# Robust Machine Learning from Weakly-Supervised, Noisy-Labeled, and Biased Data

## Masashi Sugiyama

RIKEN Center for Advanced Intelligence Project/

The University of Tokyo

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

# About Myself

■ **Masashi Sugiyama**:

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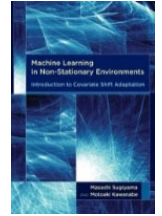- Consultant: several local startups

■ **Interests: Machine learning (ML)**

- ML theory & algorithm →
- ML applications (signal, image, language, brain, robot, mobility, advertisement, biology, medicine, education…)
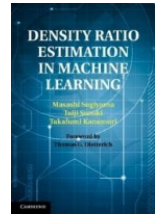
■ **Academic activities:**
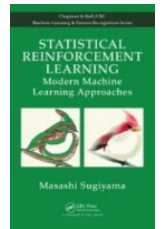
- Program Chairs for NeurIPS2015, AISTATS2019, ACML2010/2020…

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

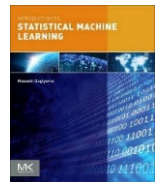Sugiyama, Suzuki & Kanamori, Density Ratio Estimation in Machine Learning, Cambridge University Press, 2012
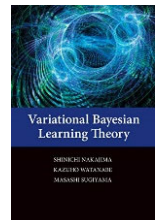
Sugiyama, Statistical Reinforcement Learning, Chapman and Hall/CRC, 2015

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.

# 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)

# What is RIKEN-AIP?

■ **MEXT Advanced Intelligence Project (2016-2025):**
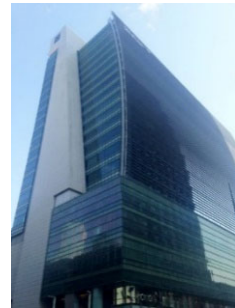
● 130 employed researchers (36% international, 23% female)

● 200 visiting researchers, 100 domestic students

● 140 international interns (total)

■ **Missions:**

● Develop new AI technology (ML, Opt, math)

● Accelerate scientific research (cancer, material, genomics)

● Solve socially critical problems (disaster, elderly healthcare)

● Study of ELSI in AI (ethical guidelines, personal data)

● Human resource development (researchers, engineers)

MEXT
MINISTRY OF EDUCATION,
CULTURE, SPORTS,
SCIENCE AND TECHNOLOGY-JAPAN

Distributed offices across Japan

Main office in the heart of Tokyo

Sendai
Kana gawa Tsukuba
Shiga
Kyoto
Tokyo
Nara Nagoya
Fukuoka

# Selected Research

## Developing New AI Technology
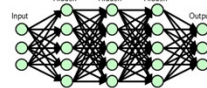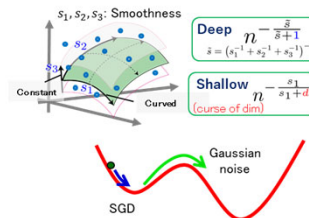
- **Theory of deep learning**:
  - Better prediction than shallow learning
  - No curse of dimensionality
  - Global optimization

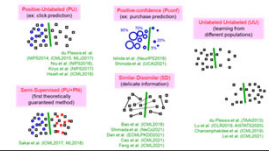$$\mathbb{E}[\|f_T - f^*\|_{L_2}^2] \le \epsilon_M + O(T^{-\frac{2r\beta}{2r\beta+1}})$$

- **Developing new methods**:
  - Weakly supervised learning
  - Noise robust learning
  - Causal inference

$s_1, s_2, s_3$: Smoothness

Deep $n^{-\frac{\bar{s}}{\bar{s}+1}}$  $\bar{s} = (s_1^{-1} + s_2^{-1} + s_3^{-1})^{-1}$

Shallow $n^{-\frac{s_1}{s_1+d}}$ (curse of dim)

Gaussian noise

SGD

**Weakly Supervised Classification**
- Various weakly supervised classification problems can be solved by risk-rewriting systematically!

**Noise Transition Correction**
- Noise transition matrix $T$:
  - 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 $T(x)$.

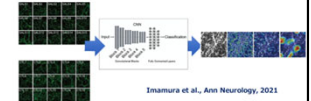**Causal Inference in the Presence of Hidden Cause**
- In causal inference, how to handle hidden cause is a big challenge!
- We developed the first method to estimate the entire structure in the presence of hidden cause:
  - Speech separation technique is employed to separate hidden cause.

Maeda & Shimizu (AISTATS2020, UAI2021)
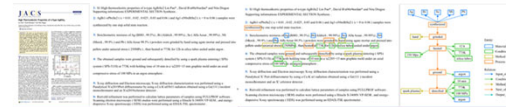
## Accelerating Scientific Research

- **Medical science**:
  - Prostate/pancreatic cancer detection
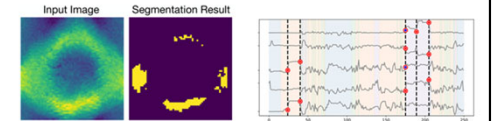  - ALS early diagnosis
  - Fetal heart screening
  - Colonoscopy

*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

- **Data-driven science**:
  - Selective inference for reliability evaluation

Input Image    Segmentation Result

## Solving Socially Critical Problems

- **Natural disaster**:
  - Fugaku-based earthquake simulation
  - Remote sensing disaster analysis

Three-dimensional low-order unstructured nonlinear finite element analysis of the ultra-large-scale and ultra-high-fidelity fault-structure model
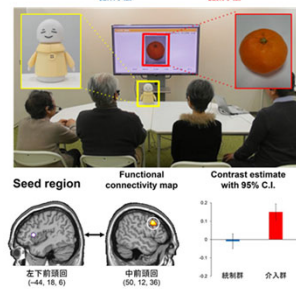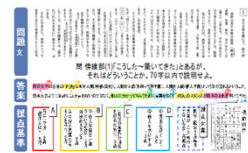
- **Elderly healthcare**:
  - Chat-robot-guided cognitive function improvement

- **Education**:
  - Automatic essay evaluation
  - Interactive essay writing support

Seed region    Functional connectivity map    Contrast estimate with 95% C.I.

## Studying AI-ELSI

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

ETHICALLY ALIGNED DESIGN
First Edition
A Vision for Prioritizing Human Well-being with Autonomous and Intelligent Systems
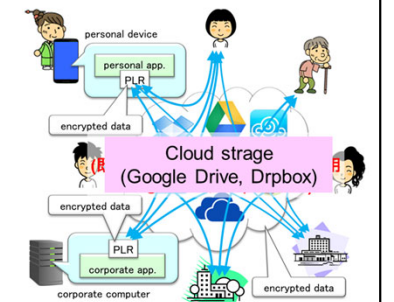#IEEEprinciples2practice    ◆IEEE

- **Personal data management**:
  - Individual-based accessibility control system

- **AI security and reliability**:
  - Adversarial attack/defense
  - Fairness faking/guarantee

Cloud strage (Google Drive, Drpbox)

# Today's Topic:
# Robust Machine Learning

■ **Goal**: Develop novel ML theories and algorithms that enable reliable learning from limited information.

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

# Contents

1. **Noisy-Label Learning**
   - A) Technical background
   - B) Single-step approach
   - C) Beyond anchor points
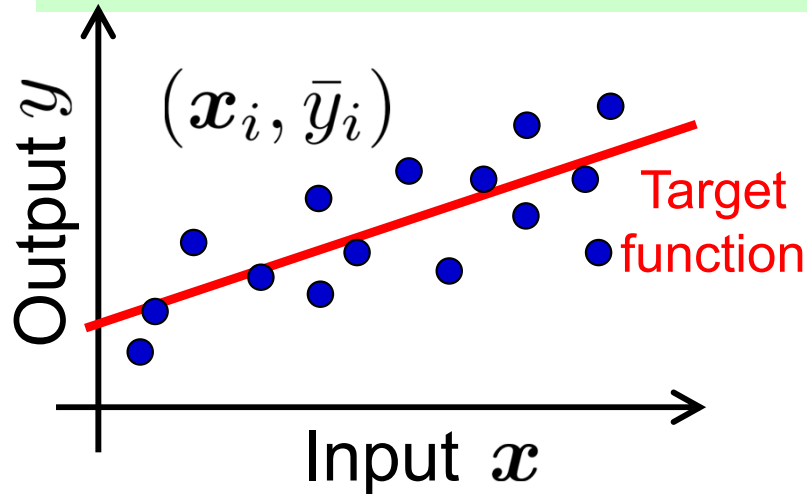   - D) Further challenges
2. Weakly Supervised Learning
3. Transfer Learning
4. Summary

Regression (additive noise)

$(\boldsymbol{x}_i, \bar{y}_i)$

Output $y$

Target function

Input $\boldsymbol{x}$

Classification (label flipping noise)

Class 1

Class 2

$(\boldsymbol{x}_i, \bar{y}_i)$

True boundary

Class 3

$$\min_{g} \sum_{i=1}^{n} \ell\Big(\bar{y}_i, g(\boldsymbol{x}_i)\Big)$$

$\ell$ : loss    $g$ : classifier

$\bar{y}$ : noisy output

■ Hasn't such a classic problem been solved?

- Regression: Yes, big data yields consistency.
- Classification: Specific noise reduction mechanism is needed to achieve consistency!

# Classical Approaches

■ **Unsupervised outlier removal**:

- Substantially more difficult than classification.

■ **Robust loss**:

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

Squared hinge

Huber · Ramp

Classification margin

Huber
Squared

Residual

■ **Regularization**:

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

$\ell_2$-regularization

Regularized
Non-regularized

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

■ **Need new approaches!**

# Contents

1. **Noisy-Label Learning**
   A) Technical background
   B) Single-step approach
   C) Beyond anchor points
   D) Further challenges
2. Weakly Supervised Learning
3. Transfer Learning
4. Summary
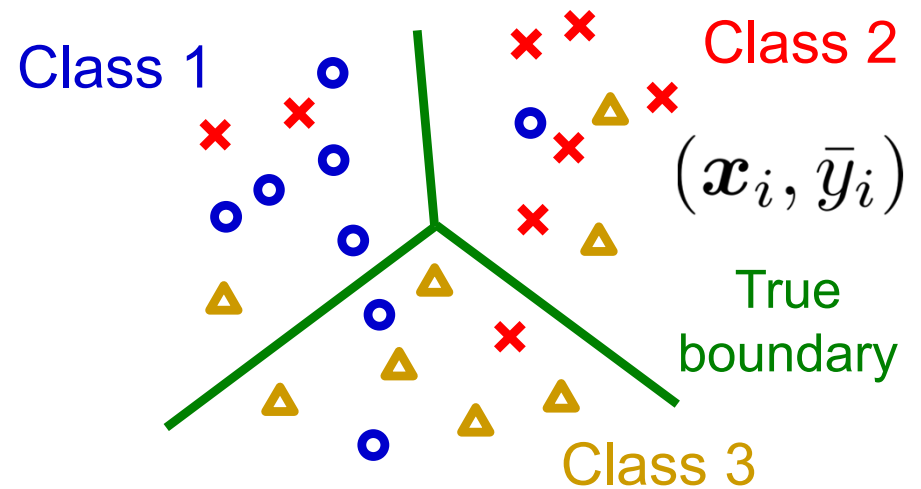
# Formulation

■ Clean training data: $\{(\boldsymbol{x}_i, y_i)\}_{i=1}^n \overset{\text{i.i.d.}}{\sim} p(\boldsymbol{x}, y)$

■ Noisy training data: $\{(\boldsymbol{x}_i, \bar{y}_i)\}_{i=1}^n \overset{\text{i.i.d.}}{\sim} \bar{p}(\boldsymbol{x}, \bar{y})$

$\boldsymbol{x} \in \mathbb{R}^d$ : Input instance

$y \in \{1, \ldots, c\}$ : Clean class label

$\bar{y} \in \{1, \ldots, c\}$ : Noisy class label

■ Probabilistic classifier in simplex: $\boldsymbol{h}(\boldsymbol{x}) \in \Delta^{c-1}$

• Each element approximates the class-posterior probability.

$$h_y(\boldsymbol{x}) \approx p(y|\boldsymbol{x})$$

■ Loss: $\ell(y, \boldsymbol{h}(\boldsymbol{x})) \in \mathbb{R}$

Class 1

Class 2

Boundary

Class 3

# Modeling Class-Conditional Noise

■ **Noise transition matrix:** $T_{y,\bar{y}} = \bar{p}(\bar{y}|y)$

| $y$ | | |
|---|---|---|
| 1 | 0 | 0 |
| 0.1 | 0.8 | 0.1 |
| 0.5 | 0.5 | 0 |

$\bar{y}$

  • Probability of flipping $y$ to $\bar{y}$.

■ We may encode human-cognitive bias:



(a) Column-diagonal      (b) Tri-diagonal      (c) Block-diagonal

Han, Yao, Niu, Zhou, Tsang,
Zhang & Sugiyama (NeurIPS2018)

■ Visualization as a simplex:     Zhang, Niu & Sugiyama (ICML2021)



Clean     Symmetric     Pairwise     General

# Loss Correction

■ **Forward correction:** Add noise by $T^{\top}$

- $\ell^{\rightarrow}(\boldsymbol{h}(\boldsymbol{x})) = \ell(T^{\top}\boldsymbol{h}(\boldsymbol{x}))$ $\quad$ $\ell_y^*(\boldsymbol{h}(\boldsymbol{x})) = \ell^*(y, \boldsymbol{h}(\boldsymbol{x}))$

<span style="color:red">Classifier-consistency</span>

$$\underset{\boldsymbol{h}}{\operatorname{argmin}}\, \mathbb{E}_{\bar{p}(\boldsymbol{x},\bar{y})}[\ell^{\rightarrow}(y, \boldsymbol{h}(\boldsymbol{x}))] = \underset{\boldsymbol{h}}{\operatorname{argmin}}\, \mathbb{E}_{p(\boldsymbol{x},y)}[\ell(y, \boldsymbol{h}(\boldsymbol{x}))]$$
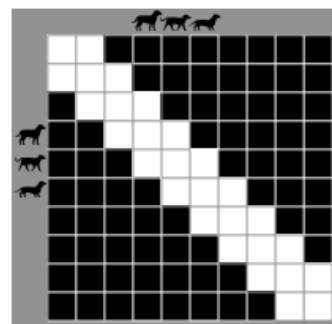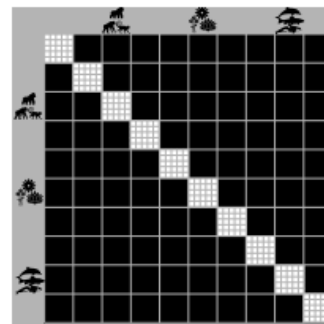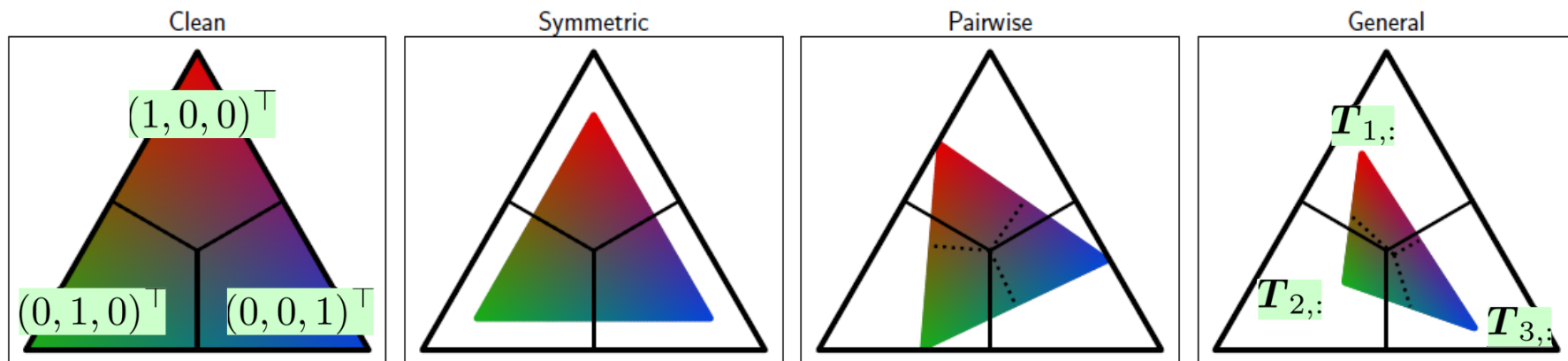
■ **Backward correction:** Remove noise by $T^{-1}$

- $\ell^{\leftarrow}(\boldsymbol{h}(\boldsymbol{x})) = T^{-1}\ell(\boldsymbol{h}(\boldsymbol{x}))$

<span style="color:red">Classifier-consistency</span>

$$\underset{\boldsymbol{h}}{\operatorname{argmin}}\, \mathbb{E}_{\bar{p}(\boldsymbol{x},\bar{y})}[\ell^{\leftarrow}(y, \boldsymbol{h}(\boldsymbol{x}))] = \underset{\boldsymbol{h}}{\operatorname{argmin}}\, \mathbb{E}_{p(\boldsymbol{x},y)}[\ell(y, \boldsymbol{h}(\boldsymbol{x}))]$$

<span style="color:red">Risk-consistency</span>

$$\forall \boldsymbol{x},\ \mathbb{E}_{\bar{p}(\bar{y}|\boldsymbol{x})}[\ell^{\leftarrow}(y, \boldsymbol{h}(\boldsymbol{x}))] = \mathbb{E}_{p(y|\boldsymbol{x})}[\ell(y, \boldsymbol{h}(\boldsymbol{x}))]$$

■ If $T$ is given, consistency can be guaranteed!

# Identifiability of Noise Transition

- In practice, we need to estimate $T$ from noisy training data $\{(\boldsymbol{x}_i, \bar{y}_i)\}_{i=1}^{n}$ .

- However, $T$ is non-identifiable in general:
  - $T$ can be decomposed as $T = UV$, where $U, V$ are some transition matrices.
  - Then $\bar{\boldsymbol{p}}_{\boldsymbol{x}} = T^{\top} \boldsymbol{p}_{\boldsymbol{x}}$
    $= V^{\top}(U^{\top} \boldsymbol{p}_{\boldsymbol{x}})$

    $T_{y,\bar{y}} = \bar{p}(\bar{y}|y)$
    $[\bar{\boldsymbol{p}}_{\boldsymbol{x}}]_{\bar{y}} = \bar{p}(\bar{y}|\boldsymbol{x})$
    $[\boldsymbol{p}_{\boldsymbol{x}}]_y = p(y|\boldsymbol{x})$

- Let's use anchor points (100%-certain samples):
  $$\{\boldsymbol{x}^y \mid p(y|\boldsymbol{x}^y) = 1\}_{y=1}^{c}$$

# Estimation of Noise Transition with Anchor Points

■ Given anchor points $\{\boldsymbol{x}^y \mid p(y|\boldsymbol{x}^y) = 1\}_{y=1}^{c}$ ,
$T_{y,\bar{y}} = \bar{p}(\bar{y}|y)$ can be naïvely estimated as

$$T_{y,\bar{y}} = \sum_{y'=1}^{c} p(\bar{y}|y')p(y'|\boldsymbol{x}^y) = \bar{p}(\bar{y}|\boldsymbol{x}^y) \approx \bar{h}_{\bar{y}}(\boldsymbol{x}^y)$$

- $\bar{\boldsymbol{h}}(\boldsymbol{x})$ is a probabilistic classifier learned from noisy training data $\{(\boldsymbol{x}_i, \bar{y}_i)\}_{i=1}^{n}$ .

■ Even if anchor points are unknown,
as long as they exist in noisy training data,
we may find them as $\boldsymbol{x}^y \leftarrow \boldsymbol{x}_i$ s.t. $\bar{h}_y(\boldsymbol{x}_i) \approx 1$.

# Further Improvements

$$x^y \leftarrow x_i \text{ s.t. } \bar{h}_y(x_i) \approx 1$$

■ We typically use deep learning to obtain $\bar{h}(x)$:

  ● Then it is often over-confident and unreliable.



Zhang, Niu & Sugiyama (ICML2021)

■ Estimated $T$ is revised during classifier training:

Xia, Liu, Wang, Han, Gong, Niu & Sugiyama (NeurIPS2019)



■ Instead of explicitly finding anchor points, latent labels are utilized: $y_i' = \text{argmax}_{y'} \bar{h}_{y'}(x_i)$

Yao, Liu, Han, Gong, Deng, Niu, Sugiyama & Tao (NeurIPS2020)

# Contents

1. **Noisy-Label Learning**
   - A) Technical background
   - B) Single-step approach
   - C) Beyond anchor points
   - D) Further challenges
2. Weakly Supervised Learning
3. Transfer Learning
4. Summary

# Challenge

■ Current approaches are in two-step:

1. Estimate transition matrix $T$.

2. Use estimated $T$ to train a classifier $h(x)$.

■ Step 1 is done without regard to Step 2:

- Estimation error of $T$ in Step 1 can be magnified in Step 2.

■ We want to estimate $T$ and $h(x)$ simultaneously in one-step.

# Naïve Solution

■ Naively, we may learn the noise transition and classifier at the same time as

$$\min_{\boldsymbol{U},\boldsymbol{h}} \mathbb{E}_{\bar{p}(\boldsymbol{x},\bar{y})}[\ell(\bar{y}, \boldsymbol{U}^\top \boldsymbol{h}(\boldsymbol{x}))]$$

■ However, the solution is not unique:

- With any invertible transition matrix $\boldsymbol{Q}$, any $(\widehat{\boldsymbol{U}}, \widehat{\boldsymbol{h}}) = (\boldsymbol{Q}^{-1}\boldsymbol{T}, \boldsymbol{Q}^\top \boldsymbol{p}_{\boldsymbol{x}})$ are solutions.

$$T_{y,\bar{y}} = \bar{p}(\bar{y}|y) \qquad [\boldsymbol{p}_{\boldsymbol{x}}]_y = p(y|\boldsymbol{x})$$

■ We need a certain constraint to obtain the right solution: $(\widehat{\boldsymbol{U}}, \widehat{\boldsymbol{h}}) = (\boldsymbol{T}, \boldsymbol{p}_{\boldsymbol{x}})$

# Total Variation Regularization

Zhang, Niu & Sugiyama (ICML2021)

🟧 Noise transition $p_x \to U^\top p_x$ is contraction in total variation distance:

$$\|U^\top p_x - U^\top p_{x'}\|_1 \leq \|p_x - p_{x'}\|_1$$

$$[p_x]_y = p(y|x)$$

- Cleaner class-posteriors have larger total variation distances!

🟧 Let's use this knowledge as a regularizer:

$$\min_{U,h} \left[ \mathbb{E}_{\bar{p}(x,\bar{y})}[\ell(\bar{y}, U^\top h(x))] - \lambda \mathbb{E}_{p(x),p(x')}\|h(x) - h(x')\|_1 \right]$$

- Under the anchor point assumption, $\lambda > 0$
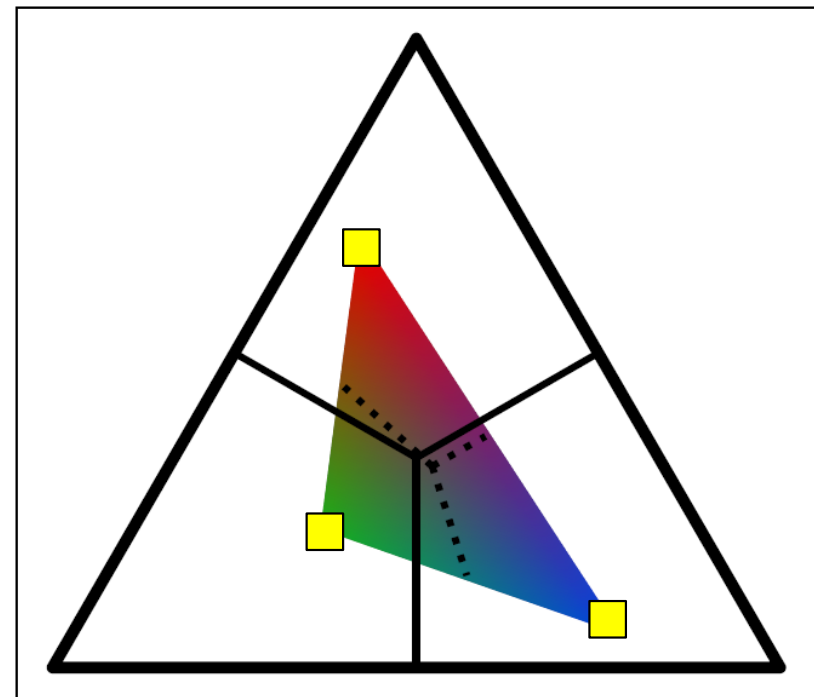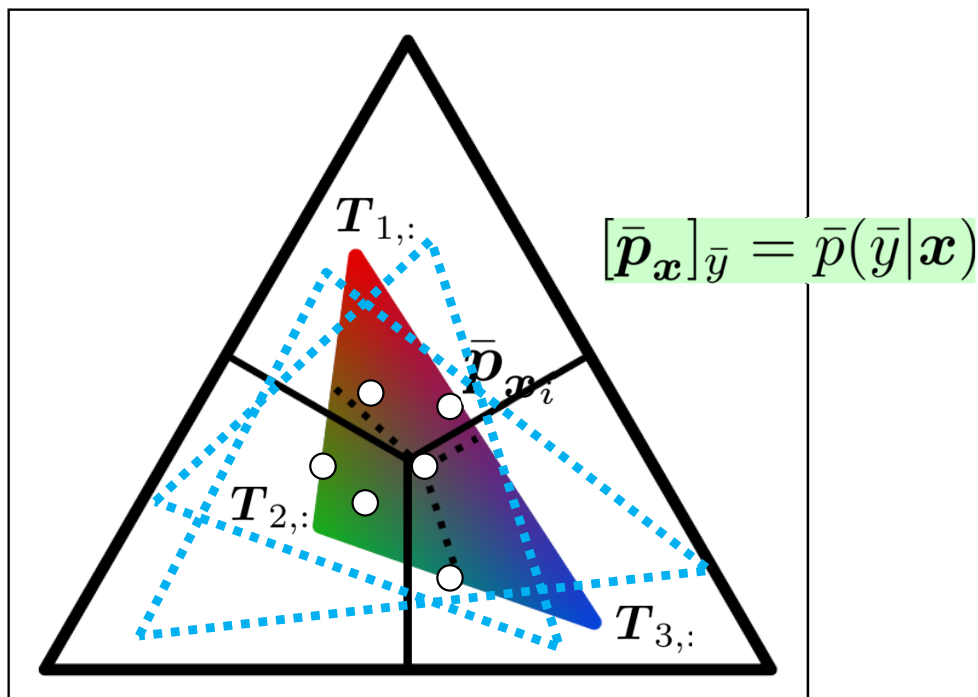  the empirical solution has statistical consistency.

# Contents

1. **Noisy-Label Learning**
   - A) Technical background
   - B) Single-step approach
   - C) Beyond anchor points
   - D) Further challenges
2. Weakly Supervised Learning
3. Transfer Learning
4. Summary

# Challenges

$$\{\boldsymbol{x}^y \mid p(y|\boldsymbol{x}^y) = 1\}_{y=1}^c$$

- To overcome the non-identifiability of $\boldsymbol{T}$ :
  - Anchor points are explicitly used.
- This condition has been relaxed to:
  - Only the existence of anchor points is assumed.

- Can we further relax this assumption?

# Non-identifiability of $T$
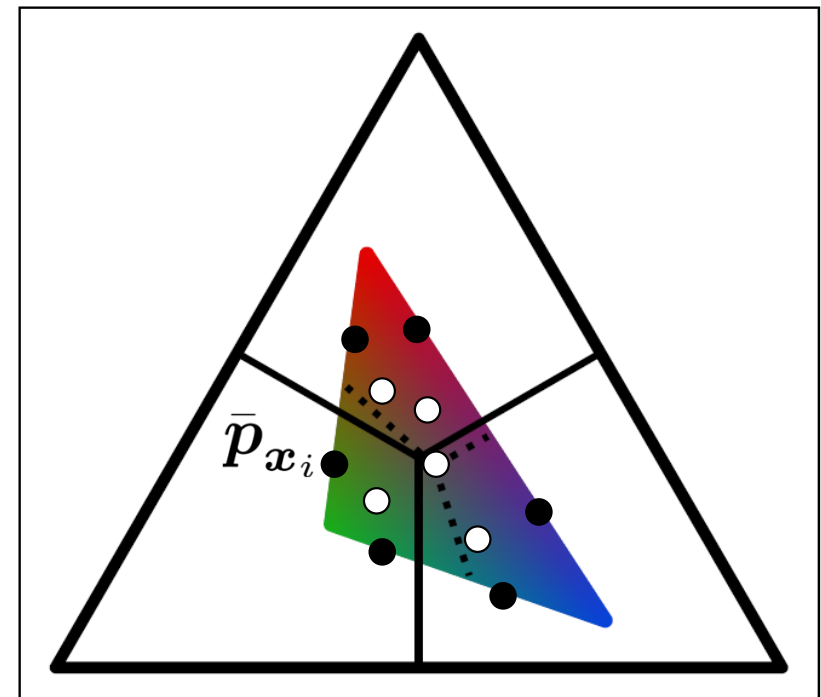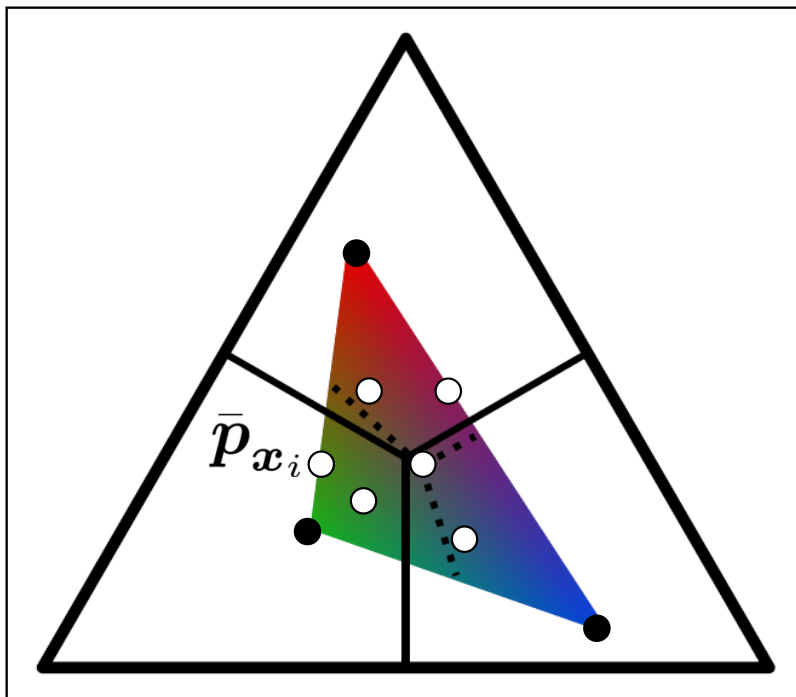
- $T$ can be visualized as a simplex, containing all training data.

- Generally, such a simplex is not unique.

- Anchor points are vertices of the true simplex.

  - Explicitly using anchor points naively recovers $T$.

$$[\bar{\boldsymbol{p}}_{\boldsymbol{x}}]_{\bar{y}} = \bar{p}(\bar{y}|\boldsymbol{x})$$

# Non-identifiability of $T$ (cont.)
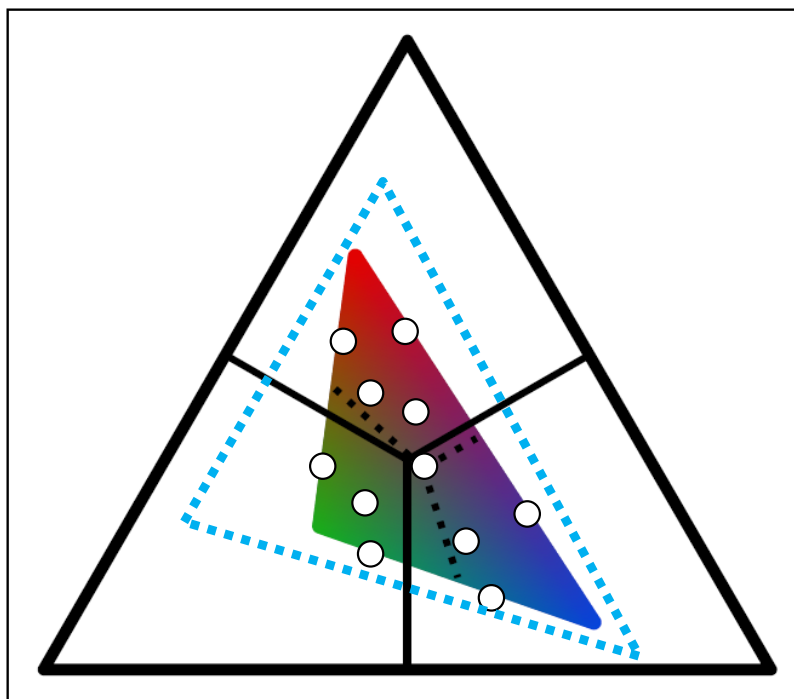
- Only the existence of anchor points still guarantees the identifiability of $T$.

- Even without anchor points, "sufficiently scattered" training data can guarantee the consistency (with the algorithm in the next page).

# Volume Minimization

- Under the "sufficiently scattered" assumption, minimizing the volume of the transition matrix guarantees consistency!

$$\min_{\boldsymbol{U}, \boldsymbol{h}} \left[ \mathbb{E}_{\bar{p}(\boldsymbol{x}, \bar{y})} [\ell(\bar{y}, \boldsymbol{U}^\top \boldsymbol{h}(\boldsymbol{x}))] + \lambda \log \det(\boldsymbol{U}) \right] \quad \lambda > 0$$
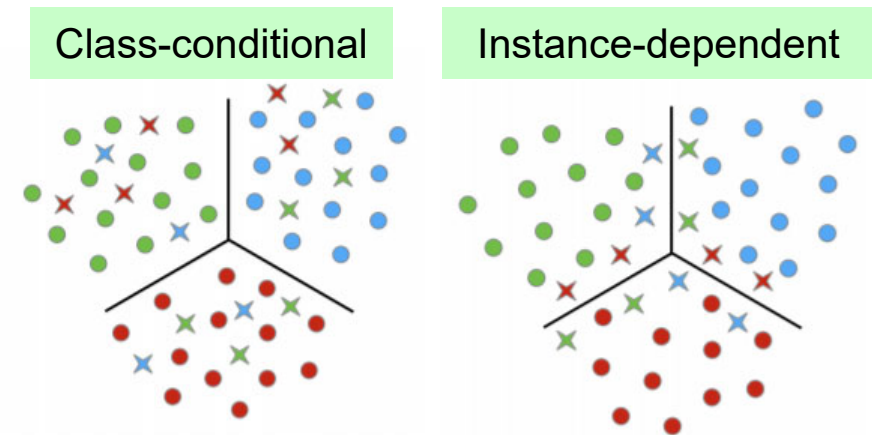
# Contents

1. **Noisy-Label Learning**
   - A) Technical background
   - B) Single-step approach
   - C) Beyond anchor points
   - D) Further challenges
2. Weakly Supervised Learning
3. Transfer Learning
4. Summary

# Beyond Class-Conditional Noise

■ Instance-independence in class-conditional noise is restrictive.



Class-conditional    Instance-dependent

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

  ● Extremely challenging problem!

■ Various heuristic solutions:

  ● Parts-based estimation

  ● Use of additional confidence scores

  ● Manifold regularization

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

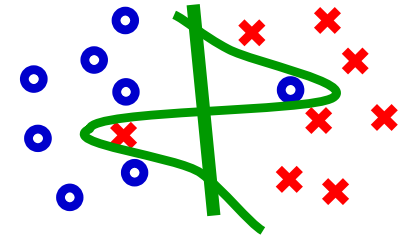Berthon, Han, Niu, Liu & Sugiyama (ICML2021)

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

# Co-teaching

■ Memorization of neural nets:

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

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



■ "Co-teaching" between two neural nets:
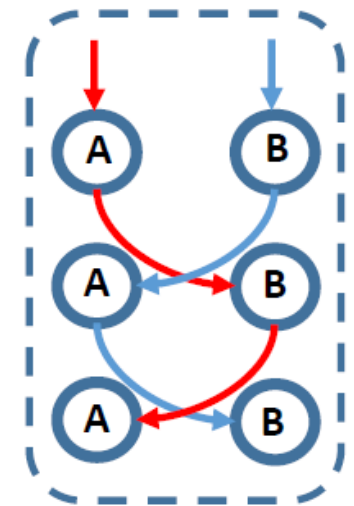
- Teach small-loss data each other.

  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!

■ **Classification requires explicit treatment of label noise:**

   ● <span style="color:red">Loss correction by noise transition</span> is promising.

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

■ **However, noise transition is generally <span style="color:red">non-identifiable</span>.**

   ● Recent development allows its consistent estimation under mild assumptions.



■ **Real-world noise is often instance-dependent:**

   ● <span style="color:red">Heuristic solutions have been developed.</span>

■ **Super-robustness by co-teaching:**

   ● <span style="color:red">Heuristic solutions have been developed.</span>

# Contents

1. Noisy-Label Learning
2. Weakly Supervised Learning
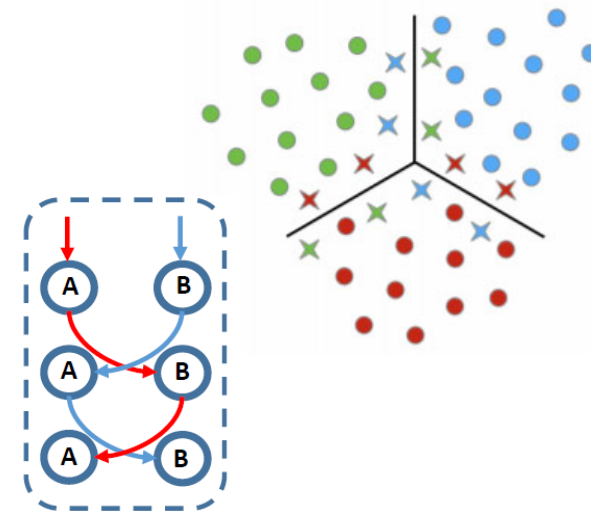3. Transfer Learning
4. Summary

# Weakly Supervised Learning

■ Fully supervised data is expensive to collect.

■ **Weakly supervised data** can be collected easily:

Positive    Negative



- Ex.) Click prediction in online ads: It is easy to automatically collect
  - ■ Clicked ads (positive),
  - ■ Unclicked ads (unlabeled).
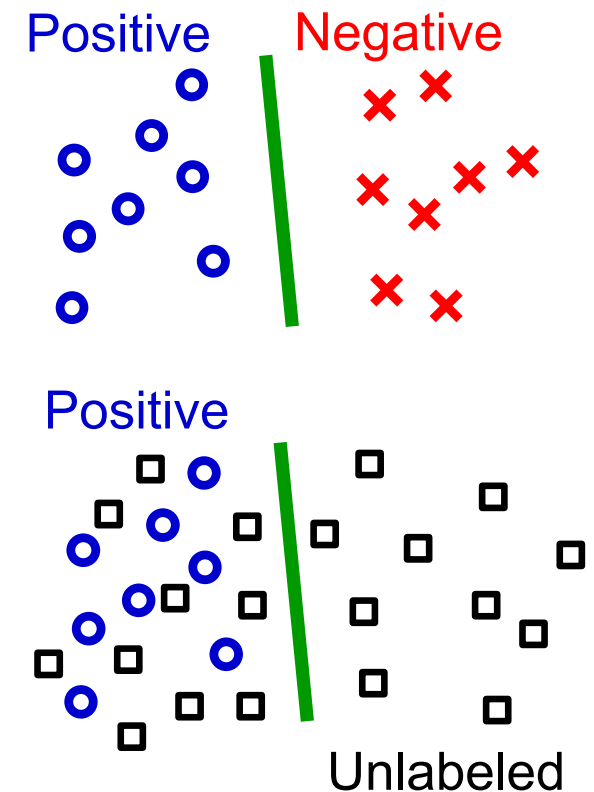
Positive



Unlabeled

■ **Learning only from P and U data is possible!**

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

- Regard U data as noisy N data and correct the loss.
- Statistically consistent.    $\mathcal{O}_p\left(1/\sqrt{n}\right)$

# Solution (Sketch)

- **Given**: Positive and unlabeled data    du Plessis, Niu & Sugiyama (NIPS2014, ICML2015)

$$\{x_i^{\mathrm{P}}\}_{i=1}^{n_{\mathrm{P}}} \overset{\text{i.i.d.}}{\sim} p(x|y=+1) \qquad \{x_j^{\mathrm{U}}\}_{j=1}^{n_{\mathrm{U}}} \overset{\text{i.i.d.}}{\sim} p(x)$$

- Decomposition of the classification risk:

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

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

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

Risk for positive data        Risk for negative data

- Eliminate the expectation over negative data as

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

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

- Unbiased risk estimation:

$$\widehat{R}_{\mathrm{PU}}(f) = \frac{\pi}{n_{\mathrm{P}}}\sum_{i=1}^{n_{\mathrm{P}}}\ell\big(f(x_i^{\mathrm{P}})\big) + \frac{1}{n_{\mathrm{U}}}\sum_{j=1}^{n_{\mathrm{U}}}\ell\big(-f(x_j^{\mathrm{U}})\big) - \frac{\pi}{n_{\mathrm{P}}}\sum_{i=1}^{n_{\mathrm{P}}}\ell\big(-f(x_i^{\mathrm{P}})\big)$$

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

$$\widehat{R}_{\mathrm{PU}}(f) = \frac{\pi}{n_{\mathrm{P}}} \sum_{i=1}^{n_{\mathrm{P}}} \ell\Big(f(\boldsymbol{x}_i^{\mathrm{P}})\Big) + \frac{1}{n_{\mathrm{U}}} \sum_{j=1}^{n_{\mathrm{U}}} \ell\Big(-f(\boldsymbol{x}_j^{\mathrm{U}})\Big) - \frac{\pi}{n_{\mathrm{P}}} \sum_{i=1}^{n_{\mathrm{P}}} \ell\Big(-f(\boldsymbol{x}_i^{\mathrm{P}})\Big)$$

■ **Optimal parametric convergence rate:**

$$R(\widehat{f}_{\mathrm{PU}}) - R(f^*) \le C(\delta) \left( \frac{2\pi}{\sqrt{n_{\mathrm{P}}}} + \frac{1}{\sqrt{n_{\mathrm{U}}}} \right) = \mathcal{O}_p\left( \frac{1}{\sqrt{n_{\mathrm{P}}}} + \frac{1}{\sqrt{n_{\mathrm{U}}}} \right)$$

$$\widehat{f}_{\mathrm{PU}} = \operatorname{argmin}_f \widehat{R}_{\mathrm{PU}}(f)$$

with probability $1 - \delta$

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

$$R(f) = \mathbb{E}_{p(\boldsymbol{x},y)}\Big[\ell\Big(yf(\boldsymbol{x})\Big)\Big]$$

■ **Risk correction further improves the performance**

Kiryo, Niu, du Plessis & Sugiyama (NIPS2017)

$$\widetilde{R}_{\mathrm{PU}}(f) = \frac{\pi}{n_{\mathrm{P}}} \sum_{i=1}^{n_{\mathrm{P}}} \ell\Big(f(\boldsymbol{x}_i^{\mathrm{P}})\Big) + \max\left\{ 0, \ \frac{1}{n_{\mathrm{U}}} \sum_{i=1}^{n_{\mathrm{U}}} \ell\Big(-f(\boldsymbol{x}_i^{\mathrm{U}})\Big) - \frac{\pi}{n_{\mathrm{P}}} \sum_{i=1}^{n_{\mathrm{P}}} \ell\Big(-f(\boldsymbol{x}_i^{\mathrm{P}})\Big) \right\}$$
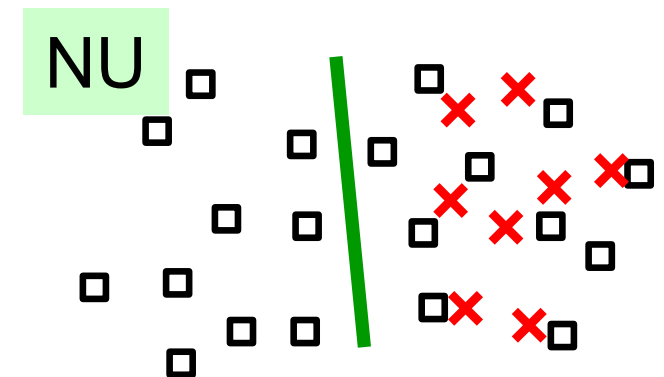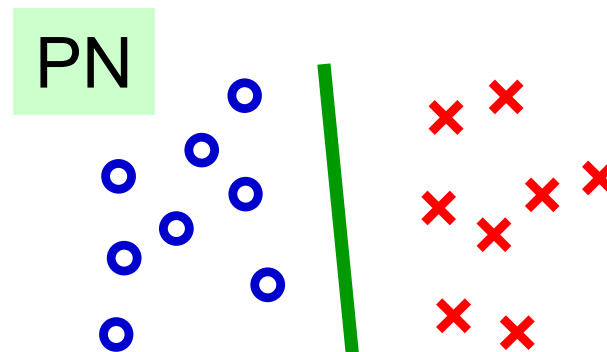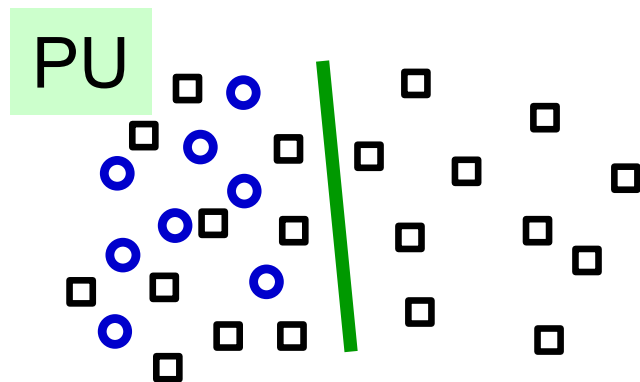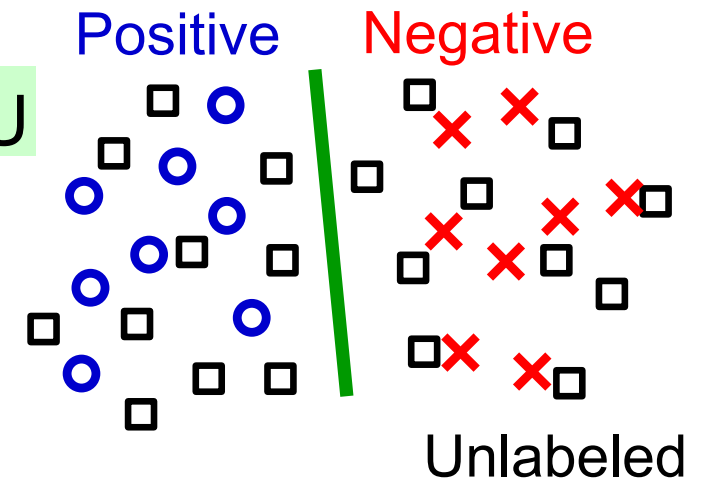
Sakai, du Plessis, Niu & Sugiyama (ICML2017)

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

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

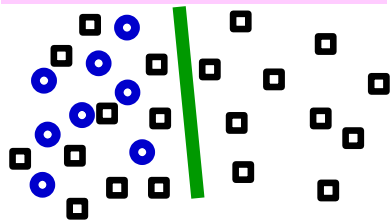■ Without cluster assumptions, PN classifiers are trainable!

$$\mathcal{O}_p\left(1/\sqrt{n_\mathrm{P}} + 1/\sqrt{n_\mathrm{N}} + 1/\sqrt{n_\mathrm{U}}\right)$$

Positive    Negative

PNU

Unlabeled

PU

PN

NU

# 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)

**Positive-confidence**

95%  70%  20%  5%

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

**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)

**Similar-Dissimilar**

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

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

# Multiclass Methods

■ 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
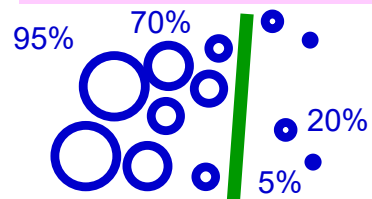  that contains the correct one ("1 or 2").

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

- Single-class confidence: Cao et al. (arXiv2021)
  One-class data with full confidence
  ("1 with 60%, 2 with 30%, and 3 with 10%")

■ Similar loss correction is possible!

Class 1    Class 2

Class 3    Boundary

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

# Summary: Weakly Supervised Learning

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



Supervised

Semi-supervised

Unsupervised

P, N, U, S, D, Pconf, Nconf, Sconf, Dconf....
Comp, Partial, SCconf…
Different weak information can be systematically combined!

High

Labeling cost

Low

Low    Classification accuracy    High

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

Machine Learning from Weak Supervision

An Empirical Risk Minimization Approach

Masashi Sugiyama, Han Bao, Takashi Ishida, Nan Lu, Tomoya Sakai, and Gang Niu

# Contents

1. Noisy-Label Learning
2. Weakly Supervised Learning
3. Transfer Learning
4. Summary

# Transfer Learning

■ **Given:**

- Training data $\{(\boldsymbol{x}_i^{\mathrm{tr}}, y_i^{\mathrm{tr}})\}_{i=1}^{n_{\mathrm{tr}}} \overset{\mathrm{i.i.d.}}{\sim} p_{\mathrm{tr}}(\boldsymbol{x}, y)$

$\boldsymbol{x}$ : Input

$y$ : Output

■ **Goal:**

- Train a predictor $y = f(\boldsymbol{x})$
  that works well in the test domain
  (with some additional data from the test domain).

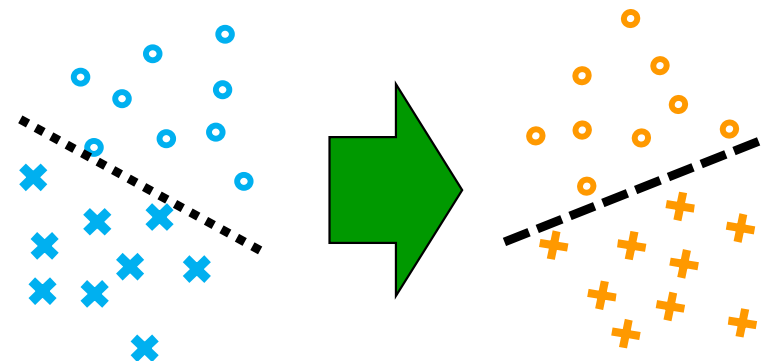$$\min_f R(f) \qquad R(f) = \mathbb{E}_{p_{\mathrm{te}}(\boldsymbol{x}, y)}[\ell(f(\boldsymbol{x}), y)]$$

$\ell$ : loss function

■ **Challenge:**

- Overcome changing distributions!

$$p_{\mathrm{tr}}(\boldsymbol{x}, y) \neq p_{\mathrm{te}}(\boldsymbol{x}, y)$$

NIPS Workshop 2006 - Whistler

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

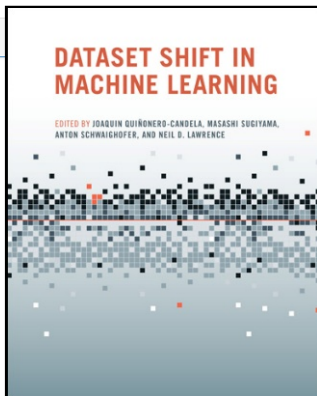*Learning when test and training inputs have different distributions*                    Workshop

Joaquin Quiñonero Candela · Masashi Sugiyama · Anton Schwaighofer · Neil D Lawrence

Sat Dec 09 05:00 PM -- 05:00 PM (JST) @ Nordic

Event URL: http://ida.first.fraunhofer.de/projects/different06/ »

Many machine learning algorithms assume that the training and the test data are drawn from the same distribution. Indeed many of the proofs of statistical consistency, etc., rely on this assumption. However, in practice we are very often faced with the situation where the training and the test data both follow the same conditional distribution, p(y|x), but the input distributions, p(x), differ. For example, principles of experimental design dictate that training data is acquired in a specific manner that bears little resemblance to the way the test inputs may later be generated. The aim of this workshop will be to try and shed light on the kind of situations where explicitly addressing the difference in the input distributions is beneficial, and on what the most sensible ways of doing this are.

Quiñonero-Candela, Sugiyama, Schwaighofer & Lawrence (Eds.), Dataset Shift in Machine Learning, MIT Press, 2009.

Learning when Training and Test Inputs Have Different Distributions
Saturday December 9, 2006
Org: Joaquin Quiñonero-Candela, Anton Schwaighofer, Neil Lawrence & Masashi Sugiyama

*Morning session: 7:30am–10:30am*

7:30am   **Opening,** *The organizers*

7:40am   **When Training and Test Distributions are Different: Characterising Learning Transfer,** *Amos Storkey, University of Edinburgh*

8:10am   **Can Adaptive Regularization Help?,** *Matthias Hein, Max Planck Institute for Biological Cybernetics*

8:40am   *coffee break*

8:50am   **Learning Classifiers in Distribution and Cost-sensitive Environments,** *Nitesh Chawla, University of Notre Dame*

9:20am   **Optimality of Bayesian Transduction - Implications for Input Non-stationarity,** *Lars Kai Hansen, Technical University of Denmark*

9:50pm   **Estimating the Joint AUC of Labelled and Unlabelled Data,** *Thomas Gärtner, Gemma Garriga, Thorsten Knopp, Peter Flach and Stefan Wrobel*

10:10am  **A Domain Adaptation Formal Framework Addressing the Training/Test Distribution Gap,** *Shai Ben-David, University of Waterloo and John Blitzer, University of Pennsylvania*

*Afternoon session: 3:30pm–6:30pm*

3:30pm   **Projection and Projectability,** *David Corfield, Max Planck Institute for Biological Cybernetics*

4:00pm   **Using features of probability distributions to achieve covariate shift,** *Arthur Gretton, MPI for Biol. Cyb. and Alex Smola, National ICT Australia*

4:20pm   **Active Learning, Model Selection and Covariate Shift,** *Masashi Sugiyama, Tokyo Institute of Technology*

4:50pm   *coffee break*

5:00pm   **Visualizing Pairwise Similarity via Semidefinite Programming,** *Amir Globerson, MIT, and Sam Roweis, University of Toronto*

5:20pm   **A Divergence Prior for Adaptive Learning,** *Xiao Li and Jeff Bilmes, University of Washington*

5:40pm   *discussion, everyone*

# Various Scenarios

$\boxed{\boldsymbol{x} : \text{Input}}$  $\boxed{y : \text{Output}}$

- **Full-distribution shift:** $p_{\mathrm{tr}}(\boldsymbol{x}, y) \neq p_{\mathrm{te}}(\boldsymbol{x}, y)$

- **Covariate shift:** $p_{\mathrm{tr}}(\boldsymbol{x}) \neq p_{\mathrm{te}}(\boldsymbol{x})$

- **Class-prior shift:** $p_{\mathrm{tr}}(y) \neq p_{\mathrm{te}}(y)$

- **Output noise:** $p_{\mathrm{tr}}(y|\boldsymbol{x}) \neq p_{\mathrm{te}}(y|\boldsymbol{x})$

- **Class-conditional shift:** $p_{\mathrm{tr}}(\boldsymbol{x}|y) \neq p_{\mathrm{te}}(\boldsymbol{x}|y)$

# Classical Approach for Transfer Learning
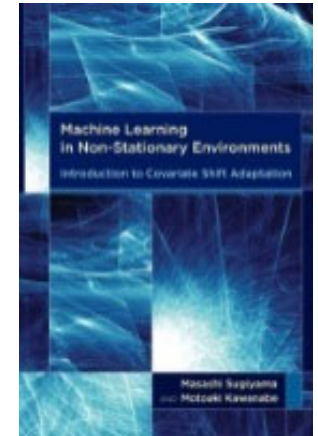
■ Two-step adaptation:

1. Importance weight estimation:

$$\widehat{w} = \underset{w}{\mathrm{argmin}} \, \widehat{\mathbb{E}}_{p_{\mathrm{tr}}(\boldsymbol{x},y)} \left[ D\left( w(\boldsymbol{x},y), \frac{p_{\mathrm{te}}(\boldsymbol{x},y)}{p_{\mathrm{tr}}(\boldsymbol{x},y)} \right) \right]$$

2. Weighted predictor training:

$$\widehat{f} = \underset{f}{\mathrm{argmin}} \, \widehat{\mathbb{E}}_{p_{\mathrm{tr}}(\boldsymbol{x},y)} \left[ \widehat{w}(\boldsymbol{x},y) \ell(f(\boldsymbol{x}),y) \right]$$

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

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

● We want to integrate these two steps!

- 🟧 <span style="color:red">Covariate shift:</span> Only input distributions change.

$$p_{\mathrm{tr}}(\boldsymbol{x}) \neq p_{\mathrm{te}}(\boldsymbol{x}) \qquad p_{\mathrm{tr}}(y|\boldsymbol{x}) = p_{\mathrm{te}}(y|\boldsymbol{x})$$

Shimodaira (JSPI2000)

- 🟧 Suppose we are given
  - Labeled training data: $\{(\boldsymbol{x}_i^{\mathrm{tr}}, y_i^{\mathrm{tr}})\}_{i=1}^{n_{\mathrm{tr}}} \overset{\mathrm{i.i.d.}}{\sim} p_{\mathrm{tr}}(\boldsymbol{x}, y)$
  - Unlabeled test data: $\{\boldsymbol{x}_i^{\mathrm{te}}\}_{i=1}^{n_{\mathrm{te}}} \overset{\mathrm{i.i.d.}}{\sim} p_{\mathrm{te}}(\boldsymbol{x})$

- 🟧 Minimize <span style="color:red">a risk upper bound</span> jointly w.r.t. weight $w$ and predictor $f$ :

Zhang et al. (ACML2020, SNCS2021)

$$J_{\ell_{\mathrm{tr}}}(f, w) \geq R_{\ell_{\mathrm{te}}}(f)^2$$

$$\widehat{f} = \operatorname*{argmin}_{f} \min_{w \geq 0} \widehat{J}_{\ell_{\mathrm{tr}}}(f, w)$$

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

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

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

  - <span style="color:red">Theoretical guarantee:</span>

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

# Dynamic Importance Weighting

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

■ Suppose we are given

- Labeled training data: $\{(\boldsymbol{x}_i^{\mathrm{tr}}, y_i^{\mathrm{tr}})\}_{i=1}^{n_{\mathrm{tr}}} \overset{\mathrm{i.i.d.}}{\sim} p_{\mathrm{tr}}(\boldsymbol{x}, y)$
- Labeled test data: $\{(\boldsymbol{x}_i^{\mathrm{te}}, y_i^{\mathrm{te}})\}_{i=1}^{n_{\mathrm{te}}} \overset{\mathrm{i.i.d.}}{\sim} p_{\mathrm{te}}(\boldsymbol{x}, y)$

■ For each mini-batch $\{(\bar{\boldsymbol{x}}_i^{\mathrm{tr}}, \bar{y}_i^{\mathrm{tr}})\}_{i=1}^{\bar{n}_{\mathrm{tr}}}, \{(\bar{\boldsymbol{x}}_i^{\mathrm{te}}, \bar{y}_i^{\mathrm{te}})\}_{i=1}^{\bar{n}_{\mathrm{te}}}$, importance weights are estimated by matching losses by kernel mean matching:

Fang et al. (NeurIPS2020)

Huang et al. (NeurIPS2007)

$$\frac{1}{\bar{n}_{\mathrm{tr}}} \sum_{i=1}^{\bar{n}_{\mathrm{tr}}} r_i \ell(f(\bar{\boldsymbol{x}}_i^{\mathrm{tr}}), \bar{y}_i^{\mathrm{tr}}) \approx \frac{1}{\bar{n}_{\mathrm{te}}} \sum_{j=1}^{\bar{n}_{\mathrm{te}}} \ell(f(\bar{\boldsymbol{x}}_j^{\mathrm{te}}), \bar{y}_j^{\mathrm{te}})$$
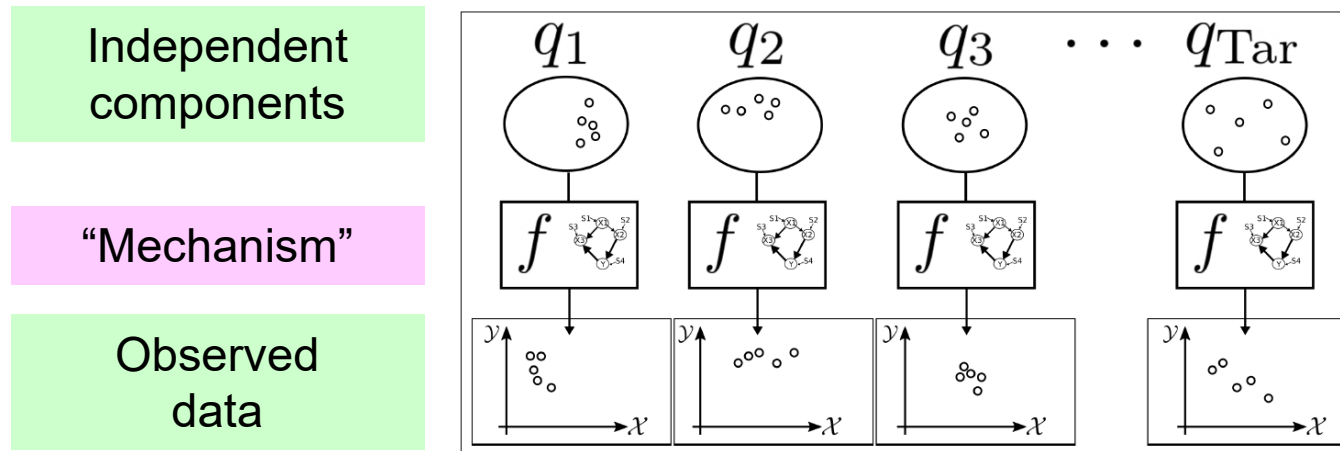
■ Extremely simple, but highly powerful!

■ In transfer learning with importance weighting, simultaneously estimation of <span style="color:red">importance</span> and <span style="color:red">predictor</span> is promising.

■ What should we do if training and test distributions look very different?

● <span style="color:red">Mechanism transfer!</span>   Teshima, Sato & Sugiyama (ICML2020)



Independent components

"Mechanism"

Observed data

■ <span style="color:red">New work:</span> Continuous distribution change.

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

# Contents

1. Noisy-Label Learning
2. Weakly Supervised Learning
3. Transfer Learning
4. Summary

# More Challenges in Reliable Machine Learning

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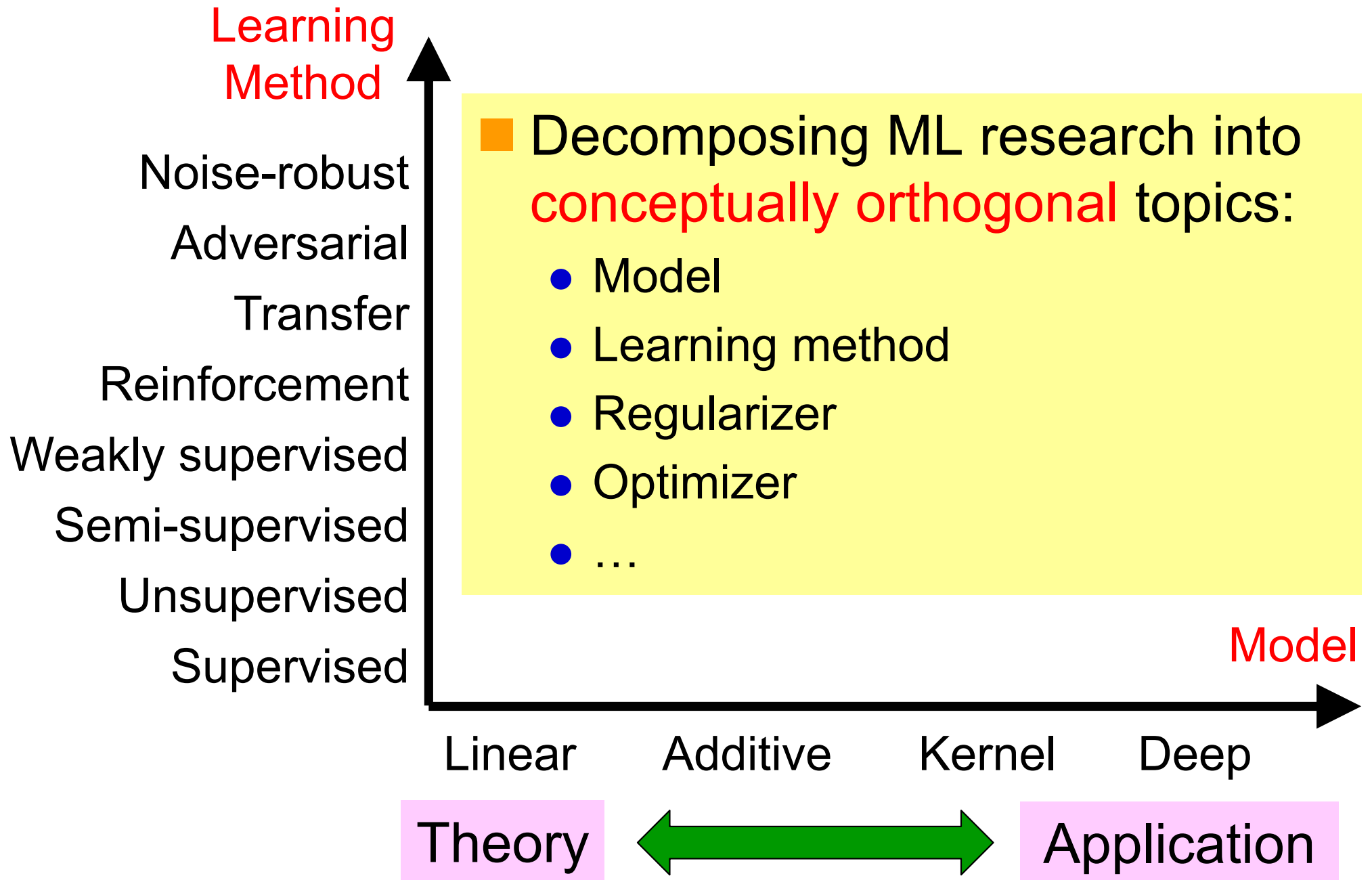        - How to handle modeling error?
- **Reliability for unexpected situations:**
    - Consider worst-case robustness ("min-max").
        - How to make it less conservative?
    - Include human support ("rejection").
        - How to handle real-time applications?
- **Exploring somewhere in the middle would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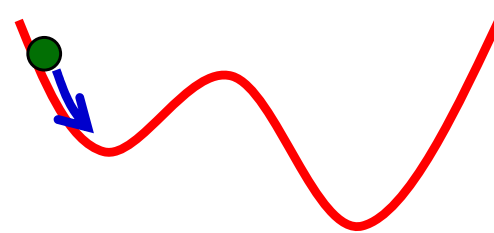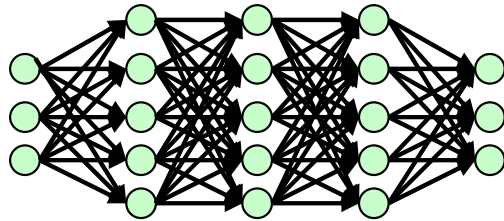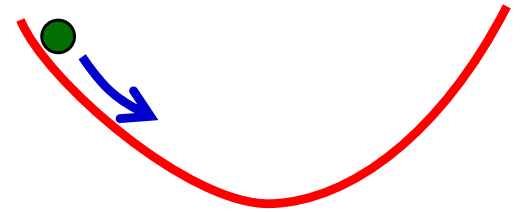    - Use partial knowledge of the corruption process.

# Axes of ML Research

# Further Investigations Needed

- Classical convex learning methods allow us to analyze the global solution.

- Since optimization in deep learning is complex, stochastic gradient descent is used.

- Thanks to the "gradual learning" nature, we can utilized intermediate learning results:
  - Strengthening supervision for weakly supervised learning.
  - Dynamic importance weighting for transfer learning.
  - Dynamic noise transition estimation for noise-robust learning.
  - Co-teaching for noise-robust learning.