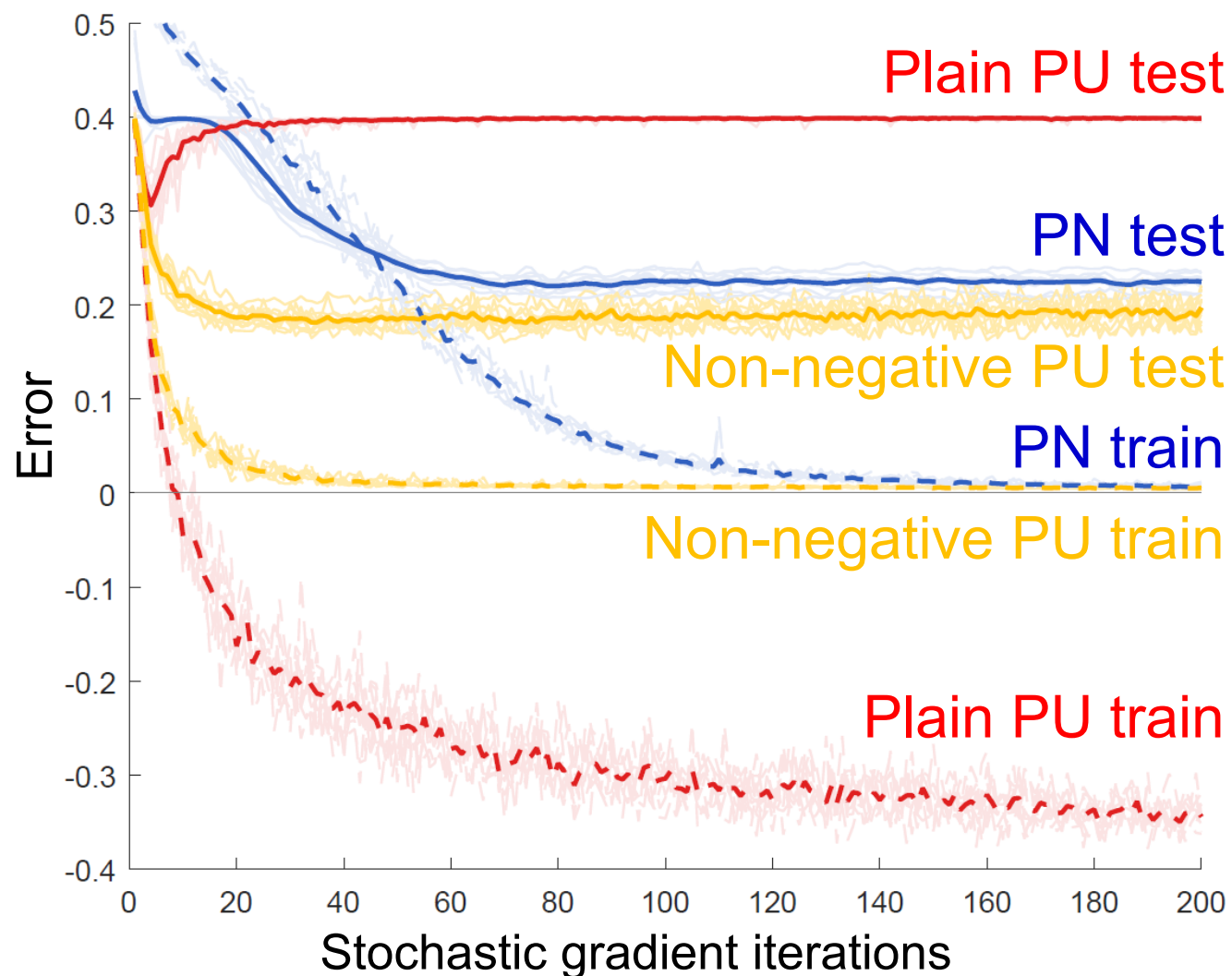
- With a large number of unlabeled data, non-negative PU can even outperform PN!

- Binary CIFAR-10:  
Positive (airplane, automobile, ship, truck)  
Negative (bird, cat, deer, dog, frog, horse)
- 13-layer CNN with ReLU

$$n_P = 1000$$

$$n_U = 50000$$

$$\pi = 0.4$$



# PU Classification: Summary

34

■ Just separating P and U is biased.

■ To be unbiased, use **composite loss**

$\tilde{\ell}(m) = \ell(m) - \ell(-m)$  for P data.

Natarajan, Dhillon, Ravikumar & Tewari (NIPS2013)

● **Optimal convergence rate achieved.**

Niu, du Plessis, Sakai, Ma & Sugiyama (NIPS2016)

■ If  $\ell(m) + \ell(-m) = \text{Const.}$ ,  
the **same loss** for P and U data.

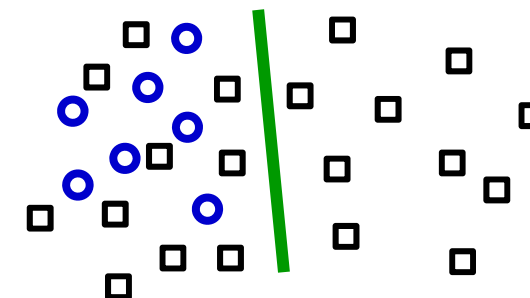
du Plessis, Niu & Sugiyama (NIPS2014)

■ If  $\tilde{\ell}(m) = am + b$ ,  
optimization becomes **convex**.

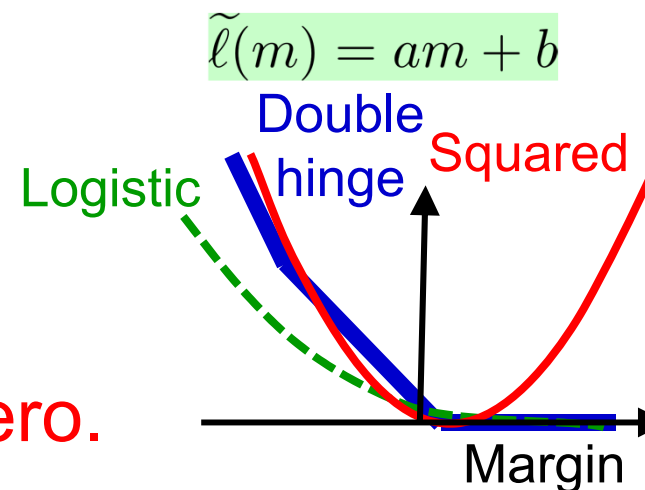
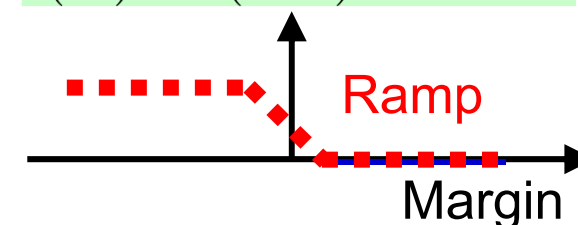
du Plessis, Niu & Sugiyama (ICML2015)

■ For deep nets, **roundup the empirical false negative error to zero.**

Kiryu, Niu, du Plessis & Sugiyama (NIPS2017)



$$\ell(m) + \ell(-m) = \text{Const.}$$





# Contents

35

1. Background
2. PN Classification
3. PU Classification
4. **PNU Classification**
5. Pconf Classification
6. UU Classification
7. SU Classification
8. Comp Classification
9. Summary

- **P**: Positive
- **N**: Negative
- **U**: Unlabeled
- **Conf**: Confidence
- **S**: Similar
- **Comp**: Complementary

**Slides:**

<http://goo.gl/meiTwY>

# PNU Classification: Setup

36

- **Given:** Positive, negative & unlabeled samples

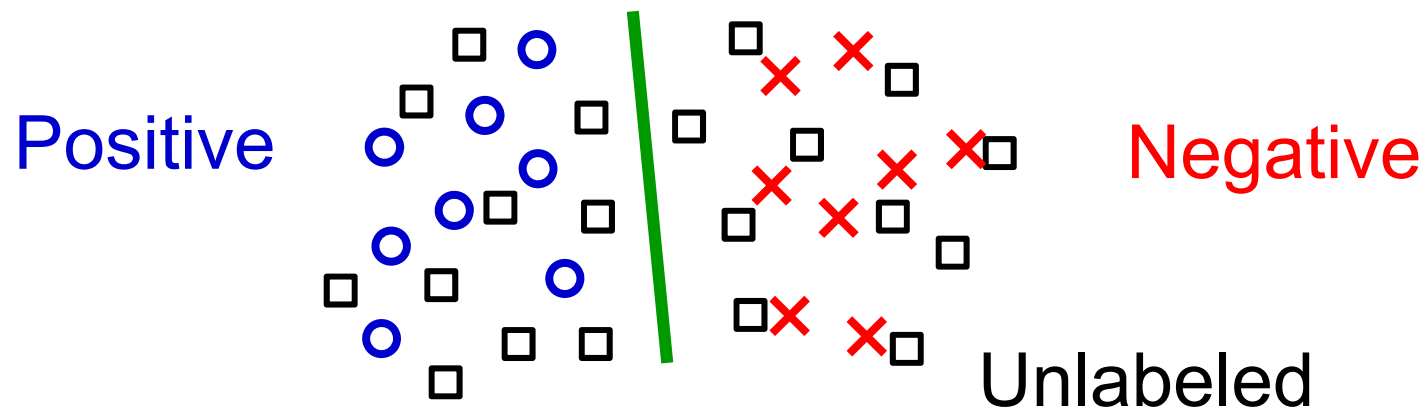
$$\{\mathbf{x}_i^P\}_{i=1}^{n_P} \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}|y = +1)$$

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

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

- **Goal:** Obtain a PN classifier

- **PNU classification** is semi-supervised learning.

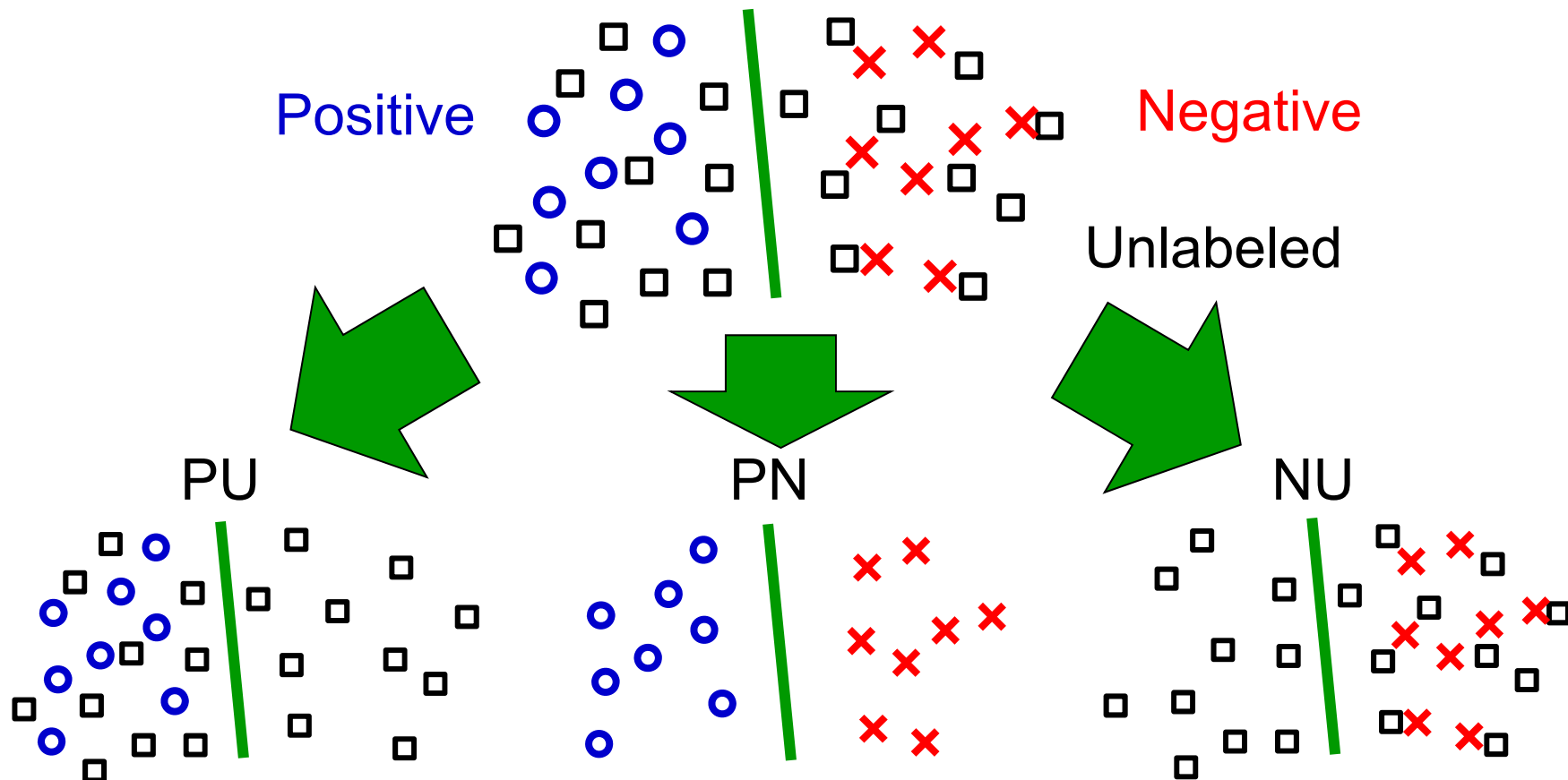


# PNU Decomposition

37

Sakai, du Plessis, Niu & Sugiyama (ICML2017)

- Let's decompose PNU into **PU**, **PN**, and **NU**:
  - Each can be solved easily.
  - Combine them!

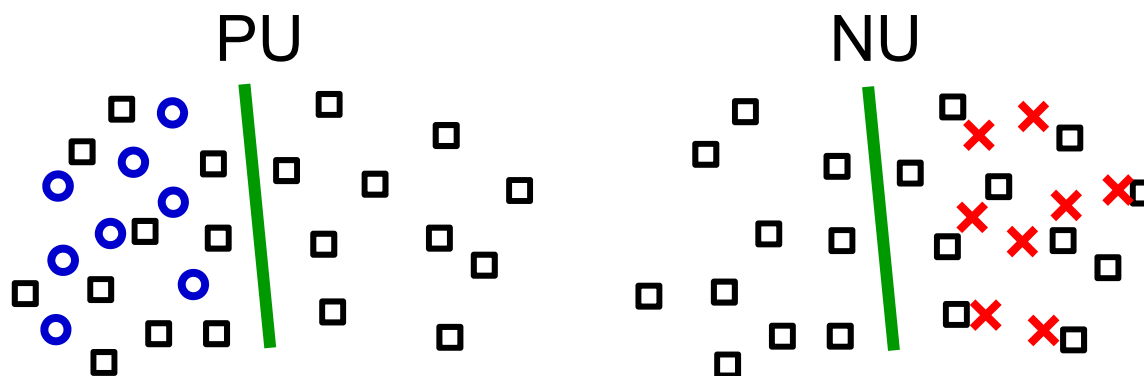


# How to Combine?

38

- **Natural choice:** Combine PU & NU (symmetric).

$$R_{\text{PU+NU}}(f) = (1 - \gamma)R_{\text{PU}}(f) + \gamma R_{\text{NU}}(f) \quad 0 \leq \gamma \leq 1$$



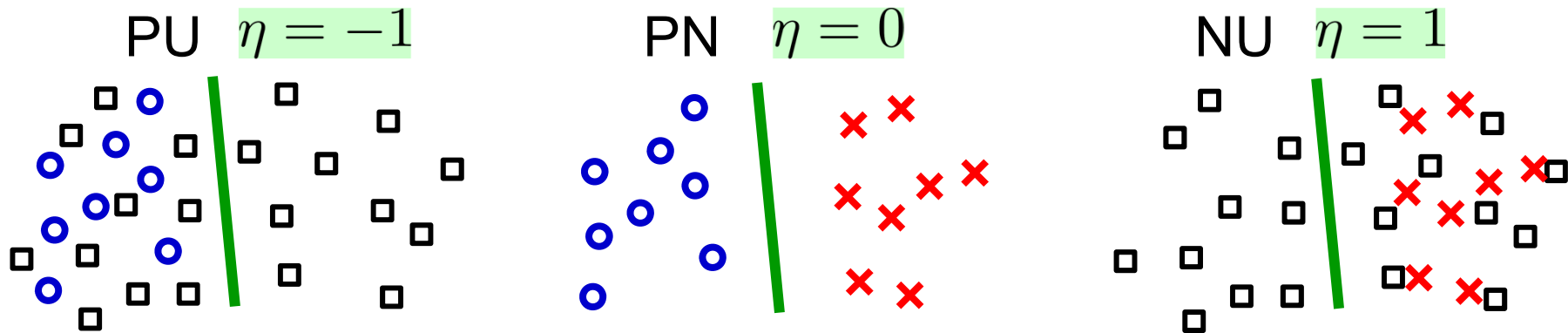
- **Theoretical risk analysis:**

Niu, du Plessis, Sakai, Ma  
& Sugiyama (NIPS2016)

- When  $\text{PU} < \text{NU}$ :  $\text{PU} < \text{PN} < \text{NU}$  or  $\text{PN} < \text{PU} < \text{NU}$ .
- When  $\text{NU} < \text{PU}$ :  $\text{NU} < \text{PN} < \text{PU}$  or  $\text{PN} < \text{NU} < \text{PU}$ .
- PU+NU is not the best possible combination.
- PU+PN & NU+PN are the best combinations.

# PN+PU & PN+NU Classification<sup>39</sup>

- **Proposed method:** Combine two best methods.



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

$$-1 \leq \eta \leq 1$$

- **PN+PU classification:**

$$R_{\text{PN+PU}}^{\gamma}(f) = (1 - \gamma)R_{\text{PN}}(f) + \gamma R_{\text{PU}}(f) \quad 0 \leq \gamma \leq 1$$

- **PN+NU classification:**

$$R_{\text{PN+NU}}^{\gamma}(f) = (1 - \gamma)R_{\text{PN}}(f) + \gamma R_{\text{NU}}(f) \quad 0 \leq \gamma \leq 1$$

# Theoretical Analysis

40

## Generalization error bound:

$$R_{0/1}(f) \leq 2\hat{R}_{\text{PNU}}^\eta(f) + \mathcal{O}\left(\frac{1}{\sqrt{n_{\text{P}}}} + \frac{1}{\sqrt{n_{\text{N}}}} + \frac{1}{\sqrt{n_{\text{U}}}}\right)$$

$$R_{0/1}(f) = \mathbb{E}_{p(\mathbf{x}, y)} \left[ \ell_{0/1}(yf(\mathbf{x})) \right]$$

$n_{\text{P}}, n_{\text{N}}, n_{\text{U}}$ : # of P, N, U samples

$\hat{R}_{\text{PNU}}^\eta(f)$ : Empirical version of  $R_{\text{PNU}}^\eta(f) = \begin{cases} R_{\text{PN+PU}}^\eta(f) & (\eta \geq 0) \\ R_{\text{PN+NU}}^{-\eta}(f) & (\eta < 0) \end{cases}$

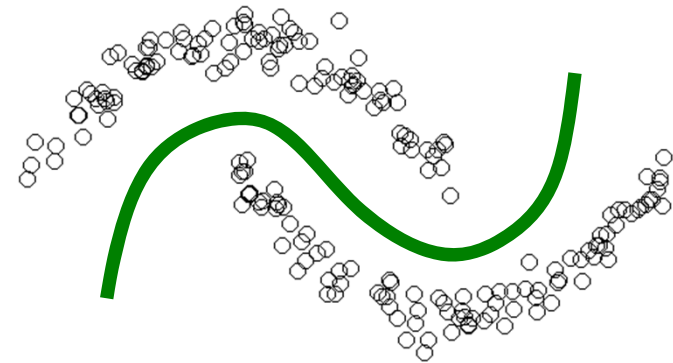
- Unlabeled data always helps without cluster assumptions!

## We use unlabeled data

for **loss evaluation**,

not for **regularization** (as manifold smoothing).

- Label information is extracted from unlabeled data!





# Experiments

Misclassification error rate: average (std)

5% t-test

(Grandvalet & Bengio, (Belkin et al., (Niu et al., (Li et al.,  
NIPS2004) JMLR2006) ICML2013) JMLR2013)

Dataset	$n_u$	$\pi$	$\hat{\pi}$	Proposed	EntReg	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)	27.5 (0.5)
	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>	<b>12.5 (0.5)</b>	<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>	24.9 (0.7)	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
Temples	1000	0.55	0.51 (0.01)	<b>43.9 (0.7)</b>	<b>43.9 (0.6)</b>	<b>43.4 (0.6)</b>	50.7 (1.6)	44.3 (0.5)
	5000	0.55	0.54 (0.01)	43.4 (0.9)	<b>43.0 (0.6)</b>	<b>43.1 (1.0)</b>	43.6 (0.7)	N/A
	10000	0.55	0.50 (0.01)	<b>45.2 (0.8)</b>	<b>44.4 (0.8)</b>	<b>44.2 (0.7)</b>	N/A	N/A

■ Proposed PN+PU & PN+NU works well!



# Contents

42

1. Background
2. PN Classification
3. PU Classification
4. PNU Classification
5. **Pconf Classification**
6. UU Classification
7. SU Classification
8. Comp Classification
9. Summary

- **P**: Positive
- **N**: Negative
- **U**: Unlabeled
- **Conf**: Confidence
- **S**: Similar
- **Comp**: Complementary

**Slides:**

<http://goo.gl/meiTwY>

# Pconf Classification: Setup

43

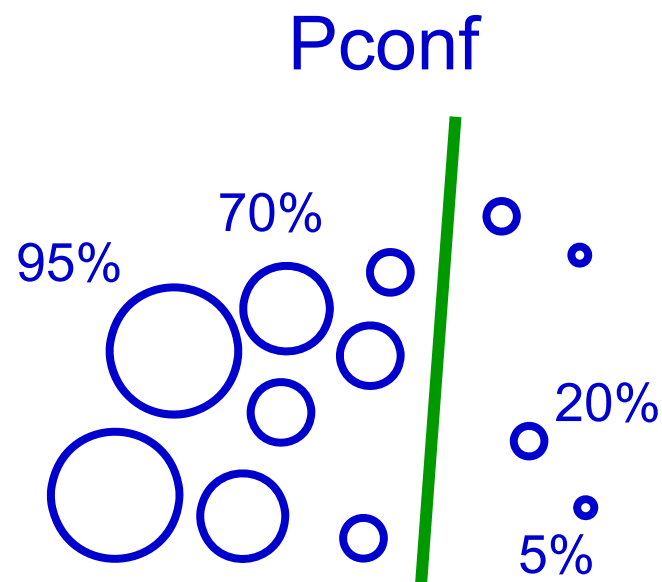
Ishida, Niu & Sugiyama (NeurIPS2018)

## ■ Given: Positive-confidence samples

$$\{(\mathbf{x}_i, r_i)\}_{i=1}^n$$

- Positive patterns:  $\{\mathbf{x}_i\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}|y = +1)$
- Their confidence:  $r_i = P(y = +1|\mathbf{x}_i)$

## ■ Goal: Obtain a PN classifier



# Pconf Risk Estimation

44

■ Classification risk:  $R(f) = \mathbb{E}_{p(\mathbf{x}, y)} \left[ \ell(y f(\mathbf{x})) \right]$

■ Naïve “confidence-weighting” is not correct.

$$R(f) \neq \mathbb{E}_{p(\mathbf{x}|y=+1)} \left[ r(\mathbf{x}) \ell(f(\mathbf{x})) + (1 - r(\mathbf{x})) \ell(-f(\mathbf{x})) \right]$$

$$r(\mathbf{x}) = P(y = +1 | \mathbf{x})$$

■ Right form is given by **importance sampling**:

$$R(f) = \pi \mathbb{E}_{p(\mathbf{x}|y=+1)} \left[ \ell(f(\mathbf{x})) + \frac{1 - r(\mathbf{x})}{r(\mathbf{x})} \ell(-f(\mathbf{x})) \right]$$

resulting in an empirical risk:

$$\hat{R}_{\text{Pconf}}(f) \propto \sum_{i=1}^n \left[ \ell(f(\mathbf{x}_i)) + \frac{1 - r_i}{r_i} \ell(-f(\mathbf{x}_i)) \right]$$

$$\{\mathbf{x}_i\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}|y = +1) \quad r_i = P(y = +1 | \mathbf{x}_i)$$

- Estimation error:

$$R(f^*) - R(\hat{f}_{\text{Pconf}}) = \mathcal{O}_p\left(\frac{1}{\sqrt{n}}\right)$$

$$f^* = \underset{f}{\operatorname{argmin}} R(f) \quad \hat{f}_{\text{Pconf}} = \underset{f}{\operatorname{argmin}} \hat{R}_{\text{Pconf}}(f)$$

$$R(f) = \mathbb{E}_{p(\mathbf{x}, y)} \left[ \ell(y f(\mathbf{x})) \right]$$

$$\hat{R}_{\text{Pconf}}(f) \propto \sum_{i=1}^n \left[ \ell(f(\mathbf{x}_i)) + \frac{1 - r_i}{r_i} \ell(-f(\mathbf{x}_i)) \right]$$

- Optimal parametric convergence rate is attained!

# Experiments

46

Correct classification rate: average (std)

5% t-test

Positive	vs.	Negative	Pconf	Weighted	Supervised
airplane	vs.	automobile	<b>84.34 ± 0.84</b>	79.32 ± 2.74	93.82 ± 0.21
airplane	vs.	bird	<b>82.50 ± 3.19</b>	<b>81.38 ± 0.48</b>	89.24 ± 0.50
airplane	vs.	cat	<b>89.10 ± 0.47</b>	86.98 ± 1.20	92.78 ± 0.49
airplane	vs.	deer	<b>87.44 ± 1.43</b>	82.00 ± 2.39	92.08 ± 0.50
airplane	vs.	dog	<b>90.24 ± 1.27</b>	86.86 ± 1.41	94.42 ± 0.89
airplane	vs.	frog	<b>91.44 ± 0.86</b>	85.12 ± 1.66	95.52 ± 0.42
airplane	vs.	horse	<b>89.26 ± 2.20</b>	<b>87.72 ± 1.99</b>	95.58 ± 0.56
airplane	vs.	ship	<b>74.36 ± 2.00</b>	70.82 ± 1.80	89.04 ± 1.06
airplane	vs.	truck	<b>84.98 ± 0.47</b>	83.22 ± 0.58	91.84 ± 1.19

■ Works better than naïve “weighted” baseline!



# Contents

1. Background
2. PN Classification
3. PU Classification
4. PNU Classification
5. Pconf Classification
6. **UU Classification**
7. SU Classification
8. Comp Classification
9. Summary

- **P**: Positive
- **N**: Negative
- **U**: Unlabeled
- **Conf**: Confidence
- **S**: Similar
- **Comp**: Complementary

**Slides:**

<http://goo.gl/meiTwY>

# UU Classification: Setup

48

du Plessis, Niu & Sugiyama (TAAI2013)  
Nan, Niu, Menon & Sugiyama (ICLR2019)

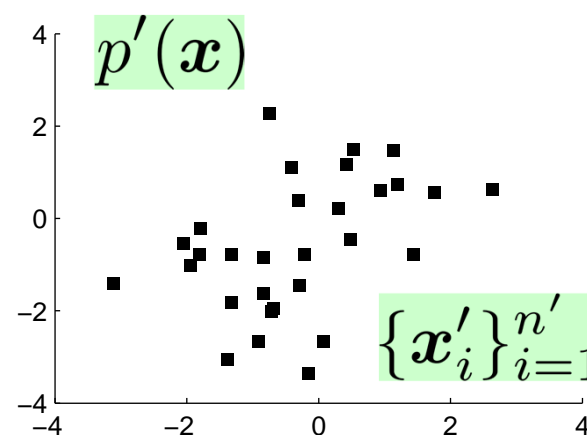
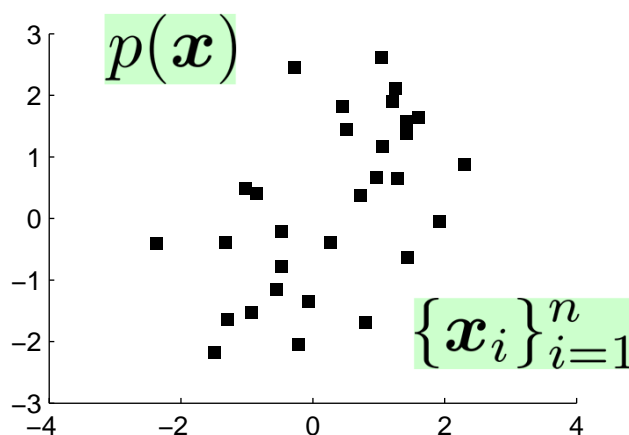
- **Given:** Two sets of unlabeled data

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

- **Assumption:** Only class-priors are different

$$p(y) \neq p'(y) \quad p(\mathbf{x}|y) = p'(\mathbf{x}|y)$$

- **Goal:** Obtain a PN classifier





# Optimal UU Classifier

49

du Plessis, Niu & Sugiyama (TAAI2013)

- Sign of the difference of class-posteriors:

$$g(\mathbf{x}) = \text{sign}[p(y = +1|\mathbf{x}) - p(y = -1|\mathbf{x})]$$

- Under **uniform** test class-prior,

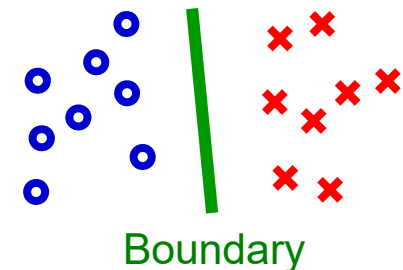
$$g(\mathbf{x}) = C \text{sign}[p(\mathbf{x}) - p'(\mathbf{x})]$$

$$C = \text{sign}[p(y = +1) - p'(y = +1)]$$

- Sign of  $C$  is unknown, but just knowing

$$\text{sign}[p(\mathbf{x}) - p'(\mathbf{x})]$$

still allows **optimal separation!**



# UU Risk Estimation

50

Nan, Niu, Menon & Sugiyama (ICLR2019)

## ■ For

- uniform test class-prior:  $\pi = 1/2$
- symmetric loss:  $\ell(m) + \ell(-m) = \text{Const.}$

the classification risk can be expressed as

$$R(f) = \mathbb{E}_{p(\mathbf{x}, y)} \left[ \ell(y f(\mathbf{x})) \right] \\ \propto \mathbb{E}_{p(\mathbf{x})} \left[ \ell(f(\mathbf{x})) \right] + \mathbb{E}_{p'(\mathbf{x}')} \left[ \ell(-f(\mathbf{x}')) \right] + \text{Const.}$$

resulting an empirical risk (up to label flip):

$$\hat{R}_{\text{UU}}(f) \propto \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i)) + \frac{1}{n'} \sum_{i=1}^{n'} \ell(-f(\mathbf{x}'_i))$$

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

- Estimation error:

$$R(f^*) - R(\hat{f}_{UU}) = \mathcal{O}_p \left( \frac{1}{\sqrt{n}} + \frac{1}{\sqrt{n'}} \right)$$

$$f^* = \operatorname{argmin}_f R(f) \quad \hat{f}_{UU} = \operatorname{argmin}_f \hat{R}_{UU}(f)$$

$$R(f) = \mathbb{E}_{p(\mathbf{x}, y)} \left[ \ell \left( y f(\mathbf{x}) \right) \right]$$

$$\hat{R}_{UU}(f) \propto \frac{1}{n} \sum_{i=1}^n \ell \left( f(\mathbf{x}_i) \right) + \frac{1}{n'} \sum_{i=1}^{n'} \ell \left( -f(\mathbf{x}'_i) \right)$$

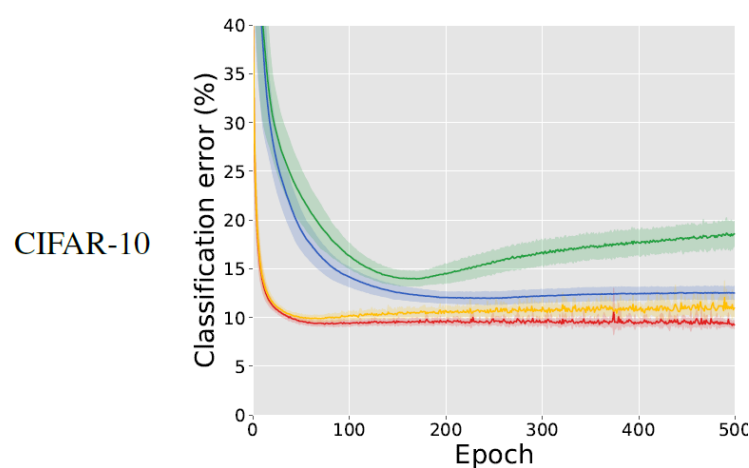
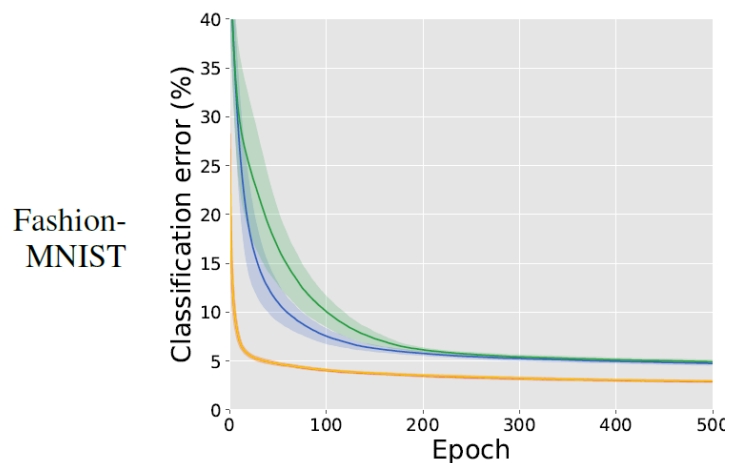
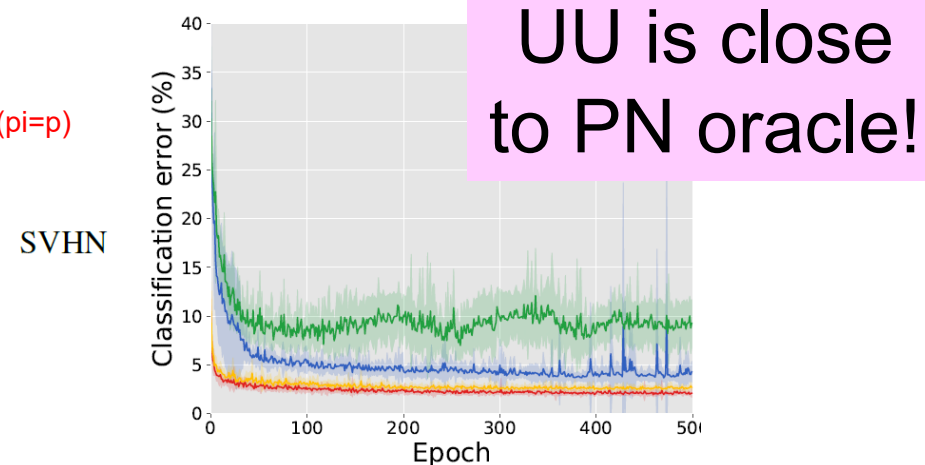
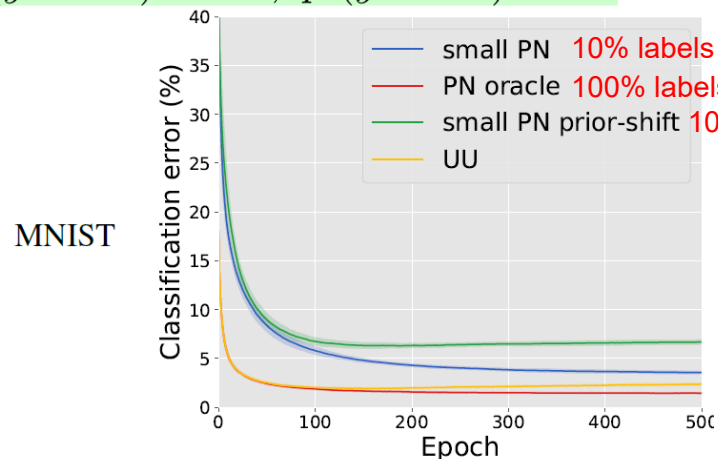
- Optimal parametric convergence rate is attained!

# Experiments

52

Dataset	# Train	# Test	# Feature	$\pi_p$	Model $g(x; \theta)$	Optimizer
MNIST	60,000	10,000	784	0.49	FC with ReLU (depth 5)	SGD
Fashion-MNIST	60,000	10,000	784	0.50	FC with ReLU (depth 5)	SGD
SVHN	100,000	26,032	3,072	0.27	AllConvNet (depth 12)	Adam
CIFAR-10	50,000	10,000	3,072	0.60	ResNet (depth 32)	Adam

$p(y = +1) = 0.9, p'(y = +1) = 0.1$





# Contents

53

1. Background
2. PN Classification
3. PU Classification
4. PNU Classification
5. Pconf Classification
6. UU Classification
7. **SU Classification**
8. Comp Classification
9. Summary

- **P**: Positive
- **N**: Negative
- **U**: Unlabeled
- **Conf**: Confidence
- **S**: Similar
- **Comp**: Complementary

**Slides:**

<http://goo.gl/meiTwY>

# SU Classification

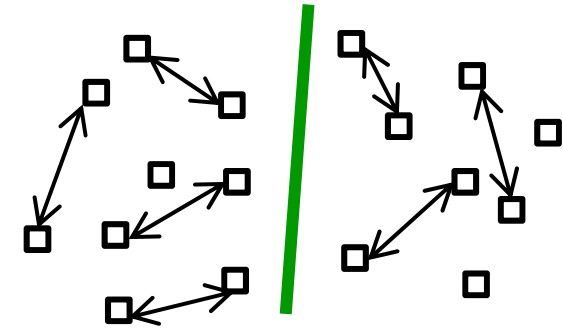
54

Bao, Niu & Sugiyama (ICML2018)

- **Given:** Similar and unlabeled samples

$$\{(\mathbf{x}_i, \mathbf{x}'_i)\}_{i=1}^{n_S} \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}, \mathbf{x}' | y = y')$$

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



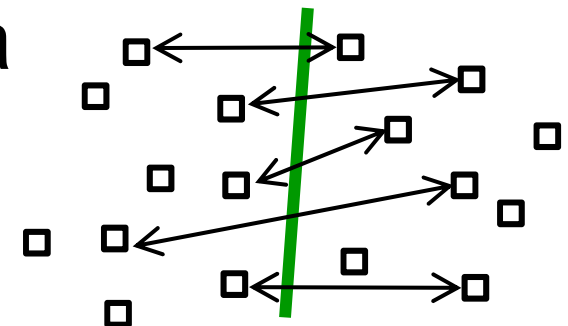
- **Goal:** Obtain a PN classifier

- This is a special case of UU classification:

$$p(y = +1) = \pi^2 / (2\pi^2 - 2\pi + 1) \quad p'(y = +1) = \pi$$

- Classification from **dissimilar** data is also possible (DU, SD, SDU)!

Shimada, Bao, Sato & Sugiyama (arXiv2019)





# Contents

55

1. Background
2. PN Classification
3. PU Classification
4. PNU Classification
5. Pconf Classification
6. UU Classification
7. SU Classification
8. **Comp Classification**
9. Summary

- **P**: Positive
- **N**: Negative
- **U**: Unlabeled
- **Conf**: Confidence
- **S**: Similar
- **Comp**: Complementary

**Slides:**

<http://goo.gl/meiTwY>

# Multiclass Labeling is Costly

56

## ■ Labeling in multi-class classification:

- What is the robot in this image?



<https://www.bostondynamics.com/atla>

1. Amazon Kiva
2. Aldebaran Nao
3. Softbank Pepper
4. Sony Aibo
5. iRobot Roomba
- ⋮
83. Boston Dynamics Atlas
- ⋮
100. Rethink Robotics Baxter

## ■ Selecting the correct class from a long list of candidates is extremely time-consuming!



# Complementary Classification <sup>57</sup>

Ishida, Niu, Hu & Sugiyama (NIPS2017)  
Ishida, Niu, Menon & Sugiyama (ICML2019)

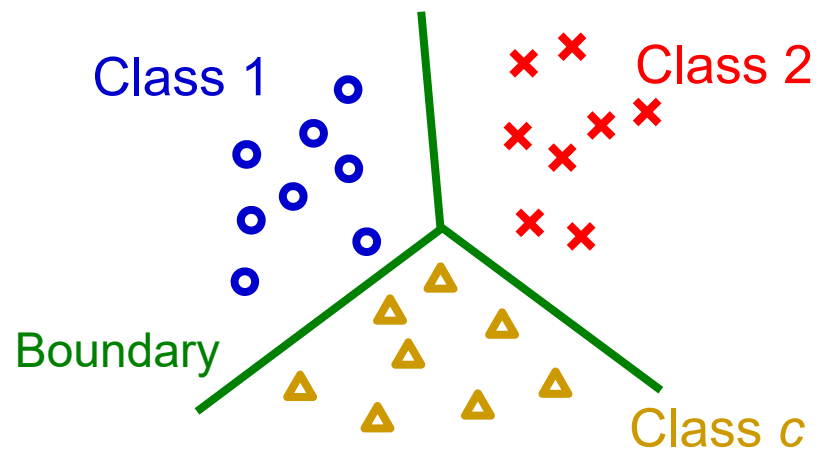
## ■ Given: Complementary labeled data

$$\{(\mathbf{x}_i, \bar{y}_i)\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} \bar{p}(\mathbf{x}, \bar{y})$$

$$\bar{p}(\mathbf{x}, \bar{y}) = \frac{1}{c-1} \sum_{y \neq \bar{y}} p(\mathbf{x}, y)$$

- Pattern  $\mathbf{x}$  does **not** belong to class  $\bar{y} \in \{1, 2, \dots, c\}$ .

## ■ Goal: Obtain a multiclass classifier



# Possible Approaches

58

$$\{(\mathbf{x}_i, \bar{y}_i)\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} \bar{p}(\mathbf{x}, \bar{y})$$

## ■ Approach 1: Classification from partial labels

Cour, Sapp & Taskar (JMLR2011)

- Multiple candidate classes are provided for each  $\mathbf{x}_i$ .
- Complementary labels are the extreme case of partial labels given to all  $c - 1$  classes other than  $\bar{y}_i$ .

## ■ Approach 2: Multi-label classification

- Each  $\mathbf{x}_i$  can belong to multiple classes.
- Negative label for  $\bar{y}_i$  and positives for the rest.

## ■ We want a more direct approach!

# Multi-Class Classification

59

■  $c$ -class classifier:  $f(\mathbf{x}) = \operatorname{argmax}_{y \in \{1, \dots, c\}} g_y(\mathbf{x})$

$g_y(\mathbf{x})$ : one-vs-rest classifier for  $y$

■  $c$ -class loss:  $L(y, \mathbf{g}(\mathbf{x}))$       $\mathbf{g}(\mathbf{x}) = (g_1(\mathbf{x}), \dots, g_c(\mathbf{x}))^\top$

● One-versus-rest:

$$L_{\text{OVR}}(y, \mathbf{g}(\mathbf{x})) = \ell(g_y(\mathbf{x})) + \frac{1}{c-1} \sum_{y' \neq y} \ell(-g_{y'}(\mathbf{x}))$$

● Pairwise comparison:

$$L_{\text{PC}}(y, \mathbf{g}(\mathbf{x})) = \sum_{y' \neq y} \ell(g_y(\mathbf{x}) - g_{y'}(\mathbf{x}))$$

■  $c$ -class classification risk:

$$R(\mathbf{g}) = \mathbb{E}_{p(\mathbf{x}, y)} \left[ L(y, \mathbf{g}(\mathbf{x})) \right]$$

# Complementary Risk Estimation<sup>60</sup>

Ishida, Niu, Menon & Sugiyama (ICML2019)

$$R(\mathbf{g}) = \mathbb{E}_{p(\mathbf{x}, y)} \left[ L(y, \mathbf{g}(\mathbf{x})) \right]$$

■ Risk can be equivalently expressed as

$$R(\mathbf{g}) = \mathbb{E}_{\bar{p}(\mathbf{x}, \bar{y})} \left[ \bar{L}(\bar{y}, \mathbf{g}(\mathbf{x})) \right]$$

● Complementary loss:

$$\bar{L}(\bar{y}, \mathbf{g}(\mathbf{x})) = -(c - 1)L(\bar{y}, \mathbf{g}(\mathbf{x})) + \sum_{y=1}^c L(y, \mathbf{g}(\mathbf{x}))$$

■ Empirical risk estimation is possible from complementary data!

$$\hat{R}_{\text{Comp}}(\mathbf{g}) = \frac{1}{n} \sum_{i=1}^n \bar{L}(\bar{y}_i, \mathbf{g}(\mathbf{x}_i)) \quad \{(\mathbf{x}_i, \bar{y}_i)\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} \bar{p}(\mathbf{x}, \bar{y})$$

# Theoretical Analysis

61

Ishida, Niu, Hu & Sugiyama (NIPS2017)

## ■ Estimation error:

$$R(\mathbf{g}^*) - R(\hat{\mathbf{g}}_{\text{Comp}}) = \mathcal{O}_p\left(\frac{1}{\sqrt{n}}\right)$$

$$\mathbf{g}^* = \underset{\mathbf{g}}{\operatorname{argmin}} R(\mathbf{g})$$

$$R(\mathbf{g}) = \mathbb{E}_{p(\mathbf{x}, y)} \left[ L(y, \mathbf{g}(\mathbf{x})) \right]$$

$$\hat{\mathbf{g}}_{\text{Comp}} = \underset{\mathbf{g}}{\operatorname{argmin}} \hat{R}_{\text{Comp}}(\mathbf{g})$$

$$\hat{R}_{\text{Comp}}(\mathbf{g}) = \frac{1}{n} \sum_{i=1}^n \bar{L}(\bar{y}_i, \mathbf{g}(\mathbf{x}_i))$$

## ■ Optimal parametric convergence rate is attained!

# Experiments

62

Correct classification rate: average (std)

5% t-test

Dataset	Class	Dim	# train	# test	Proposed	Partial-label	Multi-label
WAVEFORM1	1 ~ 3	21	1226	398	85.8(0.5)	85.7(0.9)	79.3(4.8)
WAVEFORM2	1 ~ 3	40	1227	408	84.7(1.3)	84.6(0.8)	74.9(5.2)
SATIMAGE	1 ~ 7	36	415	211	68.7(5.4)	60.7(3.7)	33.6(6.2)
PENDIGITS	1 ~ 5	16	719	336	87.0(2.9)	76.2(3.3)	44.7(9.6)
	6 ~ 10		719	335	78.4(4.6)	71.1(3.3)	38.4(9.6)
	even #		719	336	90.8(2.4)	76.8(1.6)	43.8(5.1)
	odd #		719	335	76.0(5.4)	67.4(2.6)	40.2(8.0)
	1 ~ 10		719	335	38.0(4.3)	33.2(3.8)	16.1(4.6)
DRIVE	1 ~ 5	48	3955	1326	89.1(4.0)	77.7(1.5)	31.1(3.5)
	6 ~ 10		3923	1313	88.8(1.8)	78.5(2.6)	30.4(7.2)
	even #		3925	1283	81.8(3.4)	63.9(1.8)	29.7(6.3)
	odd #		3939	1278	85.4(4.2)	74.9(3.2)	27.6(5.8)
	1 ~ 10		3925	1269	40.8(4.3)	32.0(4.1)	12.7(3.1)
LETTER	1 ~ 5	16	565	171	79.7(5.3)	75.1(4.4)	28.3(10.4)
	6 ~ 10		550	178	76.2(6.2)	66.8(2.5)	34.0(6.9)
	11 ~ 15		556	177	78.3(4.1)	67.4(3.3)	28.6(5.0)
	16 ~ 20		550	184	77.2(3.2)	68.4(2.1)	32.7(6.4)
	21 ~ 25		585	167	80.4(4.2)	75.1(1.9)	32.0(5.7)
	1 ~ 25		550	167	5.1(2.1)	5.0(1.0)	5.2(1.1)
USPS	1 ~ 5	256	652	166	79.1(3.1)	70.3(3.2)	44.4(8.9)
	6 ~ 10		542	147	69.5(6.5)	66.1(2.4)	37.3(8.8)
	even #		556	147	67.4(5.4)	66.2(2.3)	35.7(6.6)
	odd #		542	147	77.5(4.5)	69.3(3.1)	36.6(7.5)
	1 ~ 10		542	127	30.7(4.4)	26.0(3.5)	13.3(5.4)

Proposed  
method  
works  
well!

# Incorporating Ordinary Labels <sup>63</sup>

- Convert **multiclass labeling** into **yes-no labeling**:



[http://www.softbank.jp/corp/group/sbr/news/press/2014/20141029\\_01/](http://www.softbank.jp/corp/group/sbr/news/press/2014/20141029_01/)



<https://www.bostondynamics.com/atlas>

Is this Softbank Pepper?  
**Yes! (ordinary label)**

Is this iRobot Roomba?  
**No! (complementary label)**

- Use both of ordinary and complementary labels!

$$R(\mathbf{g}) = \alpha \mathbb{E}_{p(\mathbf{x}, y)} \left[ L(y, \mathbf{g}(\mathbf{x})) \right] + (1 - \alpha) \mathbb{E}_{\bar{p}(\mathbf{x}, \bar{y})} \left[ \bar{L}(\bar{y}, \mathbf{g}(\mathbf{x})) \right]$$

$\alpha \in [0, 1]$

# Experiments

$$R(f) = \alpha \mathbb{E}_{p(\mathbf{x}, y)} \left[ \mathcal{L}(f(\mathbf{x}), y) \right] + (1 - \alpha) \left\{ (c - 1) \mathbb{E}_{\bar{p}(\mathbf{x}, \bar{y})} \left[ \bar{\mathcal{L}}(f(\mathbf{x}), \bar{y}) \right] + \text{Const.} \right\}$$

Dataset	Class	Dim	# train	# test	OL ( $\alpha = 1$ )	CL ( $\alpha = 0$ )	OL & CL ( $\alpha = \frac{1}{2}$ )
WAVEFORM1	1 ~ 3	21	413/826	408	85.3(0.8)	86.0(0.4)	<b>86.9(0.5)</b>
WAVEFORM2	1 ~ 3	40	411/821	411	82.7(1.3)	82.0(1.7)	<b>84.7(0.6)</b>
SATIMAGE	1 ~ 7	36	69/346	211	74.9(4.9)	70.1(5.6)	<b>81.2(1.1)</b>
PENDIGITS	1 ~ 5	16	144/575	336	<b>91.3(2.1)</b>	84.7(3.2)	<b>93.1(2.0)</b>
	6 ~ 10		144/575	335	<b>86.3(3.5)</b>	78.3(6.2)	<b>87.8(2.8)</b>
	even #		144/575	336	94.3(1.7)	91.0(4.3)	<b>95.8(0.6)</b>
	odd #		144/575	335	<b>85.6(2.0)</b>	75.9(3.1)	<b>86.9(1.1)</b>
	1 ~ 10		72/647	335	61.7(4.3)	41.1(5.7)	<b>66.9(2.0)</b>
DRIVE	1 ~ 5	48	780/3121	1305	92.1(2.6)	89.0(2.1)	<b>94.2(1.0)</b>
	6 ~ 10		795/3180	1290	<b>87.0(3.0)</b>	86.5(3.1)	<b>89.5(2.1)</b>
	even #		657/3284	1314	<b>91.4(2.9)</b>	81.8(4.6)	<b>91.8(3.3)</b>
	odd #		790/3161	1255	91.1(1.5)	86.7(2.9)	<b>93.4(0.5)</b>
	1 ~ 10		397/3570	1292	<b>75.2(2.8)</b>	40.5(7.2)	<b>77.6(2.2)</b>
LETTER	1 ~ 5	16	113/452	171	85.2(1.3)	77.2(6.1)	<b>89.5(1.6)</b>
	6 ~ 10		110/440	178	81.0(1.7)	77.6(3.7)	<b>84.6(1.0)</b>
	11 ~ 15		111/445	177	81.1(2.7)	76.0(3.2)	<b>87.3(1.6)</b>
	16 ~ 20		110/440	184	81.3(1.8)	77.9(3.1)	<b>84.7(2.0)</b>
	21 ~ 25		117/468	167	86.8(2.7)	81.2(3.4)	<b>91.1(1.0)</b>
	1 ~ 25		22/528	167	11.9(1.7)	6.5(1.7)	<b>31.0(1.7)</b>
USPS	1 ~ 5	256	130/522	166	83.8(1.7)	76.5(5.3)	<b>89.5(1.3)</b>
	6 ~ 10		108/434	147	79.2(2.1)	67.6(4.3)	<b>85.5(2.4)</b>
	even #		108/434	166	79.6(2.7)	67.4(4.4)	<b>84.8(1.4)</b>
	odd #		111/445	147	82.7(1.9)	72.9(6.2)	<b>87.3(2.2)</b>
	1 ~ 10		54/488	147	43.7(2.6)	28.5(3.6)	<b>59.3(2.2)</b>

5% t-test

Incorporating complementary labels Improves the accuracy!





# Contents

65

1. Background
2. PN Classification
3. PU Classification
4. PNU Classification
5. Pconf Classification
6. UU Classification
7. SU Classification
8. Comp Classification
9. **Summary**

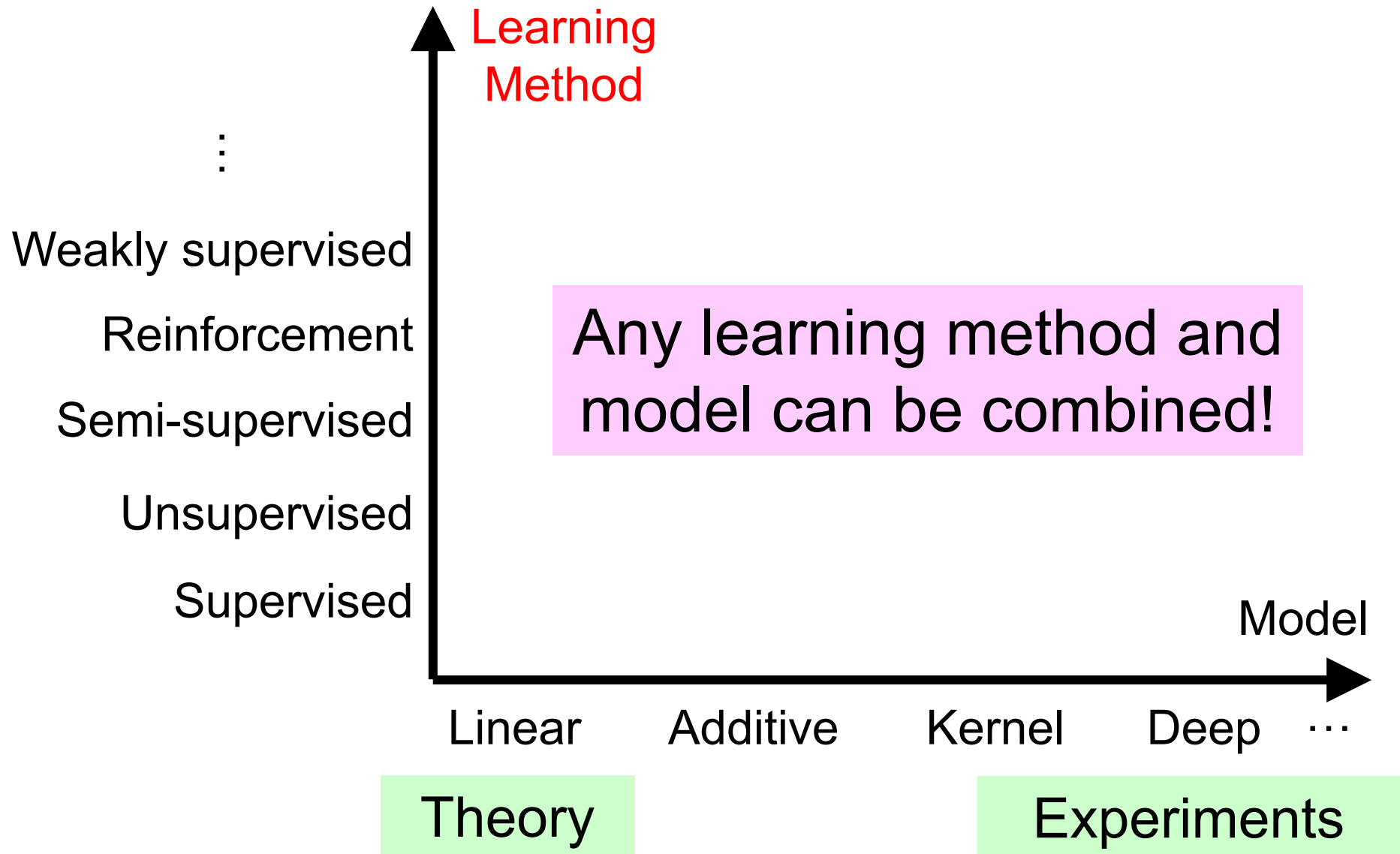
- **P**: Positive
- **N**: Negative
- **U**: Unlabeled
- **Conf**: Confidence
- **S**: Similar
- **Comp**: Complementary

**Slides:**

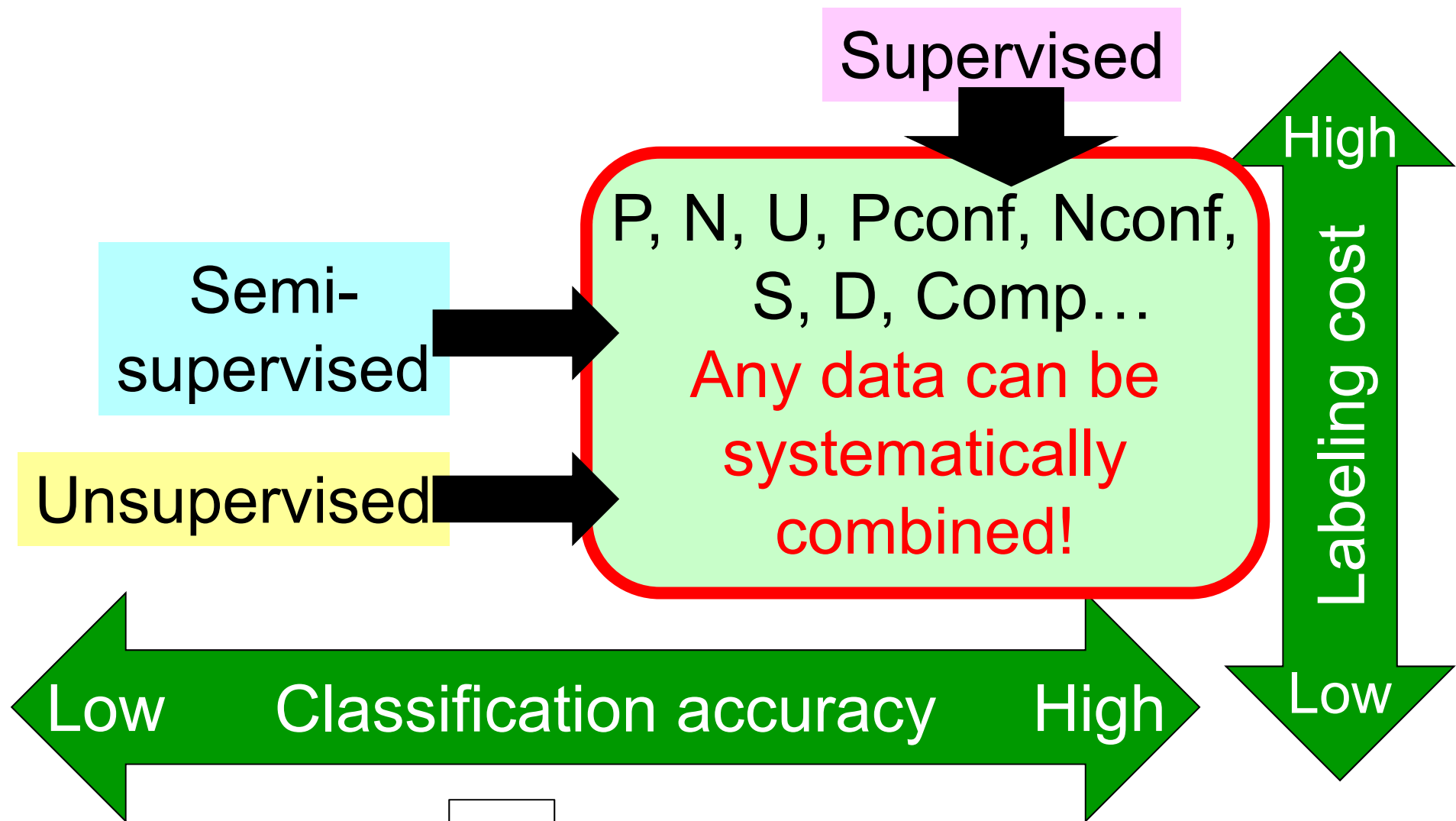
<http://goo.gl/meiTwY>

# Model vs. Learning Methods

66



# Learning from Weak Supervision<sup>67</sup>



Sugiyama, Niu, Sakai & Ishida,  
Machine Learning from Weak Supervision  
MIT Press, 2020 (?)