

Importance-Weighting Approach to Distribution Shift Adaptation

Masashi Sugiyama

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

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



Learning under Distribution Shifts

2

■ Given:

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

\mathbf{x} : Input

y : Output

■ Goal:

- Learn predictor $y = f(\mathbf{x})$ minimizing the **test risk** (with some additional data from the test domain).

$$\min_f R(f)$$

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

ℓ : loss

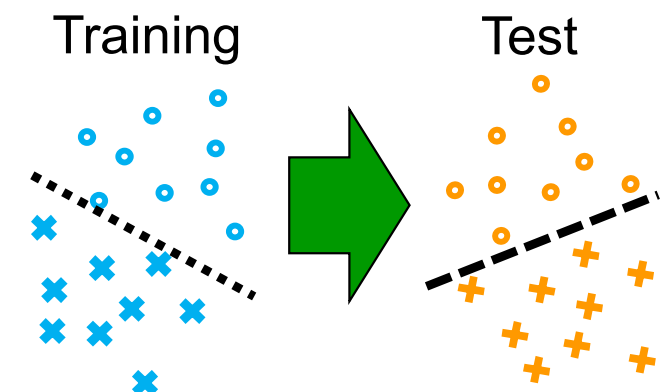
■ Challenge:

- Overcome **changing distributions!**

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

■ **Non-stationary** of the environments.

■ Sample selection bias due to **privacy** concerns.





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 Quiñero 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.

■ “Distribution Shift” workshops at NeurIPS2021, 2022, and 2023.

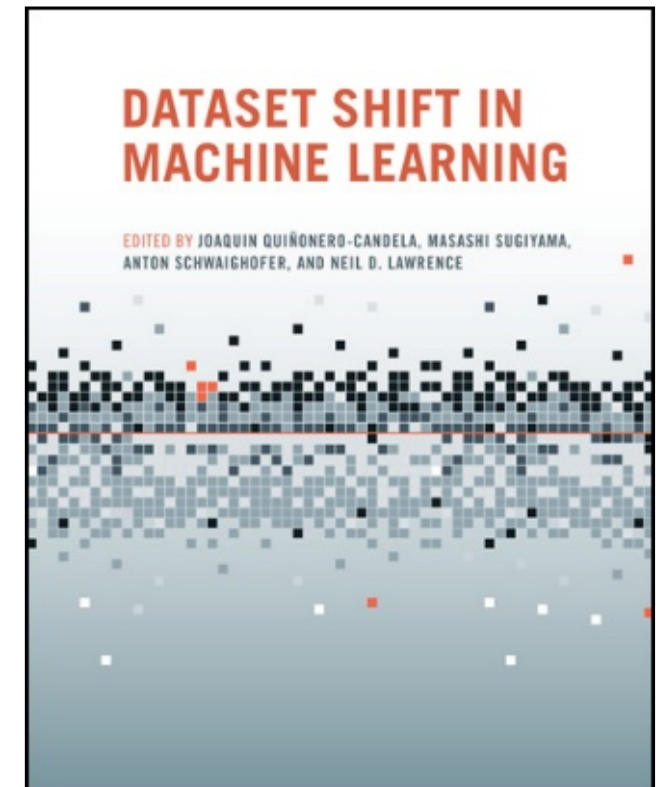
NeurIPS 2023 Workshop on Distribution Shifts (DistShift)

New Frontiers with Foundation Models

December 2023, New Orleans, USA

New Orleans Convention Center

Quiñero-Candela, Sugiyama, Schwaighofer & Lawrence (Eds.), [Dataset Shift in Machine Learning](#), MIT Press, 2009.





Contents

1. Importance Weighting
2. Continuous Shifts
3. Ongoing Challenges

Types of Distribution Shifts

\mathbf{x} : Input

y : Output

■ Joint shift:

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

■ Covariate shift:

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

■ Label shift:

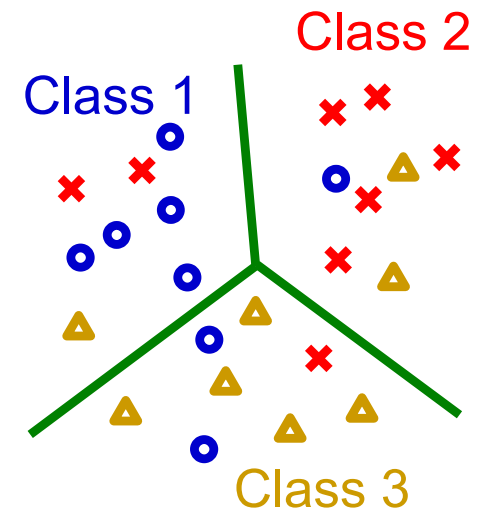
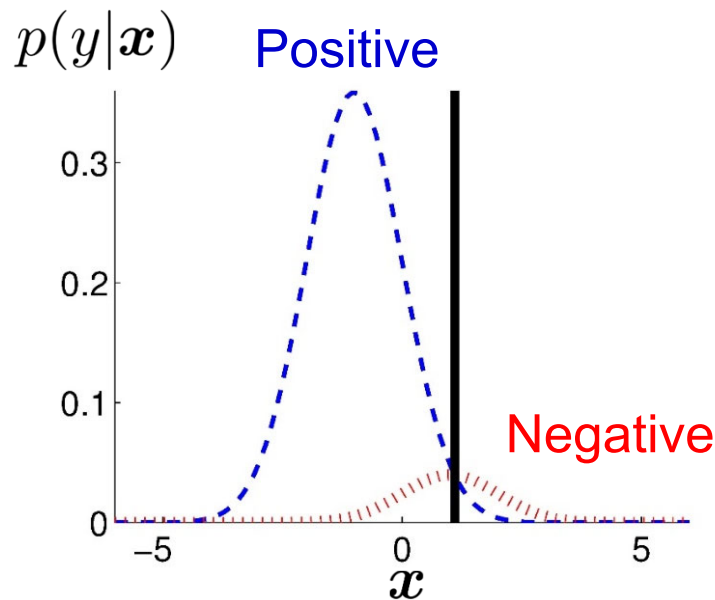
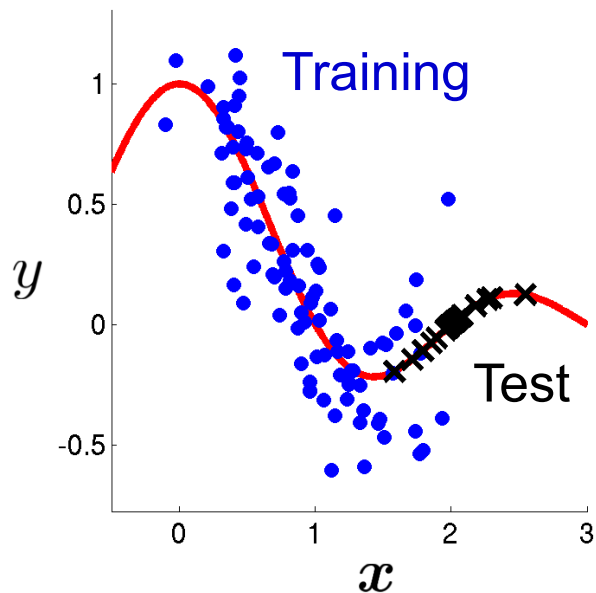
$$p_{\text{tr}}(y) \neq p_{\text{te}}(y)$$

■ Output noise:

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

■ Class-conditional shift:

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



Covariate Shift Adaptation

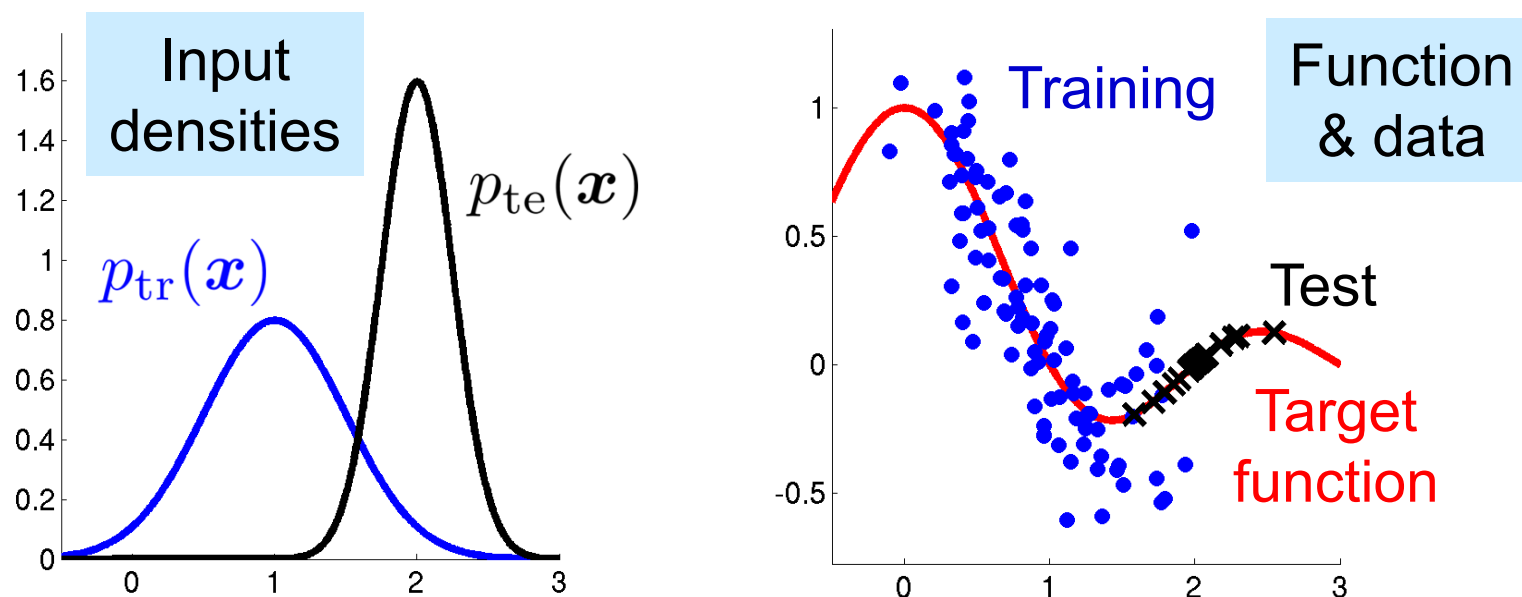
Shimodaira (JSPI2000)

- Training and test **input** distributions are different,

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

but the **output-given-input** distribution is unchanged:

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



- 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:

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

Importance-Weighted Training

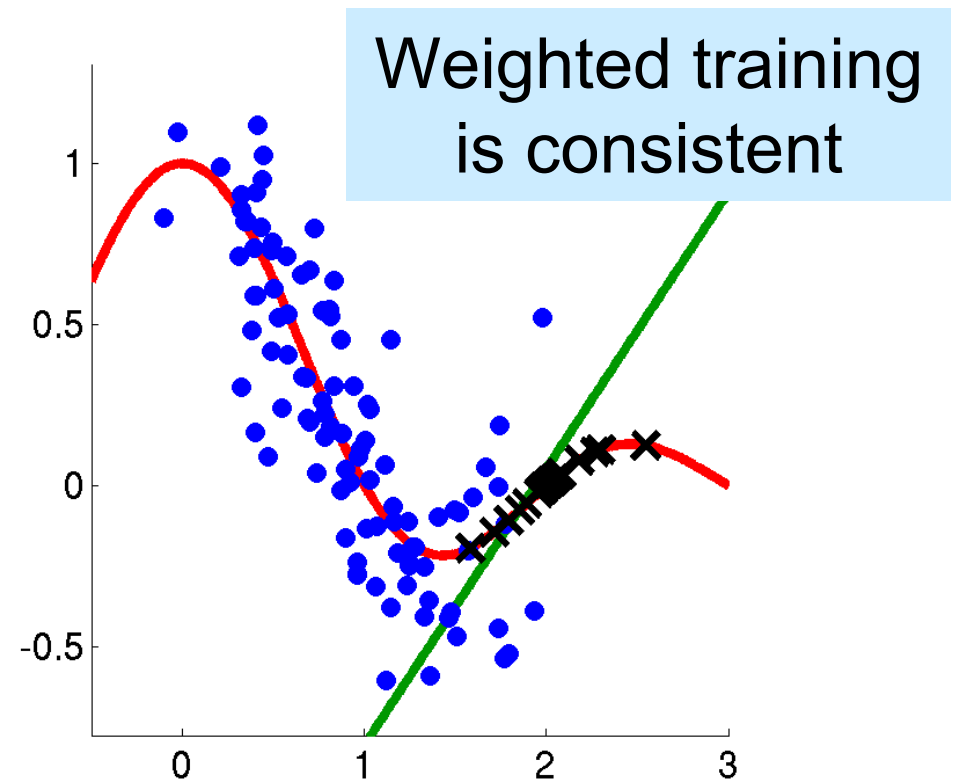
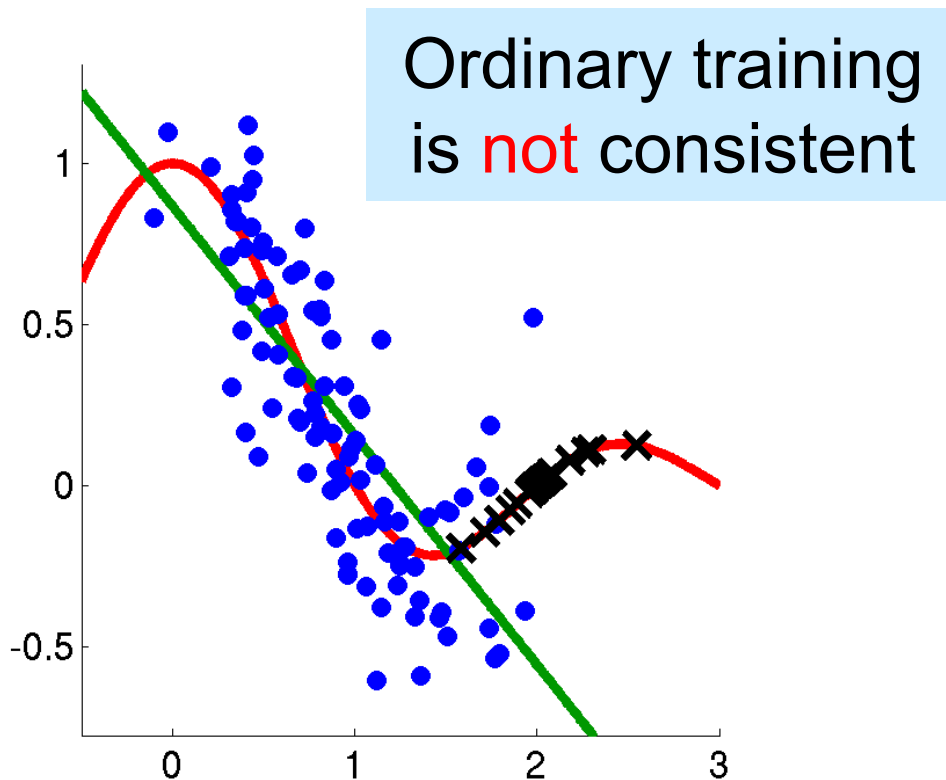
7

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

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

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

Importance



■ How do we estimate the importance?

Direct Importance Estimation

8

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

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

- **Goal:** estimate the density ratio $\frac{p_{\text{te}}(\mathbf{x})}{p_{\text{tr}}(\mathbf{x})}$

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

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

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

- **Least-squares importance fitting (LSIF):**

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

Kanamori, Hido
& Sugiyama
(NeurIPS2008)

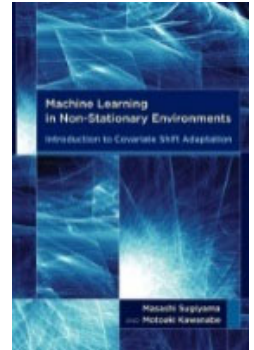
$$\operatorname{argmin}_r \mathbb{E}_{p_{\text{tr}}(\mathbf{x})} \left[\left(r(\mathbf{x}) - \frac{p_{\text{te}}(\mathbf{x})}{p_{\text{tr}}(\mathbf{x})} \right)^2 \right] = \operatorname{argmin}_r \mathbb{E}_{p_{\text{tr}}(\mathbf{x})}[r(\mathbf{x})^2] - 2\mathbb{E}_{p_{\text{te}}(\mathbf{x})}[r(\mathbf{x})]$$

- Extendable to Bregman divergences:

Sugiyama, Suzuki & Kanamori
(AISM2012)

$$\operatorname{argmin}_r \operatorname{Breg}_{\psi} \left(\frac{p_{\text{te}}(\mathbf{x})}{p_{\text{tr}}(\mathbf{x})} \parallel r \right) = \operatorname{argmin}_r \mathbb{E}_{p_{\text{tr}}(\mathbf{x})}[\partial\psi(r(\mathbf{x}))r(\mathbf{x}) - \psi(r(\mathbf{x}))] - \mathbb{E}_{p_{\text{te}}(\mathbf{x})}[\partial\psi(r(\mathbf{x}))]$$

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



Classical approaches are **two steps**:

1. Importance weight estimation (e.g., LSIF):

$$\hat{r} = \operatorname{argmin}_r \hat{J}_1(r) \quad J_1(r) = \mathbb{E}_{p_{\text{tr}}(\mathbf{x})} \left[\left(r(\mathbf{x}) - \frac{p_{\text{te}}(\mathbf{x})}{p_{\text{tr}}(\mathbf{x})} \right)^2 \right]$$

2. Importance-weighted predictor training:

$$\hat{f} = \operatorname{argmin}_f \hat{J}_2^\ell(f, \hat{r}) \quad J_2^\ell(f, r) = \mathbb{E}_{p_{\text{tr}}(\mathbf{x}, y)} [r(\mathbf{x}) \ell(f(\mathbf{x}), y)]$$

For $\ell_{\text{te}} \leq 1, \ell_{\text{tr}} \geq \ell_{\text{te}}, r \geq 0$, the test risk is bounded as

$$\frac{1}{2} R_{\ell_{\text{te}}}(f)^2 \leq J_{\ell_{\text{tr}}}(f, r) \quad J_\ell(f, r) = J_1(r) + J_2^\ell(f, r)$$

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

Zhang, Yamane, Lu & Sugiyama
(ACML2020, SNCS2021)

- **Joint upper-bound minimization:** $\hat{f} = \operatorname{argmin}_f \min_{r \geq 0} \hat{J}_{\ell_{\text{tr}}}(f, r)$



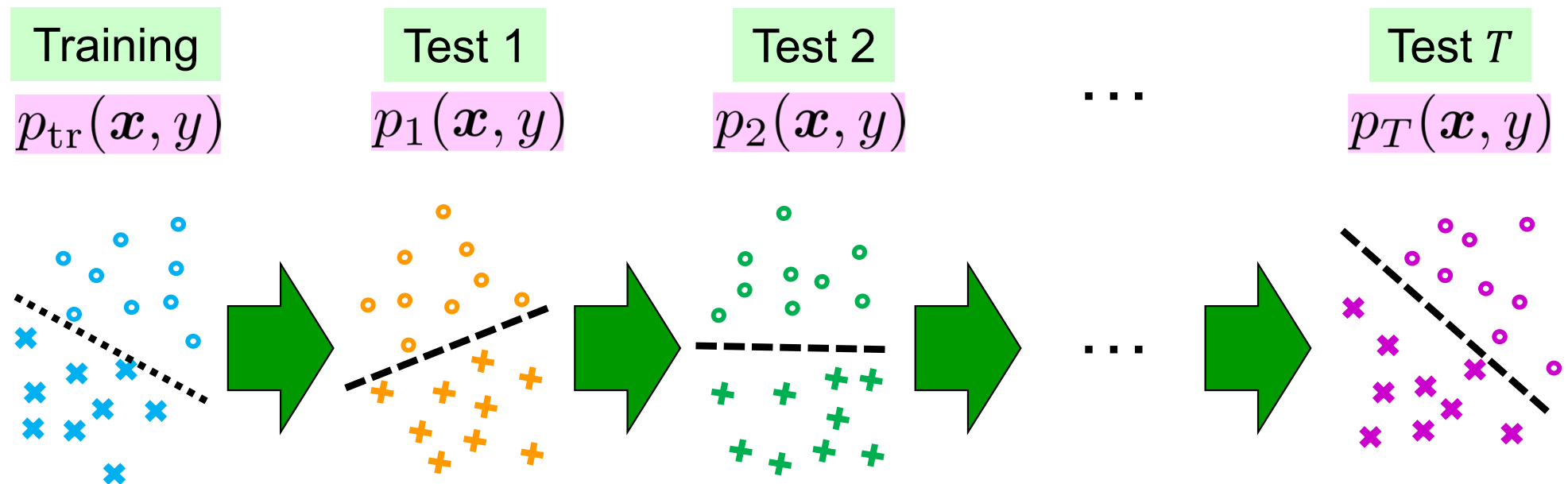
Contents

10

1. Importance Weighting
2. Continuous Shifts
3. Ongoing Challenges

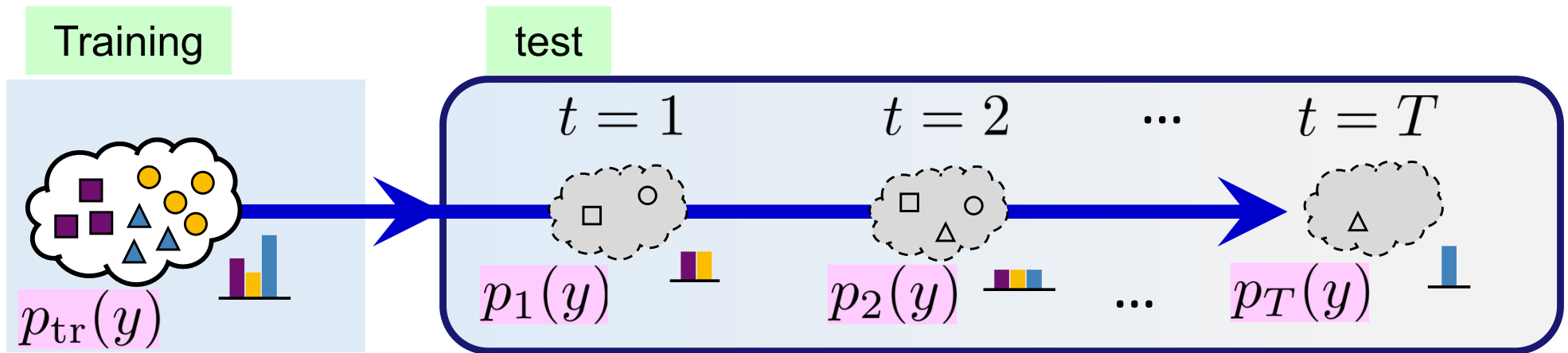
Continuous Distribution Shifts

- So far, we focused on a **fixed test domain**:
 - We trained a predictor to match the test domain.
- However, **test domains can change over time**.



- **Goal**: Obtain predictor \hat{f}_t that works well for $p_t(\mathbf{x}, y)$.

$$R_t(f) = \mathbb{E}_{p_t(\mathbf{x}, y)}[\ell(f(\mathbf{x}), y)] \quad t = 1, \dots, T$$



- **Class-priors** $p_t(y)$ **change** arbitrarily over time, but **class-conditionals stay** unchanged:

$$p_t(\mathbf{x}|y) = p_{tr}(\mathbf{x}|y) \quad t = 1, \dots, T$$

- **Given:**

- (Large) **labeled** training data: $\{(\mathbf{x}_i^{tr}, y_i^{tr})\}_{i=1}^{n_{tr}} \stackrel{\text{i.i.d.}}{\sim} p_{tr}(\mathbf{x}, y)$
- (Small) **unlabeled** test data: $\{\mathbf{x}_i^{(t)}\}_{i=1}^{n_t} \stackrel{\text{i.i.d.}}{\sim} p_t(\mathbf{x})$
 $t = 1, \dots, T$

ATLAS (Adapting To Label Shift)

13

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

$$\min_f \sum_{i=1}^{n_{\text{tr}}} \frac{\hat{p}_t(y_i^{\text{tr}})}{\hat{p}_{\text{tr}}(y_i^{\text{tr}})} \ell(f(\mathbf{x}_i^{\text{tr}}), y_i^{\text{tr}})$$

- Batch importance weighing requires **retraining** in each time step.
- Can we make it computationally more efficient?

- **Online learning!**

Hazan (2016)

- We use **online convex optimization**, assuming

- convex loss ℓ (e.g., logistic),
- linear model $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$, $\mathbf{w} \in \mathcal{W}$.

$\Pi_{\mathcal{W}}$: projection

$$\mathbf{w}_t = \Pi_{\mathcal{W}} \left[\mathbf{w}_{t-1} - \eta \nabla \hat{R}_{t-1}(\mathbf{w}_{t-1}) \right]$$

$\eta > 0$: step size

- We use **black box shift estimation** for class priors.

Lipton+ (ICML2018)

Choice of Step Size η

$$\mathbf{w}_t = \Pi_{\mathcal{W}} \left[\mathbf{w}_{t-1} - \eta \nabla \hat{R}_{t-1}(\mathbf{w}_{t-1}) \right]$$

■ If the speed of distribution shift is

- **slow**, η should be **small** to keep the previous classifier.
- **fast**, η should be **large** to quickly update the classifier.

■ How do we choose η in practice?

- **Ensemble learning!** Zhao+ (NeurIPS2020)

■ For $0 < \eta_1 < \dots < \eta_M$, we run M learners:

$$\mathbf{w}_t^{(m)} = \Pi_{\mathcal{W}} \left[\mathbf{w}_{t-1}^{(m)} - \eta_m \nabla \hat{R}_{t-1}(\mathbf{w}_{t-1}^{(m)}) \right] \quad \hat{R}_t(\mathbf{w}) = \frac{1}{n_{\text{tr}}} \sum_{i=1}^{n_{\text{tr}}} \frac{\hat{p}_t(y_i^{\text{tr}})}{\hat{p}_{\text{tr}}(y_i^{\text{tr}})} \ell(\mathbf{w}^\top \mathbf{x}_i^{\text{tr}}, y_i^{\text{tr}})$$

■ Final output is the **weighted average (cf. Hedge)**:

Freund+ (JCSS1997)

$$\mathbf{w}_t = \sum_{m=1}^M p_t^{(m)} \mathbf{w}_t^{(m)}$$

$$p_t^{(m)} \propto \exp \left(-\varepsilon \sum_{s=1}^{t-1} \hat{R}_s(\mathbf{w}_s^{(m)}) \right) \quad \varepsilon = \Theta \left(\sqrt{\frac{\ln M}{T}} \right)$$

■ **Shift intensity:** $V_T = \sum_{t=2}^T \sum_{y=1}^c |p_t(y) - p_{t-1}(y)| \geq \Theta(T^{-\frac{1}{2}})$

■ When V_T is **known:**

- **Dynamic regret** is minimized with step size $\eta = \Theta(V_T^{\frac{1}{3}} T^{-\frac{1}{3}})$:

$$\mathbb{E} \left[\underbrace{\sum_{t=1}^T R_t(\mathbf{w}_t)}_{\text{Risk of our model}} - \underbrace{\sum_{t=1}^T \min_{\mathbf{w} \in \mathcal{W}} R_t(\mathbf{w})}_{\text{Risk of the best model at each iteration}} \right] = \mathcal{O}(V_T^{\frac{1}{3}} T^{\frac{2}{3}})$$

Risk of our model

Risk of the best model at each iteration

$$R_t(\mathbf{w}) = \mathbb{E}_{p_{\text{tr}}(\mathbf{x}|y)p_t(y)} [\ell(\mathbf{w}^\top \mathbf{x}, y)]$$

■ Even when V_T is **unknown:**

- **ATLAS still achieves the same dynamic regret!**

- Number of learners: $M = 1 + \lceil \frac{1}{2} \log_2(1 + 2T) \rceil$

- Step size: $\eta_m = 2^{m-1} Z / \sqrt{T}, \quad m = 1, \dots, M$

■ If we have some **hints**, can we perform better?

- **Hint function:** $H_t(\mathbf{w}) \approx R_t(\mathbf{w}) = \mathbb{E}_{p_{\text{tr}}(\mathbf{x}|y)p_t(y)} [\ell(\mathbf{w}^\top \mathbf{x}, y)]$

■ **ATLAS-ADA:**

- $\mathbf{w}_{t-\frac{1}{2}} = \Pi_{\mathcal{W}} \left[\mathbf{w}_{t-1} - \eta \nabla \hat{R}_{t-1}(\mathbf{w}_{t-1}) \right]$ ← Same as ATLAS
- $\mathbf{w}_t = \operatorname{argmin}_{\mathbf{w} \in \mathcal{W}} \left[\eta H_t(\mathbf{w}) + \frac{1}{2} \|\mathbf{w} - \mathbf{w}_{t-\frac{1}{2}}\|^2 \right]$ ← Use the hint to match **the next loss**

■ **Dynamic regret:** $\mathbb{E} \left[\sum_{t=1}^T R_t(\mathbf{w}_t) - \sum_{t=1}^T \min_{\mathbf{w} \in \mathcal{W}} R_t(\mathbf{w}) \right] = \mathcal{O} \left(V_T^{\frac{1}{3}} G_T^{\frac{1}{3}} T^{\frac{1}{3}} \right)$

- **Reusability:** $G_T = \sum_{t=1}^T \mathbb{E} \left[\sup_{\mathbf{w} \in \mathcal{W}} \|\nabla \hat{R}_t(\mathbf{w}) - \nabla H_t(\mathbf{w})\|^2 \right] \leq \mathcal{O}(T)$

- **ATLAS-ADA is better and safe!** $\mathcal{O}(V_T^{\frac{1}{3}} T^{\frac{2}{3}})$

Continuous Covariate Shift

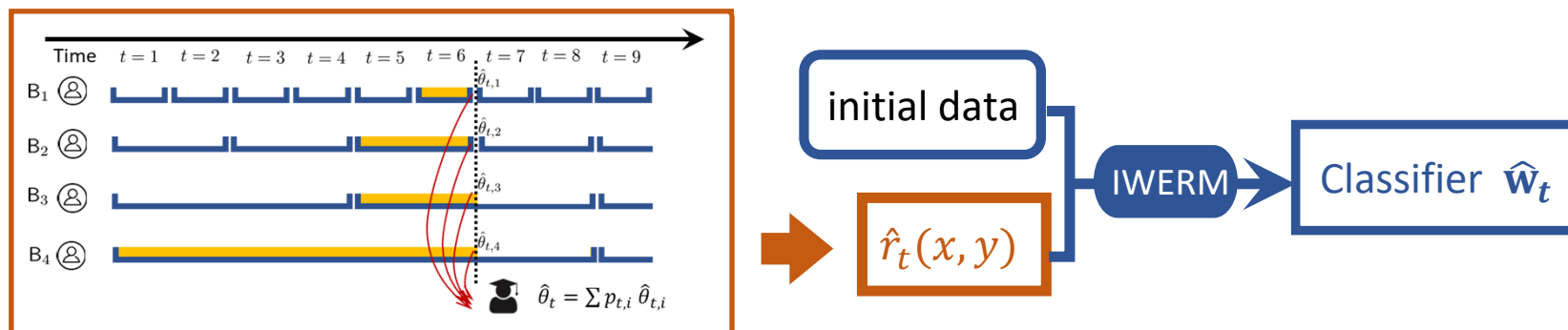
Zhang, Zhang, Zhao & Sugiyama (NeurIPS2023)

- Input density $p_t(\mathbf{x})$ changes arbitrarily over time, but output-given-input is unchanged: $p_{\text{tr}}(y|\mathbf{x}) = p_t(y|\mathbf{x})$
 $t = 1, \dots, T$

- Given:

- (Large) labeled training data: $\{(\mathbf{x}_i^{\text{tr}}, y_i^{\text{tr}})\}_{i=1}^{n_{\text{tr}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{tr}}(\mathbf{x}, y)$
- (Small) unlabeled test data: $\{\mathbf{x}_i^{(t)}\}_{i=1}^{n_t} \stackrel{\text{i.i.d.}}{\sim} p_t(\mathbf{x})$

- We use online density ratio estimation:





Contents

1. Importance Weighting
2. Continuous Shifts
3. Ongoing Challenges

- Many distribution shift works considers a specific **shift type** (e.g., covariate shift).

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

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

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

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

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

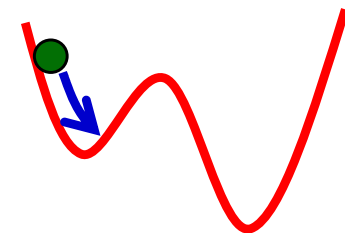
Mini-Batch-Wise Loss Matching

Fang, Lu, Niu & Sugiyama (NeurIPS2020)

Given:

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

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



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

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

Huang+
(NeurIPS2006)

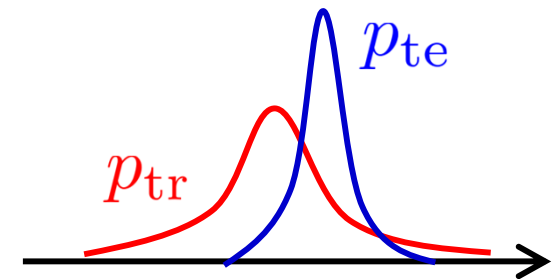
$$\frac{1}{\tilde{n}_{\text{tr}}} \sum_{i=1}^{\tilde{n}_{\text{tr}}} r_i \ell(f(\tilde{\mathbf{x}}_i^{\text{tr}}), \tilde{y}_i^{\text{tr}}) \approx \frac{1}{\tilde{n}_{\text{te}}} \sum_{j=1}^{\tilde{n}_{\text{te}}} \ell(f(\tilde{\mathbf{x}}_j^{\text{te}}), \tilde{y}_j^{\text{te}}) \quad r_i \approx \frac{p_{\text{te}}(\tilde{\mathbf{x}}_i^{\text{tr}},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

■ Limitation of importance weighting:

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

- The training domain must **cover** the test domain.

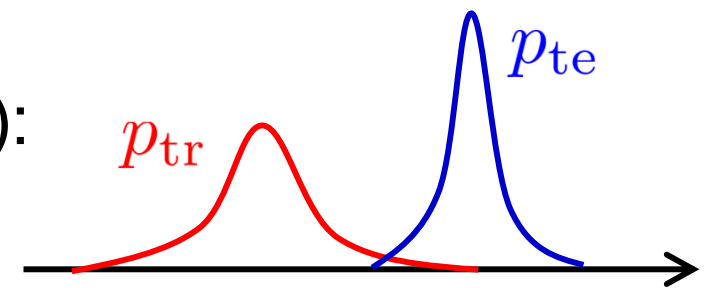
■ What if the test domain **sticks out** from the training domain?



■ Out-of-domain extension: Fang, Lu, Niu & Sugiyama (NeurIPS2023)

- Split training data into in-/out-domains by **outlier detection** (e.g., 1-class SVM):

$$\left\{ (\mathbf{x}_j^{\text{te}_{in}}, y_j^{\text{te}_{in}}) \right\}_{j=1}^{n_{\text{te}_{in}}} \quad \left\{ (\mathbf{x}_j^{\text{te}_{out}}, y_j^{\text{te}_{out}}) \right\}_{j=1}^{n_{\text{te}_{out}}}$$

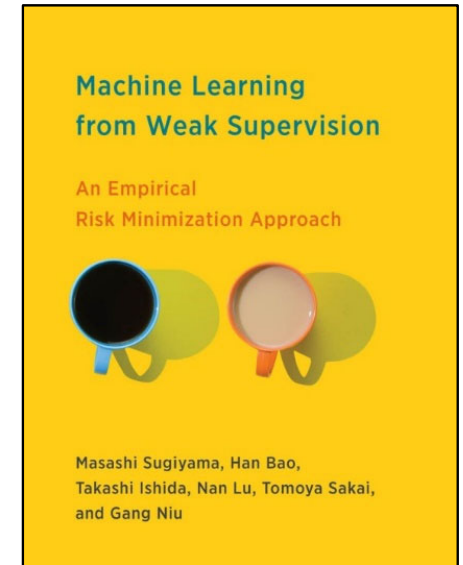


- Compute the loss separately:

$$\frac{n_{\text{te}_{in}}}{n_{\text{tr}} n_{\text{te}}} \sum_{i=1}^{n_{\text{tr}}} \frac{p_{te}(\mathbf{x}_i^{\text{tr}}, y_i^{\text{tr}})}{p_{tr}(\mathbf{x}_i^{\text{tr}}, y_i^{\text{tr}})} \ell(f(\mathbf{x}_i^{\text{tr}}), y_i^{\text{tr}}) + \frac{1}{n_{\text{te}}} \sum_{j=1}^{n_{\text{te}_{out}}} \ell(f(\mathbf{x}_j^{\text{te}_{out}}), y_j^{\text{te}_{out}})$$

Ongoing Challenges

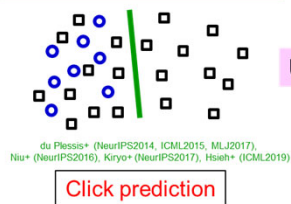
- For joint shift adaptation, requiring **labeled test data** is too strong.
- Can we use **weakly supervised learning**?
 - **Unbiased risk estimation** from weak supervision.
 - **Any loss, classifier, and optimizer** can be used.



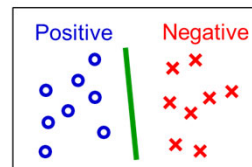
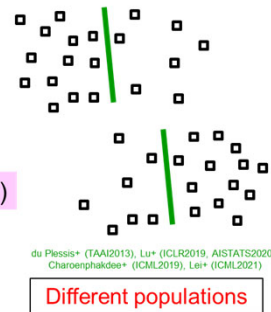
Sugiyama, Bao, Ishida, Lu, Sakai & Niu (MIT Press, 2022)

Weakly Supervised Classification (Binary)

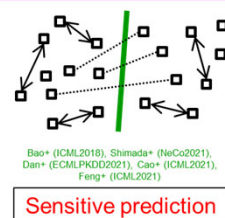
Positive-Unlabeled (PU)



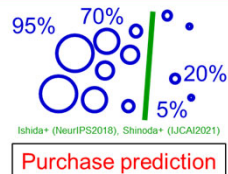
Unlabeled-Unlabeled (UU)



Similar-Dissimilar (SD)



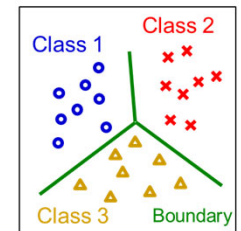
Positive-confidence (Pconf)



Weakly Supervised Classification (Multiclass)

Multi-class weak-labels:

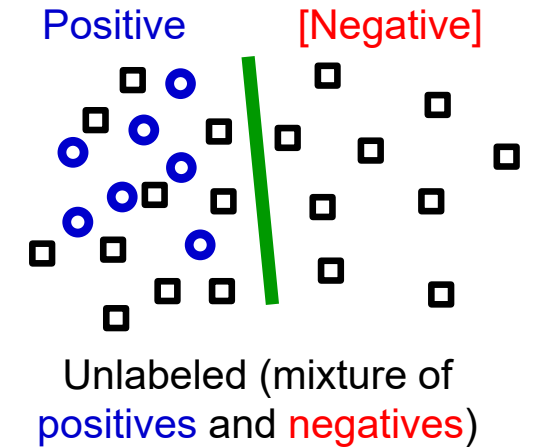
- **Complementary labels:** Specify a class that a pattern does **not** belong to ("not 1").
Ishida et al. (NIPS2017, ICML2019), Chou et al. (ICML2020)
- **Partial labels:** Specify a subset of classes that contains the correct one ("1 or 2").
Feng et al. (ICML2020, NeurIPS2020), Lv et al. (ICML2020)
- **Single-class confidence:** One-class data with full confidence ("1 with 60%, 2 with 30%, and 3 with 10%")
Cao et al. (arXiv2021)



Positive-Unlabeled (PU) Classification 23

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

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



- **Goal:** Obtain a risk-minimizing classifier

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

\mathbb{E} : expectation ℓ : loss $y = \{+1, -1\}$

- **Unbiased risk estimator:** du Plessis+ (NeurIPS2014, ICML2015)

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

$\pi = p(y = +1)$

- **Optimal convergence rate:** Niu, du Plessis, Sakai, Ma & Sugiyama (NIPS2016)

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

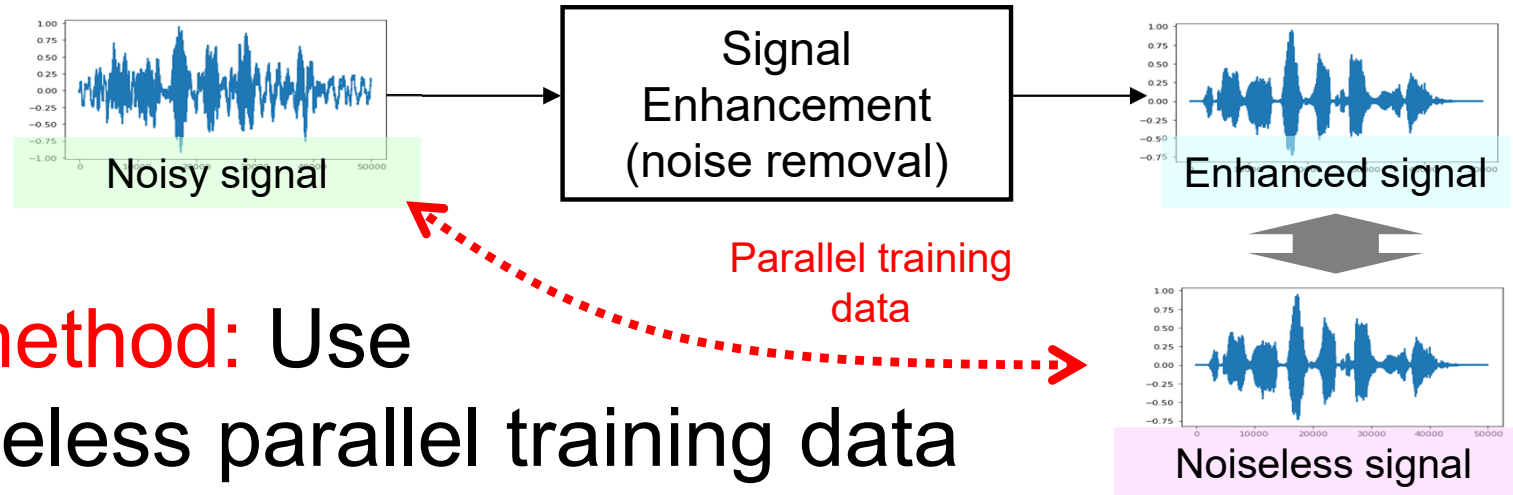
with probability $1 - \delta$

$$\hat{f}_{\text{PU}} = \operatorname{argmin}_{f \in \mathcal{F}} \hat{R}_{\text{PU}}(f)$$

$$f^* = \operatorname{argmin}_{f \in \mathcal{F}} R(f)$$

Signal Enhancement by PU Classification 24

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



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

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

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

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

