

Mixture Proportion Estimation in Weakly Supervised Learning

Masashi Sugiyama



RIKEN Center for Advanced Intelligence Project

The University of Tokyo



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



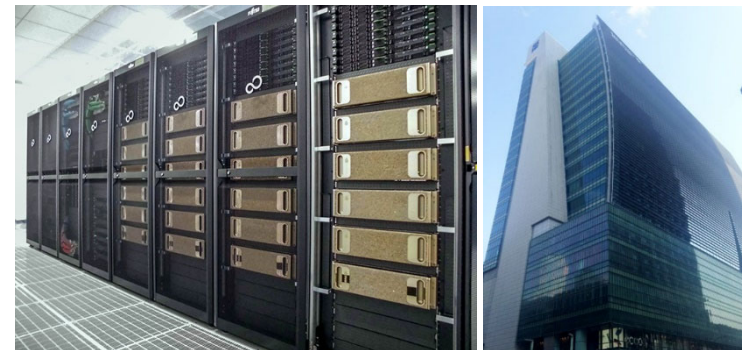
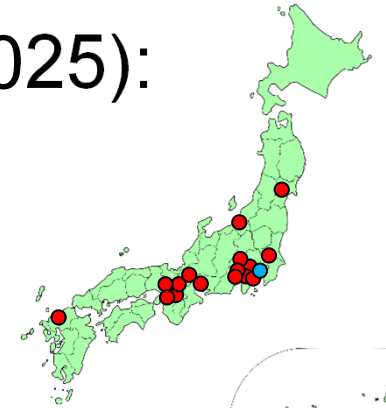
東京大学
THE UNIVERSITY OF TOKYO



RIKEN Center for Advanced Intelligence Project (AIP)

■ 10-year national project in Japan (2016-2025):

- **Develop next-generation AI technology**
(learning and optimization theory, etc.)
- **Accelerate scientific research**
(material, cancer, stem cells, genomics, etc.)
- **Solve socially critical problems**
(natural disaster, elderly healthcare, etc.)
- **Study of ethical, legal and social issues of AI**
(ethical guideline, privacy protection, etc.)
- **Human resource development**
(150+ researchers, 200+ students,
150+ interns, 300+ visiting scientists,
40+ industry projects)



My Research Interests

■ Transfer learning:

- Adaptive importance weighting

■ Density ratio estimation:

- Versatile statistical tool, where GAN is a special case.

■ Reinforcement learning:

- Sample reuse

■ Variational Bayes:

- Implicit regularization

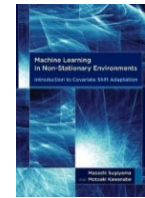
■ Weakly supervised learning:

- Empirical risk minimization approach

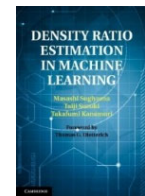
■ Noise-robust learning:

- Going beyond robust statistics and regularization

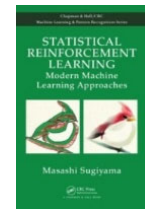
Sugiyama & Kawanabe, *Machine Learning in Non-Stationary Environments*, MIT Press, 2012



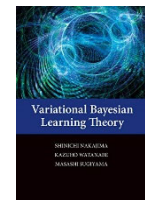
Sugiyama, Suzuki & Kanamori, *Density Ratio Estimation in Machine Learning*, Cambridge University Press, 2012



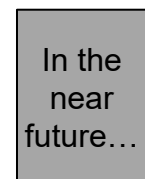
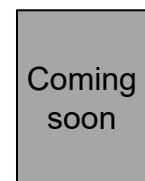
Sugiyama, *Statistical Reinforcement Learning*, Chapman and Hall/CRC, 2015



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



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



Today's Topic:

Mixture Proportion Estimation

- **Goal:** Find a mixture proportion of unknown probability distributions.
- From some data, find $\theta_1, \dots, \theta_c$ such that

$$p_0 = \sum_{y=1}^c \theta_y p_y \quad \sum_{y=1}^c \theta_y = 1 \quad \theta_1, \dots, \theta_c \geq 0$$

p_0, p_1, \dots, p_c : Unknown probability distributions

- Various applications in machine learning:
 - **Class-prior shift adaptation:** Importance weight estimation
 - **Positive-unlabeled classification:** Class-prior estimation
 - **Noisy label classification:** Noise transition estimation

Contents

5

1. Semi-supervised class-prior shift adaptation
 - A) Basic solution
 - B) Distribution matching
 - C) Summary
2. Positive-unlabeled classification
3. Conclusions



Semi-Supervised Classification with Class-Prior Shift

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

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

$\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^d$: Input pattern

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

$y \in \mathcal{Y} = \{1, \dots, c\}$: Class label

- **Goal:** Train a classifier $y = f(\mathbf{x})$ that works well in the test domain.

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

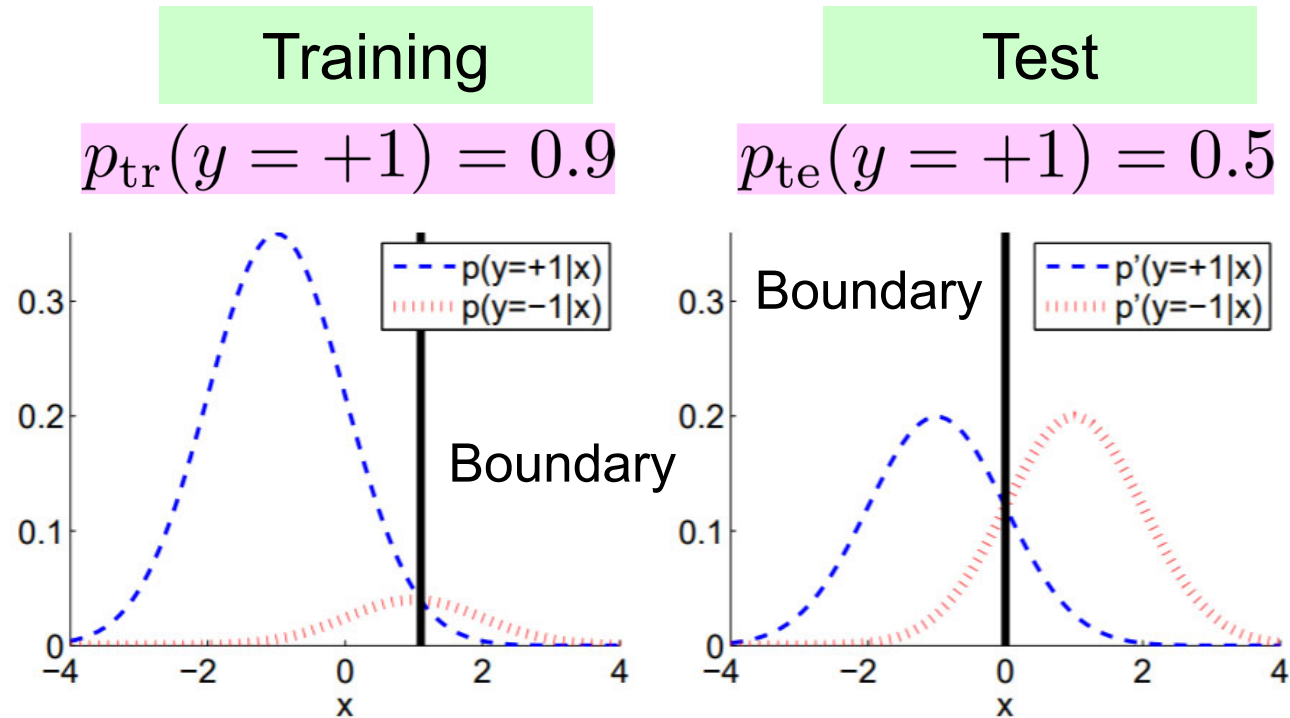
ℓ : loss function

- **Challenge:** Overcome the class-prior shift!

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

Illustration of Class-Prior Shift

7



- Class-prior shift changes the optimal boundary.
- Adaptation is needed!

Contents

1. Semi-supervised class-prior shift adaptation
 - A) Basic solution
 - B) Distribution matching
 - C) Summary
2. Positive-unlabeled classification
3. Conclusions



Empirical Risk Minimization (ERM)

9

$$\operatorname{argmin}_{f \in \mathcal{F}} \left[\sum_{i=1}^{n_{\text{tr}}} \ell(f(\mathbf{x}_i^{\text{tr}}), y_i^{\text{tr}}) \right] \quad \{(\mathbf{x}_i^{\text{tr}}, y_i^{\text{tr}})\}_{i=1}^{n_{\text{tr}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{tr}}(\mathbf{x}, y)$$

■ Generally, ERM is **consistent**:

- Learned function converges to the optimal solution when $n_{\text{tr}} \rightarrow \infty$.

■ However, class-prior shift makes ERM **inconsistent**:

$$\operatorname{argmin}_{f \in \mathcal{F}} \left[\frac{1}{n_{\text{tr}}} \sum_{i=1}^{n_{\text{tr}}} \ell(f(\mathbf{x}_i^{\text{tr}}), y_i^{\text{tr}}) \right] \xrightarrow{n_{\text{tr}} \rightarrow \infty} \operatorname{argmin}_{f \in \mathcal{F}} \left[\mathbb{E}_{p(\mathbf{x}|y) p_{\text{tr}}(y)} [\ell(f(\mathbf{x}), y)] \right] \neq R(f)$$

$$R(f) = \mathbb{E}_{p(\mathbf{x}|y) p_{\text{te}}(y)} [\ell(f(\mathbf{x}), y)]$$
$$p_{\text{tr}}(y) \neq p_{\text{te}}(y)$$

Importance-Weighted ERM (IWERM) 10

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

- IWERM is consistent even under class-prior shift.

$$\begin{aligned} \operatorname{argmin}_{f \in \mathcal{F}} \left[\frac{1}{n_{\text{tr}}} \sum_{i=1}^{n_{\text{tr}}} \frac{p_{\text{te}}(y_i^{\text{tr}})}{p_{\text{tr}}(y_i^{\text{tr}})} \ell(f(\mathbf{x}_i^{\text{tr}}), y_i^{\text{tr}}) \right] \\ \xrightarrow[n_{\text{tr}} \rightarrow \infty]{} \operatorname{argmin}_{f \in \mathcal{F}} \left[\mathbb{E}_{p(\mathbf{x}|y)p_{\text{tr}}(y)} \left[\frac{p_{\text{te}}(y)}{p_{\text{tr}}(y)} \ell(f(\mathbf{x}), y) \right] \right] \\ = \operatorname{argmin}_{f \in \mathcal{F}} \left[\underbrace{\mathbb{E}_{p(\mathbf{x}|y)p_{\text{te}}(y)} [\ell(f(\mathbf{x}), y)]}_{= R(f)} \right] \end{aligned}$$

- How can we know the importance weight?

Class-Prior Estimation by the EM Algorithm

Saerens et al. (NeCo2001)

1. Obtain a training class-posterior estimator $\hat{p}_{\text{tr}}(y|\mathbf{x})$
from $\{(\mathbf{x}_i^{\text{tr}}, y_i^{\text{tr}})\}_{i=1}^{n_{\text{tr}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{tr}}(\mathbf{x}, y)$.

2. Estimate the training class-prior by $\hat{p}_{\text{tr}}(y) \propto n_y$.

n_y : Number of training samples in class y

3. Set $\hat{p}_{\text{te}}(y|\mathbf{x}) = \hat{p}_{\text{tr}}(y|\mathbf{x})$ and $\hat{p}_{\text{te}}(y) = \hat{p}_{\text{tr}}(y)$.

4. Repeat until convergence:

i. Update the test class-posterior as $\hat{p}_{\text{te}}(y|\mathbf{x}) \propto \frac{\hat{p}_{\text{te}}(y)}{\hat{p}_{\text{tr}}(y)} \hat{p}_{\text{tr}}(y|\mathbf{x})$.

ii. Update the test class-prior as $\hat{p}_{\text{te}}(y) = \frac{1}{n_{\text{te}}} \sum_{j=1}^{n_{\text{te}}} \hat{p}_{\text{te}}(y|\mathbf{x}_j^{\text{te}})$.

■ Can we avoid using $\hat{p}_{\text{tr}}(y|\mathbf{x})$? $\{\mathbf{x}_j^{\text{te}}\}_{j=1}^{n_{\text{te}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{te}}(\mathbf{x})$

1. Semi-supervised class-prior shift adaptation
 - A) Basic solution
 - B) Distribution matching
 - C) Summary
2. Positive-unlabeled classification
3. Conclusions



EM Method as Distribution Matching under KL Divergence

13

du Plessis et al. (NN2014)

■ Let $q(\mathbf{x}) = \sum_{y=1}^c \theta_y p_{\text{tr}}(\mathbf{x}|y)$. $\theta_1, \dots, \theta_c \geq 0$ $\sum_{y=1}^c \theta_y = 1$

■ Fit $q(\mathbf{x})$ to $p_{\text{te}}(\mathbf{x})$ under KL divergence:

$$\operatorname{argmin}_{\theta_1, \dots, \theta_c} \text{KL}[p_{\text{te}}(\mathbf{x}) \| q(\mathbf{x})] = \operatorname{argmin}_{\theta_1, \dots, \theta_c} \int p_{\text{te}}(\mathbf{x}) \log \frac{p_{\text{te}}(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x}$$

$$\approx \operatorname{argmin}_{\theta_1, \dots, \theta_c} \frac{1}{n_{\text{te}}} \sum_{j=1}^{n_{\text{te}}} \log \frac{p_{\text{te}}(\mathbf{x}_j^{\text{te}})}{q(\mathbf{x}_j^{\text{te}})} \quad \{\mathbf{x}_j^{\text{te}}\}_{j=1}^{n_{\text{te}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{te}}(\mathbf{x})$$

■ Fixed-point iteration to solve the KKT condition recovers the EM approach!

■ Without estimating $\hat{p}_{\text{tr}}(y|\mathbf{x})$, can we directly minimize the KL divergence?

Direct KL-Divergence Approximation by Density Ratio Estimation

Keziou (2003), Nguyen et al. (NIPS2007), Sugiyama et al. (NIPS2007)

Identity (from **Fenchel's inequality**):

$$\text{KL}[p_{\text{te}}(\mathbf{x}) \| q(\mathbf{x})] = \sup_s \left\{ - \int p_{\text{te}}(\mathbf{x}) s(\mathbf{x}) d\mathbf{x} + \int q(\mathbf{x}) \log s(\mathbf{x}) d\mathbf{x} \right\} + 1$$

- Maximizer is $s(\mathbf{x}) = q(\mathbf{x}) / p_{\text{te}}(\mathbf{x})$.

$$q(\mathbf{x}) = \sum_{y=1}^c \theta_y p_{\text{tr}}(\mathbf{x}|y)$$

Empirical approximation:

$$\widehat{\text{KL}}[p_{\text{te}}(\mathbf{x}) \| q(\mathbf{x})] = \max_s \left\{ - \frac{1}{n_{\text{te}}} \sum_{j=1}^{n_{\text{te}}} s(\mathbf{x}_j^{\text{te}}) + \sum_{y=1}^c \frac{\theta_y}{n_y} \sum_{i:y_i=y} \log s(\mathbf{x}_i^{\text{tr}}) \right\} + 1$$

- Maximization corresponds to **estimating density ratio** $s(\mathbf{x})$.

Then we can directly estimate the test class-prior as

$$\underset{\theta_1, \dots, \theta_c}{\text{argmin}} \text{KL}[p_{\text{te}}(\mathbf{x}) \| q(\mathbf{x})] \approx \underset{\theta_1, \dots, \theta_c}{\text{argmin}} \widehat{\text{KL}}[p_{\text{te}}(\mathbf{x}) \| q(\mathbf{x})]$$

Distribution Matching under the f -Divergence

du Plessis et al. (NN2014)

- We don't have to stick to the KL divergence.
 - We can use any divergence such as the f -divergence:

For convex f such that $f(1) = 0$, Ali & Slivey (1966), Csiszár (1967)

$$\text{Div}_f[p_{\text{te}}(\mathbf{x}) \| q(\mathbf{x})] = \int p_{\text{te}}(\mathbf{x}) f\left(\frac{q(\mathbf{x})}{p_{\text{te}}(\mathbf{x})}\right) d\mathbf{x} \quad q(\mathbf{x}) = \sum_{y=1}^c \theta_y p_{\text{tr}}(\mathbf{x}|y)$$

- Directly estimate the f -divergence from data:

$$\widehat{\text{Div}}_f[p_{\text{te}}(\mathbf{x}) \| q(\mathbf{x})] \quad \{\mathbf{x}_i^{\text{tr}}\}_{i=1}^{n_{\text{tr}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{tr}}(\mathbf{x}) \quad \{\mathbf{x}_j^{\text{te}}\}_{j=1}^{n_{\text{te}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{te}}(\mathbf{x})$$

- Estimate the test class-prior as

$$\underset{\theta_1, \dots, \theta_c}{\text{argmin}} \text{Div}_f[p_{\text{te}}(\mathbf{x}) \| q(\mathbf{x})] \approx \underset{\theta_1, \dots, \theta_c}{\text{argmin}} \widehat{\text{Div}}_f[p_{\text{te}}(\mathbf{x}) \| q(\mathbf{x})]$$

- How do we estimate the f -divergence from data?

Direct f -Divergence Approximation by Density Ratio Estimation

Keziou (2003), Nguyen et al. (NIPS2007), Sugiyama et al. (AISM2012)

Identity (from **Fenchel's inequality**):

$$\text{Div}_f[p_{\text{te}}(\mathbf{x})||q(\mathbf{x})] \\ = - \inf_s \left\{ \int p_{\text{te}}(\mathbf{x}) \left(\partial f(s(\mathbf{x}))s(\mathbf{x}) - f(s(\mathbf{x})) \right) d\mathbf{x} - \int q(\mathbf{x}) \partial f(s(\mathbf{x})) d\mathbf{x} \right\}$$

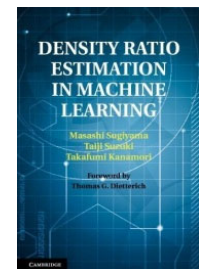
- Equality holds when $s(\mathbf{x}) = q(\mathbf{x})/p_{\text{te}}(\mathbf{x})$.

Empirical approximation:

$$\widehat{\text{Div}}_f[p_{\text{te}}(\mathbf{x})||q(\mathbf{x})] \\ = - \min_s \frac{1}{n_{\text{te}}} \sum_{j=1}^{n_{\text{te}}} \left(\partial f(s(\mathbf{x}_j^{\text{te}}))s(\mathbf{x}_j^{\text{te}}) - f(s(\mathbf{x}_j^{\text{te}})) \right) - \sum_{y=1}^c \frac{\theta_y}{n_y} \sum_{i:y_i=y} f'(s(\mathbf{x}_i^{\text{tr}}))$$

- Minimization corresponds to **density ratio matching under the Bregman divergence.**

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



Various Choices of Function f

For convex f such that $f(1) = 0$,

$$\text{Div}_f[p_{\text{te}}(\mathbf{x}) \| q(\mathbf{x})] = \int p_{\text{te}}(\mathbf{x}) f\left(\frac{q(\mathbf{x})}{p_{\text{te}}(\mathbf{x})}\right) d\mathbf{x}$$

■ Kullback-Leibler (KL) divergence: $f(t) = -\log t$

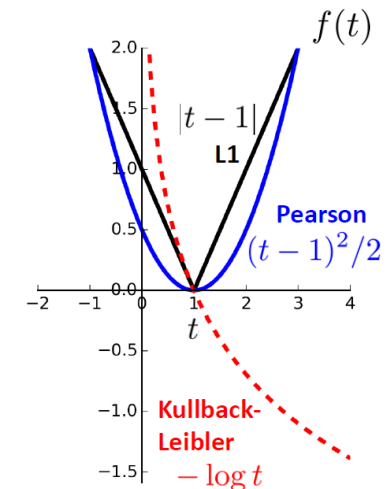
- Popular choice, but sensitive to outliers.
- Optimization is convex if $s(\mathbf{x})$ is a linear model.

■ Pearson (PE) divergence: $f(t) = (t - 1)^2 / 2$

- Robust to outliers.
- Optimization is analytic if $s(\mathbf{x})$ is a linear model.

■ Power divergence: $f(t) = (t^\alpha - 1)t / \alpha$ for $\alpha > 0$

- Generalization of KL ($\alpha \rightarrow 0$) and PE ($\alpha = 1$).
- More robust for $\alpha > 1$, but optimization becomes non-convex.



1. Semi-supervised class-prior shift adaptation
 - A) Basic solution
 - B) Distribution matching
 - C) Summary
2. Positive-unlabeled classification
3. Conclusions



Summary: Semi-Supervised Classification with Class-Prior Shift

■ Importance-weighted empirical risk minimization.

- Estimation of the test class-prior $p_{te}(y)$ is needed.

■ EM is seminal, but requires $\hat{p}_{tr}(y|\mathbf{x})$.

- EM is KL-div minimization with fix-point iteration.
- Can we directly minimize KL-div without $\hat{p}_{tr}(y|\mathbf{x})$?

■ KL-div approximation with density ratio estimation.

- Can we use another divergence?

■ Various divergences/distances can be used.

- f -div approximation by density ratio estimation.
- L2-distance approximation by **density difference estimation**.

1. Semi-supervised class-prior shift adaptation
2. **Positive-unlabeled classification**
 - A) Basic solution
 - B) Identifiability
 - C) Density ratio estimation with anchor points
 - D) Partial distribution matching with irreducibility
 - E) Regrouping without irreducibility
 - F) Summary
3. Conclusions



Positive-Unlabeled (PU) Classification ²¹

- **Given:** Positive and unlabeled samples

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

$\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^d$: Input pattern

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

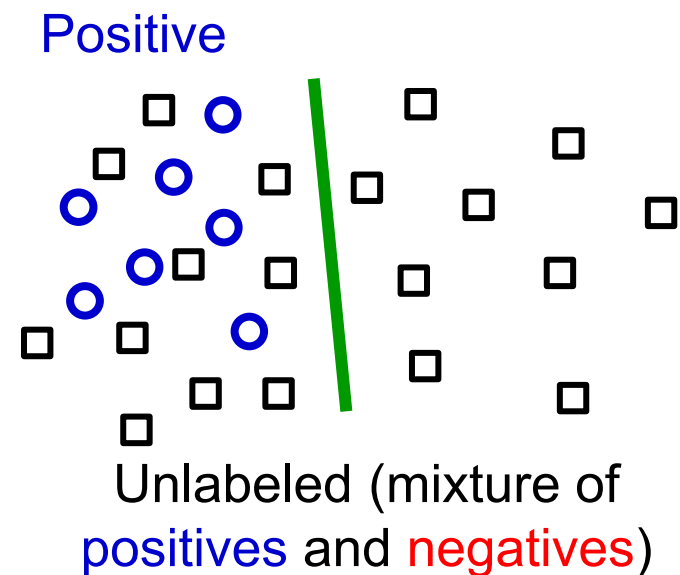
$y \in \mathcal{Y} = \{+1, -1\}$: Class label

- No negative data

- **Goal:** Obtain a positive-negative (PN) classifier

- **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)



1. Semi-supervised class-prior shift adaptation
2. **Positive-unlabeled classification**
 - A) **Basic solution**
 - B) Identifiability
 - C) Density ratio estimation with anchor points
 - D) Partial distribution matching with irreducibility
 - E) Regrouping without irreducibility
 - F) Summary
3. Conclusions



PN Risk Decomposition

- Risk of classifier f :

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

$\pi = p(y = +1)$: Class-prior probability
(for the moment, assume it is known)

- Since we do not have N data in the PU setting, the risk cannot be directly estimated.
 - How can we overcome this problem?

du Plessis et al. (ICML2015)

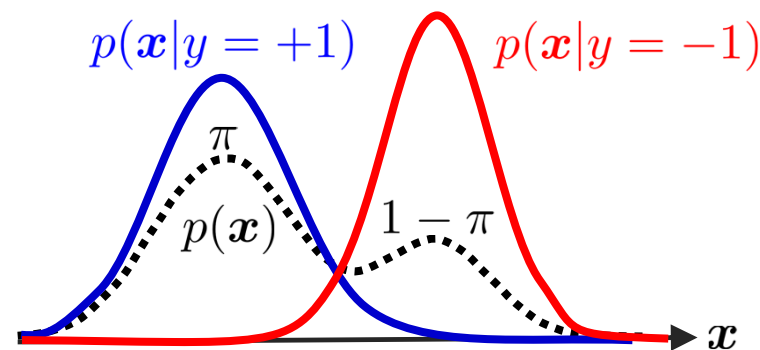
$$R(f) = \pi \mathbb{E}_{p(\mathbf{x}|y=+1)} \left[\ell(f(\mathbf{x})) \right] + (1 - \pi) \mathbb{E}_{p(\mathbf{x}|y=-1)} \left[\ell(-f(\mathbf{x})) \right]$$

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

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

- This allows us to eliminate the N-density as

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



PU Empirical Risk Minimization

25

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

- Replacing expectations by sample averages gives an empirical risk:

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

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

- Optimal convergence rate is attained: Niu et al. (NIPS2016)

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

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

with probability $1 - \delta$

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

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

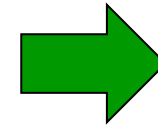
- But, in practice, $\pi = p(y = +1)$ is unknown!

Class-Prior Estimation with Non-Traditional Classification

Elkan & Noto (KDD2008)

- Consider **PU label** $s \in \{0, 1\}$:

- If \mathbf{x} is P (or U), $s = 1$ (or $s = 0$).



$$\pi = \frac{p(s = 1)}{p(s = 1|y = +1)}$$

- Train a **“non-traditional” classifier** $\hat{p}(s|\mathbf{x})$ from PU data.

- Usual supervised classification from $\{\mathbf{x}_i^P\}_{i=1}^{n_P}, \{\mathbf{x}_i^U\}_{i=1}^{n_U}$
(Assume P is labeled from U when $s = 1$.)

- Obtain $\hat{\pi}$ with

\mathcal{P} : Set of validation P data

$$\hat{p}(s = 1) = \frac{n_P}{n_P + n_U} \quad \hat{p}(s = 1|y = +1) = \frac{1}{|\mathcal{P}|} \sum_{\mathbf{x} \in \mathcal{P}} \hat{p}(s = 1|\mathbf{x})$$

- **Can we avoid training a non-traditional classifier?**

cf. Original paper solves
PU classification by

$$\hat{p}(y = +1|\mathbf{x}) = \frac{\hat{p}(s = 1|\mathbf{x})}{\hat{p}(s = 1|y = +1)}$$

1. Semi-supervised class-prior shift adaptation
2. **Positive-unlabeled classification**
 - A) Basic solution
 - B) **Identifiability**
 - C) Density ratio estimation with anchor points
 - D) Partial distribution matching with irreducibility
 - E) Regrouping without irreducibility
 - F) Summary
3. Conclusions



Non-Traditional Classification as Partial Distribution Matching

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

- PN classification: (Full) distribution matching

$$\min_{\theta \in [0,1]} \text{Div}[p(\mathbf{x}) || q(\mathbf{x})] \quad q(\mathbf{x}) = \theta p(\mathbf{x}|y = +1) + (1 - \theta)p(\mathbf{x}|y = -1)$$

- PU classification: **Partial distribution matching**

$$\min_{\theta \in [0,1]} \text{Div}[p(\mathbf{x}) || q'(\mathbf{x})] \quad q'(\mathbf{x}) = \theta p(\mathbf{x}|y = +1)$$

du Plessis & Sugiyama (IEICE2014)

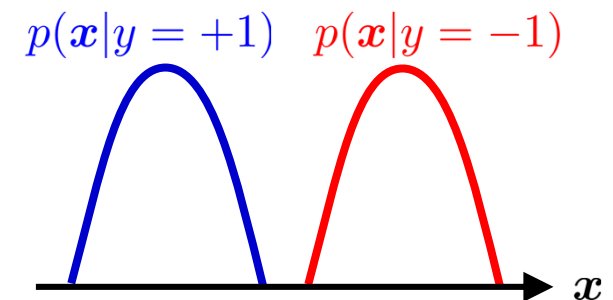
- **Class-prior estimation by non-traditional classification** can be interpreted as **partial matching with Pearson divergence.**

$$\frac{1}{2} \int p(\mathbf{x}) \left(\frac{q'(\mathbf{x})}{p(\mathbf{x})} - 1 \right)^2 d\mathbf{x}$$

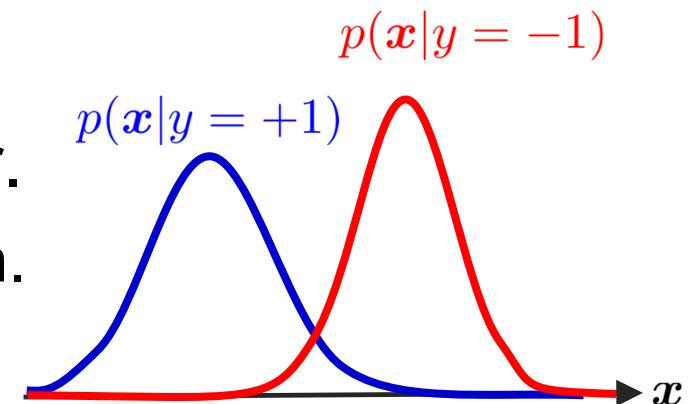
Behaviors of Partial Matching

$$\min_{\theta \in [0,1]} \text{Div}_f[p(\mathbf{x}) \| q'(\mathbf{x})] \quad q'(\mathbf{x}) = \theta p(\mathbf{x}|y = +1)$$

- If **two classes have no overlap**, naïve partial matching works.
 - Just fitting $p(\mathbf{x}|y = +1)$ is sufficient.



- If **two classes are overlapped**, partial matching generally **over-estimates** the true class-prior.
 - Tails of $p(\mathbf{x}|y = -1)$ affect the solution.



Non-Identifiability of the Class-Prior

30

Blanchard et al. (JMLR2010)

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

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

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

Non-estimable

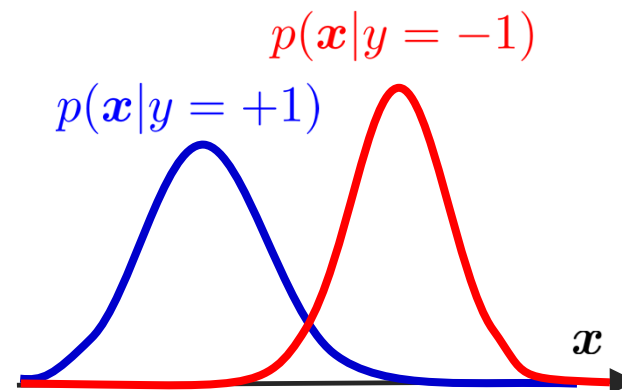
■ When $p(\mathbf{x}|y = +1)$ and $p(\mathbf{x}|y = -1)$ are overlapped, they may share **some common component**.

● Its proportion can be **arbitrarily changed**.

● Indeed, any $\theta \in \left\{ \exists p'(\mathbf{x}), p(\mathbf{x}) = \theta p(\mathbf{x}|y = +1) + (1 - \theta)p'(\mathbf{x}) \right\}$

can be a valid solution.

■ We need a **reasonable assumption** to obtain a meaningful solution!



1. Semi-supervised class-prior shift adaptation
2. **Positive-unlabeled classification**
 - A) Basic solution
 - B) Identifiability
 - C) **Density ratio estimation with anchor points**
 - D) Partial distribution matching with irreducibility
 - E) Regrouping without irreducibility
 - F) Summary
3. Conclusions



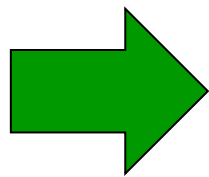
Class-Prior Estimation under Anchor Point Assumption

Sugiyama et al. (MIT Press, in press)

- Assume there exists an **anchor point** in $\{\mathbf{x}_i^P\}_{i=1}^{n_P}$:

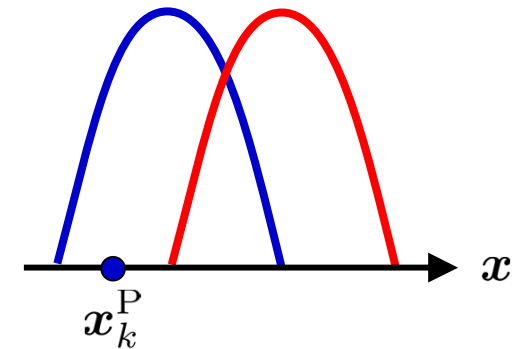
- For some $k \in \{1, \dots, n_P\}$,

$$p(\mathbf{x}_k^P | y = +1) > 0 \text{ and } p(\mathbf{x}_k^P | y = -1) = 0.$$



$$\pi = \frac{p(\mathbf{x}_k^P)}{p(\mathbf{x}_k^P | y = +1)}$$

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



- Density ratio estimation gives

$$\hat{\pi} = \min_{i \in \{1, \dots, n_P\}} \hat{r}^{-1}(\mathbf{x}_i^P)$$

$$\hat{r}(\mathbf{x}) \approx \frac{p(\mathbf{x} | y = +1)}{p(\mathbf{x})}$$

- Simple and nice!

- But the anchor point assumption may be too strong.

1. Semi-supervised class-prior shift adaptation
2. **Positive-unlabeled classification**
 - A) Basic solution
 - B) Identifiability
 - C) Density ratio estimation with anchor points
 - D) **Partial distribution matching with irreducibility**
 - E) Regrouping without irreducibility
 - F) Summary
3. Conclusions



Partial Matching with Differentiable f -Divergence

du Plessis et al.
(ACML2015, MLJ2017)

$$\text{div}(\theta) = \text{Div}_f[p(\mathbf{x})||q'(\mathbf{x})] = \int p(\mathbf{x}) f\left(\frac{q'(\mathbf{x})}{p(\mathbf{x})}\right) d\mathbf{x}$$

$$q'(\mathbf{x}) = \theta p(\mathbf{x}|y = +1)$$

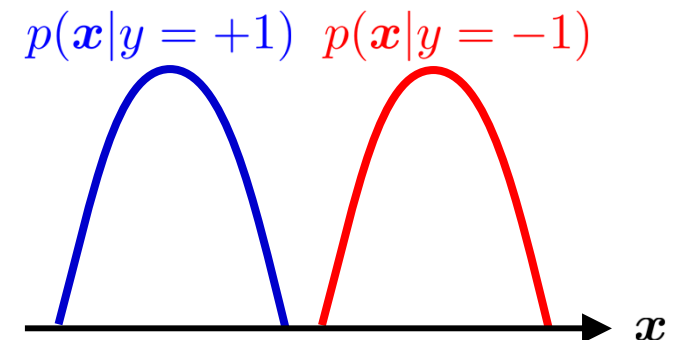
Suppose $f(t)$ has the minimum at $t \geq 1$.

- When $f(t)$ is **differentiable**,

$\text{div}'(\pi) = 0$ is necessary for $\pi = \underset{\theta \in [0,1]}{\text{argmin}} \text{div}(\theta)$.

$$\text{div}'(\pi) = \int f'(p(y = +1|\mathbf{x})) p(\mathbf{x}|y = +1) d\mathbf{x}$$

- $\text{div}'(\pi) = 0$ if and only if
 - Two classes are non-overlapped,
 - and $f'(1) = 0$ (e.g., Pearson div).



With Non-Differentiable f -Divergence

35

$$\text{div}(\theta) = \text{Div}_f[p(\mathbf{x}) \| q'(\mathbf{x})] = \int p(\mathbf{x}) f\left(\frac{q'(\mathbf{x})}{p(\mathbf{x})}\right) d\mathbf{x}$$

Suppose $f(t)$ has the minimum at $t \geq 1$.

- When $f(t)$ is **non-differentiable** at $t = 1$,
 $0 \in \partial \text{div}(\pi)$ is necessary for $\pi = \underset{\theta \in [0,1]}{\text{argmin}} \text{div}(\theta)$.

∂ : subdifferential

$$\partial \text{div}(\pi) = \int \partial f(p(y = +1 | \mathbf{x})) p(\mathbf{x} | y = +1) d\mathbf{x}$$

- $0 \in \partial \text{div}(\pi)$ even when two classes are overlapped, if

- $f(t)$ is **penalized** as $f(t) \leftarrow \begin{cases} f(t) & (t \leq 1) \\ \infty & (t > 1) \end{cases}$,

- and the **irreducibility assumption** holds: Blanchard et al. (JMLR2010)

- $p(\mathbf{x} | y = +1)$ is not a component of $p(\mathbf{x} | y = -1)$.

Irreducibility and Anchor Points

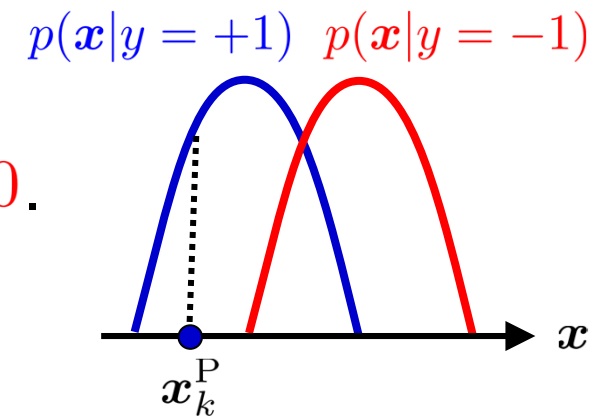
■ Irreducibility: Blanchard et al. (JMLR2010)

- $p(\mathbf{x}|y = +1)$ is not a component of $p(\mathbf{x}|y = -1)$.

$$\pi = \sup \left\{ \pi' \mid \exists p'(\mathbf{x}), p(\mathbf{x}) = \pi' p(\mathbf{x}|y = +1) + (1 - \pi') p'(\mathbf{x}) \right\}$$

■ Anchor points: Liu & Tao (IEEE-TPAMI2015)

- For some $k \in \{1, \dots, n_P\}$,
 $p(\mathbf{x}_k^P | y = +1) > 0$ and $p(\mathbf{x}_k^P | y = -1) = 0$.



■ Irreducibility holds if and only if at least one anchor point exists:

- Density ratio based method uses the anchor point **explicitly**.
- Partial matching only assumes its existence **implicitly**.

■ Therefore, the required assumption is weaker!

Practical Choice of f :

Penalized L1-Distance

du Plessis et al. (ACML2015, MLJ2017)

$$\text{pen}L_1(\theta) = \int p(\mathbf{x}) f\left(\frac{q'(\mathbf{x})}{p(\mathbf{x})}\right) d\mathbf{x}$$

$$f(t) = \begin{cases} 1 - t & (t \leq 1) \\ \infty & (t > 1) \end{cases}$$

- Regularized least-squares density ratio estimation gives a divergence approximator **analytically** as

$$\widehat{\text{pen}L_1}(\theta) = \frac{1}{\lambda} \sum_{b=1}^B \max(0, \beta_b) \beta_b - \theta + 1$$

$\lambda > 0$: Regularization parameter

$$\beta_b = \frac{\theta}{n_P} \sum_{i=1}^{n_P} \varphi_b(\mathbf{x}_i^P) - \frac{1}{n_U} \sum_{j=1}^{n_U} \varphi_b(\mathbf{x}_j^U)$$

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

- Model of density ratio $q'(\mathbf{x})/p(\mathbf{x})$:

$$s(\mathbf{x}) = \sum_{b=1}^B \alpha_b \varphi_b(\mathbf{x}) + 1$$

$\varphi_b(\mathbf{x}) \geq 0$: Basis function

$\alpha_b \geq 0$: Parameter

Implementation and Analysis

$$\widehat{\text{pen}L_1}(\theta) = \frac{1}{\lambda} \sum_{b=1}^B \max(0, \beta_b) \beta_b - \theta + 1$$

- **Algorithm:** Find a minimizer w.r.t. $\theta \in [0, 1]$.
 - **Computationally very efficient!**
- **Optimal convergence rate is achieved!**

$$\text{pen}L_1(\hat{\pi}) - \text{pen}L_1(\pi) = \mathcal{O}_p(1/\sqrt{n_P} + 1/\sqrt{n_U})$$

$$\hat{\pi} = \underset{0 \leq \theta \leq 1}{\text{argmin}} \widehat{\text{pen}L_1}(\theta)$$

$$\beta_b = \frac{\theta}{n_P} \sum_{i=1}^{n_P} \varphi_b(\mathbf{x}_i^P) - \frac{1}{n_U} \sum_{j=1}^{n_U} \varphi_b(\mathbf{x}_j^U) \quad \begin{matrix} \{\mathbf{x}_i^P\}_{i=1}^{n_P} \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}|y=+1) \\ \{\mathbf{x}_j^U\}_{j=1}^{n_U} \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}) \end{matrix}$$

$$\text{pen}L_1(\theta) = \int p(\mathbf{x}) f\left(\frac{q'(\mathbf{x})}{p(\mathbf{x})}\right) d\mathbf{x}$$

$$f(t) = \begin{cases} 1-t & (t \leq 1) \\ \infty & (t > 1) \end{cases} \quad \begin{matrix} p(\mathbf{x}) = \pi p(\mathbf{x}|y=+1) + (1-\pi)p(\mathbf{x}|y=-1) \\ q'(\mathbf{x}) = \theta p(\mathbf{x}|y=+1) \end{matrix}$$

- **However, there is no way to assess irreducibility in practice.**

1. Semi-supervised class-prior shift adaptation
2. **Positive-unlabeled classification**
 - A) Basic solution
 - B) Identifiability
 - C) Density ratio estimation with anchor points
 - D) Partial distribution matching with irreducibility
 - E) **Regrouping without irreducibility**
 - F) Summary
3. Conclusions



Class-Prior Estimation without Irreducibility

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

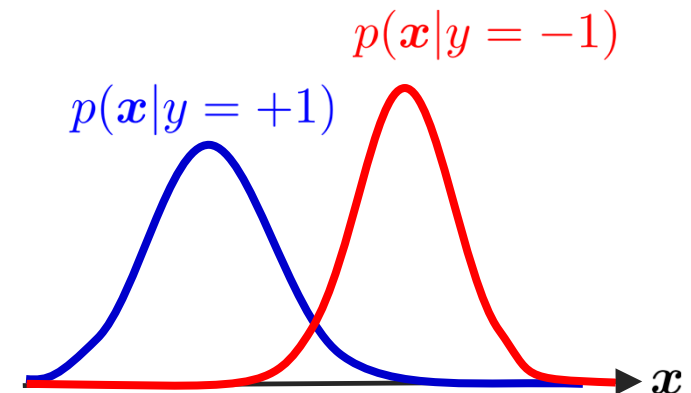
- Without irreducibility, any

$$\theta \in \left\{ \exists p'(\mathbf{x}), p(\mathbf{x}) = \theta p(\mathbf{x}|y = +1) + (1 - \theta)p'(\mathbf{x}) \right\}$$

can be a valid solution, due to common components.

- Partial matching actually gives its maximum value.

- Can we mitigate the positive bias in the absence of irreducibility?



Regrouping

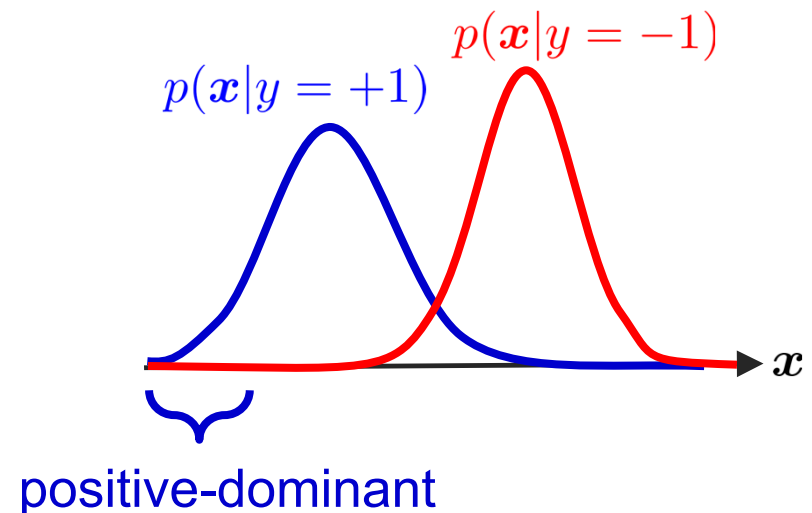
Yao et al. (arXiv2020)

41

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

- **Idea:** Regroup a small **positive-dominant region** to be **fully positive**.

- By this regrouping,
 - π is slightly increased,
 - but **irreducibility is satisfied!**



- How can we find a positive-dominant region?

Implementation

- Consider **PU label** $s \in \{0, 1\}$:
 - If x is P (or U), $s = 1$ (or $s = 0$).
- Train a **“non-traditional” classifier** $\hat{p}(s|\mathbf{x})$:
 - Usual supervised classification from $\{\mathbf{x}_i^P\}_{i=1}^{n_P}, \{\mathbf{x}_i^U\}_{i=1}^{n_U}$.
- Select some unlabeled samples that have the highest positive-confidence: $\hat{p}(s = 1|\mathbf{x}_j^U)$
- Copy them and give positive labels.
- Solve the converted class-prior estimation problem.

1. Semi-supervised class-prior shift adaptation
2. **Positive-unlabeled classification**
 - A) Basic solution
 - B) Identifiability
 - C) Density ratio estimation with anchor points
 - D) Partial distribution matching with irreducibility
 - E) Regrouping without irreducibility
 - F) **Summary**
3. Conclusions



Summary: PU Classification

- There is a nice **empirical risk minimization** method, given class-prior $\pi = p(y = +1)$ can be estimated.
- However, the class-prior is **not identifiable** in general.
- Simple **density ratio estimation** solution:
 - Use anchor points (i.e., 100% positive), which may be strong.
- Computationally efficient **penL₁-div partial matching**.
 - Without **irreducibility** (P-density is not part of N-density), its solution is positively biased.
 - Existence of anchor points is sufficient, but not assessable.
- **Regrouping**:
 - By preprocessing of data, the positive bias can be reduced.

Contents

1. Semi-supervised class-prior shift adaptation
2. Positive-unlabeled classification
3. **Conclusions**



Summary:

Mixture Proportion Estimation

- Many applications in machine learning:
 - **Class-prior shift adaptation**: Importance weight $p_{te}(y)$
 - Identifiability allows naïve distribution matching to solve.
 - **Positive-unlabeled (PU) classification**: Class-prior $p(y)$
 - Non-identifiability posed significant challenges.
 - **Noisy label classification**: Noise transition $p(\bar{y}|y)$

$$\underbrace{p(\bar{y}|\mathbf{x})}_{\text{Observed}} = \sum_y \underbrace{p(\bar{y}|y)}_{\text{Non-observed}} \underbrace{p(y|\mathbf{x})}_{\text{Non-observed}}$$

y : Clean class label

\bar{y} : Noisy class label

- Multiple non-identifiability is even more challenging!

Challenge:

Overcoming Non-Identifiability

■ **Identifiability conditions** have been investigated:

- Irreducibility, anchor set, anchor points...

Blanchard et al. (JMLR2010)

Liu & Tao (IEEE-TPAMI2015)

■ However, these identifiability conditions may not be satisfied in practice.

■ Even without identifiability, it is promising to

- Reduce estimation bias by **regrouping** (in PU classification). Yao et al. (arXiv2020)
- Use a weaker “**sufficiently scattered**” assumption (in noisy-label classification). Li et al. (ICML2021)

Challenge: Towards Better Machine Learning (ML)

■ The estimated proportion is later used in ML tasks.

■ Current approach is **two-step**:

● Estimate the mixture proportion.

$$\hat{\pi} = \operatorname{argmin}_{\pi} \operatorname{MPE}(\pi)$$

● Use the estimated proportion to solve the target ML problem.

$$\hat{f} = \operatorname{argmin}_f \operatorname{ML}(f, \hat{\pi})$$

■ **1st step is preformed without regards to 2nd step.**

■ Combining them into **one-step** is more promising:

$$\hat{f} = \operatorname{argmin}_f \min_{\pi} \operatorname{ML} \& \operatorname{MPE}(f, \pi)$$

● Alternate optimization.

Kato et al. (arXiv2018)

● Joint upper-bound optimization.

Zhang et al. (ACML2020, LNCS2021)

● Dynamic stochastic optimization.

Xia et al. (NeurIPS2019)

Fang et al. (NeurIPS2020)

Zhang et al. (ICML2021)

Grateful to Great Collaborators!



| | | |
|--|--------|---|
| Team leader Masashi Sugiyama | PAGE > | Research scientist Gang Niu |
| Research scientist Fumiko Kawasaki | | Research scientist Takahiro Mimori |
| Postdoctoral researcher Voot Tangkaratt | | Postdoctoral researcher Jingfeng Zhang |
| Postdoctoral researcher Jiaqi Lyu | | Postdoctoral researcher Shuo Chen |
| Technical scientist Yuka Mori | | Technical staff I Masashi Ugawa |
| Senior visiting scientist Shinichi Nakajima | | Visiting scientist Yoshihiro Nagano |
| Visiting scientist Florian Yger | | Visiting scientist Takashi Ishida |
| Visiting scientist Miao Xu | | Visiting scientist Takayuki Osa |
| Visiting scientist Bo Han | | Visiting scientist Daichi Noborio |
| Visiting scientist Yuko Kuroki | | Visiting scientist Hisashi Yoshida |
| Visiting scientist Ryohei Kasai | | Visiting scientist Feng Liu |
| Visiting scientist Lei Feng | | Visiting scientist Tongliang Liu |
| Junior research associate Takeshi Teshima | | Junior research associate Yifan Zhang |
| Part-time worker I Zhenghang Cui | | Part-time worker I Han Bao |
| Part-time worker I Masahiro Fujisawa | | |

and many interns
over the world!

- Professor
 - [Masashi Sugiyama](#) (Complexity, Computer, Information, RIKEN)
- Lecturer
 - [Naoto Yokoya](#) (Complexity, Computer, Information, RIKEN)
 - [Takashi Ishida](#) (Complexity, Computer, Information)
- Project Lecturer
 - [Nobutaka Ito](#) (Complexity)
- Associate professor (to [Sato Lab](#) from April 2020)
 - [Issai Sato](#) (Computer, Information, Complexity, RIKEN)
- Project Assistant Professor
 - Chao-Kai Chiang (Complexity)
- Project Researcher (Postdoctoral Researcher)
 - [Yoshihiro Nagano](#) (Complexity)
 - Dongxian Wu (Complexity)
- Project Specialist
 - Yuko Kawashima (Complexity)
 - Soma Yokoi (Complexity)
 - Fumi Sato (Complexity)
- Doctor Student
 - Seiya Tokui (Computer) * [Sato lab](#).
 - Shinji Nakadai (Computer)
 - [Kento Nozawa](#) (Complexity) * [Sato lab](#).
 - Kento Suzuki (Complexity)
 - [Han Bao](#) (Computer)
 - [Zhenghang Cui](#) (Computer) * [Sato lab](#).
 - [Liyuan Xu](#) (Computer)
 - [Takeshi Teshima](#) (Complexity)
 - Ryuichi Kiryo (Computer)
 - Masahiro Fujisawa (Complexity) * [Sato lab](#).
 - [Jongyeong Lee](#) (Computer)
 - Tianyi Zhang (Complexity)
 - [Yivan Zhang](#) (Computer)
 - Riou Charles (Computer)
 - [Vallappa Chockalingam](#) (Computer)
 - Tongtong Fang (Complexity)
 - Boyo Chen (Complexity)
 - Xiaoyu Dong (Complexity)
 - Yujie Zhang (Complexity)
 - Xinqiang Cai (Complexity)
 - Jian Song (Complexity)
 - Wanshui Gan (Complexity)
- Master Student
 - Atsushi Ito (Complexity)
 - Shinji Kawakami (Complexity) * [Sato lab](#).
 - Wataru Ohtori (Computer)
 - Tokio Kajitsuka (Computer)
 - Shota Nakajima (Computer)
 - Takahiro Suzuki (Computer) * [Sato lab](#).
 - Kei Mukaiyama (Computer) * [Sato lab](#).
 - Mingcheng Hou (Computer) * [Sato lab](#).
 - Hyunggyu Park (Complexity) * [Sato lab](#).
 - Yuting Tang (Complexity)
 - Shintaro Nakamura (Complexity)
 - Xujie Wang (Complexity)
 - Toshiki Kodera (Computer)
 - Yuma Aoki (Computer)
 - Jiahuan Li (Computer)
 - Chengwei Liang (Computer)
 - Huanjian Zhou (Complexity)
 - Kun Yang (Complexity)
 - Takanori Shirasaka (Complexity)
 - Reo Iizuka (Complexity)
 - Xiaomou Hou (Complexity)
- Bachelor Student
 - Kanma Noda (Information Science)
 - Kento Yamamoto (Information Science)
 - Kazuki Ota (Information Science)
 - Kazutaka Yahiro (Information Science)
 - Hikaru Fujita (Information Science)
- Research Student
 - [Kenny Song](#) (Computer)
 - Or Raveh (Computer)
 - Johannes Ackermann (Computer)
 - Anan Methasate (Computer)
 - Cemal Erat (Computer)

