

Importance-Weighting Approach to Distribution Shift Adaptation

Masashi Sugiyama

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



Slides →

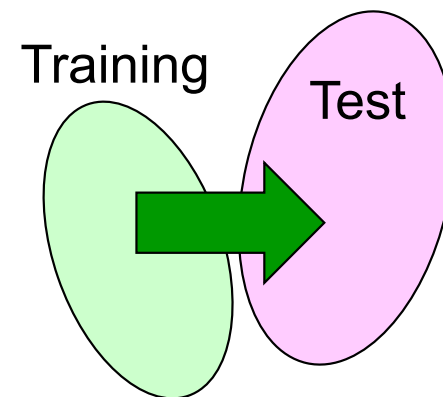
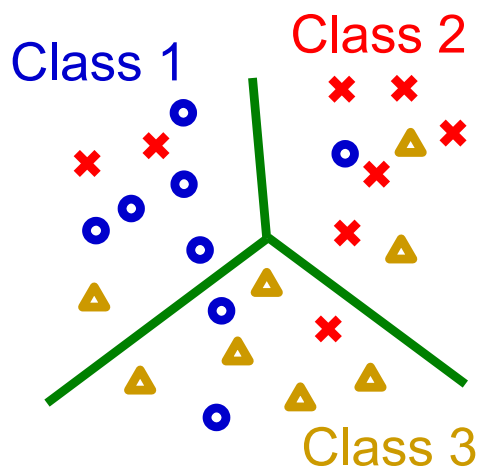
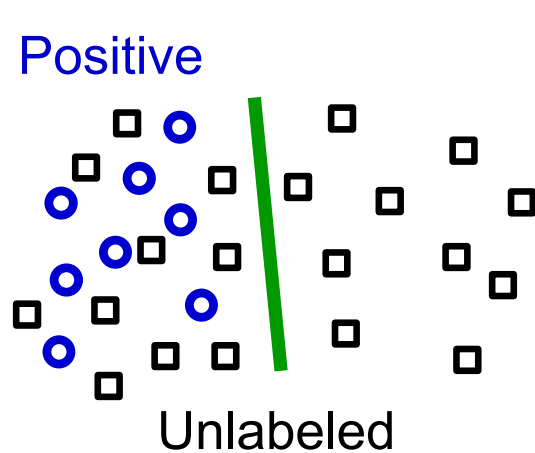


Reliable Machine Learning

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

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

■ Improving the reliability is an urgent challenge!





Contents

1. Weakly Supervised Learning

- A) Positive-Unlabeled Classification
- B) Various Extensions

2. Noisy-Label Learning

3. Transfer Learning

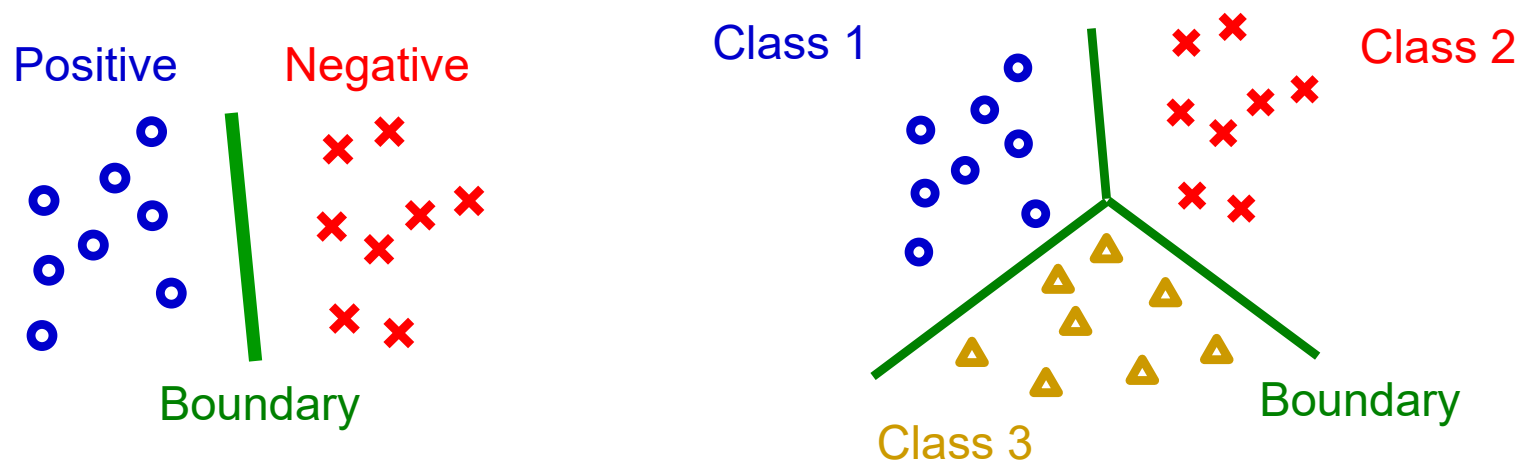
4. Towards More Reliable Learning

Slides →



Weakly Supervised Classification

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



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



Contents

1. Weakly Supervised Learning
 - A) Positive-Unlabeled Classification
 - B) Various Extensions
2. Noisy-Label Learning
3. Transfer Learning
4. Towards More Reliable Learning

Slides →



Positive-Unlabeled (PU) Classification

6

Li+ (IJCAI2003)

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

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

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

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

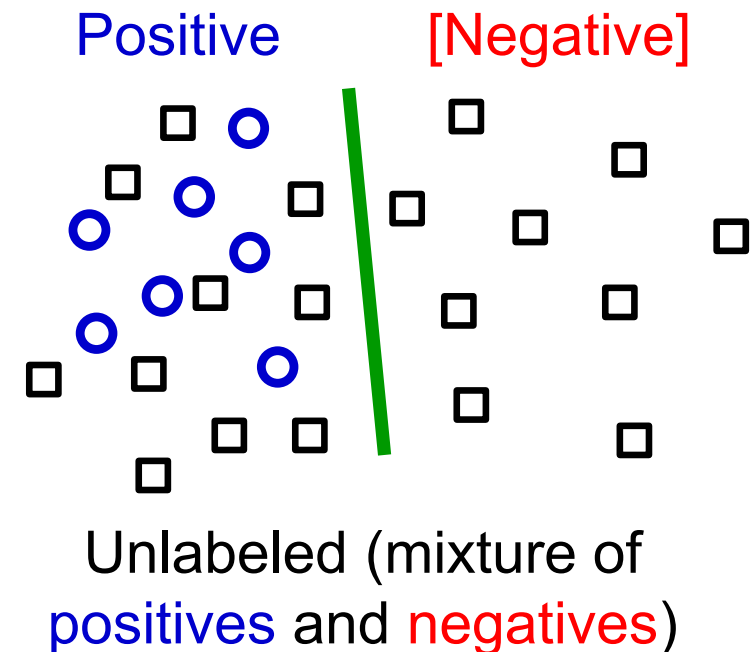
\mathbb{E} : expectation

ℓ : loss

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

Example: Ad click prediction

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



PU Unbiased Risk Estimation

7

du Plessis+ (NeurIPS2014, ICML2015)

Decompose the risk:

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

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

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

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

- Eliminate the expectation over N data as

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

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

Unbiased risk estimator:

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

Non-Negative Risk Correction

8

Kiryu+ (NeurIPS2017) , Lu+ (AISTATS2020)

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

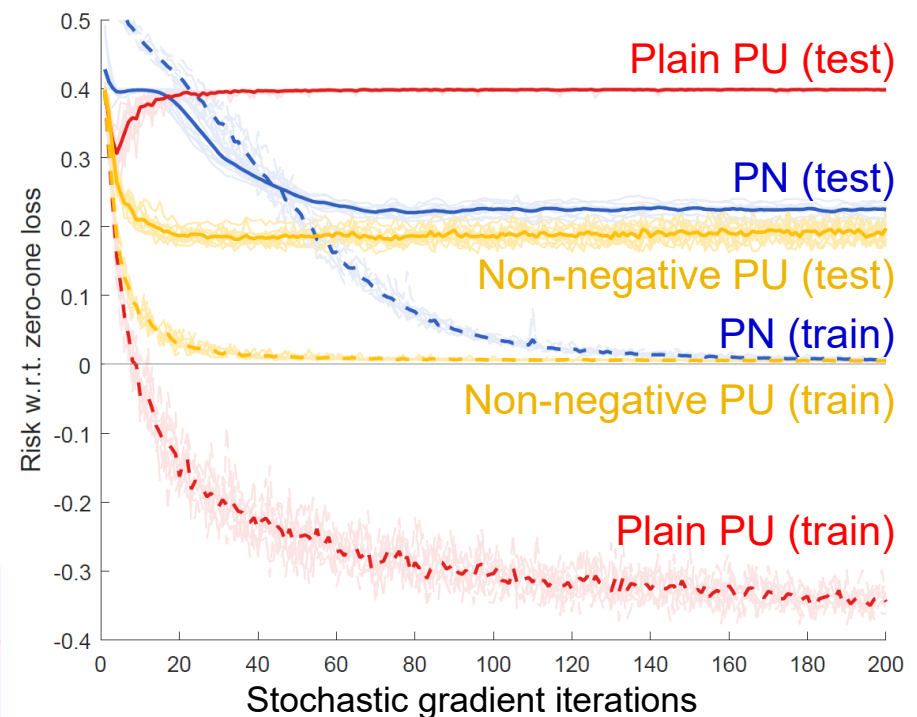
■ Risk for N data: $R^-(f) = \mathbb{E}_{p(\mathbf{x})} \left[\ell(-1, f(\mathbf{x})) \right] - \pi \mathbb{E}_{p(\mathbf{x}|y=+1)} \left[\ell(-1, f(\mathbf{x})) \right]$

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

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

■ **Non-negative correction:**

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





Contents

1. **Weakly Supervised Learning**
 - A) Positive-Unlabeled Classification
 - B) **Various Extensions**
2. Noisy-Label Learning
3. Transfer Learning
4. Towards More Reliable Learning

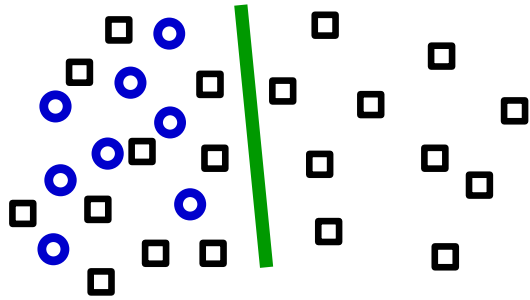
Slides →



Various Extensions (Binary)

■ Similar unbiased risk estimation is possible!

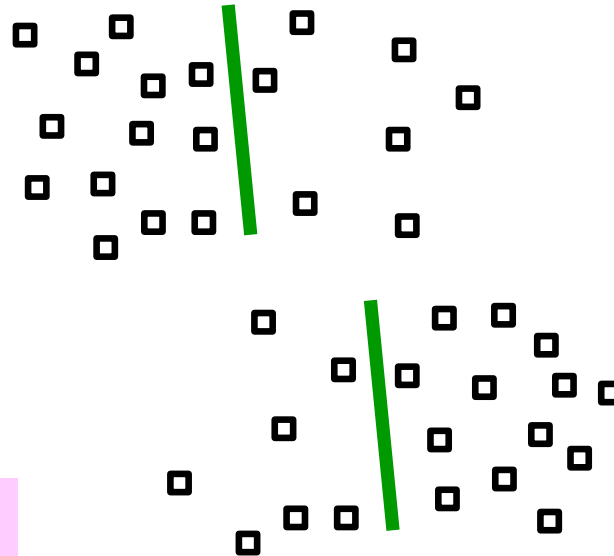
Positive-Unlabeled (PU)



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

Click prediction

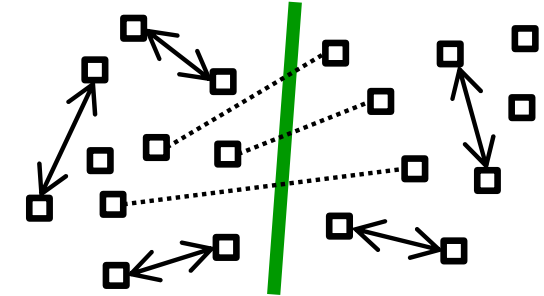
Unlabeled-Unlabeled (UU)



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

Different populations

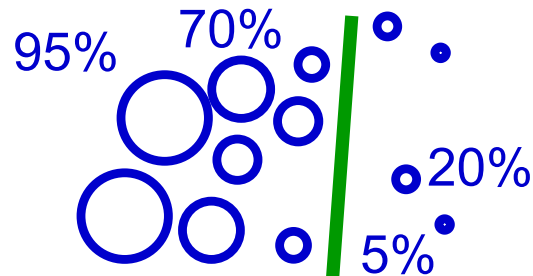
Similar-Dissimilar (SD)



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

Sensitive prediction

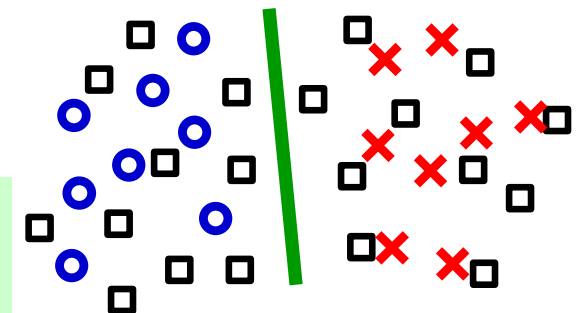
Positive-confidence (Pconf)



Ishida+ (NeurIPS2018), Shinoda+ (IJCAI2021)

Purchase prediction

Positive-Negative-Unlabeled (PNU)



Sakai+ (ICML2017, ML2018)

Semi-supervised classification
without manifold/clusters →

Various Extensions (Multiclass)

11

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

- **Multi-class weak-labels:**

- **Complementary label:**

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

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

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

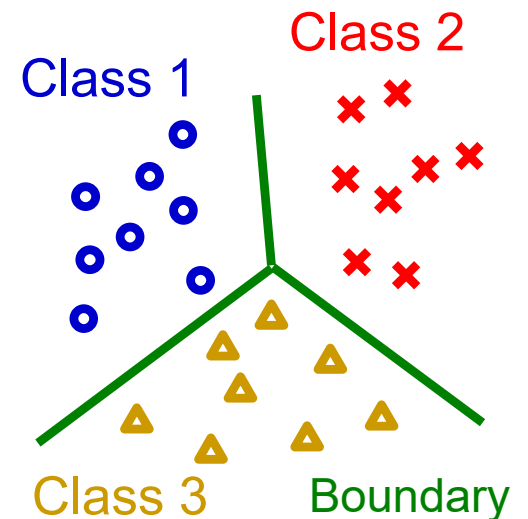
Feng+ (ICML2020,
NeurIPS2020),
Lv+ (ICML2020)

- **Single-class confidence:**

Cao+ (arXiv2021)

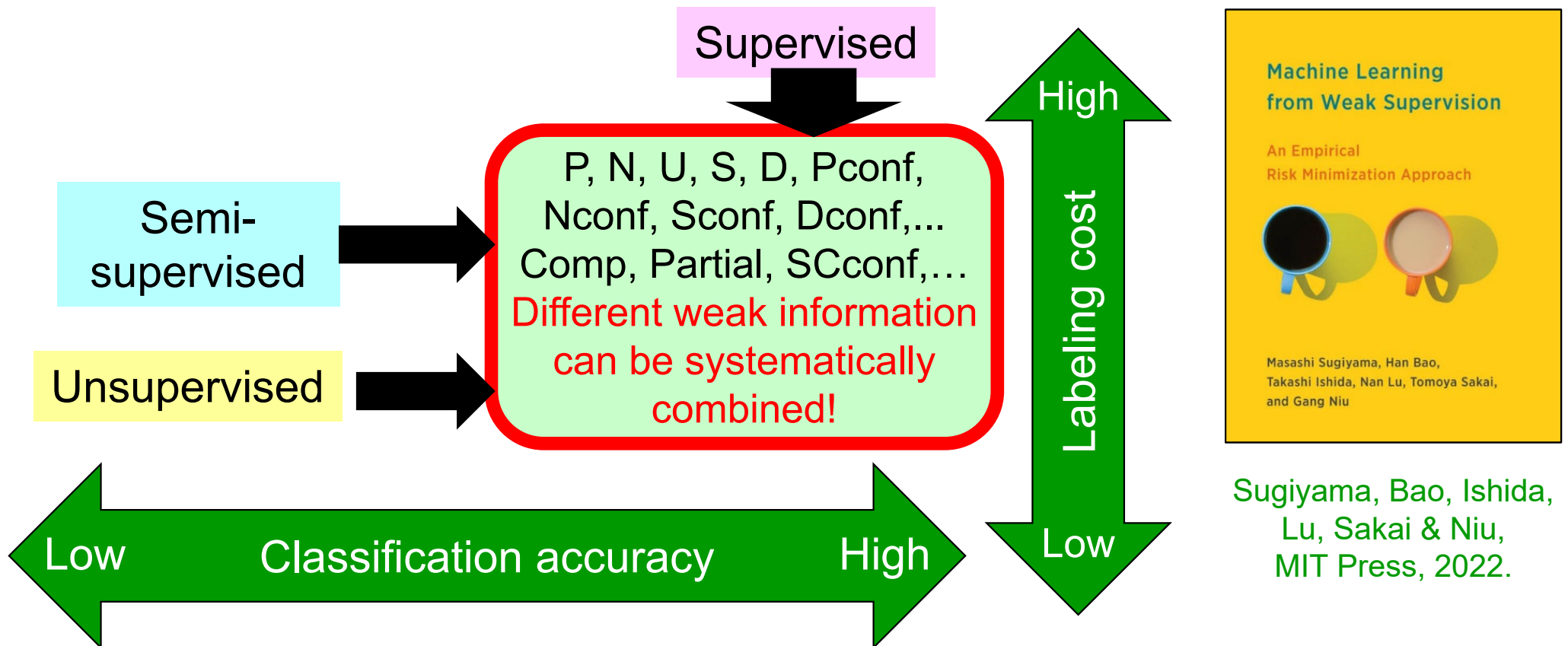
- One-class data with full confidence

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



- Similar unbiased risk estimation is possible!

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





Contents

13

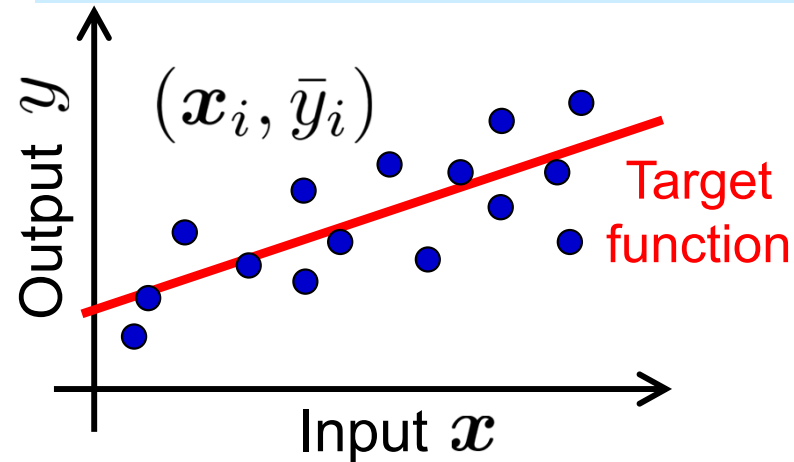
1. Weakly Supervised Learning
2. Noisy-Label Learning
 - A) Noise Transition
 - B) Algorithms
3. Transfer Learning
4. Towards More Reliable Learning

Slides →



Supervised Learning with Noisy Output 14

Regression (additive noise)

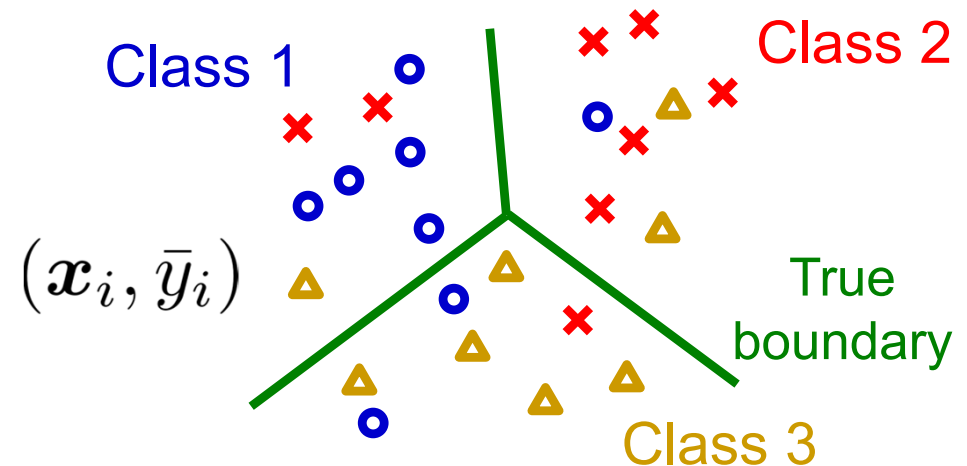


$$\min_g \sum_{i=1}^n \ell(\bar{y}_i, g(\mathbf{x}_i))$$

ℓ : loss

\bar{y} : noisy output

Classification (label flipping noise)



g : probabilistic classifier

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



Contents

15

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

Slides →



Modeling Input-Independent Noise

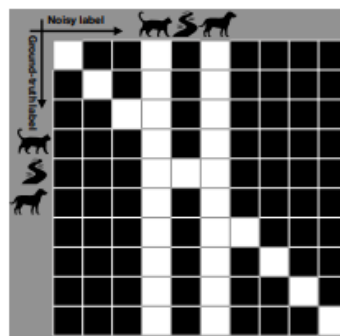
■ Noise transition matrix: $T_{y,\bar{y}} = \bar{p}(\bar{y}|y)$

- Probability of flipping y to \bar{y} .

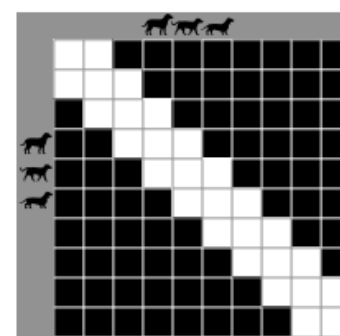
$$T = \begin{matrix} & \bar{y} \\ \begin{matrix} y \\ y \\ y \end{matrix} & \begin{bmatrix} 1 & 0 & 0 \\ 0.1 & 0.8 & 0.1 \\ 0.5 & 0.5 & 0 \end{bmatrix} \end{matrix}$$

■ Human-cognitive bias can be encoded in T .

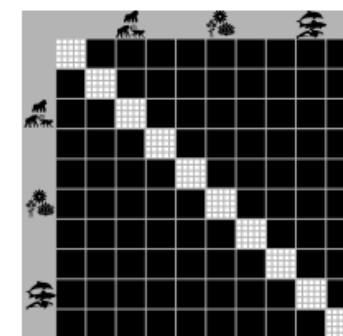
Han+ (NeurIPS2018)



(a) Column-diagonal

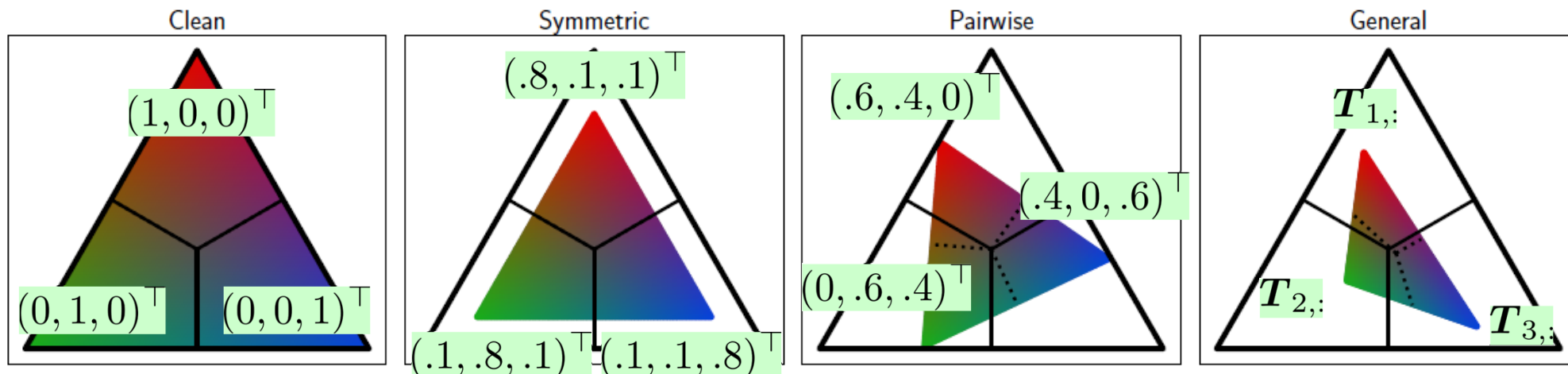


(b) Tri-diagonal



(c) Block-diagonal

■ T can be visualized in a simplex.



Loss Correction with Noise Transition 17

Patrini+ (CVPR2017)

- **Forward correction:** $\ell(\mathbf{T}^\top \mathbf{g}(\mathbf{x}))$ ℓ : vectorized loss
 - Add noise by \mathbf{T}^\top . $\ell_y(\mathbf{g}(\mathbf{x})) = \ell(y, \mathbf{g}(\mathbf{x}))$
- **Backward correction:** $\mathbf{T}^{-1} \ell(\mathbf{g}(\mathbf{x}))$
 - Remove noise by \mathbf{T}^{-1} .
- If \mathbf{T} is given, **consistency** can be guaranteed!
- If \mathbf{T} is unknown, how is it estimated?



Contents

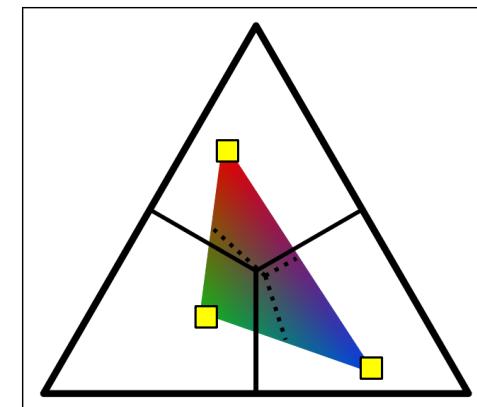
18

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

Slides →



- With **noiseless labels**, T can be obtained naively:
 - We know the **vertices** of the triangle.



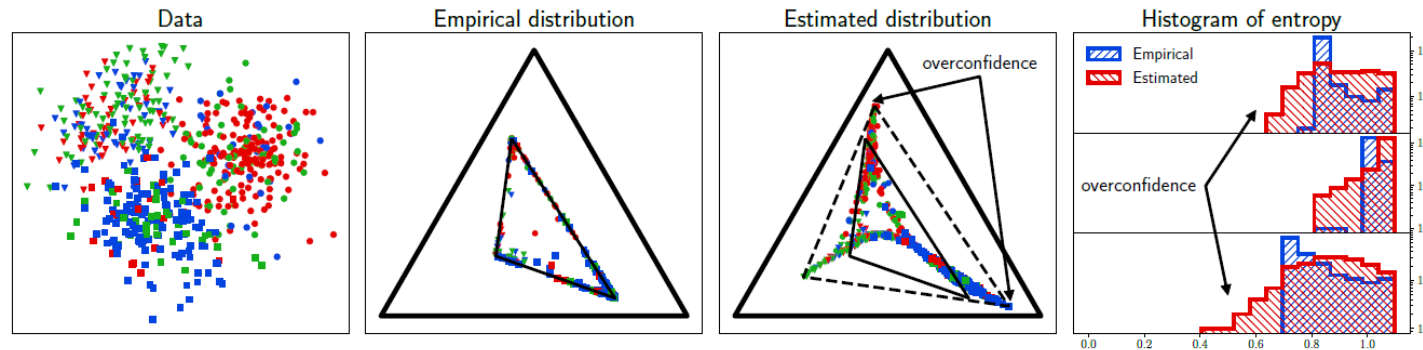
- Can we estimate T from noisy labels?
 - Generally, T is **non-identifiable**:

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

- Assume **noiseless labels** exist in the training set:
 - Select the **most confident data** as noiseless ones.

$$\mathbf{x}^y \leftarrow \mathbf{x}_i \text{ s.t. } \hat{g}_y(\mathbf{x}_i) \approx 1 \quad \hat{g} = \operatorname{argmin}_g \sum_{i=1}^n \ell(\bar{y}_i, g(\mathbf{x}_i))$$

- **Over-confidence** of neural networks is harmful.



Zhang+ (ICML2021)

- The **two-step nature** magnifies the estimation error:

1. Noise transition estimation: $\hat{\mathbf{T}}$
2. Classifier training with estimated $\hat{\mathbf{T}}$:
 - Naïve simultaneous estimation suffers **non-identifiability**.

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

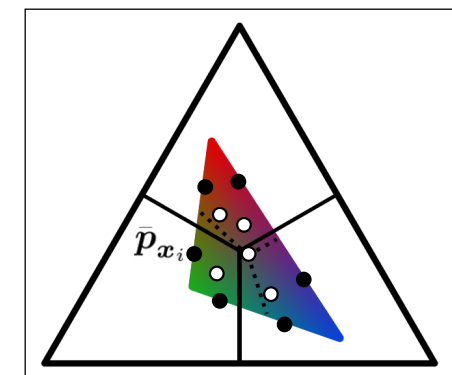
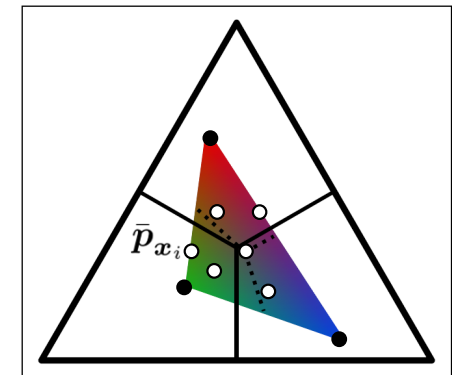
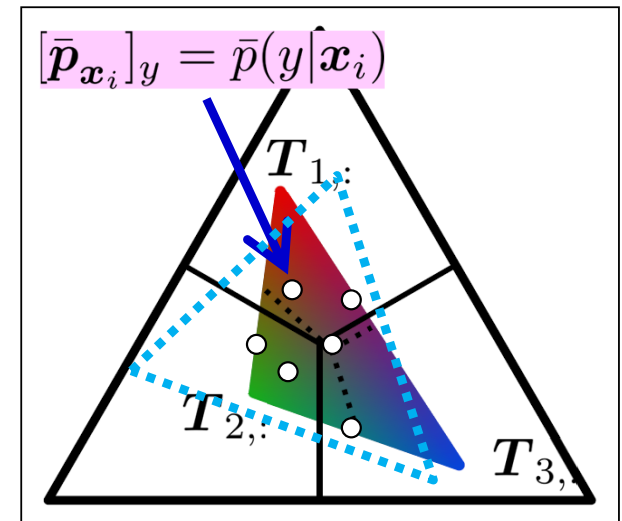
$$\min_{\mathbf{T}', \mathbf{g}} \sum_{i=1}^n \ell(\bar{y}_i, \mathbf{T}'^\top \mathbf{g}(\mathbf{x}_i))$$

- Assumption of having **noiseless labels** is too strong.

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

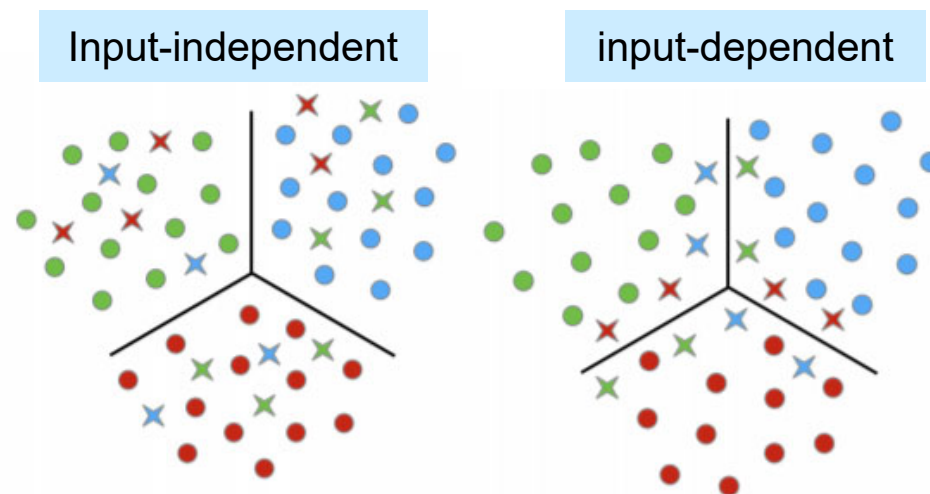
$$\min_{T', g} \sum_{i=1}^n \ell(\bar{y}_i, T'^{\top} g(\mathbf{x}_i)) + \lambda \log \det(T') \quad \lambda > 0$$

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



- Real-world noise is often **input-dependent**.

- E.g., more noise near the boundary.



- **Noise transition function:**

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

- Extremely challenging to estimate it!

- **Heuristics:**

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

Xia+ (NeurIPS2020)

Berthon+ (ICML2021)

Cheng+ (CVPR2022)

Summary: Noisy-Label Learning

23

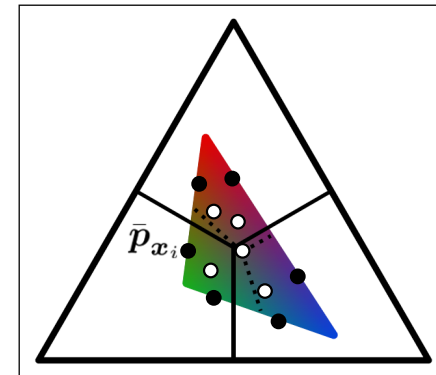
■ Explicit treatment of label noise is necessary:

- Loss correction by noise transition is promising.

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

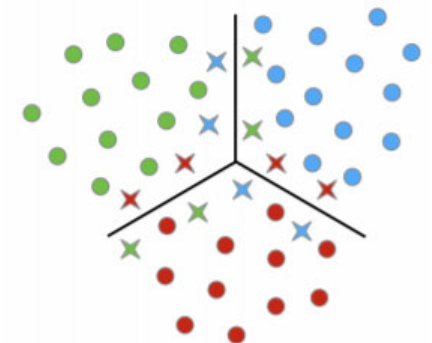
■ However, noise transition is generally **non-identifiable**:

- Recent development allows consistent estimation under mild assumptions.



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

- Heuristic solutions have been developed.
- Further **theoretical development** is needed.





Contents

24

1. Weakly Supervised Learning
2. Noisy-Label Learning
3. **Transfer Learning**
 - A) Importance Weighting
 - B) Continuous Distribution Shift
4. Towards More Reliable Learning

Slides →



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

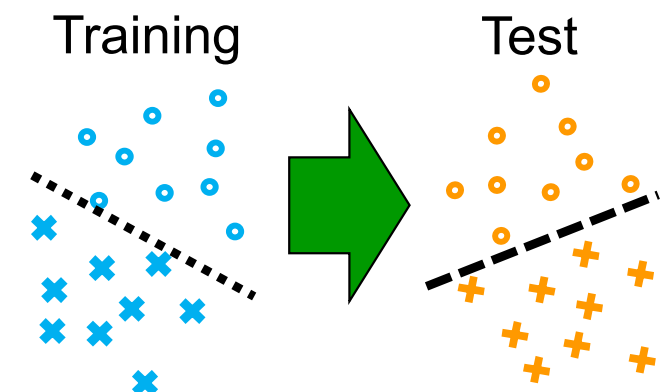
Challenge:

- 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.



Types of Distribution Shift

26

\mathbf{x} : Input

y : Output

■ Joint shift:

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

■ Covariate shift:

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

■ Class-prior shift:

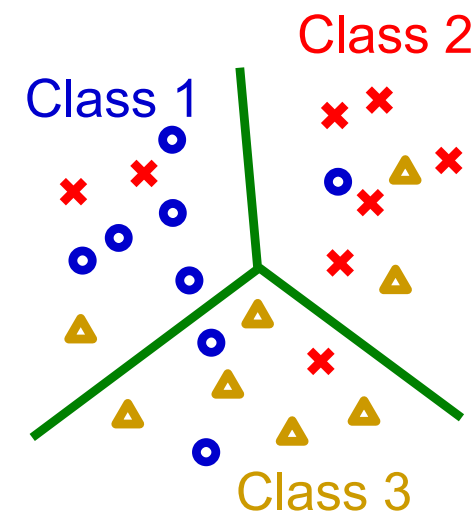
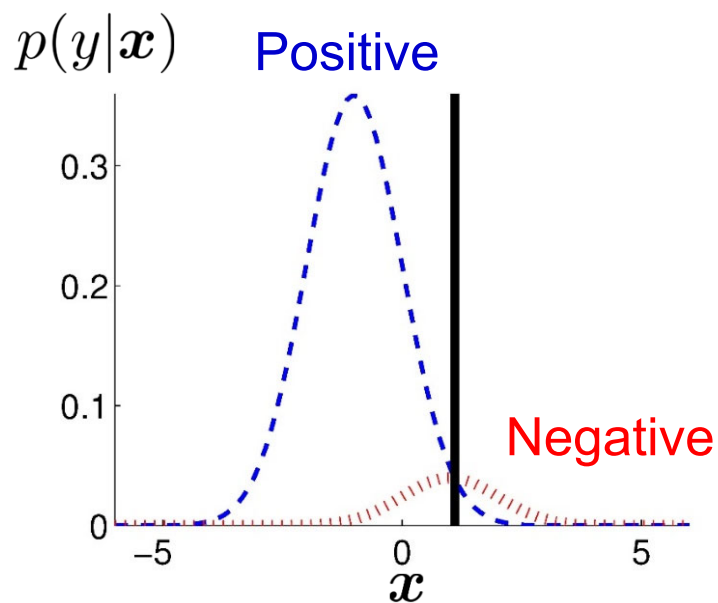
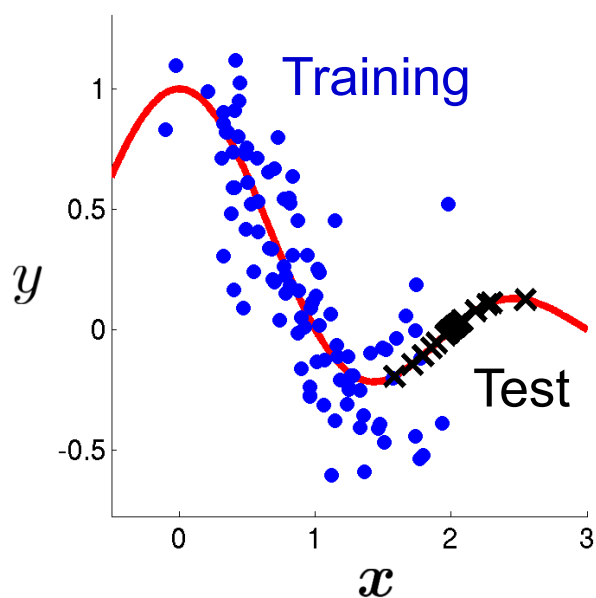
$$p_{\text{tr}}(y) \neq p_{\text{te}}(y)$$

■ Output noise:

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

■ Class-conditional shift:

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





Contents

27

1. Weakly Supervised Learning
2. Noisy-Label Learning
3. Transfer Learning
 - A) Importance Weighting
 - B) Continuous Distribution Shift
4. Towards More Reliable Learning

Slides →

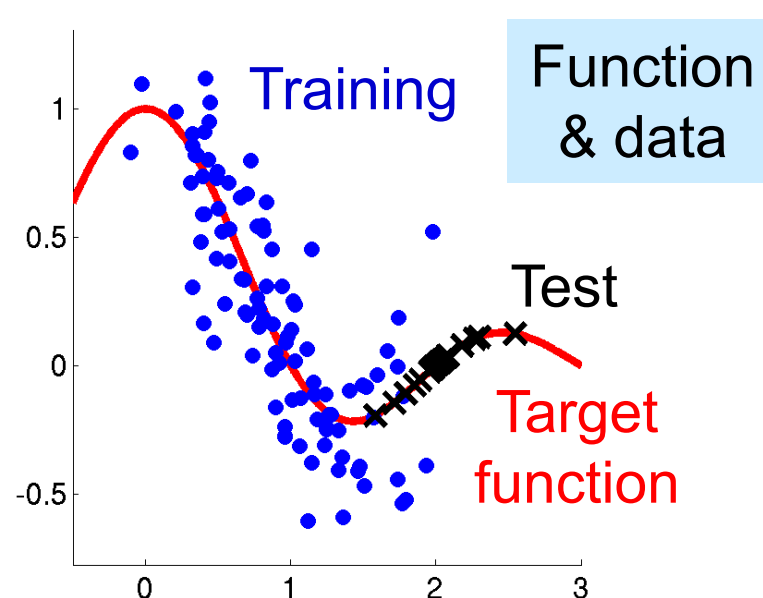
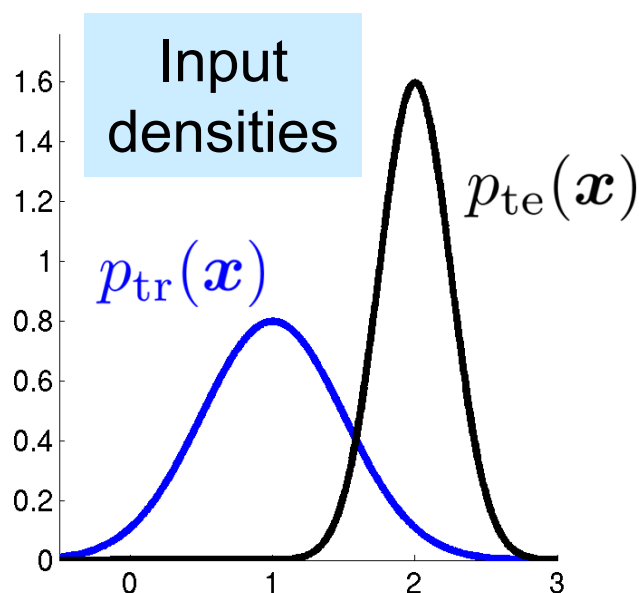


- Training and test **input** distributions are different,

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

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:

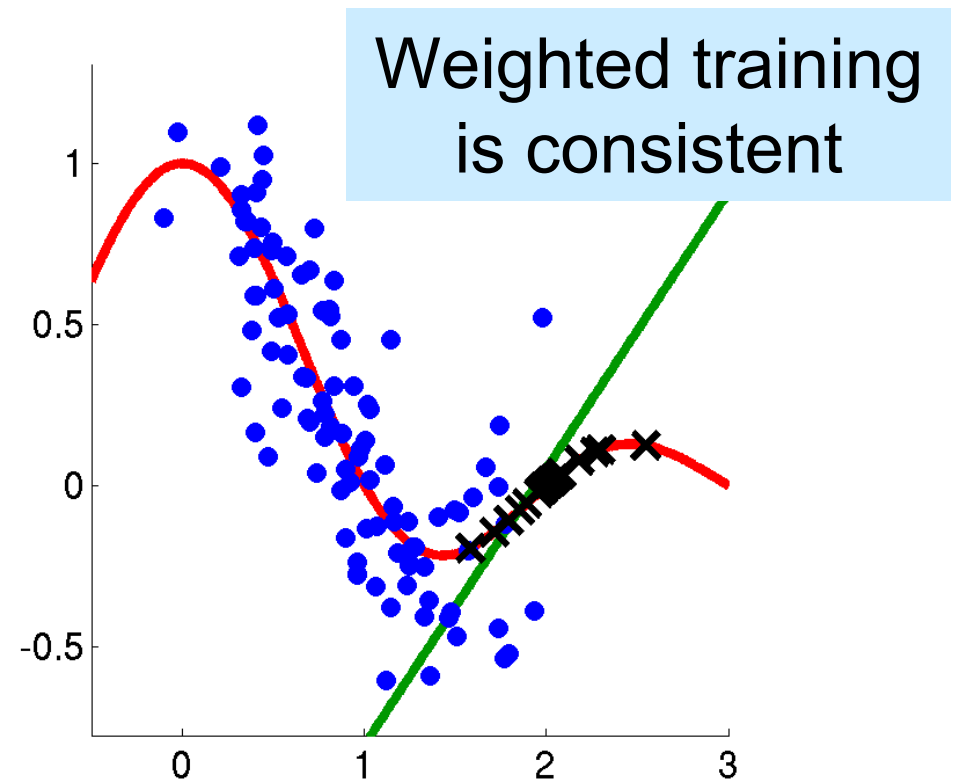
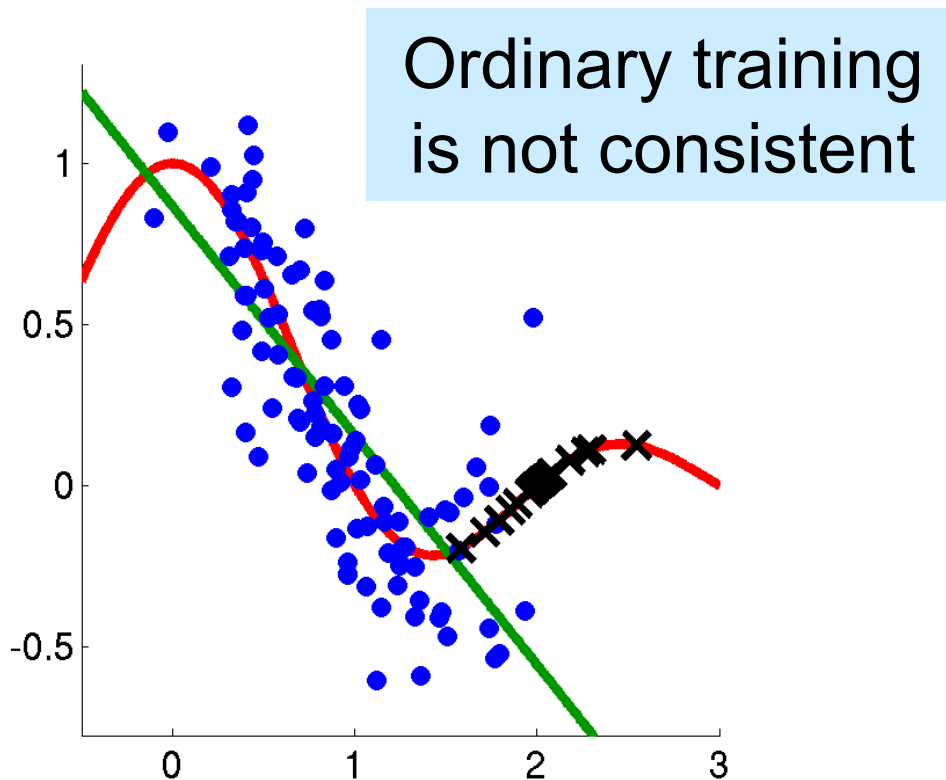
$$\{\mathbf{x}_j^{\text{te}}\}_{j=1}^{n_{\text{te}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{te}}(\mathbf{x})$$

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

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

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

Importance



■ How do we estimate the importance?

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

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

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

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

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

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

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

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

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

Joint Importance-Predictor Estimation 31

Zhang+
(ACML2020,
SNCS2021)

- The classical approaches are **two steps**:

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

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

2. Importance-weighted predictor training:

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

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

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

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

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



Contents

32

1. Weakly Supervised Learning
2. Noisy-Label Learning
3. Transfer Learning
 - A) Importance Weighting
 - B) Continuous Distribution Shift
4. Towards More Reliable Learning

Slides →



Bai+ (NeurIPS2022)

- **Class-priors** $p_t(y)$ **change** arbitrarily over time, but **class-conditional** is unchanged: $p_{\text{tr}}(\mathbf{x}|y) = p_t(\mathbf{x}|y)$
 $t = 1, \dots, T$

- **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 convex optimization**: Hazan (2016)

- convex loss ℓ (e.g., logistic),
- linear model $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$, $\mathbf{w} \in \mathcal{W}$.
- $p_{t-1}(y)$ is estimated by **black box shift estimation**. Lipton+ (ICML2018)

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

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

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

 $\eta > 0$: step size

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:

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

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

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

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

Freund+ (JCSS1997)

■ **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**:

- Online learning with step size $\eta = \Theta(V_T^{\frac{1}{3}} T^{-\frac{1}{3}})$ achieves the **optimal dynamic regret**:

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

Risk of our model

Risk of the best model at each iteration

■ Even when V_T is **unknown**:

- Our method **still achieves the optimal dynamic regret!**

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

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

Zhang+ (arXiv2023)

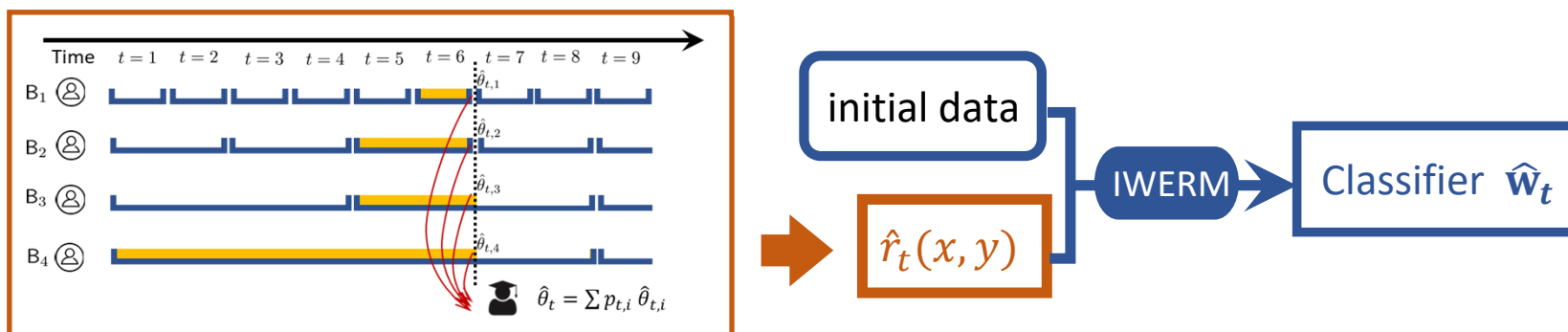
- Input density $p_t(\mathbf{x})$ change arbitrarily over time, but output-given-input is unchanged: $p_{\text{tr}}(y|\mathbf{x}) = p_t(y|\mathbf{x})$

$$t = 1, \dots, T$$

- 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:



Stay tuned!



Contents

37

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

Slides →



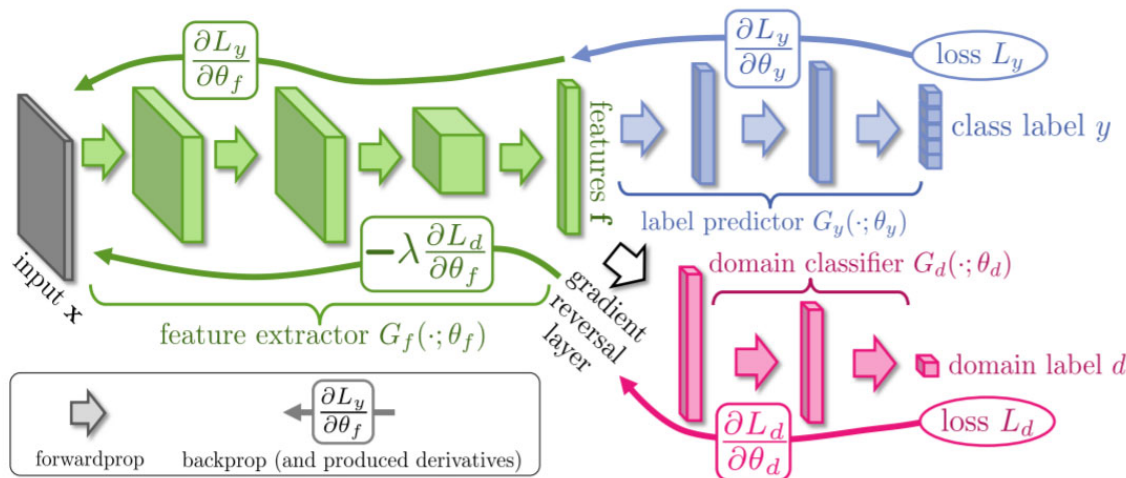
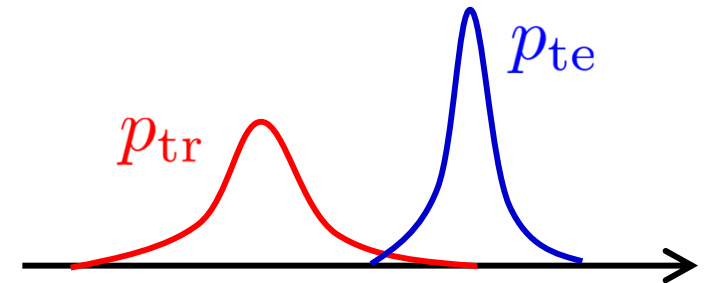
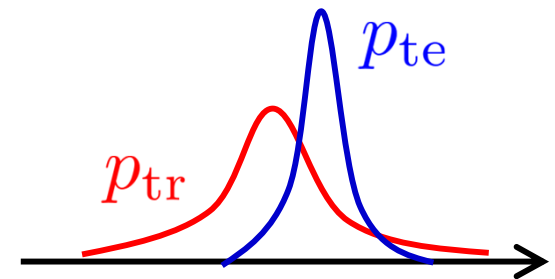
■ Limitation of importance weighting:

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

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

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

- **Feature matching**



Ben-David+ (NeurIPS2006), Ganin+ (ICML2015)

- However, considering **covariate shift** is still essential.

Joint Shift

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

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

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

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

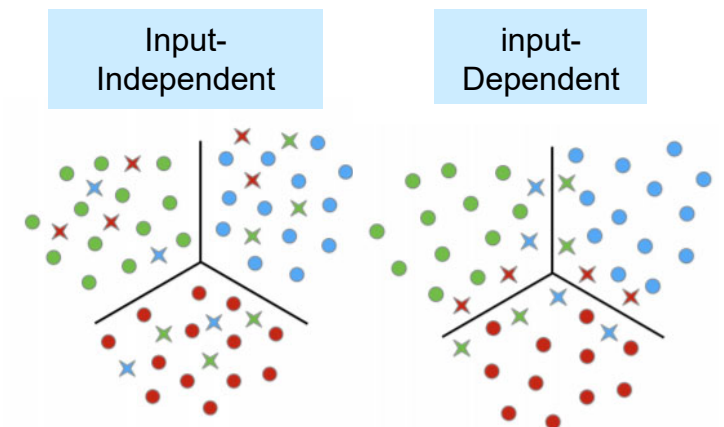
$$p_{\text{tr}}(\bar{y}|\mathbf{x}) = \sum_y \underbrace{p(\bar{y}|y, \mathbf{x})}_{\text{Noise transition}} p_{\text{te}}(y|\mathbf{x})$$

\bar{y} : Noisy class label

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

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

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



Mini-Batch-Wise Loss Matching

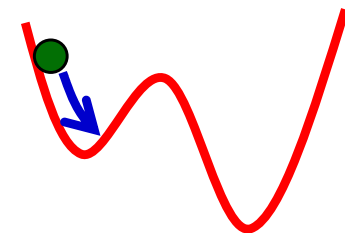
40

Fang+ (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)

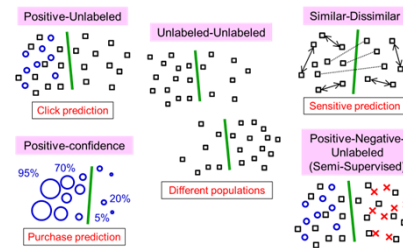
$$\frac{1}{\tilde{n}_{\text{tr}}} \sum_{i=1}^{\tilde{n}_{\text{tr}}} r_i \ell(f(\tilde{\mathbf{x}}_i^{\text{tr}}), \tilde{y}_i^{\text{tr}}) \approx \frac{1}{\tilde{n}_{\text{te}}} \sum_{j=1}^{\tilde{n}_{\text{te}}} \ell(f(\tilde{\mathbf{x}}_j^{\text{te}}), \tilde{y}_j^{\text{te}}) \quad r_i \approx \frac{p_{\text{tr}}(\tilde{\mathbf{x}}_i^{\text{tr}},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

Future Challenges

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

- Can we perform joint shift adaptation from **weak supervision**?
- Can we extend it to **continuous joint shift**?
- Can we extend it to a **limited-memory setting**?

Weakly Supervised Classification (Binary)

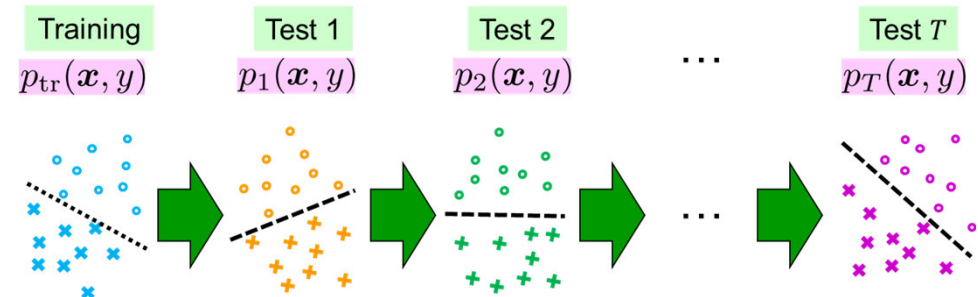
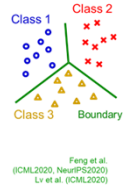


Weakly Supervised Classification (Multiclass)

■ 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").
- **Partial labels:** Specify a subset of classes that contains the correct one ("1 or 2").
- **Single-class confidence:** One-class data with full confidence ("1 with 60%, 2 with 30%, and 3 with 10%")



■ In real-world application, updating the system online is **dangerous** because new data can be **malicious**:

- Updating the system **periodically** (daily, etc.) is practical.
- But we want the system to **reflect the latest data**.
- Can we systematically use a **buffer** for temporary update?



Grateful to Collaborators!



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

and many great interns!

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

