

# Recent Advances in Robust Machine Learning



Masashi Sugiyama

RIKEN Center for Advanced Intelligence Project/  
The University of Tokyo



<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

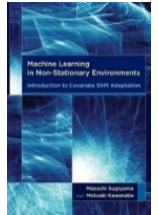
## ■ Interests: Machine learning (ML)

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

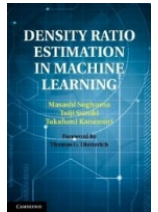
## ■ Academic activities:

- Program Chairs for NeurIPS2015, AISTATS2019, ACML2010/2020...

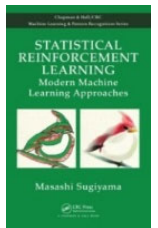
Sugiyama & Kawanabe, **Machine Learning in Non-Stationary Environments**, MIT Press, 2012



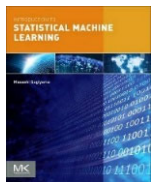
Sugiyama, Suzuki & Kanamori, **Density Ratio Estimation in Machine Learning**, Cambridge University Press, 2012



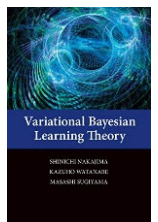
Sugiyama, **Statistical Reinforcement Learning**, Chapman and Hall/CRC, 2015



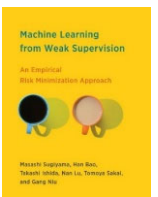
Sugiyama, **Introduction to Statistical Machine Learning**, Morgan Kaufmann, 2015



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.





# Contents

1. Introduction of RIKEN-AIP
2. Robust Machine Learning
  - A) Weakly Supervised Learning
  - B) Transfer Learning
  - C) Noise-Robust Learning
3. Summary

# 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



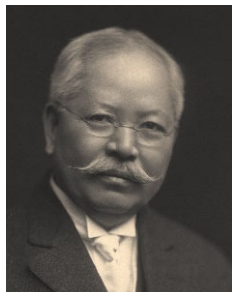
Suzuki Umetaro



Yukawa Hideki



Noyori Ryoji



Takamine Jokichi

**1917–**  
RIKEN  
private  
foundation

**1948–**  
Scientific  
Research  
Institute Ltd.  
(KAKEN)

**1958–**  
RIKEN  
public  
corporation

**2003–**  
RIKEN  
Independent  
Administrative  
Institution

**2015–**  
RIKEN  
National Research  
and  
Development  
Institute

**RIKEN Konzern**  
In 1939, 63 companies  
and 121 plants

1967  
moved to Wako

**2017** The centennial celebration

- 1984 Tsukuba
- 1990 Sendai
- 1993 Nagoya
- 1997 Harima
- 1998 Yokohama
- 2000 Kobe
- 2011 Osaka
- 2016 Keihanna

**2017**  
Wako: 50<sup>th</sup> Anniversary  
Harima: 20<sup>th</sup> Anniversary

# Office and Research

## Headquarters

RIKEN Information R&D and Strategy Headquarters



Dr. Mino Michihiko

RIKEN Cluster for Pioneering Research



Dr. Koyasu Shigeo

RIKEN Cluster for Science, Technology and Innovation Hub



Director  
Dr. Koyasu Shigeo



Gonokami Makoto

## Executive Directors



Dr. Koyasu Shigeo



Dr. Miyazono Kohei



Mr. Kagaya Satoru



Dr. Naka Makiko



Mr. Matsuo Hiromichi

Program for Drug Discovery and Medical Technology Platforms



Mr. Matsuo Hiromichi

RIKEN Baton Zone Program/RIKEN Industrial Co-Creation Program



Spring-8 Center



Physics

Dr. Ishikawa Tetsuya

Life

Harima  
Kobe

Keihan'na

Sendai

Wako

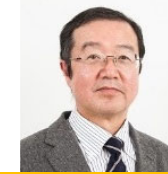
Interdisciplinary Theoretical and Mathematical Sciences Program

Physics

Yasunobu

Physics

Center for Advanced Photonics



Katsumi

Physics

Center for Brain Science



Ryuichiro

Physics

Physics

Life

Center for Computational Science



Informatics

Dr. Matsuoka Satoshi

Center for Biosystems Dynamics Research



Life

Dr. Nishida Eisuke

Center for Integrative Medical Sciences



Life

Dr. Yamamoto Kazuhiko

Center for Sustainable Resource Science



Life

Dr. Saito Kazuki

Center for Advanced Intelligence Project



Informatics

Dr. Sugiyama Masashi

BioResource Research Center



Life

Dr. Sniroishi Toshihiko

<Overseas>  
RAL-RIKEN (UK)  
BNL-RIKEN (USA)  
Beijing Representative Office (China)  
Singapore Representative Office (Singapore)  
European Representative Office (Belgium)

Osaka

Osaka

Yokohama

Tokyo

Tsukuba

# International Researchers

## Latin America 16, 2%

Argentine, Brazil, Columbia, Costa Rica, Mexico, Peru

## Europe 190, 24%

Austria, Belarus, Belgium, Bulgaria, Croatia, Denmark, Finland, France, Germany, Greece, Hungary, Ireland, Italy, Moldova, Netherlands, Poland, Romania, Russia, Slovakia, Slovenia, Spain, Sweden, Switzerland, Ukraine, United Kingdom, Uzbekistan

## Oceania 15, 2%

Australia, Fiji, New Zealand

## Africa 18, 2%

Algeria, Cameroon, Egypt, Ghana, Tunisia, South Africa, Senegal, Zambia

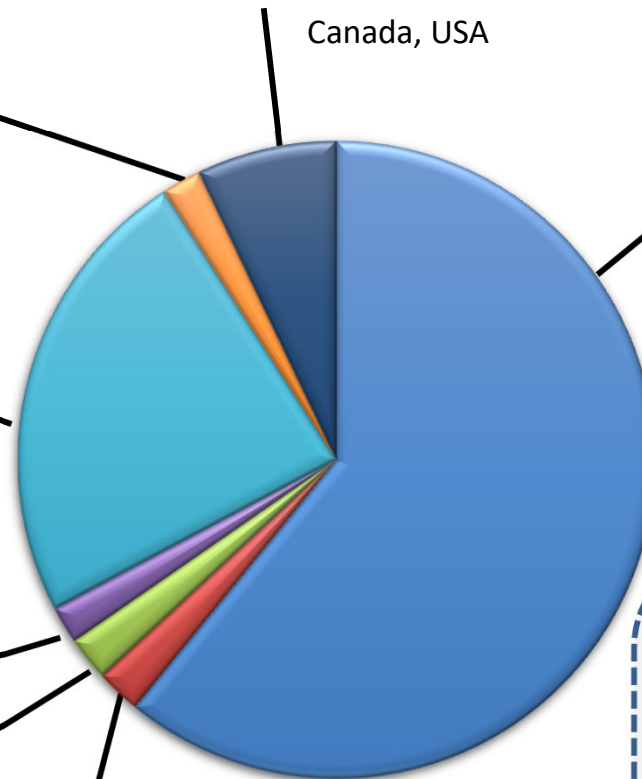
## Middle East 18, 2%

Iran, Israel, Lebanon, Turkey, Yemen

Canada, USA

## Asia 489, 61%

Bangladesh, China, Hong Kong, India, Indonesia, Korea, Malaysia, Nepal, Pakistan, Philippines, Singapore, Sri Lanka, Taiwan, Thailand, Vietnam



### In order of representation at RIKEN

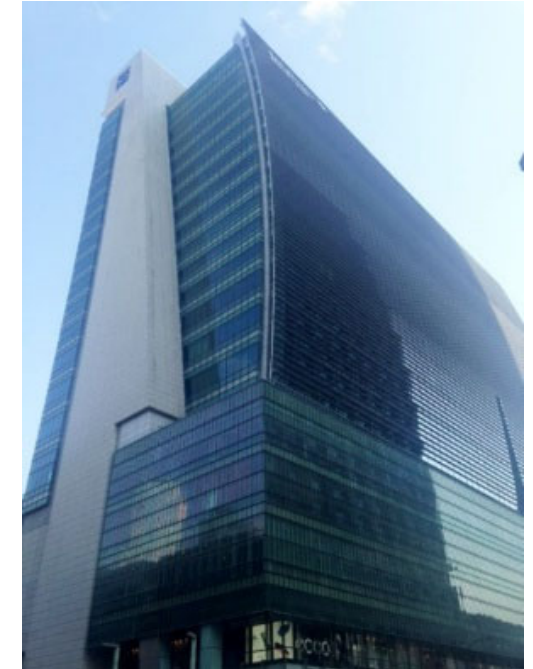
1. China	260
2. India	58
3. Korea	54
4. USA	43
5. Germany	37
6. UK	28
7. Taiwan	27
France	27
9. Indonesia	17
10. Russia	15

**TOTAL 803 people, 65 countries and regions**

# What is RIKEN-AIP?

- RIKEN founded **Center for Advanced Intelligence Project (AIP)** in 2016, under Ministry of Education, Culture, Sports, Science and Technology (MEXT).

Main office located in the heart of Tokyo



Distributed office across Japan



In-house GPU servers



Open discussion space





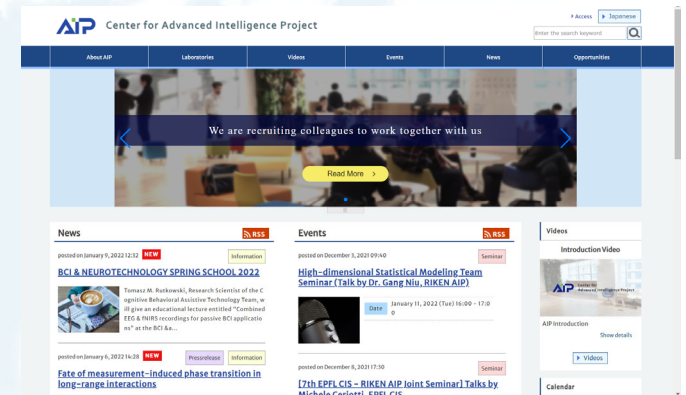
As of Apr. 1, 2022

## ■ Diverse research staffs:

- 130 employed researchers (36% international, 23% female)
- 210 visiting researchers
- 100 domestic students
- 140 international interns (total)

## ■ Extensive collaboration:

- 40+ international collaboration partners
- 40+ industry projects



# AIP's 5 Missions

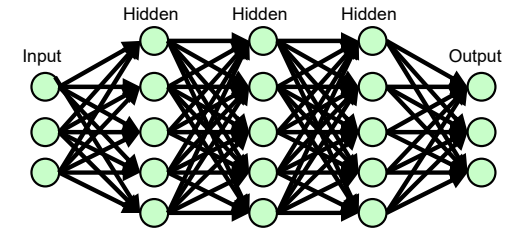
- **Develop next-generation AI technology:**
  - machine learning & optimization theory, etc.
- **Accelerate scientific research:**
  - cancer, material, genomics, etc.
- **Solve socially critical problems:**
  - natural disaster, elderly healthcare, etc.
- **Study of ethical, legal and social issues of AI:**
  - ethical guidelines, personal data, etc.
- **Human resource development:**
  - researchers, engineers, etc.



# Developing New AI Technology

## Theory of deep learning:

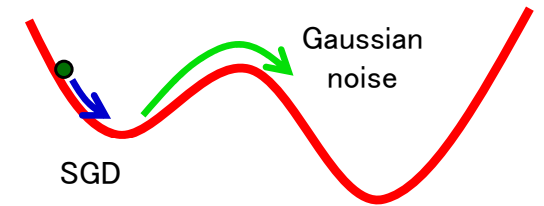
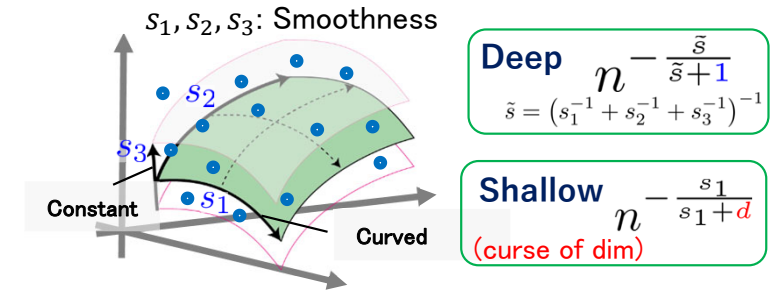
- Better prediction than shallow learning
- No curse of dimensionality
- Global optimization



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

## 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!**

**Positive-Unlabeled (PU)**  
(ex: click prediction)



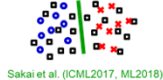
**Positive-confidence (Pconf)**  
(ex: purchase prediction)



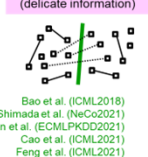
**Unlabeled-Unlabeled (UU)**  
(learning from different populations)



**Semi-Supervised (PU+PN)**  
(first theoretically guaranteed method)



**Similar-Dissimilar (SD)**  
(delicate information)



du Plessis et al. (TAAI2013), Lu et al. (ICLR2019, AISTATS2020), Charoenphakdee et al. (ICML2019), Lei et al. (ICML2021)

### Noise Transition Correction

Noise transition matrix  $T$ :

$$T^T = \begin{bmatrix} 1 & 0.1 & 0.5 \\ 0 & 0.8 & 0.5 \\ 0 & 0.1 & 0 \end{bmatrix}$$

- Clean-to-noisy flipping probability.

Major approaches: Patirini et al. (CVPR2017)

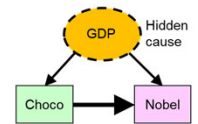
- **Loss correction** by  $T^{-1}$  to eliminate noise.
- **Classifier adjustment** by  $T^T$  to simulate noise.

We want to estimate  $T$  **only from noisy data:**

- Use human cognition as a "mask" for  $T$ . Han et al. (NeurIPS2018)
- Learn  $T$  and a classifier dynamically. Xia et al. (NeurIPS2019)
- Decompose  $T$  into simpler components. Yao et al. (NeurIPS2020)
- Regularize  $T$  to be estimable. Zhang et al. (ICML2021), Li et al. (ICML2021)
- Extension to input-dependent noise  $T(x)$ . Xia et al. (NeurIPS2020), Berthon et al. (ICML2021)

### 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.

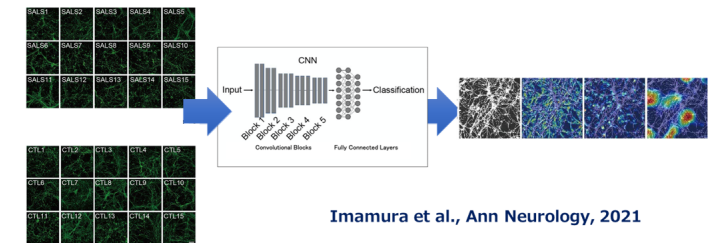
# Accelerating Scientific Research

## Medical science:

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



\*Yoichiro Yamamoto, et al. Automated acquisition of explainable knowledge from unannotated histopathology images. *Nature Communications* 10:5642, 2019.



Imamura et al., *Ann Neurology*, 2021

## Material science:

- Database creation with text mining

1: Si High thermoelectric properties of a-type Ag<sub>10</sub>Sn<sub>2</sub>Li<sub>2</sub>Pan<sup>+</sup>, David Brath<sup>†</sup> and Nita Dogue Supporting information EXPERIMENTAL SECTION Synthesis

2: Ag<sub>10</sub>-aPnSe<sub>2</sub> (x = 0.01, 0.02, 0.025, 0.03 and 0.06) and Ag<sub>10</sub>-nNnSe<sub>2</sub> (x = 0 to 0.06) samples were synthesized by one step solid state reaction.

3: Stoichiometric mixtures of Ag (99.99%), Bi (99.99%), Sb (99.99%), Ni (99.99%), Pb (99.99%) and Pt (99.99%) powders were ground by hand using agate mortar and pressed into pellets under uniaxial stress (250MPa), then heated at 773K for 12h in silica tubes sealed under argon.

4: The obtained samples were ground and subsequently densified by using a spark plasma sintering (SPS) system (SPS-515) at 773K with holding time of 10 min in a 1x20x15 mm graphite mold under an axial compressive stress of 100 MPa in an argon atmosphere.

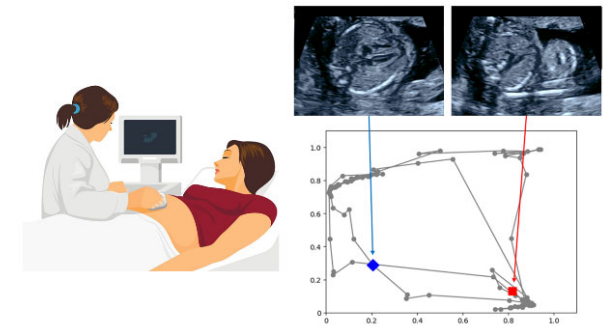
5: X-ray diffraction and Electron microscopy X-ray diffraction characterization was performed using a Panalytical X'Pert diffractometer by using a Cu K $\alpha$  radiation obtained using a Ge(111) incident monochromator and an X'celerator detector.

6: Rietveld refinement was performed to calculate lattice parameters of samples using FULLPROF software. Scanning electron microscopy (SEM) analyses were performed using a Hitachi S-3600N VP-SEM, and energy-dispersive X-ray spectroscopy (EDX) was performed using an EDAX-TSL spectrometer.

```

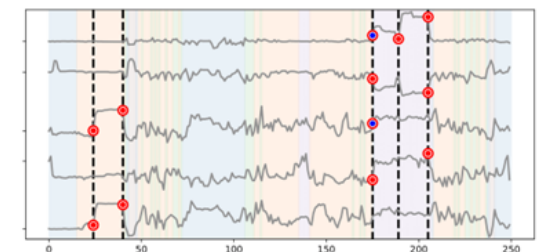
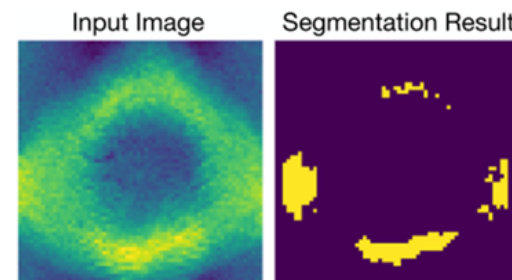
    graph TD
      A[Ag] --> B[hand]
      B --> C[grinded]
      C --> D[250 MPa]
      D --> E[heated]
      E --> F[12h]
      F --> G[773K]
      G --> H[silica tubes]
      H --> I[ground]
      I --> J[spark plasma]
      J --> K[10 min]
      K --> L[100 MPa]
      L --> M[argon]
      M --> N[diffraction]
      N --> O[SEM]
      O --> P[EDX]
      
```

Relationship graph showing values on the x-axis (0 to 1.0) and y-axis (0.0 to 1.0). The graph features several data series with different markers and colors, showing various trends and peaks.



## Data-driven science:

- Selective inference for reliability evaluation



# Solving Socially Critical Problems

## Natural disaster:

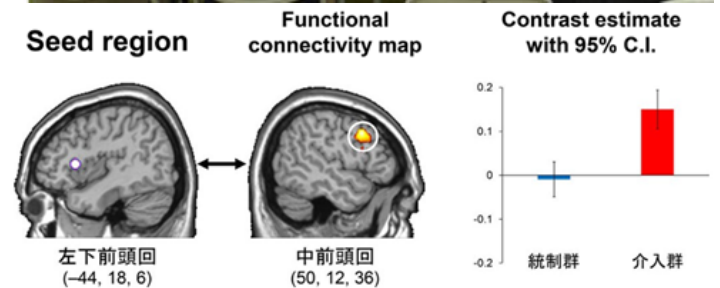
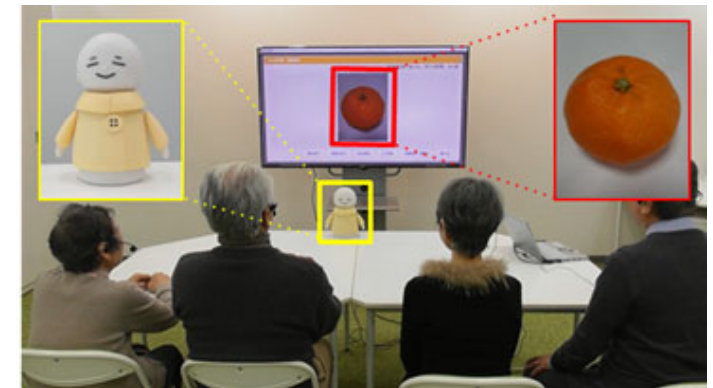
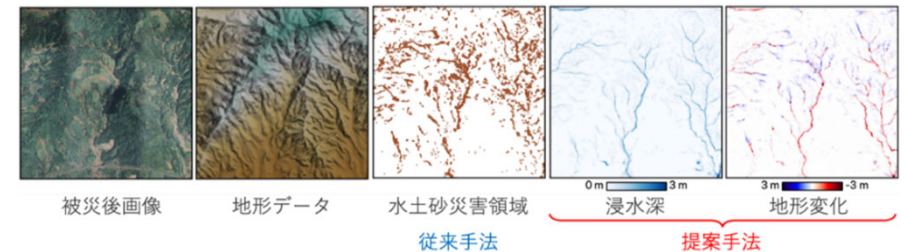
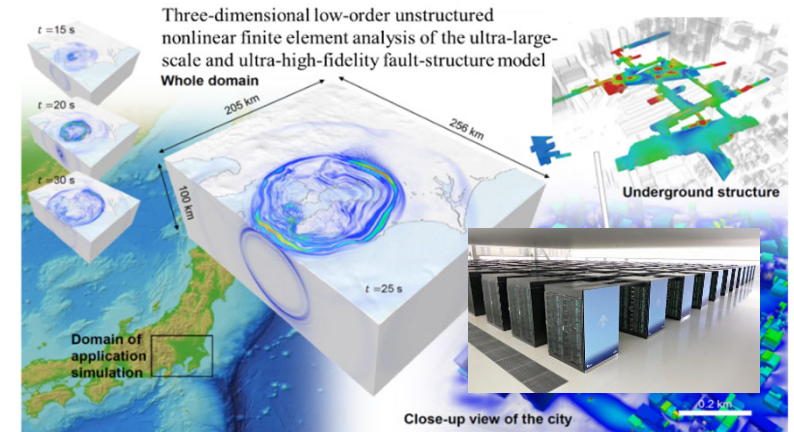
- Fugaku-based earthquake simulation
- Remote sensing disaster analysis

## Elderly healthcare:

- Chat-robot-guided cognitive function improvement

## Education:

- Automatic essay evaluation
- Interactive essay writing support



問題文

解答

採点基準

問 傍線部(1)「こうした～葉ってきた」とあるが、それはどういふことか。70字以内で説明せよ。

西洋文化の基礎の対決は神対人間(宗教=契約)、人間対自然(科学=合理主義)、人間対人間(個人主義)という形で現れるという。日本人のように余計なことを言わないのではなく、他人に分かってもらうために言葉を尽くして対決のスタンスが西洋の文化を帯びてきた。

採点対象

減点対象

模範回答

Smoking **should be banned!** It causes **bad breath** and makes your **clothes stinky!**

While your point is **logically valid**, it **does not justify banning** the substance. Can you provide a more relevant reason?

Oh..I know that **many people die of second-hand smoke** each day!

Indeed, **people dying from second-hand smoke** is much more convincing than **causing bad breath**. *Where did you learn about this information?*

# Studying AI-ELSI

## AI ethical guidelines:

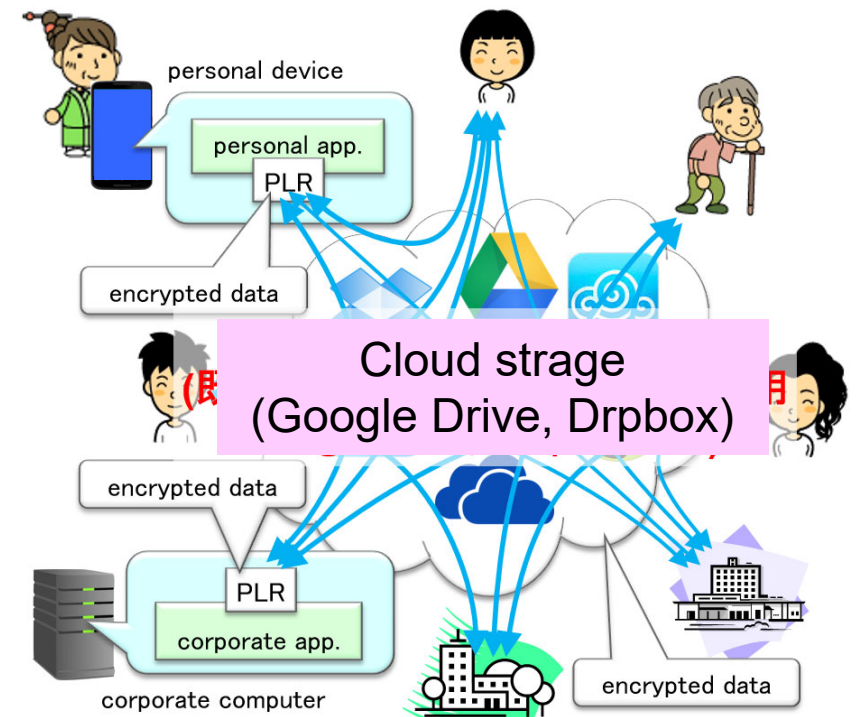
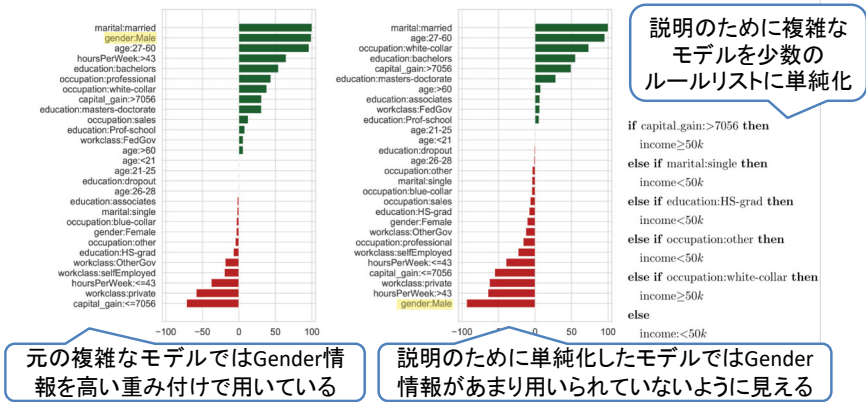
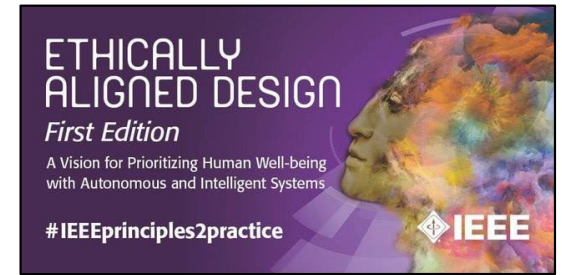
- Japanese Society for AI, Ministry of Internal Affairs and Communications, Cabinet Office
- IEEE, G20, OECD

## Personal data management:

- Individual-based accessibility control system

## AI security and reliability:

- Adversarial attack/defense
- Fairness faking/guarantee





# Contents

1. Introduction of RIKEN-AIP
2. Robust Machine Learning
  - A) Weakly Supervised Learning
  - B) Transfer Learning
  - C) Noise-Robust Learning
3. Summary

- **Goal:** Develop novel ML theories and algorithms that enable reliable learning from limited information.
  - **Insufficient information:** weak supervision.
  - **Data bias:** changing environments, privacy.
  - **Label noise:** human error, sensor error.
  - **Attack:** adversarial noise, distribution shift.



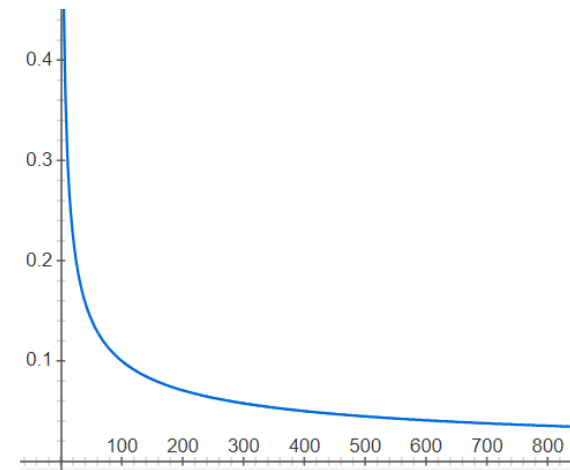
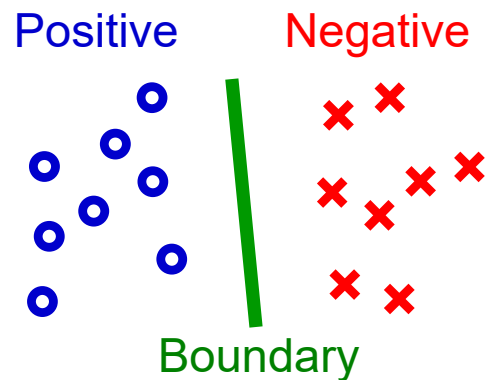


# Contents

1. Introduction of RIKEN-AIP
2. Robust Machine Learning
  - A) Weakly Supervised Learning
  - B) Transfer Learning
  - C) Noise-Robust Learning
3. Summary

- **ML from big labeled data** is successful.
  - Speech, image, language, advertisement,...
  - Estimation error of the boundary decreases in order  $1/\sqrt{n}$  .

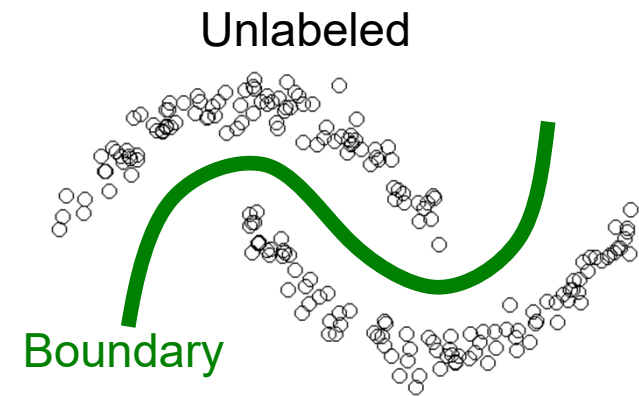
$n$  : Number of labeled samples



- However, there are various applications where **big labeled data is not available**.
  - Medicine, disaster, robots, brain, ...

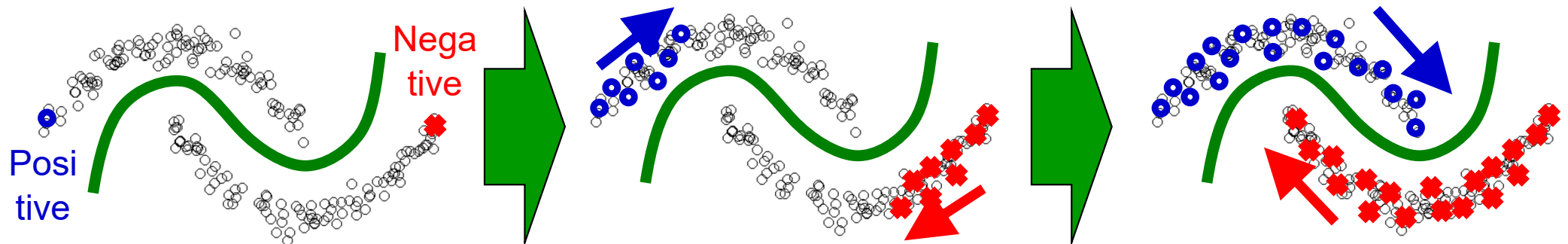
## ■ Unsupervised classification:

- No label is used.
- Essentially clustering.
- No guarantee for prediction.



## ■ Semi-supervised classification:

- Additionally use a small amount of labeled data.
- Propagate labels along clusters.
- No guarantee for prediction.

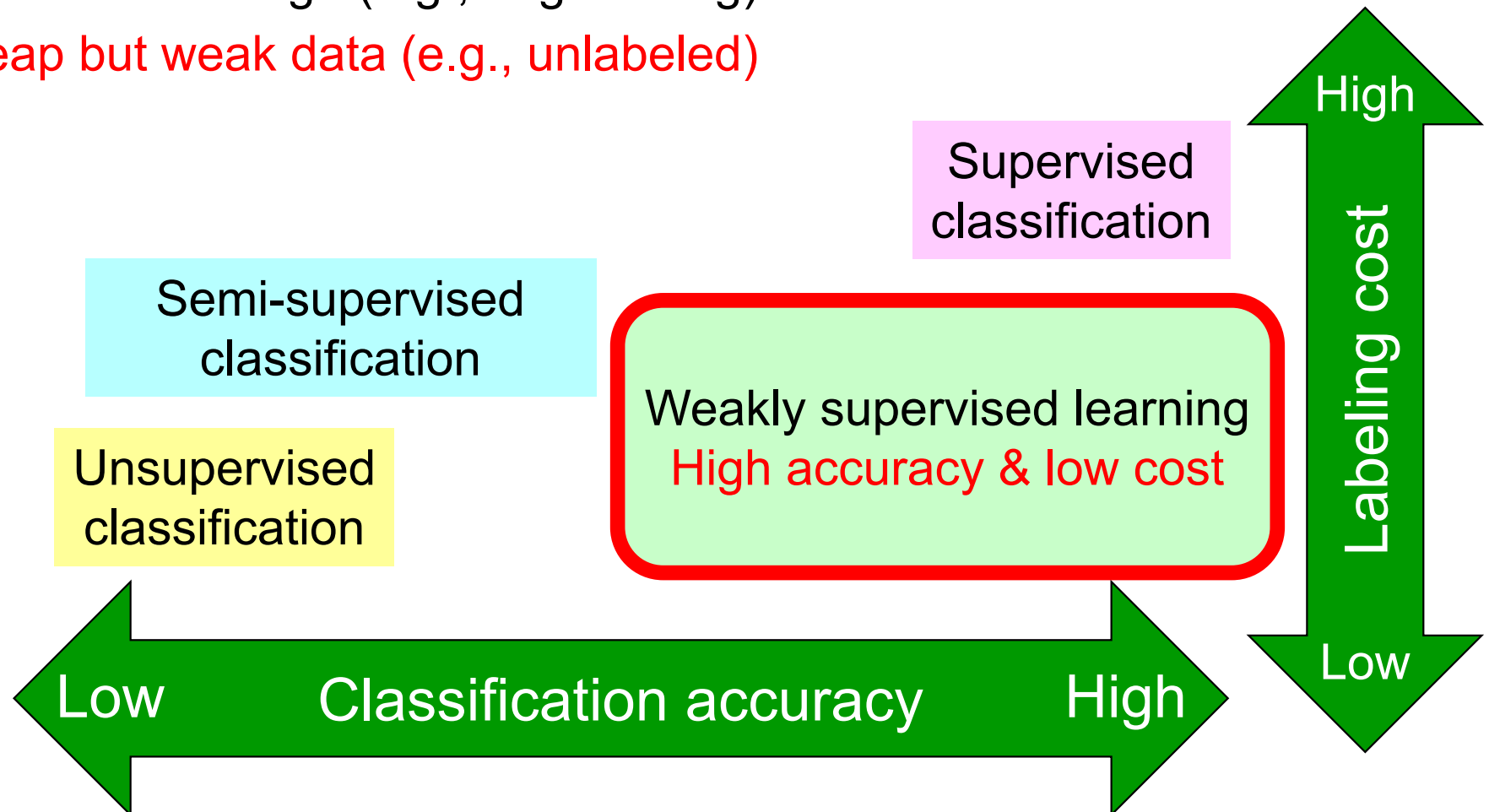


# Weakly Supervised Learning

20

## ■ Coping with labeling cost:

- Improve data collection (e.g., crowdsourcing)
- Use a simulator to generate pseudo data (e.g., physics, chemistry, robotics, etc.)
- Use domain knowledge (e.g., engineering)
- Use cheap but weak data (e.g., unlabeled)



- **Given:** Positive and unlabeled samples

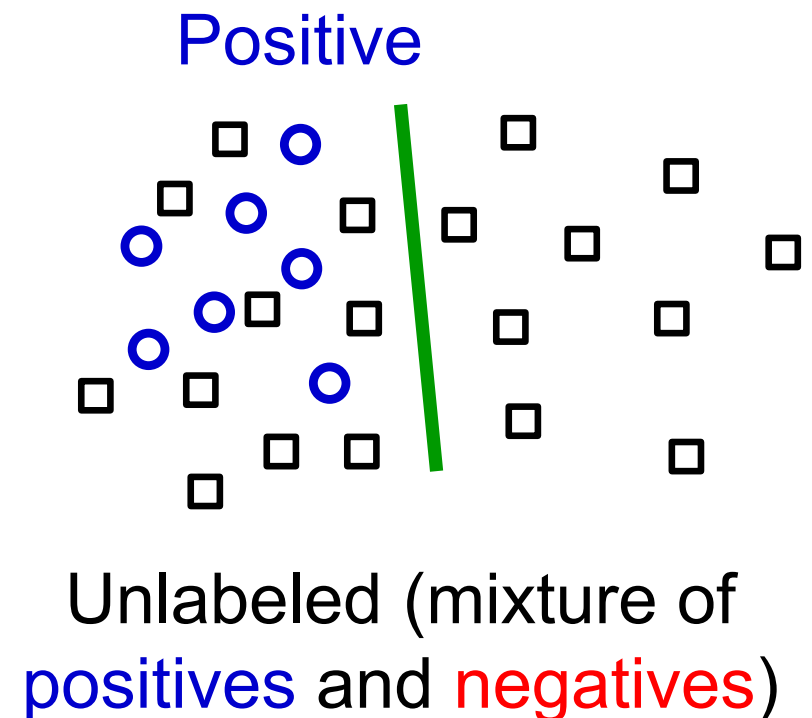
$$\{\mathbf{x}_i^P\}_{i=1}^{n_P} \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}|y = +1)$$

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

- **Goal:** Obtain a PN classifier

Example: Ad-click prediction

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



- **Given:** Positive and unlabeled data

du Plessis, Niu & Sugiyama  
(NIPS2014, ICML2015)

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

- Decomposition of the classification risk:

$$R(f) = \mathbb{E}_{p(\mathbf{x}, y)} \left[ \ell(y f(\mathbf{x})) \right] \quad \ell : \text{loss} \quad \pi = p(y = +1) : \text{Class prior (assumed known)}$$

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

Risk for positive data

Risk for negative data

- Eliminate the expectation over negative data as

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

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

- Unbiased risk estimation:

$\mathcal{O}_p\left(1/\sqrt{n_P} + 1/\sqrt{n_U}\right)$

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

# Positive-Negative-Unlabeled Classification 23 (Semi-Supervised Classification)

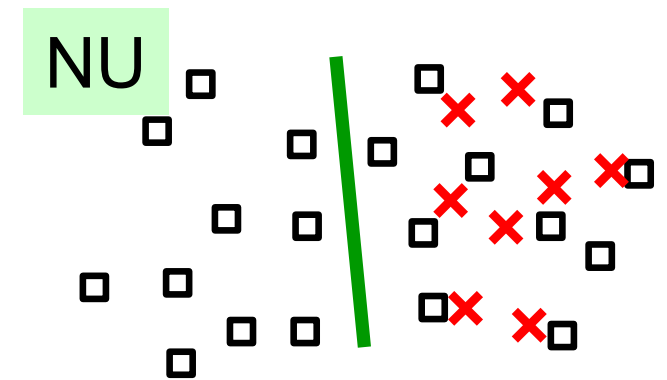
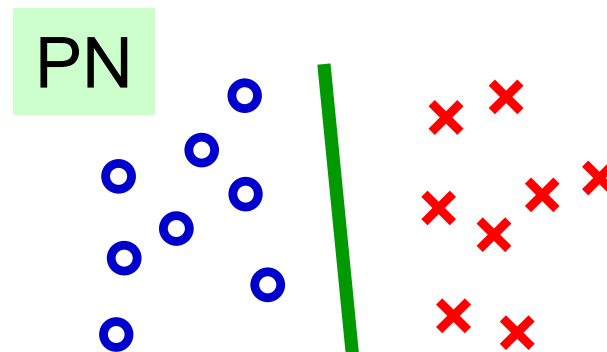
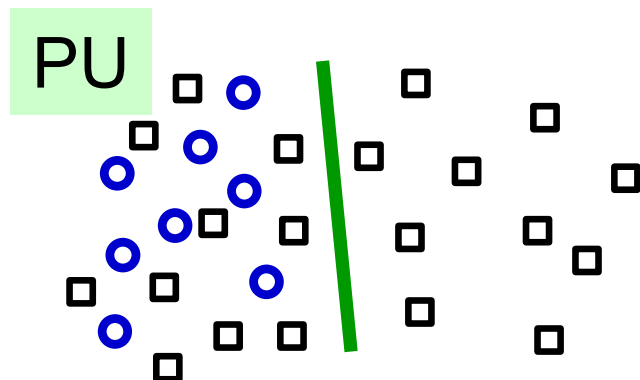
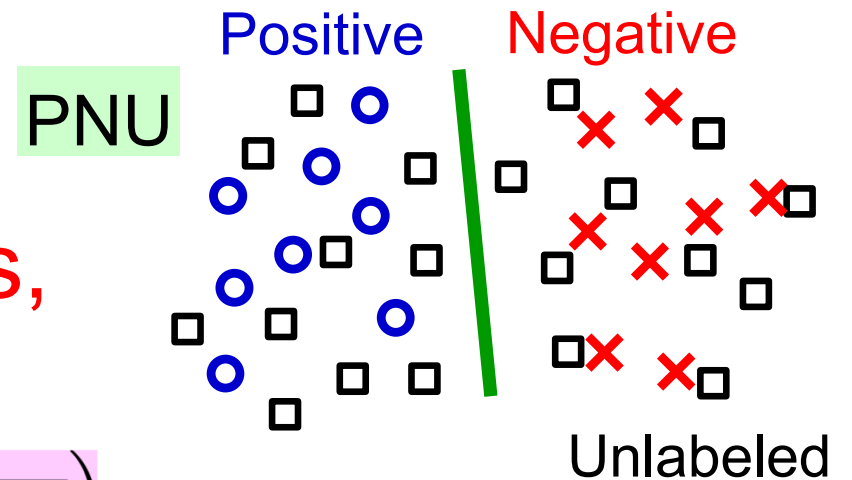
Sakai, du Plessis, Niu & Sugiyama (ICML2017)

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

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

## Without cluster assumptions, PN classifiers are trainable!

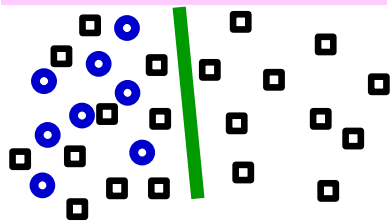
$$\mathcal{O}_p\left(1/\sqrt{n_P} + 1/\sqrt{n_N} + 1/\sqrt{n_U}\right)$$



# Various Extensions

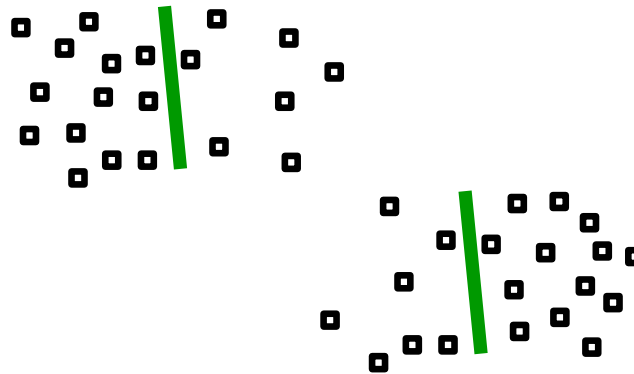
- Learning from weakly supervised data is possible in many different forms!

Positive-Unlabeled



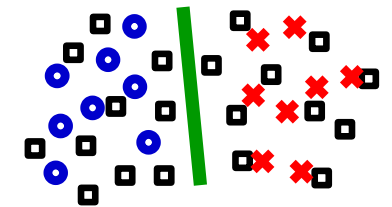
du Plessis et al. (NIPS2014, ICML2015, MLJ2017)  
Niu et al. (NIPS2016), Kiryo et al. (NIPS2017)  
Hsieh et al. (ICML2019)

Unlabeled-Unlabeled



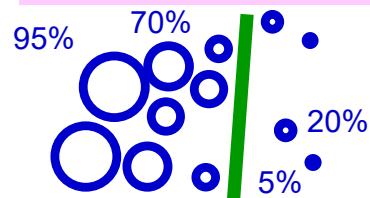
du Plessis et al., (TAAI2013)  
Lu et al. (ICLR2019, AISTATS2020)  
Charoenphakdee et al. (ICML2019)  
Lei et al. (ICML2021)

Semi-Supervised



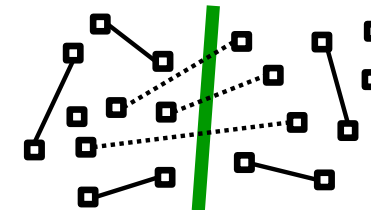
Sakai et al. (ICML2017, ML2018)

Positive-confidence



Ishida et al. (NeurIPS2018)  
Shinoda et al. (IJCAI2021)

Similar-Dissimilar



Bao et al. (ICML2018)  
Shimada et al. (NeCo2021)  
Dan et al. (ECMLPKDD2021)  
Cao et al. (ICML2021)  
Feng et al. (ICML2021)

- All are loss-correction based and consistent.
- Any loss, classifier, and optimizer can be used.

$$\mathcal{O}_p\left(1/\sqrt{n}\right)$$



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

■ **Multi-class weak-labels:**

- **Complementary labels:** Specify a class that a pattern does **not** belong to (“not 1”).

Ishida et al.  
(NIPS2017, ICML2019)  
Chou et al. (ICML2020)

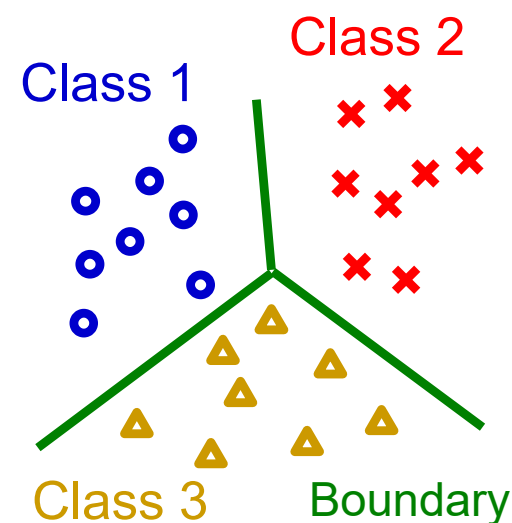
- **Partial labels:** Specify a subset of classes that contains the correct one (“1 or 2”).

Feng et al.  
(ICML2020, NeurIPS2020)  
Lv et al. (ICML2020)

- **Single-class confidence:** One-class data with full confidence (“1 with 60%, 2 with 30%, and 3 with 10%”)

■ **Systematic loss correction is possible!**

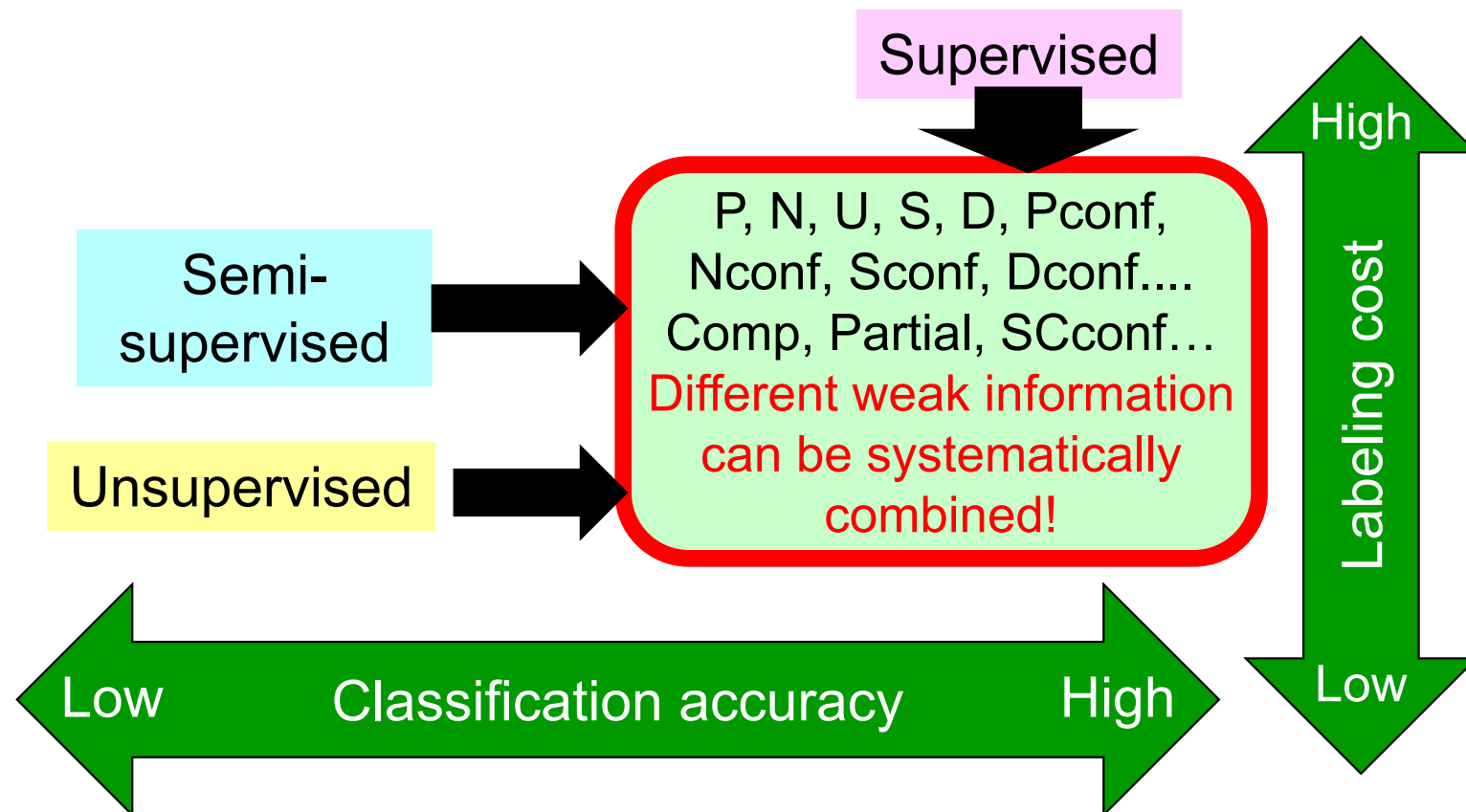
$$\mathcal{O}_p\left(1/\sqrt{n}\right)$$



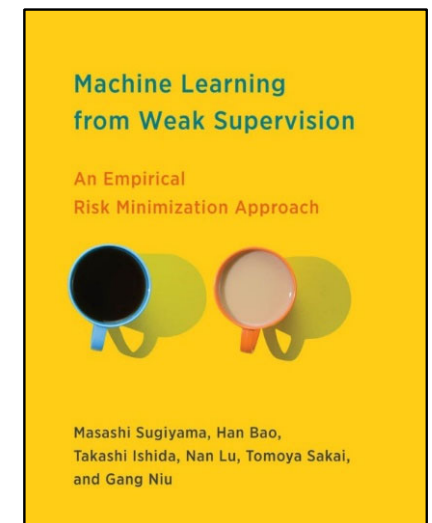
# Summary: Weakly Supervised Learning

26

- We developed an empirical risk minimization framework for weakly supervised learning:
  - Any loss, classifier, and optimizer can be used.
  - Statistical consistency with optimal convergence.



Sugiyama, Bao, Ishida, Lu, Sakai & Niu,  
*Machine Learning from Weak Supervision: An Empirical Risk Minimization Approach.*  
MIT Press, August 2022.





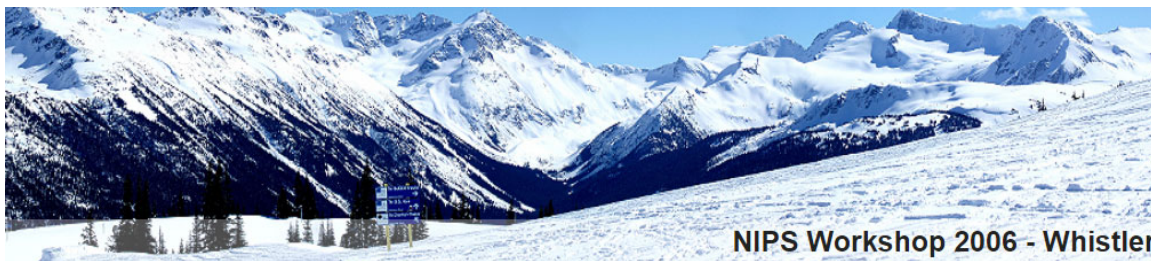
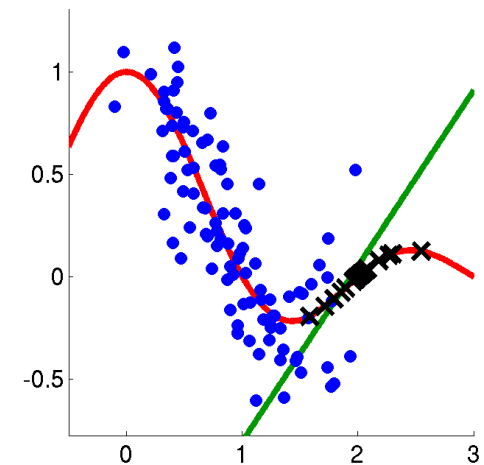
# Contents

1. Introduction of RIKEN-AIP
2. Robust Machine Learning
  - A) Weakly Supervised Learning
  - B) Transfer Learning
  - C) Noise-Robust Learning
3. Summary

# Transfer Learning

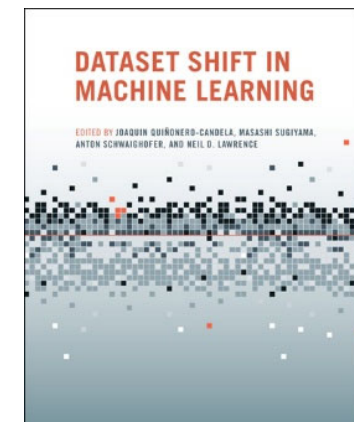
28

- Training and test data often have different distributions, due to
  - changing environments,
  - sample selection bias (privacy).
- **Transfer learning (domain adaptation):**
  - 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 et al. (MIT Press 2009)



# Classical Approach for Transfer Learning

## ■ Two-step adaptation:

1. Importance weight estimation:

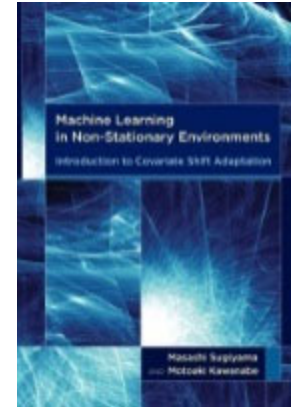
$$\hat{w} = \operatorname{argmin}_w \hat{\mathbb{E}}_{p_{\text{tr}}(\mathbf{x}, y)} \left[ D \left( w(\mathbf{x}, y), \frac{p_{\text{te}}(\mathbf{x}, y)}{p_{\text{tr}}(\mathbf{x}, y)} \right) \right]$$

2. Weighted predictor training:

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

- 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 30

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

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

Shimodaira (JSPI2000)

- Suppose we are given

- Labeled training data:  $\{(\mathbf{x}_i^{\text{tr}}, y_i^{\text{tr}})\}_{i=1}^{n_{\text{tr}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{tr}}(\mathbf{x}, y)$

- Unlabeled test data:  $\{\mathbf{x}_i^{\text{te}}\}_{i=1}^{n_{\text{te}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{te}}(\mathbf{x})$

- Minimize **a risk upper bound** jointly

Zhang et al.  
(ACML2020, SNCS2021)

w.r.t. weight  $w$  and predictor  $f$ :  $J_{\ell_{\text{tr}}}(f, w) \geq R_{\ell_{\text{te}}}(f)^2$

$$\hat{f} = \underset{f}{\operatorname{argmin}} \min_{w \geq 0} \hat{J}_{\ell_{\text{tr}}}(f, w)$$

$$R_{\ell}(f) = \mathbb{E}_{p_{\text{te}}(\mathbf{x}, y)}[\ell(f(\mathbf{x}), y)]$$

$$\ell_{\text{te}} \leq 1, \ell_{\text{tr}} \geq \ell_{\text{te}}$$

$\hat{J}_{\ell}$  : Empirical approximation of  $J_{\ell}$

- **Theoretical guarantee:**

$$R_{\ell_{\text{te}}}(\hat{f}) \leq \sqrt{2} \min_f R_{\ell_{\text{te}}}(f) + \mathcal{O}_p(n_{\text{tr}}^{-1/4} + n_{\text{te}}^{-1/4})$$

# Dynamic Importance Weighting 31

■ General changing distributions:  $p_{\text{tr}}(\mathbf{x}, y) \neq p_{\text{te}}(\mathbf{x}, y)$

■ Suppose we are given

• Labeled training data:  $\{(\mathbf{x}_i^{\text{tr}}, y_i^{\text{tr}})\}_{i=1}^{n_{\text{tr}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{tr}}(\mathbf{x}, y)$

• Labeled test data:  $\{(\mathbf{x}_i^{\text{te}}, y_i^{\text{te}})\}_{i=1}^{n_{\text{te}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{te}}(\mathbf{x}, y)$

■ For **each mini-batch**  $\{(\bar{\mathbf{x}}_i^{\text{tr}}, \bar{y}_i^{\text{tr}})\}_{i=1}^{\bar{n}_{\text{tr}}}, \{(\bar{\mathbf{x}}_i^{\text{te}}, \bar{y}_i^{\text{te}})\}_{i=1}^{\bar{n}_{\text{te}}}$ ,

importance weights are estimated by

Fang et al.  
(NeurIPS2020)

matching **losses** by **kernel mean matching**:

Huang et al. (NeurIPS2007)

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

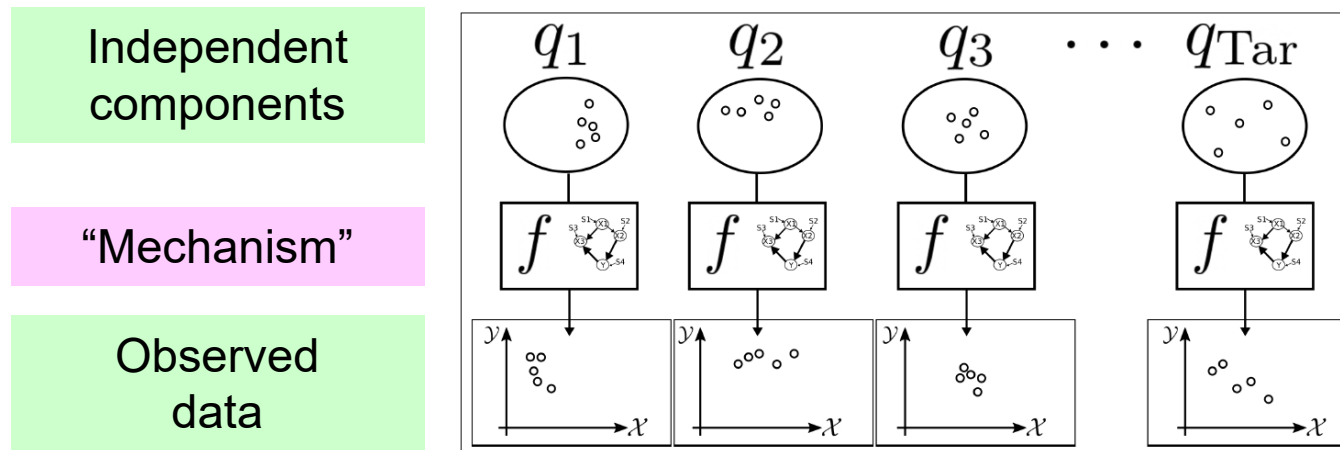
■ **Extremely simple, but highly powerful!**

# Summary

- In transfer learning with importance weighting, simultaneously performing **importance estimation** and **predictor training** is promising.
- What should we do if training and test distributions look very different?

- **Mechanism transfer!**

Teshima, Sato & Sugiyama (ICML2020)



Bai, Zhang, Zhao, Sugiyama & Zhou (NeurIPS2022)

- **Current challenge:** Continuous distribution change



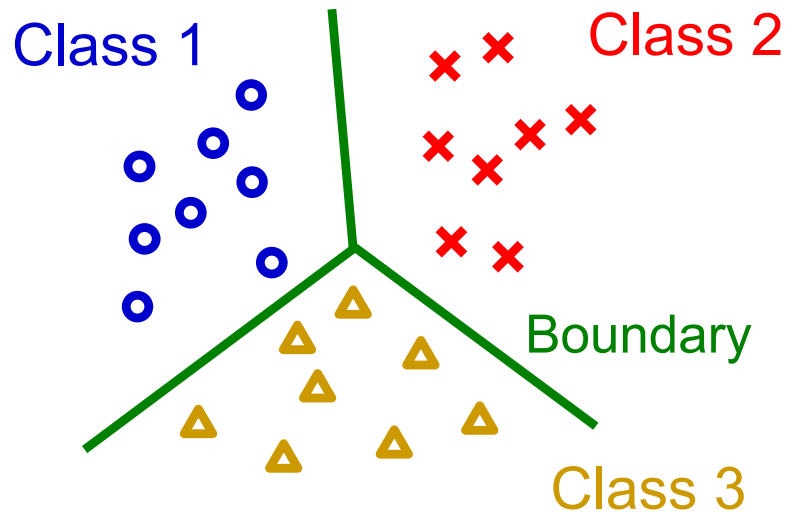


# Contents

1. Introduction of RIKEN-AIP
2. Robust Machine Learning
  - A) Weakly Supervised Learning
  - B) Transfer Learning
  - C) Noise-Robust Learning
3. Summary

# Supervised Classification

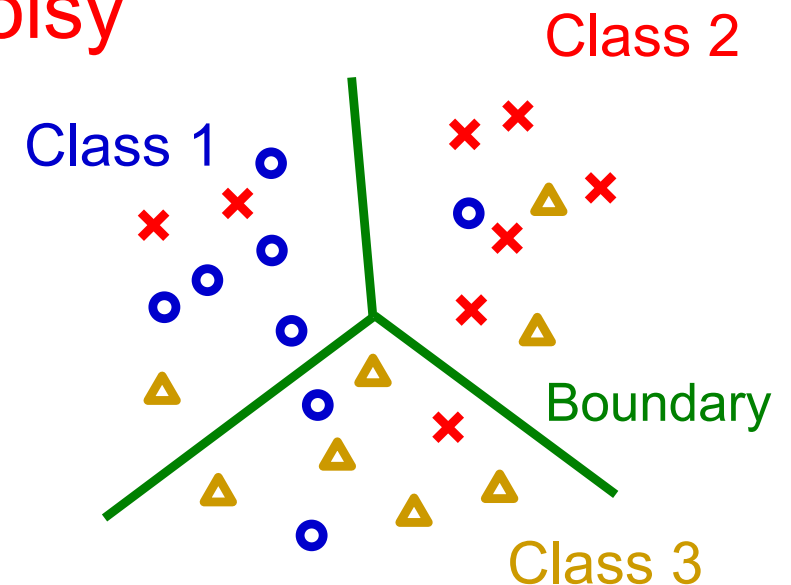
- Supervised classification with **clean** labels:



Training error minimization is **statistically consistent** and work well in practice.

- However, real-world labels are **noisy** possibly due to human error:

Training error minimization is **no longer consistent** and does not work well in practice.

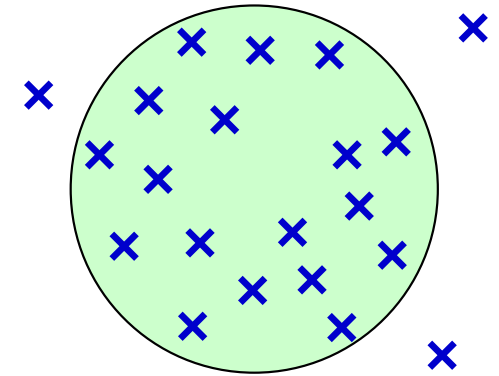


# Classical Approaches

35

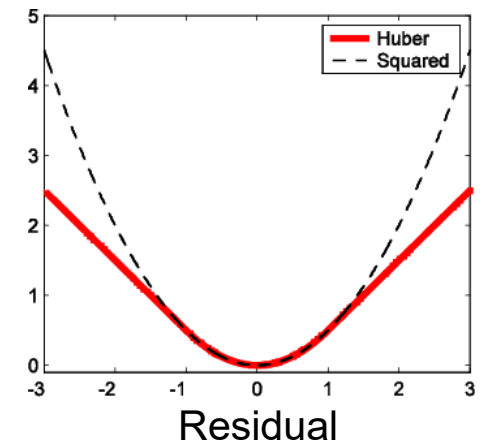
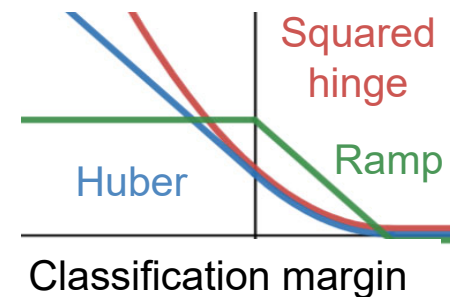
## ■ Unsupervised outlier removal:

- Substantially more difficult than classification.



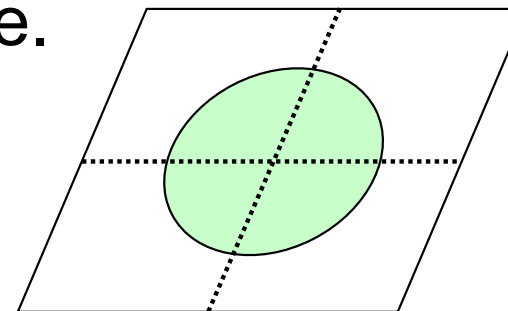
## ■ Robust loss:

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

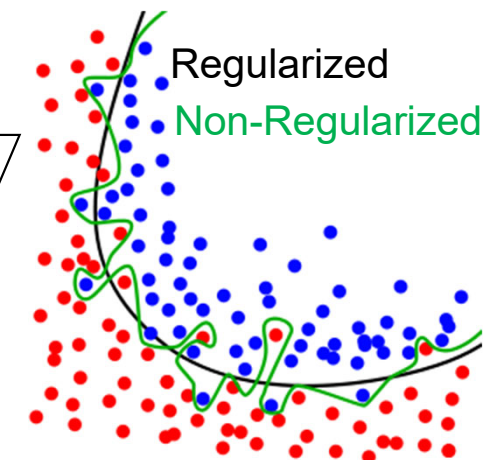


## ■ Regularization:

- 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 et al. (CVPR2017)

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

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

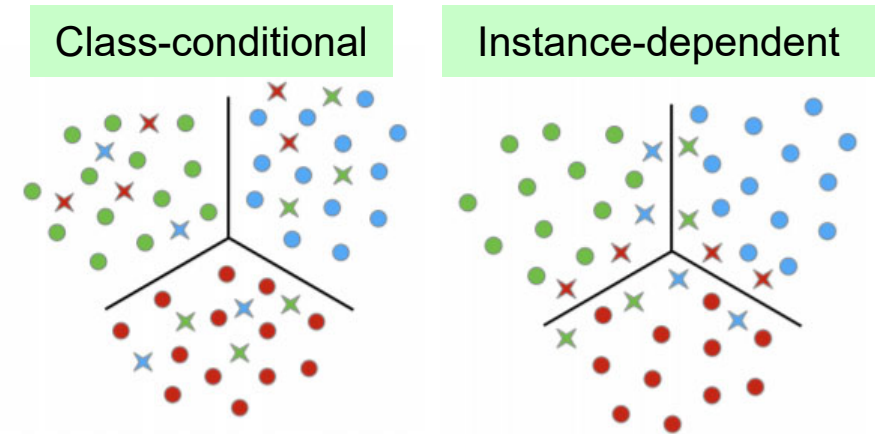
- Use human cognition as a “mask” for  $T$ . Han, Yao, Niu, Zhou, Tsang, Zhang & Sugiyama (NeurIPS2018)
- Reduce estimation error of  $T$ . Xia, Liu, Wang, Han, Gong, Niu & Sugiyama (NeurIPS2019)  
Yao, Liu, Han, Gong, Deng, Niu, Sugiyama & Tao (NeurIPS2020)
- Learn  $T$  and classifier simultaneously. Zhang, Niu & Sugiyama (ICML2021)
- Estimate  $T$  under weaker conditions. Li, Liu, Han, Niu & Sugiyama (ICML2021)

# Beyond Class-Conditional Noise

37

- Real-world noise may be instance-dependent:

- Ex.: Noise is large near the boundary.



- Instance-dependent noise:**  $T_{y, \bar{y}}(\mathbf{x}) = \bar{p}(\bar{y} | y, \mathbf{x})$

- Extremely challenging to estimate the noise transition matrix function!

- Various heuristic solutions:**

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

Xia, Liu, Han, Wang,  
Gong, Liu, Niu, Tao  
& Sugiyama (NeurIPS2020)

Berthon, Han, Niu, Liu  
& Sugiyama (ICML2021)

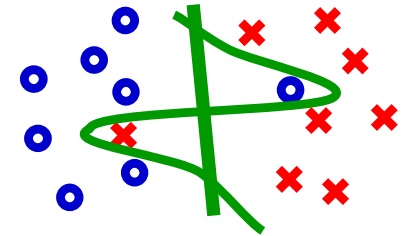
Cheng, Liu, Ning, Wang, Han, Niu,  
Gao & Sugiyama (CVPR2022)

# Co-teaching

## Memorization of neural nets:

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

Arpit et al. (ICML2017)  
Zhang et al. (ICLR2017)



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

- Teach small-loss data each other.

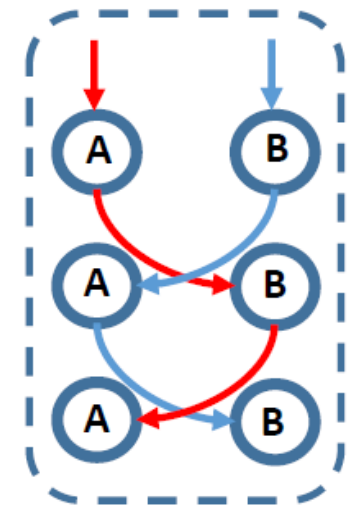
Han, Yao, Yu, Niu, Xu, Hu, Tsang & Sugiyama (NeurIPS2018)

- Teach only disagreed data.

Yu, Han, Yao, Niu, Tsang & Sugiyama (ICML2019)

- Gradient ascent for large-loss data.

Han, Niu, Yu, Yao, Xu, Tsang & Sugiyama (ICML2020)



## No theory but very robust in experiments:

- Works well even if 50% random label flipping!





# Contents

1. Introduction of RIKEN-AIP
2. Robust Machine Learning
  - A) Weakly Supervised Learning
  - B) Transfer Learning
  - C) Noise-Robust Learning
3. Summary

# Challenges in Reliable ML

40

- Reliability for expectable situations:
  - Model the corruption process explicitly and correct the solution.
    - How to handle modeling error?
  
- Reliability for unexpected situations:
  - Consider worst-case robustness (“min-max”).
    - How to make it less conservative?
  - Include human support (“rejection”).
    - How to handle real-time applications?
  
- Exploring somewhere in the middle would be practically more useful:
  - Use partial knowledge of the corruption process.



# History of AI and Future

41

## ■ Classic AI:

- 1960s:  
symbolic, logical AI
- 1980s:  
Expert systems

## ■ Neuro-inspired AI:

- 1960s:  
1-layer perceptrons
- 1980s:  
Multilayer perceptrons

## ■ Statistical machine learning:

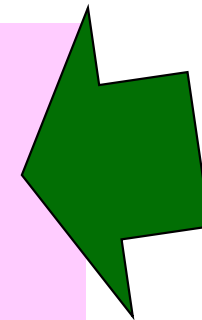
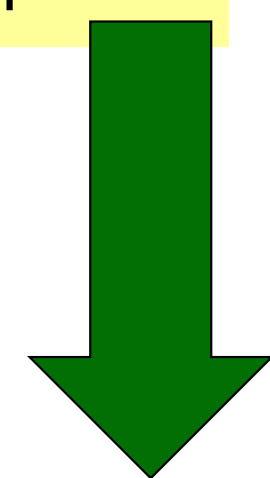
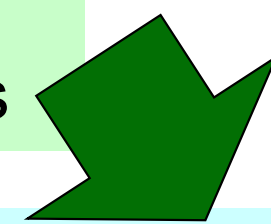
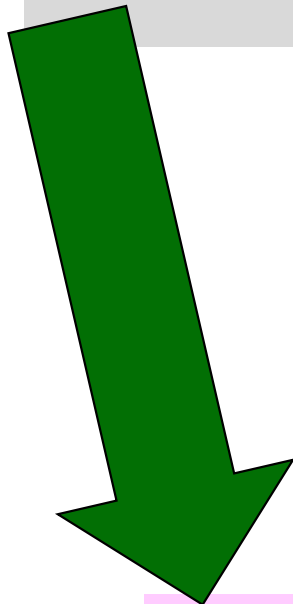
- 2000s: Statistics, Bayes,  
convex optimization, kernels

## ■ Deep learning:

- 2010s: Stochastic  
gradient, gigantic  
deep models

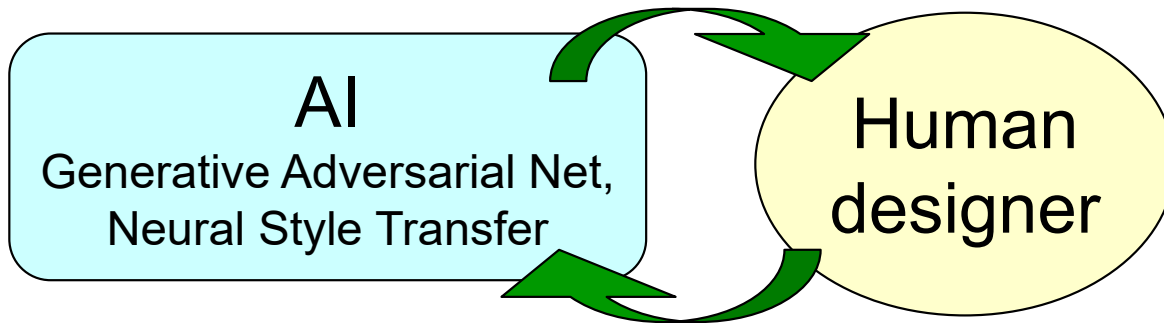
## ■ Next-generation AI:

- Integration of elements
- Human-like AI?



## ■ Is Human-like AI ultimate?

- Future AI needs not be autonomous.
- Future AI may learn **together with humans**.



Fashion show at UTokyo in Mar. 2019  
(with Prof. Aihara and Emarie)



## ■ AI needs to be inclusive to human society:

- **Technology**

**X**

**Human creativity,  
culture, and ethics.**