

# Machine Learning from Weak, Noisy, and Biased Supervision

Masashi Sugiyama

RIKEN Center for Advanced Intelligence Project/  
The University of Tokyo, Japan



<http://www.ms.k.u-tokyo.ac.jp/sugi/>



東京大学  
THE UNIVERSITY OF TOKYO

# About Myself



## ■ Masashi Sugiyama:

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

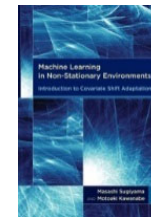
## ■ Interests: Machine learning (ML)

- ML theory & algorithm →
- ML applications (signal, image, language, brain, robot, mobility, advertisement, biology, medicine, education...)

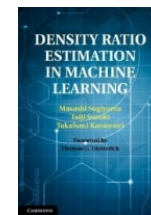
## ■ Academic activities:

- Program Chairs for NeurIPS2015, AISTATS2019, ACML2010/2020...
- Keynote speaker at ICLR2023.

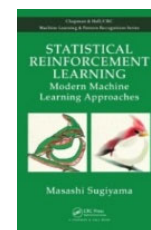
Sugiyama & Kawanabe,  
**Machine Learning  
in Non-Stationary  
Environments**,  
MIT Press, 2012



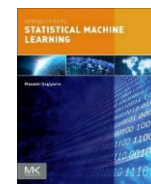
Sugiyama, Suzuki &  
Kanamori, **Density Ratio  
Estimation in Machine  
Learning**, Cambridge  
University Press, 2012



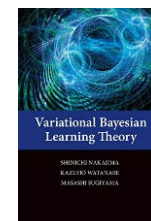
Sugiyama,  
**Statistical Reinforcement  
Learning**,  
Chapman and Hall/CRC,  
2015



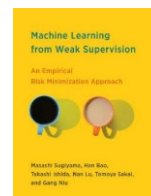
Sugiyama,  
**Introduction to Statistical  
Machine Learning**,  
Morgan Kaufmann,  
2015



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



Sugiyama, Bao, Ishida,  
Lu, Sakai & Niu.  
**Machine Learning from  
Weak Supervision**,  
MIT Press, 2022.



# What is “RIKEN”?

## ■ Name in Japanese:

理化学研究所



- Pronounced as: rikagaku kenkyusho
- Meaning: Physics and Chemistry Research Institute

## ■ Acronym in Japanese: 理研 (RIKEN)

# Brief History

Shibusawa Eiichi



Okochi Masatoshi



Nishina Yoshio



Tomonaga Shinichiro



Matsumoto Hiroshi



Gonokami Makoto



Kikuchi Dairoku



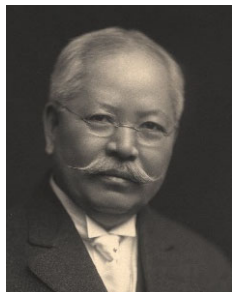
Suzuki Umetaro



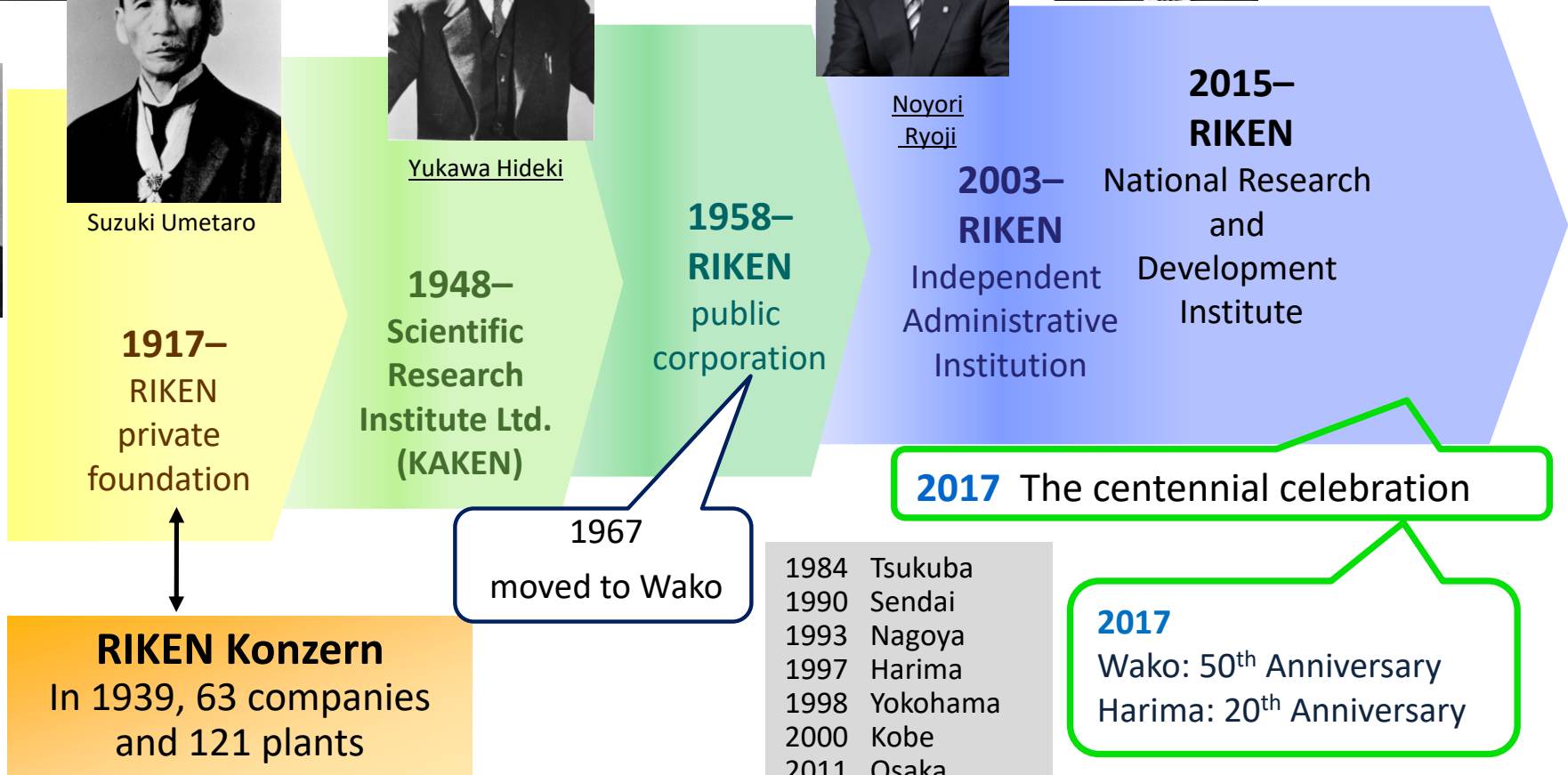
Yukawa Hideki



Noyori Ryoji



Takamine Jokichi



# Office and Research

2900 researchers  
500 admin staffs

RIKEN Information R&D and Strategy Headquarters



Informatics

Spring-8 Center

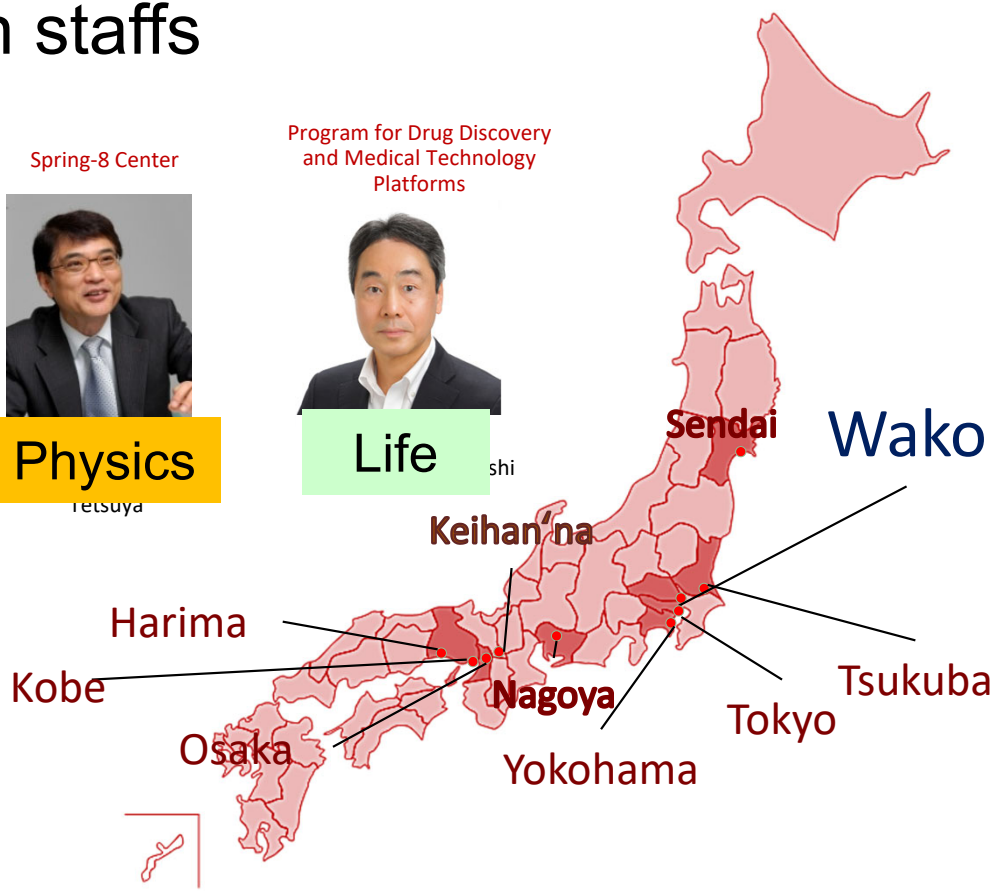


Physics  
Tetsuya

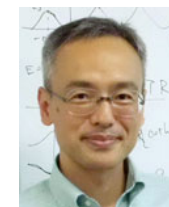
Program for Drug Discovery and Medical Technology Platforms



Life  
Shi



Center for Quantum Computing



Physics  
Yasunobu

Center for Emergent Matter Science



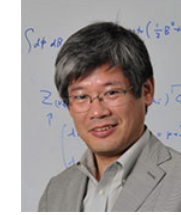
Physics

Nishina Center for Accelerator-Based Science



Physics

Interdisciplinary Theoretical and Mathematical Sciences Program



Physics

Center for Advanced Photonics



Physics

Center for Brain Science



Life

Center for Computational Science



Informatics  
Director  
Dr. Matsuoka Satoshi

Center for Biosystems Dynamics Research



Life  
Director  
Dr. Nishida Eisuke

Center for Integrative Medical Sciences



Life  
Director  
Dr. Yamamoto Kazuhiko

Center for Sustainable Resource Science



Life  
Director  
Dr. Saito Kazuki

Center for Advanced Intelligence Project



Informatics  
Director  
Dr. Sugiyama Masashi

BioResource Research Center



Life  
Director  
Dr. Shiroishi Toshihiko

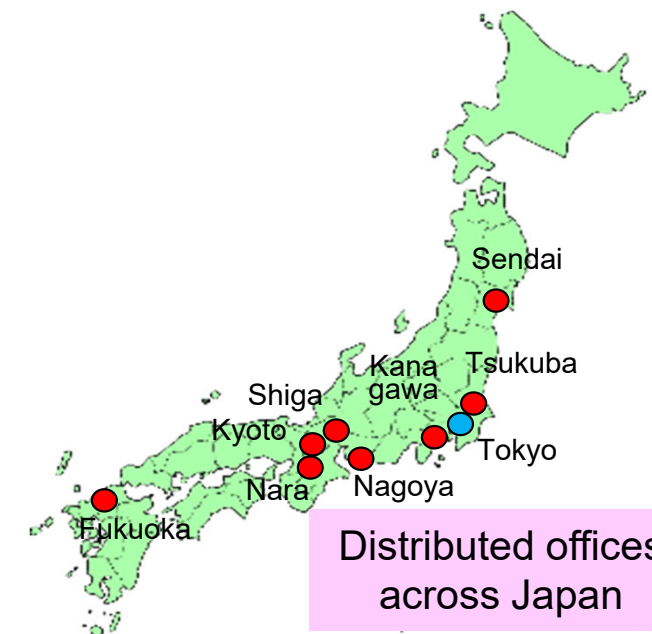


# What is RIKEN-AIP?

- RIKEN founded **Center for Advanced Intelligence Project (AIP)** in 2016, under Ministry of Education, Culture, Sports, Science and Technology (MEXT):
  - 130 employed researchers (40% international, 25% female)
  - 250 visiting researchers
  - 130 domestic students
  - 140 international interns (total)
  - 40+ international collaboration partners
  - 40+ industry projects



Main office  
in the heart of Tokyo

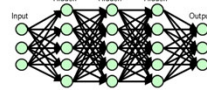


Distributed offices  
across Japan

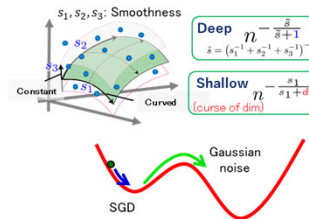
## Developing New AI Technology

### Theory of deep learning:

- Better prediction than shallow learning
- No curse of dimensionality
- Global optimization



$$\mathbb{E}[\|f_T - f^*\|_{L_2}^2] \leq \epsilon_M + O(T^{-\frac{2r\beta}{2r\beta+1}})$$

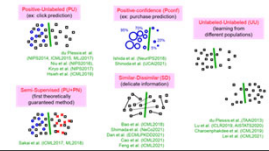


### Developing new methods:

- Weakly supervised learning
- Noise robust learning
- Causal inference

#### Weakly Supervised Classification

Various weakly supervised classification problems can be solved by risk-rewriting systematically!



#### Noise Transition Correction

Noise transition matrix  $T$ :

$$T = \begin{bmatrix} 0.7 & 0.05 \\ 0 & 0.95 \end{bmatrix}$$

- Clean-to-noisy flipping probability.
- Major approaches:
  - Loss correction by  $T^{-1}$  to eliminate noise.
  - Classifier adjustment by  $T$  to simulate noise.
- We want to estimate  $T$  only from noisy data:
  - Use human cognition as a "mask" for  $T$ .
  - Learn  $T$  and a classifier dynamically.
  - Decompose  $T$  into simpler components.
  - Regularize  $T$  to be estimable.
  - Extension to input-dependent noise  $T(x)$ .

#### Causal Inference in the Presence of Hidden Cause

In causal inference, how to handle hidden cause is a big challenge!

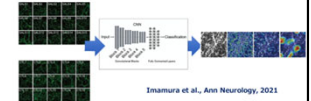
- We developed the first method to estimate the entire structure in the presence of hidden cause:
  - Speech separation technique is employed to separate hidden cause.



## Accelerating Scientific Research

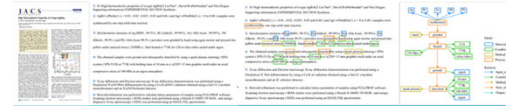
### Medical science:

- Prostate/pancreatic cancer detection
- ALS early diagnosis
- Fetal heart screening
- Colonoscopy



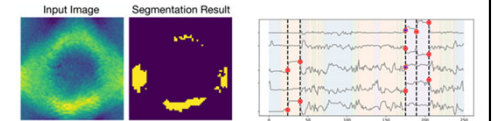
### Material science:

- Database creation with text mining



### Data-driven science:

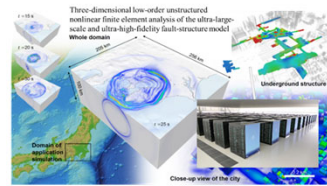
- Selective inference for reliability evaluation



## Solving Socially Critical Problems

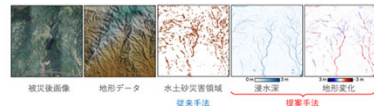
### Natural disaster:

- Fugaku-based earthquake simulation
- Remote sensing disaster analysis



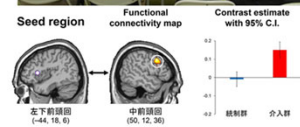
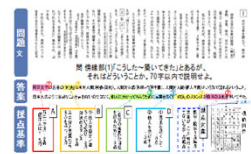
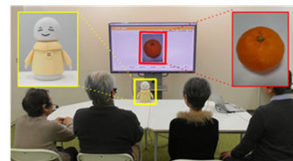
### Elderly healthcare:

- Chat-robot-guided cognitive function improvement



### Education:

- Automatic essay evaluation
- Interactive essay writing support



## Studying AI-ELSI

### AI Ethical guidelines:

- Japanese Society for AI, Ministry of Internal Affairs and Communications, Cabinet Office
- IEEE, G20, OECD



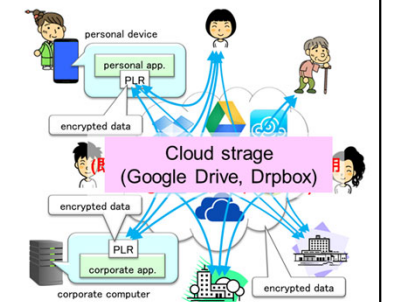
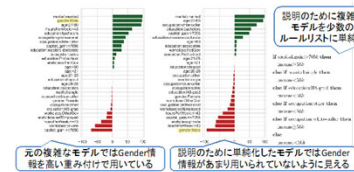
### Personal data management:

- Individual-based accessibility control system



### AI security and reliability:

- Adversarial attack/defense
- Fairness faking/guarantee

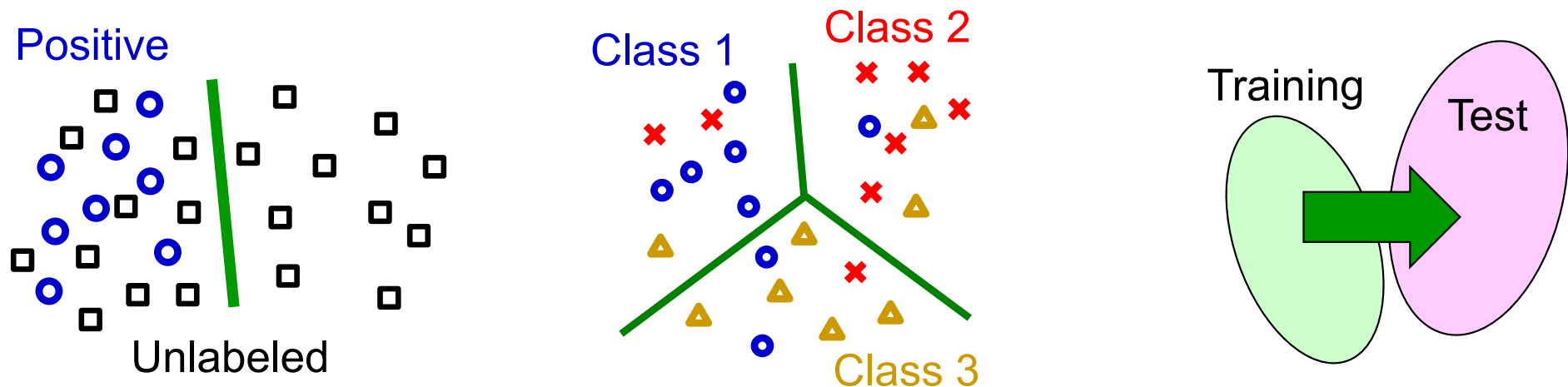


# Reliable Machine Learning

■ **Reliability** of machine learning systems can be degraded by various factors:

- **Insufficient information:** weak supervision.
- **Label noise:** human error, sensor error.
- **Data bias:** changing environments, privacy.

■ Improving the reliability is an urgent challenge!







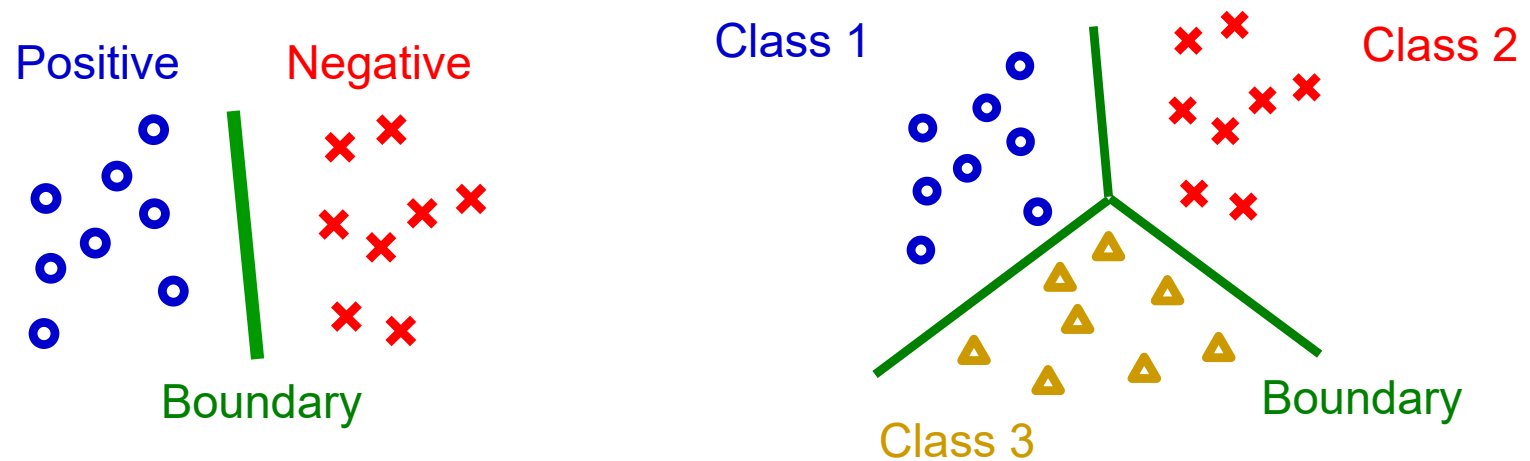
# Contents

1. Weakly Supervised Learning
2. Noisy-Label Learning
3. Transfer Learning
4. Towards More Reliable Learning

# Weakly Supervised Classification

10

- Supervised classification from big labeled data is successful: speech, image, language, ...



- However, there are many applications where **big labeled data is not available**:
  - Medicine, disaster, robot, brain, ...
- We want to utilize “**weak**” supervision that can be collected easily!

# Positive-Unlabeled (PU) Classification 11

Li+ (IJCAI2003)

- **Given:** PU samples (no N samples).

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

- **Goal:** Obtain a classifier minimizing the PN risk.

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

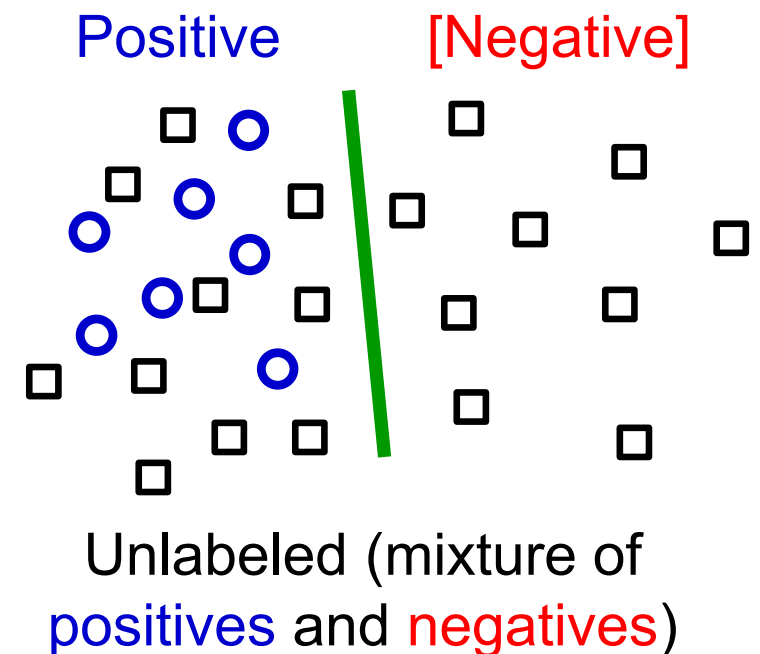
$\mathbb{E}$  : expectation

$\ell$  : loss

$y = \{+1, -1\}$

**Example:** Ad click prediction

- **Clicked ad:** User likes it  $\rightarrow$  P
- **Unclicked ad:** User dislikes it or User likes it but doesn't have time to click it  $\rightarrow$  U (=P or N)



## Decompose the risk:

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

$\pi = p(y = +1)$  : Class prior (assumed known)  $\rightarrow$

Scott+ (AISTATS2009)  
Ramaswamy+ (ICML2016)  
du Plessis+ (MLJ2017)  
Yao+ (ICLR2022)

## Without N data, $R^-(f)$ can not be estimated directly:

- Eliminate the expectation over N data as

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

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

## Unbiased risk estimator:

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



# Non-Negative Risk Correction

Kiryō+ (NeurIPS2017) , Lu+ (AISTATS2020)

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

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

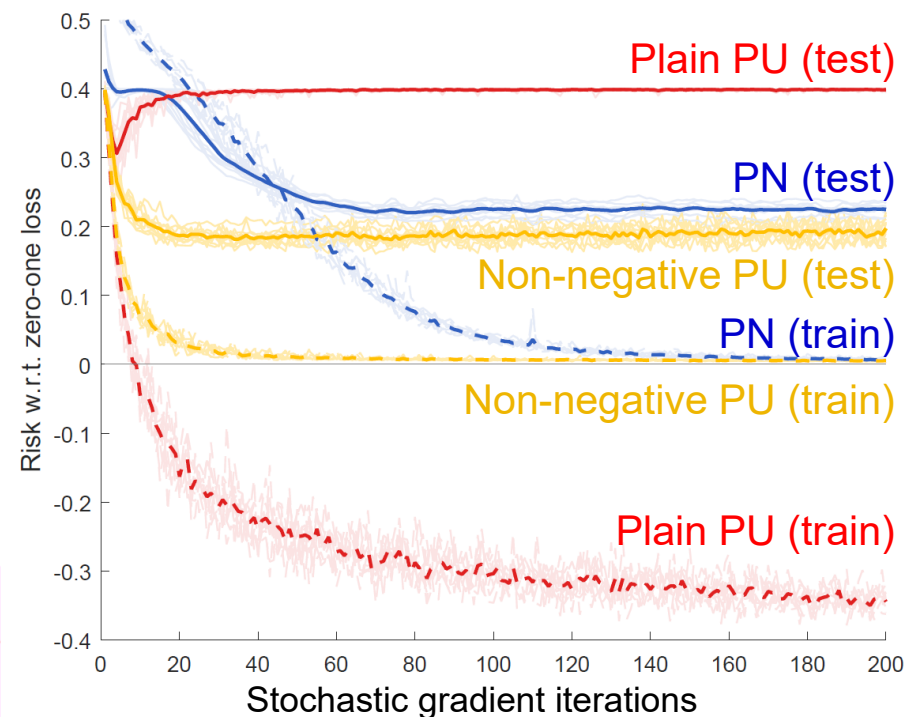
■ Empirical estimate:  $\hat{R}_{\text{PU}}^-(f) = \frac{1}{n_U} \sum_{i=1}^{n_U} \ell(-1, f(\mathbf{x}_i^U)) - \frac{\pi}{n_P} \sum_{i=1}^{n_P} \ell(-1, f(\mathbf{x}_i^P))$

■ When **loss is non-negative**:

- True  $R^-(f)$  is non-negative.
- But empirical estimate can be **negative**!

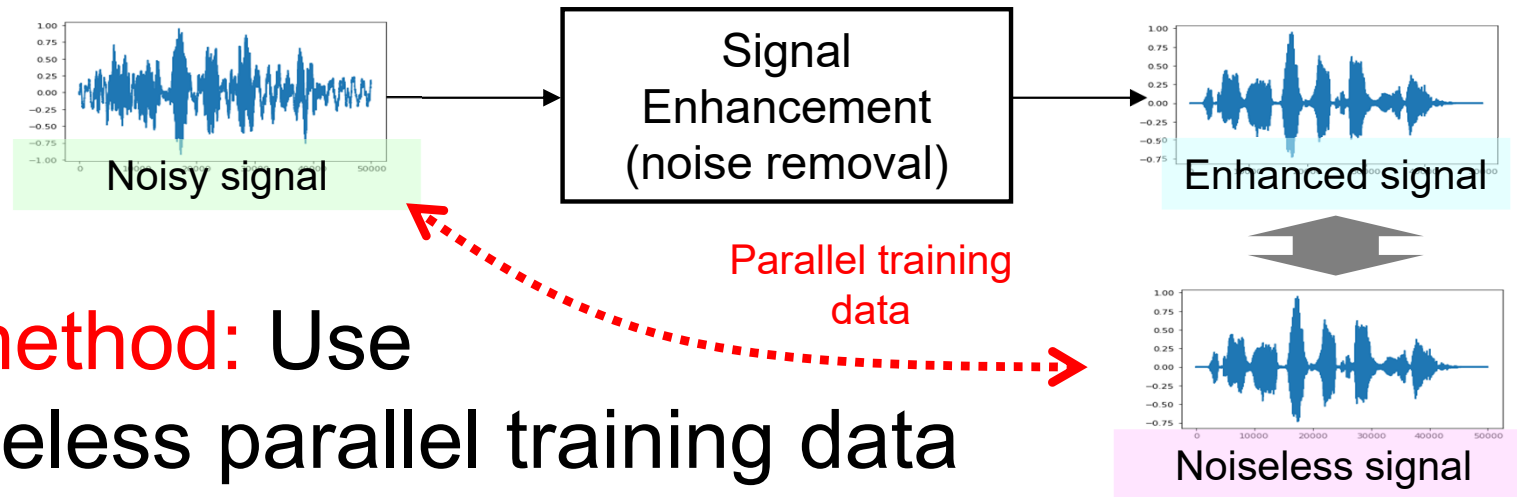
■ **Non-negative correction**:

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



# Signal Enhancement by PU Classification 14

Ito & Sugiyama (ICASSP2023, [Best Paper Award](#))

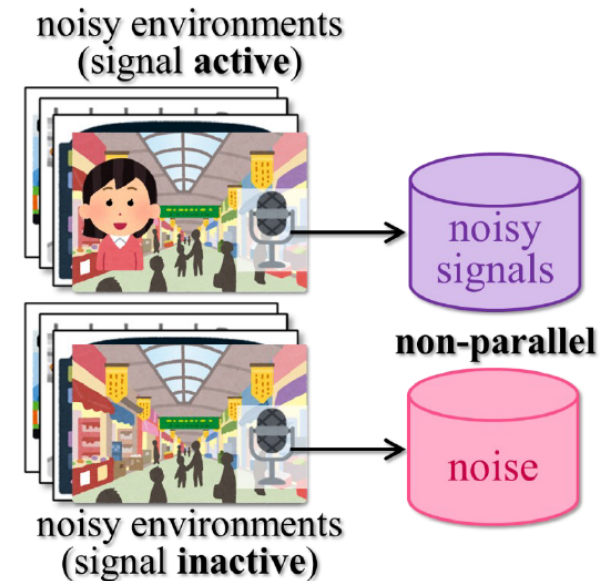


■ **Existing method:** Use noisy/noiseless parallel training data

- In practice, use synthetic data  
→ Do not generalize well in reality.

■ **Proposed method:** Use non-parallel noisy signal and noise.

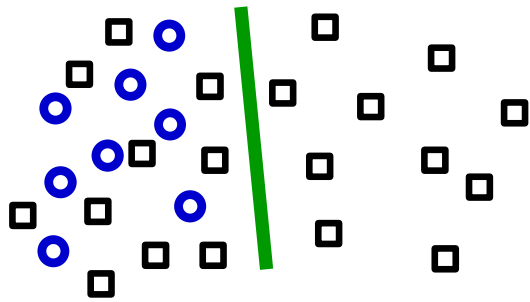
	Methods	SI-SNRi [dB]
Non-parallel	Proposed	14.62 (0.20)
	MixIT <small>Wisdom+ (NeurIPS2020)</small>	12.19 (4.50)
Parallel	Supervised	15.86 (1.28)



# Various Extensions (Binary)

■ Similar unbiased risk estimation is possible!

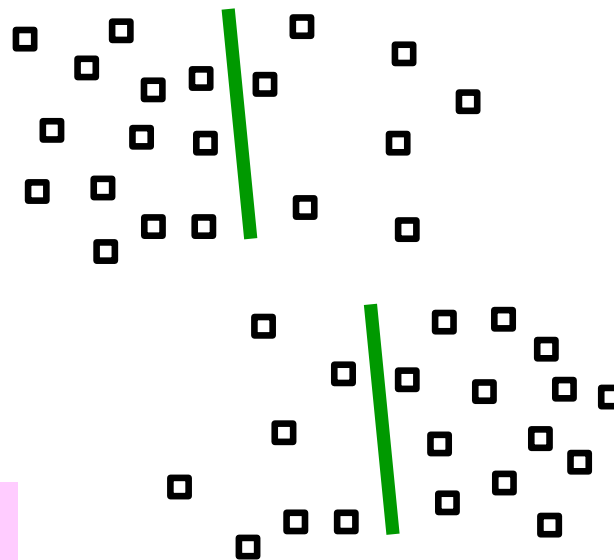
## Positive-Unlabeled (PU)



du Plessis+ (NeurIPS2014, ICML2015, MLJ2017),  
Niu+ (NeurIPS2016), Kiryo+ (NeurIPS2017), Hsieh+ (ICML2019)

Click prediction

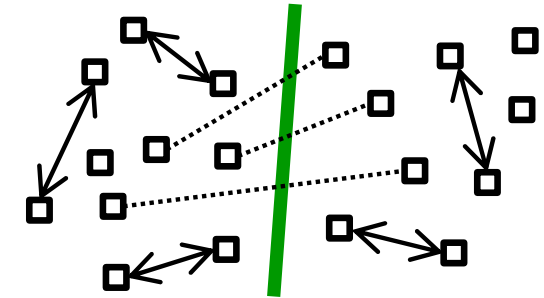
## Unlabeled-Unlabeled (UU)



du Plessis+ (TAAI2013), Lu+ (ICLR2019, AISTATS2020),  
Charoenphakdee+ (ICML2019), Lei+ (ICML2021)

Different populations

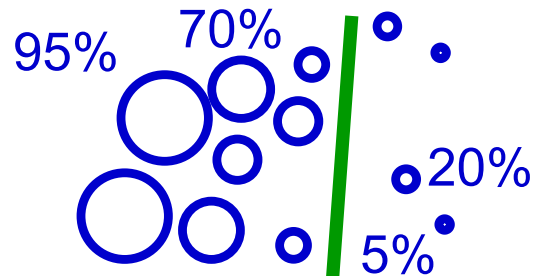
## Similar-Dissimilar (SD)



Bao+ (ICML2018), Shimada+ (NeCo2021),  
Dan+ (ECMLPKDD2021), Cao+ (ICML2021),  
Feng+ (ICML2021)

Sensitive prediction

## Positive-confidence (Pconf)

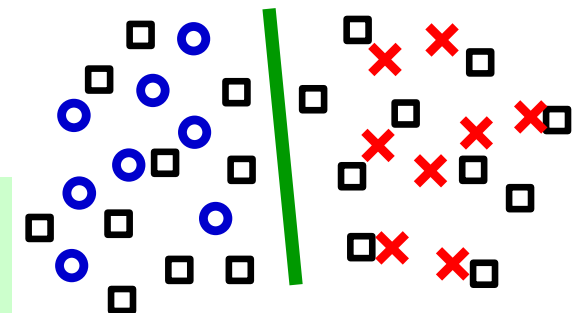


Ishida+ (NeurIPS2018), Shinoda+ (IJCAI2021)

Purchase prediction

Semi-supervised classification  
without manifold/clusters →

## Positive-Negative-Unlabeled (PNU)



Sakai+ (ICML2017, ML2018)

# Various Extensions (Multiclass)

16

■ Labeling patterns in **multi-class** problems is even more painful.

■ **Multi-class weak-labels:**

● **Complementary label:**

Specifies a class that a pattern does **not** belong to (“not 1”).

● **Partial label:** Specifies a subset of classes that contains the correct one (“1 or 2”).

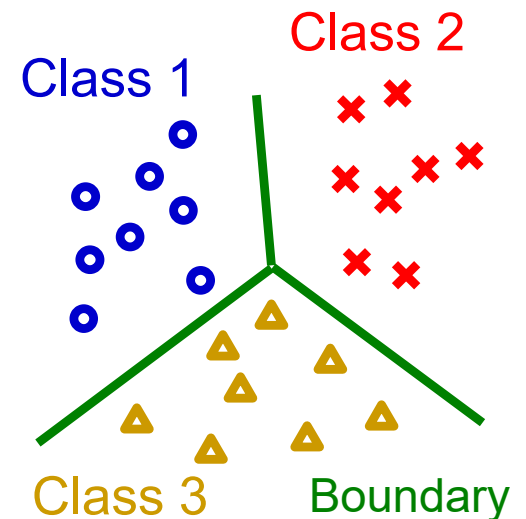
● **Single-class confidence:**

One-class data with full confidence

(“1 with 60%, 2 with 30%, and 3 with 10%”)

Ishida+ (NeurIPS2017,  
ICML2019),  
Chou+ (ICML2020)

Cao+ (arXiv2021)

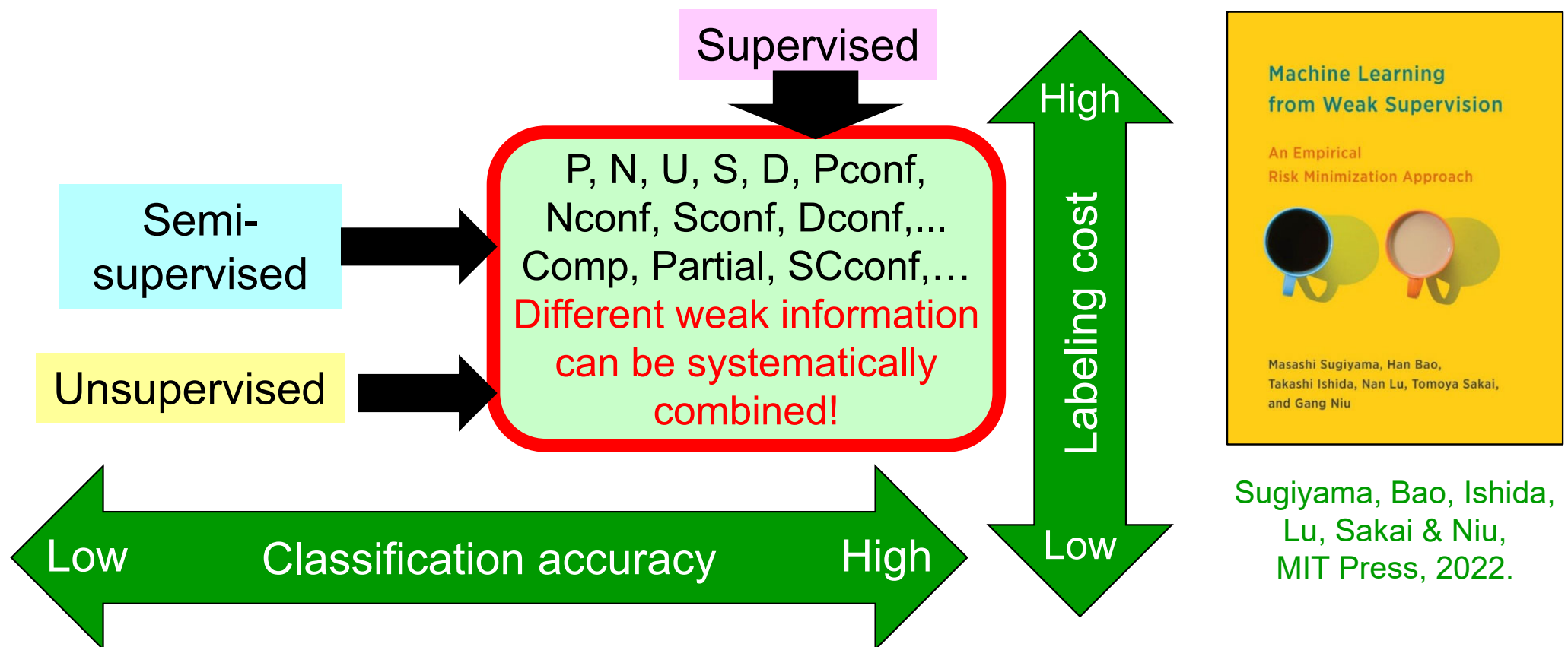


Feng+ (ICML2020,  
NeurIPS2020),  
Lv+ (ICML2020)

■ Similar unbiased risk estimation is possible!



- Empirical risk minimization framework for weakly supervised learning:
  - Any loss, classifier, and optimizer can be used.



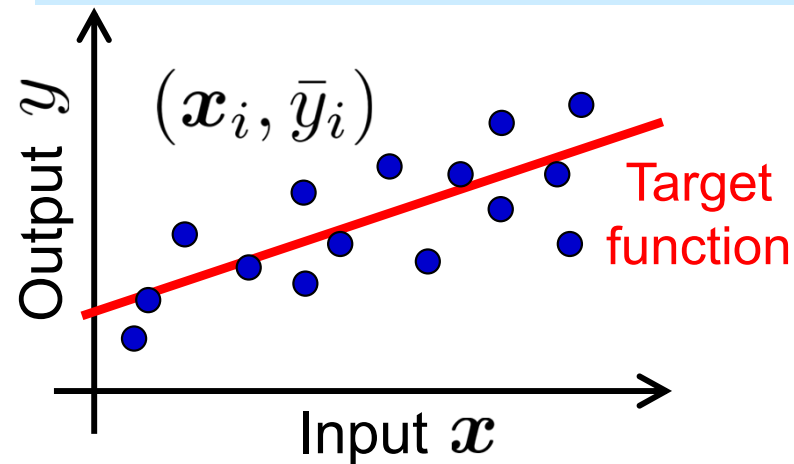


# Contents

1. Weakly Supervised Learning
2. Noisy-Label Learning
  - A) Noise Transition
  - B) Algorithms
3. Transfer Learning
4. Towards More Reliable Learning

# Supervised Learning with Noisy Output 19

Regression (additive noise)

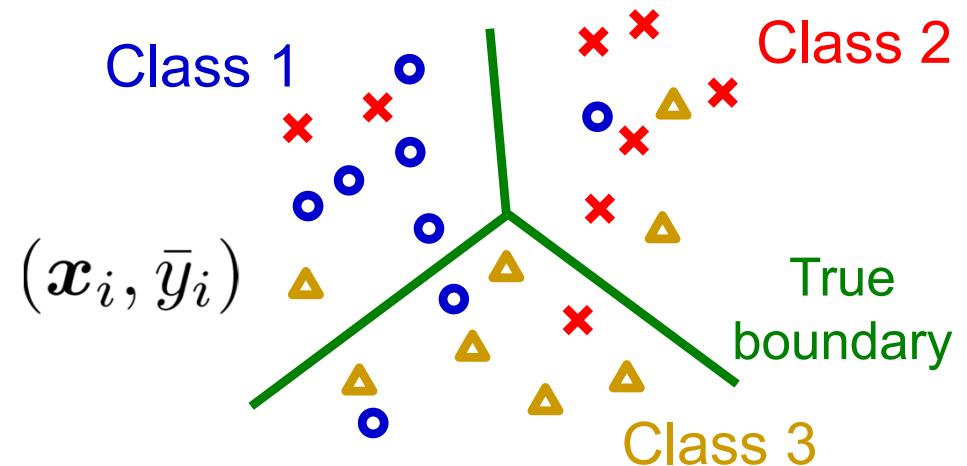


$$\min_g \sum_{i=1}^n \ell(\bar{y}_i, g(\mathbf{x}_i))$$

$\ell$ : loss

$\bar{y}$ : noisy output

Classification (label flipping noise)



$g$ : probabilistic classifier

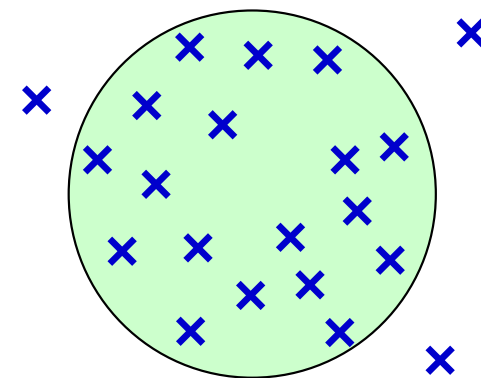
- Hasn't such a classic problem been solved?
  - **Regression**: Yes, noisy big data yield consistency.
  - **Classification**: Specific noise reduction mechanism is needed to achieve consistency!

# Classical Approaches

20

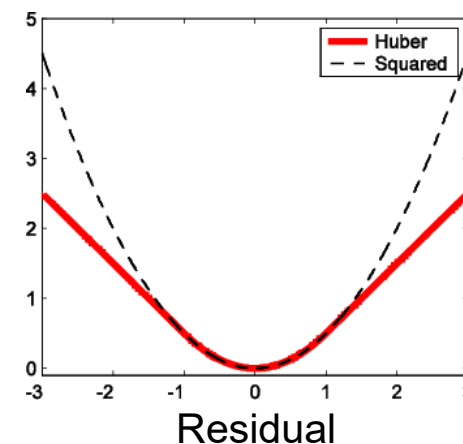
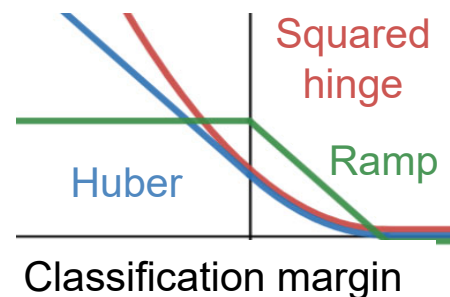
## ■ Unsupervised outlier removal:

- Substantially more difficult than classification.



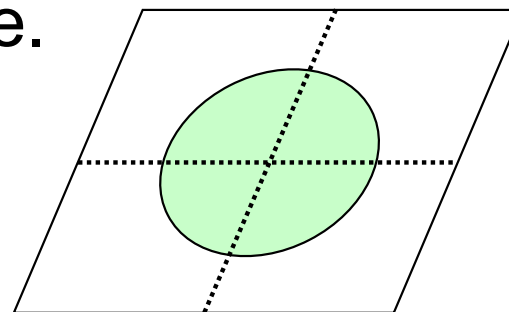
## ■ Robust loss:

- Works well for regression, but limited effectiveness for classification.

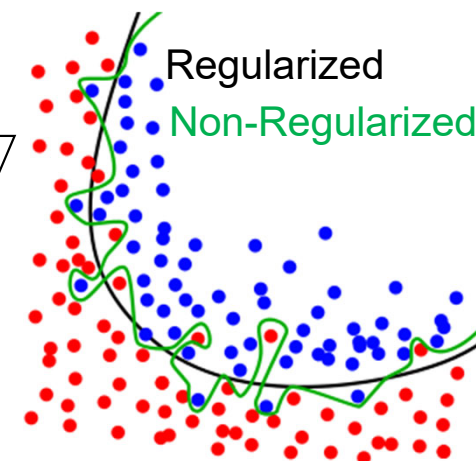


## ■ Regularization:

- Effective in suppressing overfitting, but too smooth for strong noise.



$\ell_2$ -regularization



<https://en.wikipedia.org/wiki/Overfitting>

## ■ Need new approaches!



## ■ Noise transition matrix $T$ :

- Clean-to-noisy flipping probability.

$$T = \begin{array}{|c|c|c|} \hline 1 & 0 & 0 \\ \hline 0.1 & 0.8 & 0.1 \\ \hline 0.5 & 0.5 & 0 \\ \hline \end{array}$$

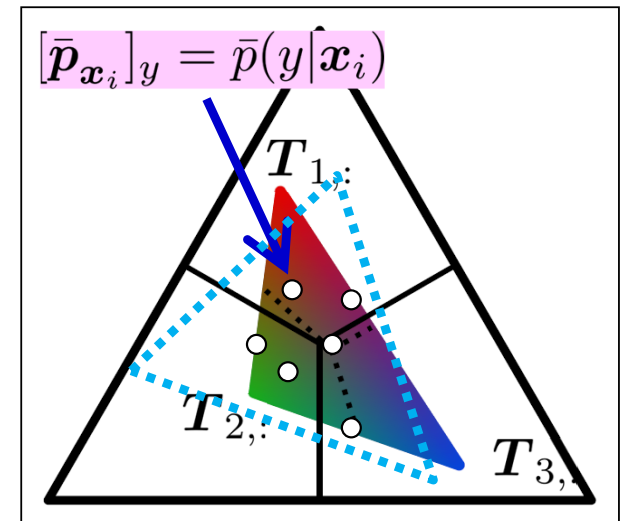
## ■ Major approaches: Patrini+ (CVPR2017)

- Classifier adjustment by  $T^{-1}$  to simulate noise.
- Loss correction by  $T^{-1}$  to eliminate noise.

## ■ We want to estimate $T$ only from noisy data:

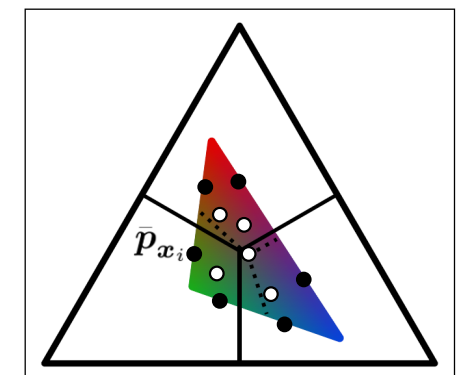
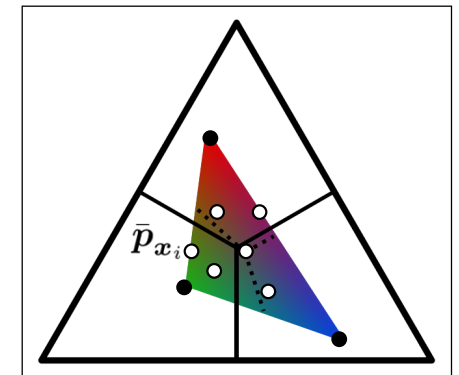
- Use human cognition as a “mask” for  $T$ . Han+ (NeurIPS2018)
- Reduce estimation error of  $T$ . Xia+ (NeurIPS2019)
- Learn  $T$  and classifier simultaneously. Yao+ (NeurIPS2020)
- Estimate  $T$  under weaker conditions. Zhang+ (ICML2021)
- Li+ (ICML2021)

- Noisy training data  $\{(\mathbf{x}_i, \bar{y}_i)\}_{i=1}^n$  can be **mapped in the simplex** formed by noise transition matrix  $T$ .
- Minimizing the **volume** of the simplex can give a solution:



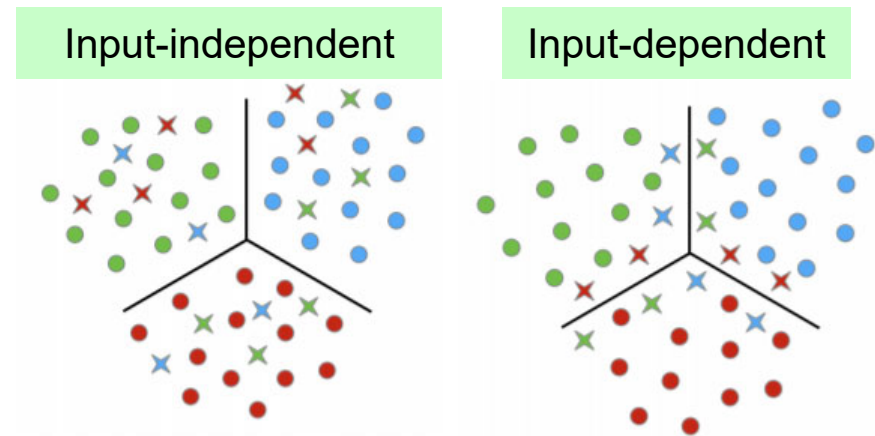
$$\min_{\mathbf{T}', \mathbf{g}} \sum_{i=1}^n \ell(\bar{y}_i, \mathbf{T}'^\top \mathbf{g}(\mathbf{x}_i)) + \lambda \log \det(\mathbf{T}') \quad \lambda > 0$$

- With noiseless labels, we can find the true  $T$ .
- Even without noiseless labels, “**sufficiently scattered**” training data allow identification of the true  $T$ !



■ Real-world noise may be **input-dependent**:

- E.g., noise level is high **near the boundary**.



■ Modeling input-dependent noise:  $T_{y, \bar{y}}(\mathbf{x}) = \bar{p}(\bar{y} | y, \mathbf{x})$

- Extremely challenging to estimate the noise transition matrix **function**!

■ **Exploring heuristic solutions:**

- Parts-based estimation.
- Use of additional confidence scores.
- Manifold regularization.

Xia+ (NeurIPS2020)

Berthon+ (ICML2021)

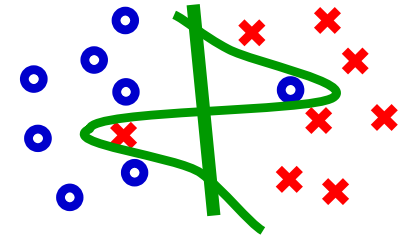
Cheng+ (CVPR2022)

# Co-teaching

## Memorization of neural nets:

- Stochastic gradient descent fits clean data faster.
- However, naïve early stopping does not work well.

Arpit+ (ICML2017)  
Zhang+ (ICLR2017)



## “Co-teaching” between two neural nets:

- Teach small-loss data each other.

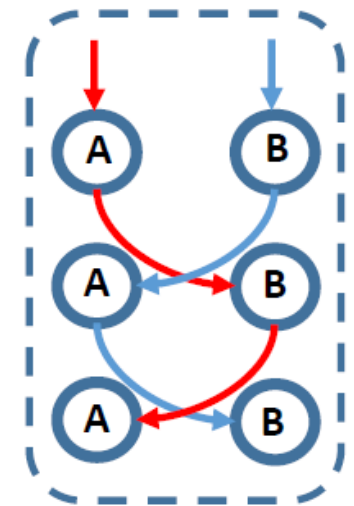
Han+ (NeurIPS2018)

- Teach only disagreed data.

Yu+ (ICML2019)

- Gradient ascent for large-loss data.

Han+ (ICML2020)



## No theory but very robust in experiments:

- Works well even if 50% random label flipping!



# Summary: Noisy-Label Learning

25

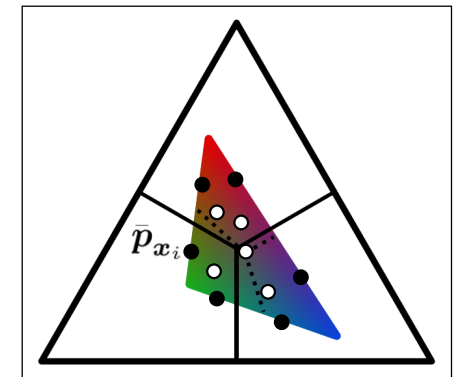
- Explicit treatment of label noise is necessary:
  - Loss correction by noise transition is promising.

$$T_{y, \bar{y}} = \bar{p}(\bar{y} | y)$$

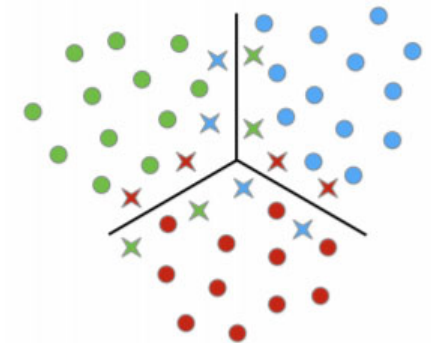
- However, noise transition is generally non-identifiable:

$$T^\top p = T_2^\top (T_1^\top p) \quad T = T_1 T_2$$

- Recent development allows consistent estimation under mild assumptions.



- Real-world noise is often input-dependent:
  - Heuristic solutions have been developed.
  - Further theoretical development is needed.





# Contents

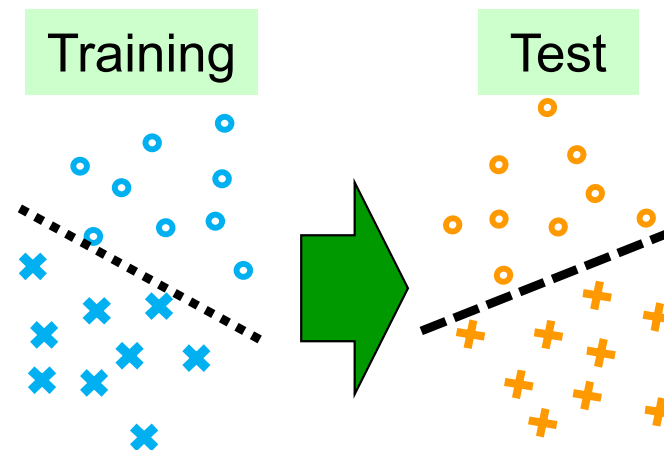
26

1. Weakly Supervised Learning
2. Noisy-Label Learning
3. **Transfer Learning**
4. Towards More Reliable Learning



■ Training and test data often follow **different distributions**, due to

- changing environments,
- sample selection bias (privacy).

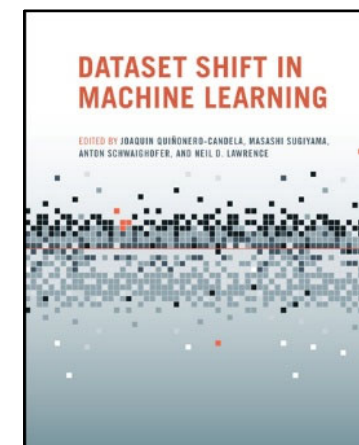


■ **Transfer learning:**

- Train a test-domain predictor using training data from different domains.



NIPS Workshop on Learning when Test and Training Inputs Have Different Distributions, Whistler 2006



Quiñonero-Candela, Sugiyama, Schwaighofer & Lawrence (MIT Press 2009)

# Basics: Importance-Weighted Training 28

- **Covariate shift:** Only input distributions change.

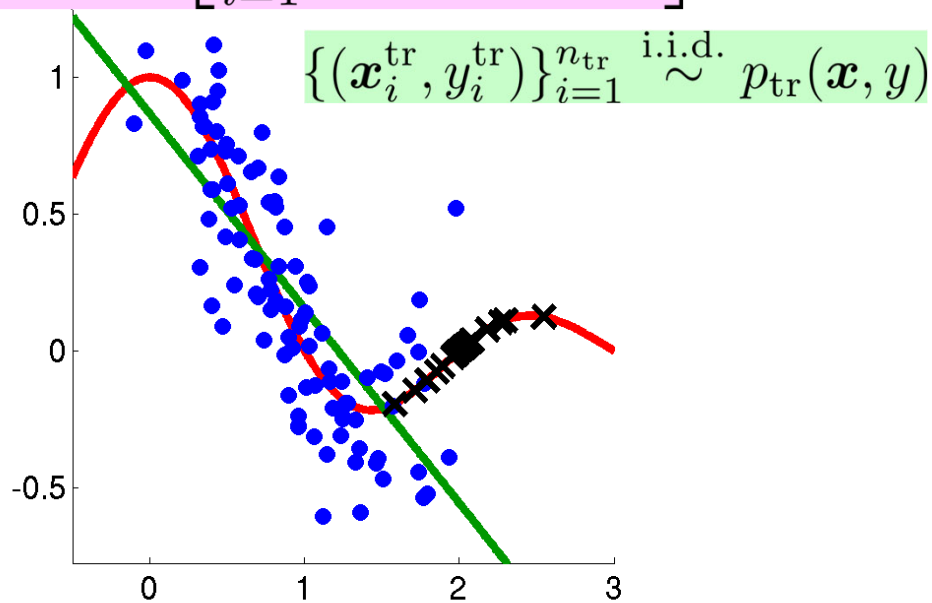
$$p_{\text{tr}}(\mathbf{x}) \neq p_{\text{te}}(\mathbf{x})$$
$$p_{\text{tr}}(y|\mathbf{x}) = p_{\text{te}}(y|\mathbf{x})$$

$\mathbf{x}$ : Input

$y$ : Output

Shimodaira (JSPI2000)

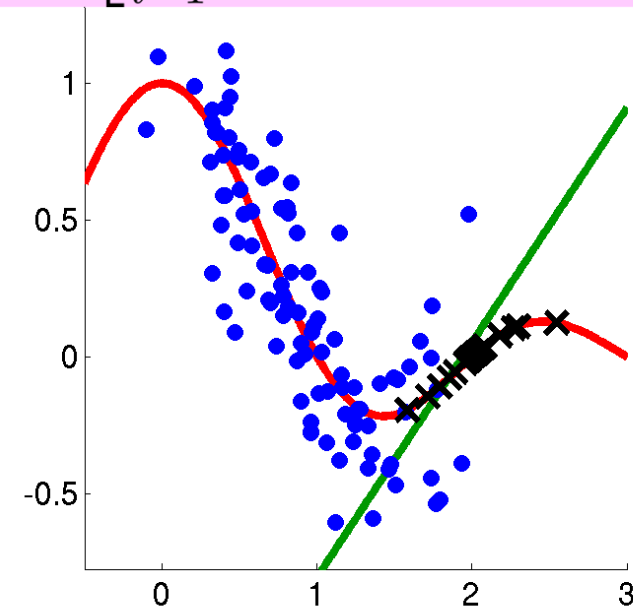
$$\operatorname{argmin}_f \left[ \sum_{i=1}^{n_{\text{tr}}} \ell(f(\mathbf{x}_i^{\text{tr}}), y_i^{\text{tr}}) \right]$$



Ordinary training is not consistent

Importance

$$\operatorname{argmin}_f \left[ \sum_{i=1}^{n_{\text{tr}}} \frac{p_{\text{te}}(\mathbf{x}_i^{\text{tr}})}{p_{\text{tr}}(\mathbf{x}_i^{\text{tr}})} \ell(f(\mathbf{x}_i^{\text{tr}}), y_i^{\text{tr}}) \right]$$



Importance-weighted training is consistent

- **Given:** training and test input data

$$\{\mathbf{x}_i^{\text{tr}}\}_{i=1}^{n_{\text{tr}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{tr}}(\mathbf{x}) \quad \{\mathbf{x}_j^{\text{te}}\}_{j=1}^{n_{\text{te}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{te}}(\mathbf{x})$$

- **Kernel mean matching:** Huang+ (NeurIPS2006)

- Match the means of  $r(\mathbf{x})p_{\text{tr}}(\mathbf{x})$  and  $p_{\text{te}}(\mathbf{x})$  in RKHS  $\mathcal{H}$ .

$$\min_{r \in \mathcal{H}} \left\| \int K(\mathbf{x}, \cdot) p_{\text{te}}(\mathbf{x}) d\mathbf{x} - \int K(\mathbf{x}, \cdot) r(\mathbf{x}) p_{\text{tr}}(\mathbf{x}) d\mathbf{x} \right\|_{\mathcal{H}}^2 \quad K(\mathbf{x}, \cdot) : \text{kernel}$$

- **Least-squares importance fitting (LSIF):** Kanamori+ (NeurIPS2008)

- Fit a model  $r(\mathbf{x})$  to  $\frac{p_{\text{te}}(\mathbf{x})}{p_{\text{tr}}(\mathbf{x})}$  by least squares:

$$\begin{aligned} \operatorname{argmin}_r \left[ \int \left( r(\mathbf{x}) - \frac{p_{\text{te}}(\mathbf{x})}{p_{\text{tr}}(\mathbf{x})} \right)^2 p_{\text{tr}}(\mathbf{x}) d\mathbf{x} \right] \\ = \operatorname{argmin}_r \left[ \int r(\mathbf{x})^2 p_{\text{tr}}(\mathbf{x}) d\mathbf{x} - 2 \int r(\mathbf{x}) p_{\text{te}}(\mathbf{x}) d\mathbf{x} \right] \end{aligned}$$

- They do **not** estimate  $p_{\text{tr}}(\mathbf{x}), p_{\text{te}}(\mathbf{x})$ , but  $\frac{p_{\text{te}}(\mathbf{x})}{p_{\text{tr}}(\mathbf{x})}$  **directly!**

## 1. Importance weight estimation

(e.g., least-squares importance fitting):

Kanamori+  
(JMLR2009)

$$\hat{w} = \operatorname{argmin}_w \hat{\mathbb{E}}_{p_{\text{tr}}(\mathbf{x})} \left[ \left( w(\mathbf{x}) - \frac{p_{\text{te}}(\mathbf{x})}{p_{\text{tr}}(\mathbf{x})} \right)^2 \right]$$

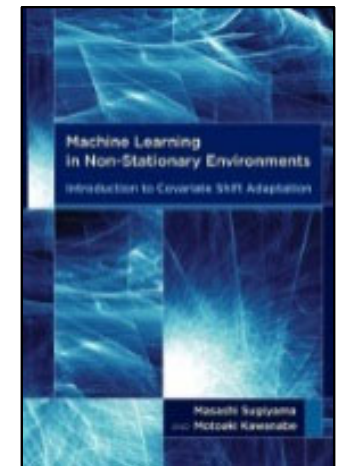
## 2. Weighted predictor training:

$$\hat{f} = \operatorname{argmin}_f \hat{\mathbb{E}}_{p_{\text{tr}}(\mathbf{x}, y)} [\hat{w}(\mathbf{x}) \ell(f(\mathbf{x}), y)]$$

- However, estimation error in Step 1 is not taken into account in Step 2.

- We want to integrate these two steps!

Sugiyama & Kawanabe  
(MIT Press 2012)



# Joint Weight-Predictor Optimization 31

Zhang+ (ACML2020, SNCS2021)

- **Given:** Labeled training data and unlabeled test data

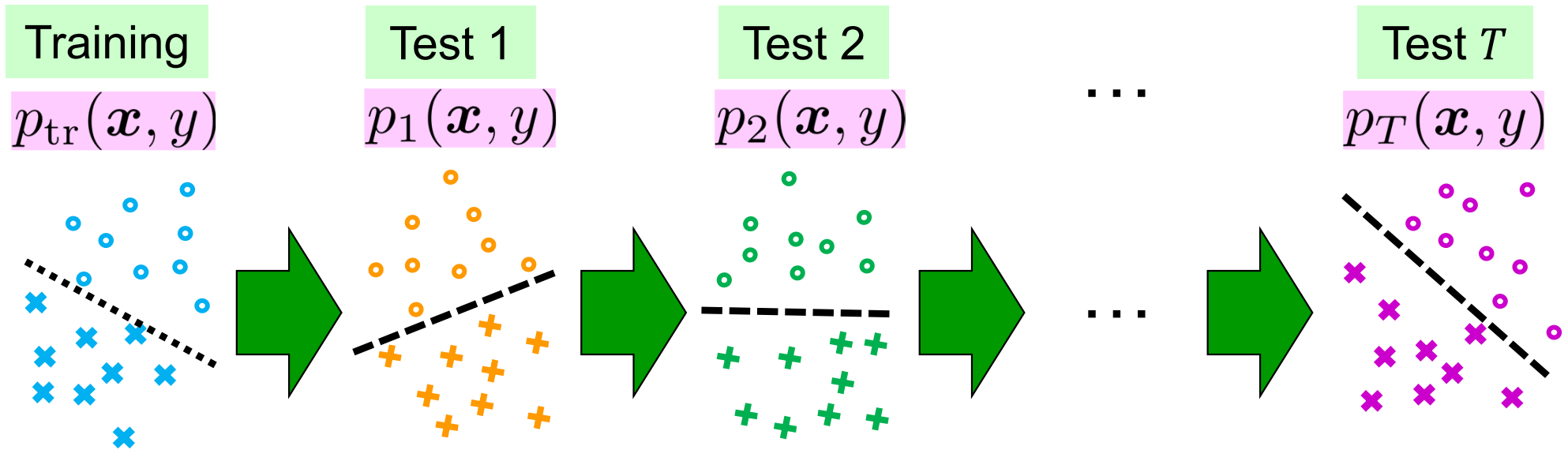
$$\{(\mathbf{x}_i^{\text{tr}}, y_i^{\text{tr}})\}_{i=1}^{n_{\text{tr}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{tr}}(\mathbf{x}, y) \quad \{\mathbf{x}_j^{\text{te}}\}_{j=1}^{n_{\text{te}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{te}}(\mathbf{x})$$

- **Joint minimization of a risk upper bound:**

$$\min_{w \geq 0, f \in \mathcal{F}} J_{\ell'}(w, f) \quad \frac{1}{2} R_{\ell}(f)^2 \leq J_{\ell'}(w, f) \quad \ell \leq 1, \ell' \geq \ell$$
$$R_{\ell}(f) = \mathbb{E}_{p_{\text{te}}(\mathbf{x}, y)}[\ell(f(\mathbf{x}), y)]$$

$$J_{\ell'}(w, f) = \mathbb{E}_{p_{\text{tr}}(\mathbf{x})} \left[ \left( w(\mathbf{x}) - \frac{p_{\text{te}}(\mathbf{x})}{p_{\text{tr}}(\mathbf{x})} \right)^2 \right] \leftarrow \text{1st step}$$
$$+ \left( \mathbb{E}_{p_{\text{tr}}(\mathbf{x}, y)} [w(\mathbf{x}) \ell'(f(\mathbf{x}), y)] \right)^2 \leftarrow \text{2nd step}$$

- Classic approach corresponds to 2-step minimization.



■ **Continuous label shift:** Bai+ (NeurIPS2022)

- Only class-prior  $p_t(y)$  changes.

■ **Continuous covariate shift:** Zhang+ (arXiv2023)

- Only input density  $p_t(\mathbf{x})$  changes.

■ Without knowing the shift intensity, our methods achieve the **same dynamic regret** as the case with known shift intensity.

$$\mathbb{E} \left[ \sum_{t=1}^T R_t(f_t) - \sum_{t=1}^T \min_{f \in \mathcal{F}} R_t(f) \right]$$





# Contents

33

1. Weakly Supervised Learning
2. Noisy-Label Learning
3. Transfer Learning
4. Towards More Reliable Learning

# Joint Shift

- Many distribution shift works focus on a particular **shift type** (e.g., covariate shift):

$$p_{\text{tr}}(\mathbf{x}) \neq p_{\text{te}}(\mathbf{x}) \quad p_{\text{tr}}(y|\mathbf{x}) = p_{\text{te}}(y|\mathbf{x})$$

- However, **identification** of the shift type is challenging.

- Label noise** is also a type of distribution shift:

$$p_{\text{tr}}(\bar{y}|\mathbf{x}) = \sum_y \underbrace{p(\bar{y}|y, \mathbf{x})}_{\text{Noise transition}} p_{\text{te}}(y|\mathbf{x})$$

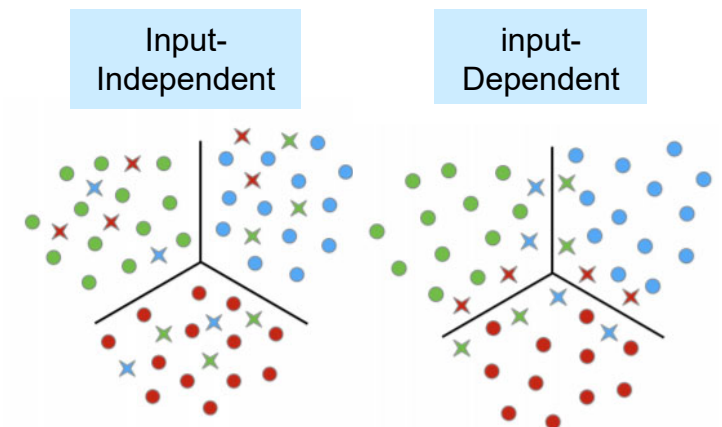
$\bar{y}$  : Noisy class label

Noise transition

- Nice theory for input-independent noise.
- But **input-dependent noise** is hard.

- Let's consider **joint shift**:

$$p_{\text{tr}}(\mathbf{x}, y) \neq p_{\text{te}}(\mathbf{x}, y)$$



# Mini-Batch-Wise Loss Matching

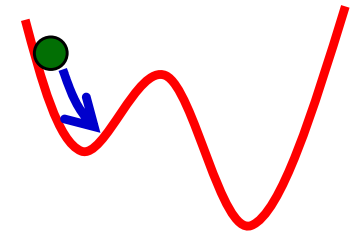
35

## Given:

Fang+ (NeurIPS2020)

- (Large) labeled training data:  $\{(\mathbf{x}_i^{\text{tr}}, y_i^{\text{tr}})\}_{i=1}^{n_{\text{tr}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{tr}}(\mathbf{x}, y)$
- (Small) **labeled** test data:  $\{(\mathbf{x}_j^{\text{te}}, y_j^{\text{te}})\}_{j=1}^{n_{\text{te}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{te}}(\mathbf{x}, y)$

- We try to learn the importance weight **dynamically** in the **mini-batch-wise** manner.



$$f \leftarrow f - \eta \nabla \hat{R}(f) \quad \eta > 0 : \text{step size}$$

- For **each mini-batch**  $\{(\tilde{\mathbf{x}}_i^{\text{tr}}, \tilde{y}_i^{\text{tr}})\}_{i=1}^{\tilde{n}_{\text{tr}}}, \{(\tilde{\mathbf{x}}_j^{\text{te}}, \tilde{y}_j^{\text{te}})\}_{j=1}^{\tilde{n}_{\text{te}}}$ , importance weights are estimated by **kernel mean matching** for **loss values**:

Huang+ (NeurIPS2006)

$$\frac{1}{\tilde{n}_{\text{tr}}} \sum_{i=1}^{\tilde{n}_{\text{tr}}} r_i \ell(f(\tilde{\mathbf{x}}_i^{\text{tr}}), \tilde{y}_i^{\text{tr}}) \approx \frac{1}{\tilde{n}_{\text{te}}} \sum_{j=1}^{\tilde{n}_{\text{te}}} \ell(f(\tilde{\mathbf{x}}_j^{\text{te}}), \tilde{y}_j^{\text{te}}) \quad r_i \approx \frac{p_{\text{te}}(\tilde{\mathbf{x}}_i^{\text{tr}}, \tilde{y}_i^{\text{tr}})}{p_{\text{tr}}(\tilde{\mathbf{x}}_i^{\text{tr}}, \tilde{y}_i^{\text{tr}})}$$

# Current Challenges

■ For joint shift adaptation, requiring **labeled test data** is too strong.

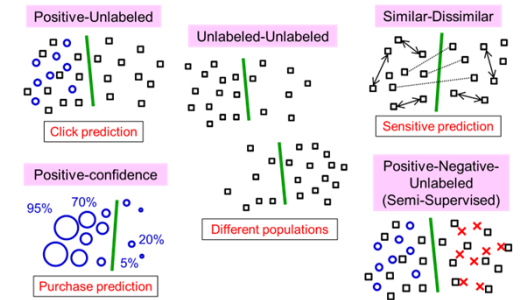
- Can we use **weakly supervised learning**?

■ Importance weighting requires the test domain to be included in the training domain.

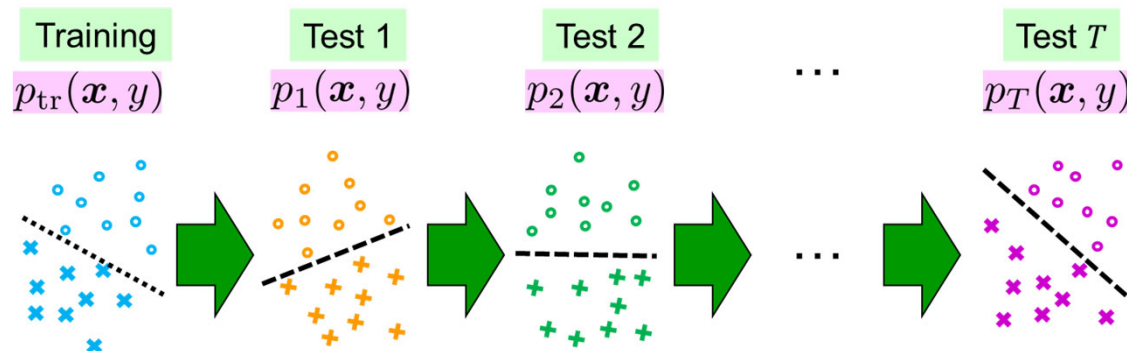
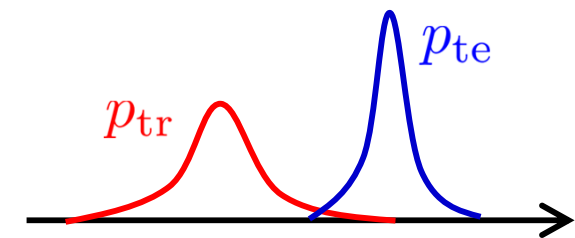
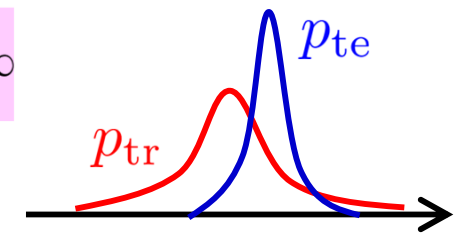
- Can we properly handle **out-of-training-domain test data**?

■ Can we handle **continuous joint shift**?

Weakly Supervised Classification (Binary)



$$\frac{p_{te}(\mathbf{x}, y)}{p_{tr}(\mathbf{x}, y)} < \infty$$



# Grateful to Collaborators!



Team leader Masashi Sugiyama	Research scientist Gang Niu
Postdoctoral researcher Jingfeng Zhang	Postdoctoral researcher Jiaqi Lyu
Postdoctoral researcher Shuo Chen	Senior visiting scientist Shinichi Nakajima
Visiting scientist Futoshi Futami	Visiting scientist Florian Yger
Visiting scientist Takashi Ishida	Visiting scientist Miao Xu
Visiting scientist Takayuki Osa	Visiting scientist Bo Han
Visiting scientist Takahiro Mimori	Visiting scientist Feng Liu
Visiting scientist Lei Feng	Visiting scientist Tongliang Liu
Part-time worker I Masahiro Fujisawa	Part-time worker I Yifan Zhang

and many great interns!

- Professor
  - [Masashi Sugiyama](#) (Complexity, Computer, Information, RIKEN)
- Associate Professor
  - [Naoto Yokoya](#) (Complexity, Computer, Information, RIKEN)
- Lecturer
  - [Takashi Ishida](#) (Complexity, Computer, Information)
- Project Lecturer
  - [Nobutaka Ito](#) (Complexity)
- Professor (to [Sato Lab](#) from April 2022)
  - [Issei Sato](#) (Computer, Informati)
- Project Assistant Professor
  - Chao-Kai Chiang (Complexity)
- Project Researcher (Postdoctoral Rese)
  - [Dongxian Wu](#) (Complexity)
- Project Specialist
  - Yuko Kawashima (Complexity)
  - Soma Yokoi (Complexity)
  - Fumi Sato (Complexity)
- Doctoral Student
  - Shinji Nakadai (Computer)
  - Ryuichi Kiryo (Computer)
  - [Jongyeong Lee](#) (Computer)
  - Tianyi Zhang (Complexity)
  - [Yivan Zhang](#) (Computer)
  - Riou Charles (Computer)
  - [Valliappa Chockalingam](#) (Comput
  - Tongtong Fang (Complexity)
  - Boyo Chen (Complexity)
  - Xiaoyu Dong (Complexity)
  - Yujie Zhang (Complexity)
  - [Xinqiang Cai](#) (Complexity)
  - Jian Song (Complexity)
  - Wanshui Gan (Complexity)
  - Yuting Tang (Complexity)
  - Shintaro Nakamura (Complexity)
  - Or Raveh (Complexity)
  - [Johannes Ackermann](#) (Computer
  - [Wei Wang](#) (Complexity)
  - Hongruixuan Chen (Complexity)
  - Huanjian Zhou (Complexity)
  - Zhiyuan Zhan (Complexity)
  - Zhihao Liu (Complexity)
- Master Student
  - Hyunggyu Park (Complexity)\* [Sato lab.](#)
  - Jiahuan Li (Computer)
  - Kun Yang (Complexity)
  - Xiaomou Hou (Complexity)
  - Anan Methasate (Computer)
  - Cemal Erat (Computer)
  - Kento Yamamoto (Computer)
  - Kazuki Ota (Computer)
  - Iu Yahiro (Computer)
  - Hikaru Fujita (Computer)
  - Yu Yao (Complexity)
  - Yoshifumi Nakano (Complexity)
  - Soichiro Nishimori (Complexity)
  - Ryota Ushio (Complexity)
  - Tiankui Xian (Complexity)
  - Thanawat Lodkaew (Computer)
  - Masahiro Negishi (Computer)
  - Yuto Nozaki (Computer)
  - Kanta Shimizu (Computer)
  - Zhongle Zhu (Computer)
  - Fang Liu (Computer)
  - Ming Li (Complexity)
  - Luheng Wang (Complexity)
  - Liuzhuozheng Li (Complexity)
- Research Student
  - Meike Tütken (Computer)
  - Serhii Khomenko (Information Science)
  - Artem Lubkivskyi (Information Science)

