

Machine Learning from Weak Supervision: An Empirical Risk Minimization Approach

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



Today's Topic:

Robust Machine Learning

- In real-world applications, it becomes increasingly important to consider **robustness** against various factors:
 - **Data bias**: changing environments, privacy.
 - **Insufficient information**: weak supervision.
 - **Label noise**: human error, sensor error.
 - **Attack**: adversarial noise, distribution shift.
- In this talk, I will give an overview of our recent advances in robust machine learning.

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

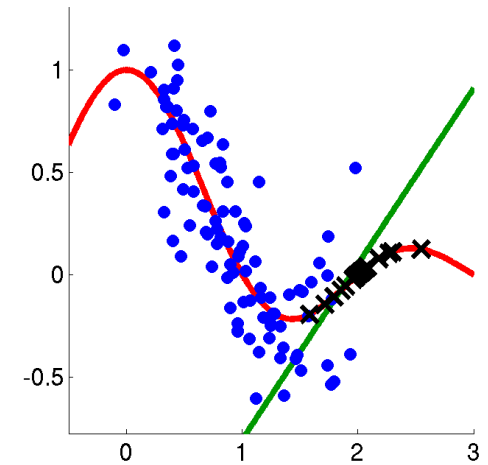


Contents

1. Transfer learning
2. Weakly supervised classification
3. Future outlook

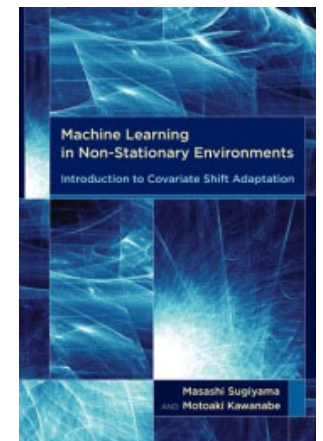
Transfer Learning

- Training and test data often have different distributions, due to
 - changing environments,
 - sample selection bias (privacy).
- **Transfer learning (domain adaptation):**
 - Train a test-domain predictor using training data from different domains.



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

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



Problem Setup

- **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

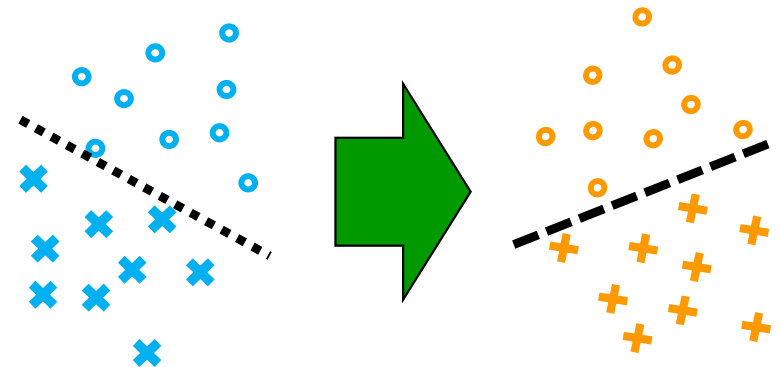
- **Goal:** Train a predictor $y = f(\mathbf{x})$ that works well in the test domain.

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

ℓ : loss function

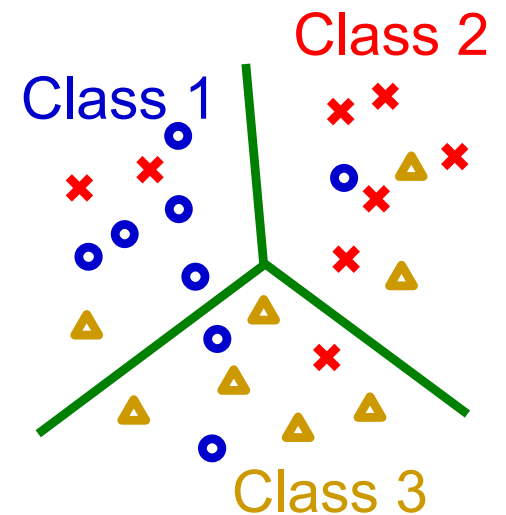
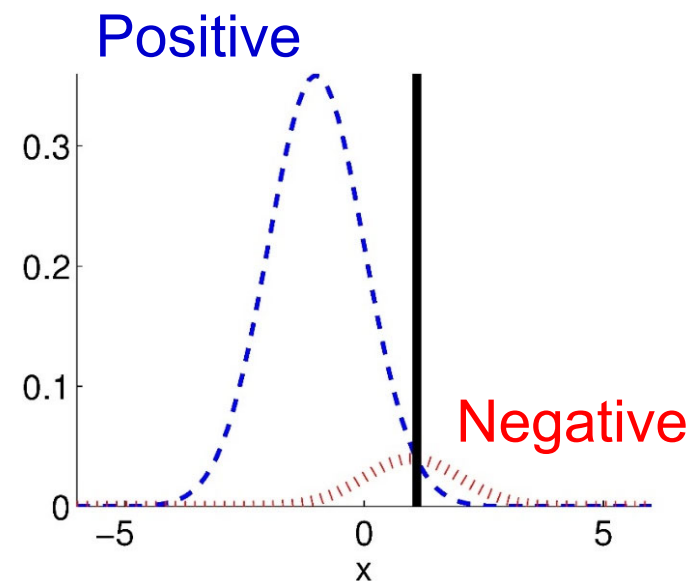
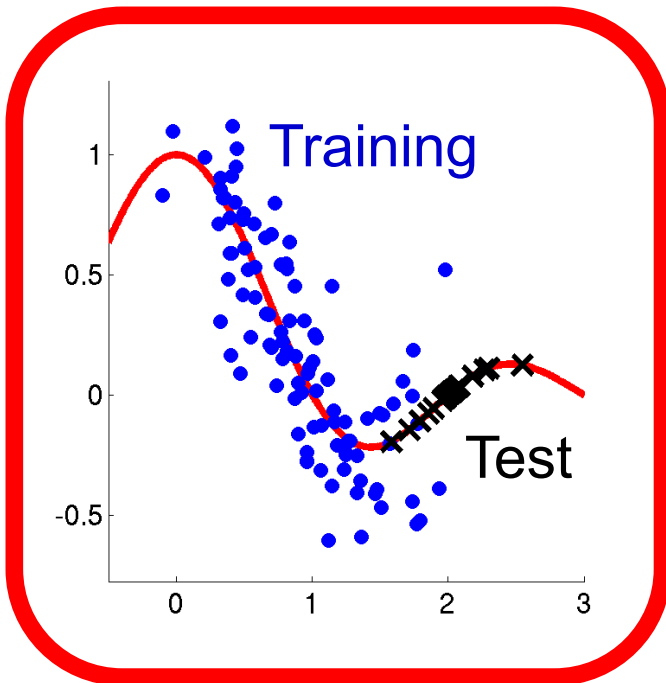
- **Challenge:** Overcome changing distributions!

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



Various Scenarios

- Full-distribution 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/target 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: $p_{\text{tr}}(\mathbf{x}|y) \neq p_{\text{te}}(\mathbf{x}|y)$



Regression under Covariate Shift

7

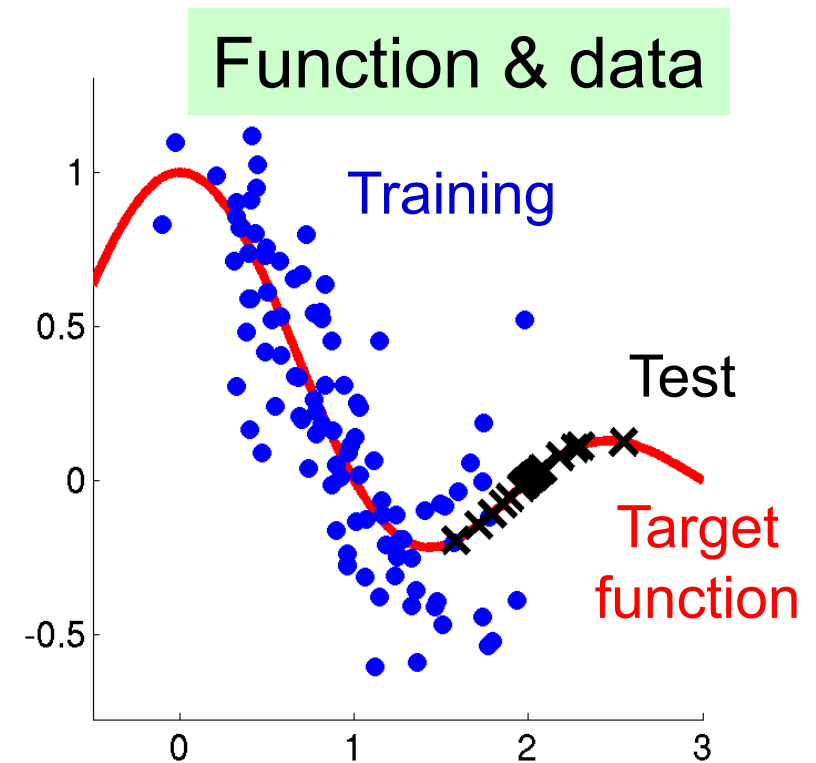
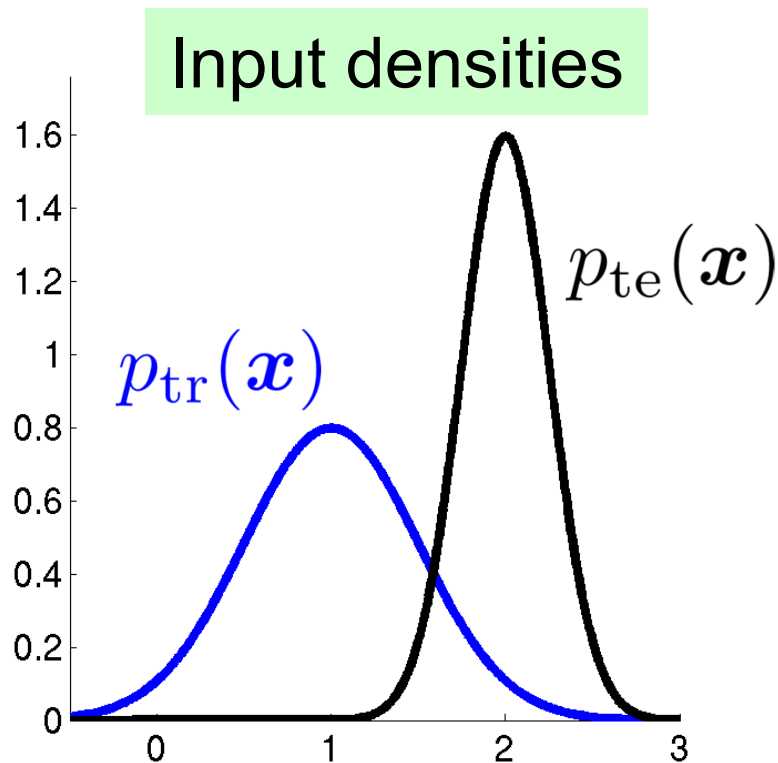
■ Covariate shift: Shimodaira (JSPI2000)

- Training and test input distributions are different:

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

- But the output-given-input distribution remains unchanged:

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



Empirical Risk Minimization (ERM)

8

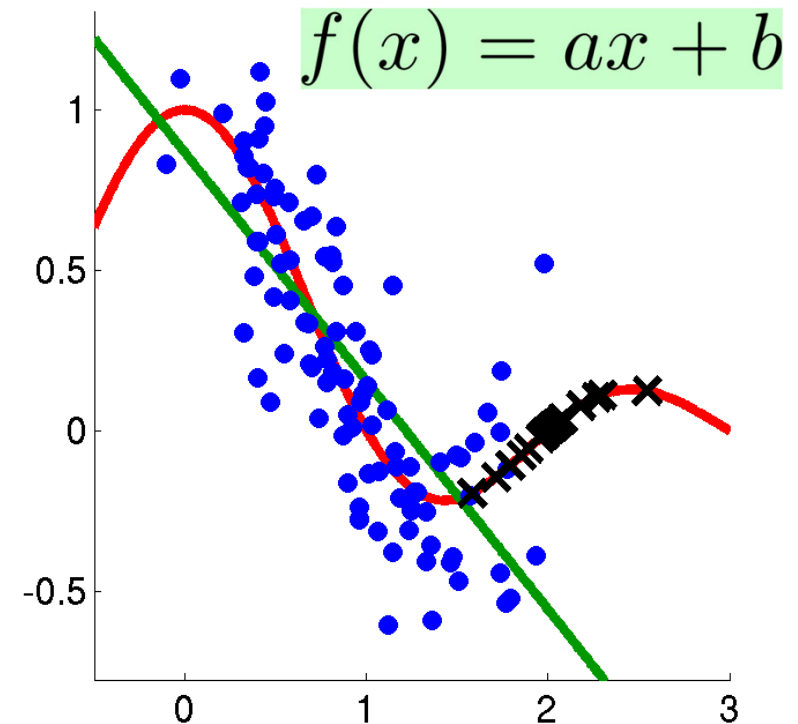
$$\min_f \left[\sum_{i=1}^{n_{\text{tr}}} \ell(f(\mathbf{x}_i^{\text{tr}}), y_i^{\text{tr}}) \right]$$

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

■ Generally, ERM is **consistent**:

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

■ However, covariate shift makes ERM **inconsistent**:



$$\frac{1}{n_{\text{tr}}} \sum_{i=1}^{n_{\text{tr}}} \ell(f(\mathbf{x}_i^{\text{tr}}), y_i^{\text{tr}}) \xrightarrow{n_{\text{tr}} \rightarrow \infty} \mathbb{E}_{p_{\text{tr}}(\mathbf{x}, y)} [\ell(f(\mathbf{x}), y)] \neq R(f)$$

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

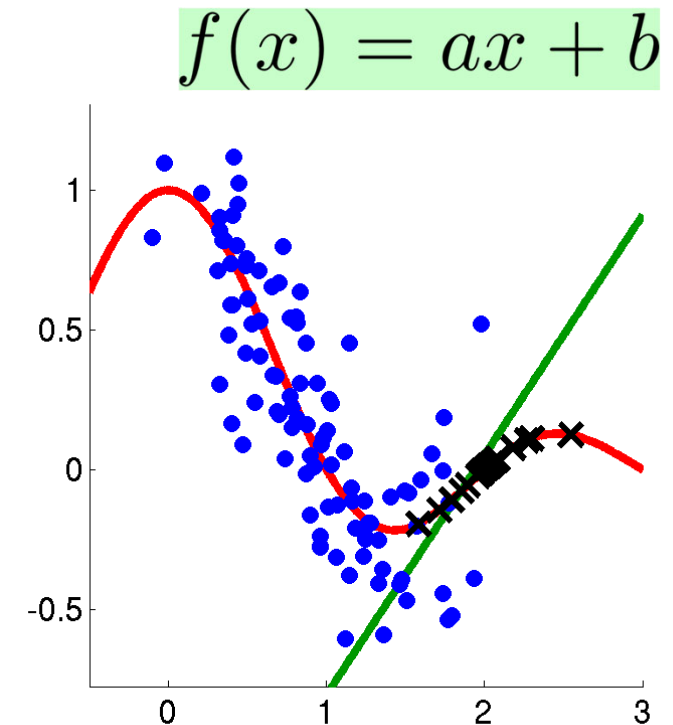
Importance-Weighted ERM (IWERM) 9

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

- IWERM is **consistent** even under **covariate shift**.

$$\frac{1}{n_{\text{tr}}} \sum_{i=1}^{n_{\text{tr}}} \frac{p_{\text{te}}(\mathbf{x}_i^{\text{tr}})}{p_{\text{tr}}(\mathbf{x}_i^{\text{tr}})} \ell(f(\mathbf{x}_i^{\text{tr}}), y_i^{\text{tr}})$$

$$\begin{aligned} & \xrightarrow[n_{\text{tr}} \rightarrow \infty]{} \mathbb{E}_{p_{\text{tr}}(\mathbf{x}, y)} \left[\frac{p_{\text{te}}(\mathbf{x})}{p_{\text{tr}}(\mathbf{x})} \ell(f(\mathbf{x}), y) \right] \\ & = \mathbb{E}_{p_{\text{te}}(\mathbf{x}, y)} [\ell(f(\mathbf{x}), y)] = R(f) \end{aligned}$$



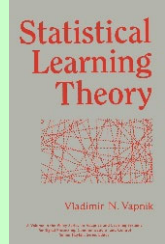
- How can we know the importance weight?



Vapnik's principle:

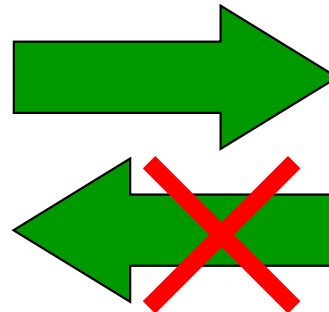
Vapnik (Wiley, 1998)

When solving a problem of interest,
one should not solve a more general problem
as an intermediate step



Knowing densities

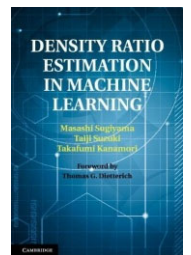
$$p_{\text{te}}(\mathbf{x}), p_{\text{tr}}(\mathbf{x})$$



Knowing ratio

$$r^*(\mathbf{x}) = \frac{p_{\text{te}}(\mathbf{x})}{p_{\text{tr}}(\mathbf{x})}$$

- Estimating the density ratio is substantially easier than estimating both the densities!
- Various direct density-ratio estimators were developed.



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

Least-Squares Importance Fitting (LSIF)

11

Kanamori et al. (JMLR2009)

- Given training and test input data:

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

- Directly fit a model r to $r^*(\mathbf{x}) = \frac{p_{\text{te}}(\mathbf{x})}{p_{\text{tr}}(\mathbf{x})}$ by LS:

$$\min_r Q(r) \quad Q(r) = \int \left(r(\mathbf{x}) - r^*(\mathbf{x}) \right)^2 p_{\text{tr}}(\mathbf{x}) d\mathbf{x}$$

- Empirical approximation:

$$Q(r) = \int r(\mathbf{x})^2 p_{\text{tr}}(\mathbf{x}) d\mathbf{x} - 2 \int r(\mathbf{x}) p_{\text{te}}(\mathbf{x}) d\mathbf{x} + C$$

$$\approx \frac{1}{n_{\text{tr}}} \sum_{i=1}^{n_{\text{tr}}} r(\mathbf{x}_i^{\text{tr}})^2 - \frac{2}{n_{\text{te}}} \sum_{j=1}^{n_{\text{te}}} r(\mathbf{x}_j^{\text{te}}) + C$$

From Two-Step Adaptation to One-Step Adaptation

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

1. Weight estimation (e.g., LSIF):

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

2. Weighted predictor training (e.g., IWERM):

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

- Can we integrate these two steps?



Contents

1. Transfer learning
 - A) Joint upper-bound minimization
 - B) Dynamic importance weighting
2. Weakly supervised classification
3. Future outlook

Joint Upper-Bound Minimization

14

Zhang et al. (ACML2020, SNCS2021)

■ Suppose we are given

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

■ **Goal:** We want to minimize the test risk.

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

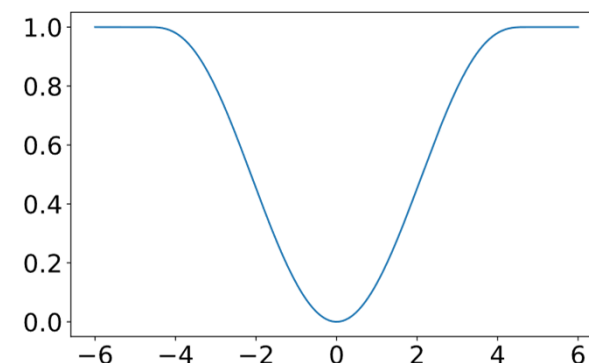
 ℓ : loss function

■ We use two losses $\ell \leq 1, \ell' \geq \ell$.

For example:

- ℓ : 0/1, ℓ' : hinge or softmax cross-entropy (classification)
- ℓ : Tukey, ℓ' : squared (regression)

Tukey loss



- For $\ell \leq 1, \ell' \geq \ell, r \geq 0$,
the test risk is upper-bounded as

$$\frac{1}{2} R_{\ell}(f)^2 \leq J_{\ell'}(r, f)$$

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

$$J_{\ell'}(r, f) = \left(\mathbb{E}_{p_{\text{tr}}(\mathbf{x}, y)} [r(\mathbf{x}) \ell'(f(\mathbf{x}), y)] \right)^2 \leftarrow \text{IWERM}$$
$$+ \mathbb{E}_{p_{\text{tr}}(\mathbf{x})} [(r(\mathbf{x}) - r^*(\mathbf{x}))^2] \leftarrow \text{LSIF}$$

- In terms of this upper-bound minimization,
2-step (LSIF followed by IWERM) is not optimal:
 - Let's directly minimize the upper bound w.r.t. r, f !

- Under some mild conditions, the test risk of the empirical solution $\hat{f} = \operatorname{argmin}_f \min_r \hat{J}_{\ell'}(r, f)$ is upper-bounded as

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

$$\hat{J}_{\ell'}(r, f) = \left(\frac{1}{n_{\text{tr}}} \sum_{i=1}^{n_{\text{tr}}} r(\mathbf{x}_i^{\text{tr}}) \ell'(f(\mathbf{x}_i^{\text{tr}}), y_i^{\text{tr}}) \right)^2 + \left(\frac{1}{n_{\text{tr}}} \sum_{i=1}^{n_{\text{tr}}} r(\mathbf{x}_i^{\text{tr}})^2 - \frac{2}{n_{\text{te}}} \sum_{j=1}^{n_{\text{te}}} r(\mathbf{x}_j^{\text{tr}}) + C \right)$$

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

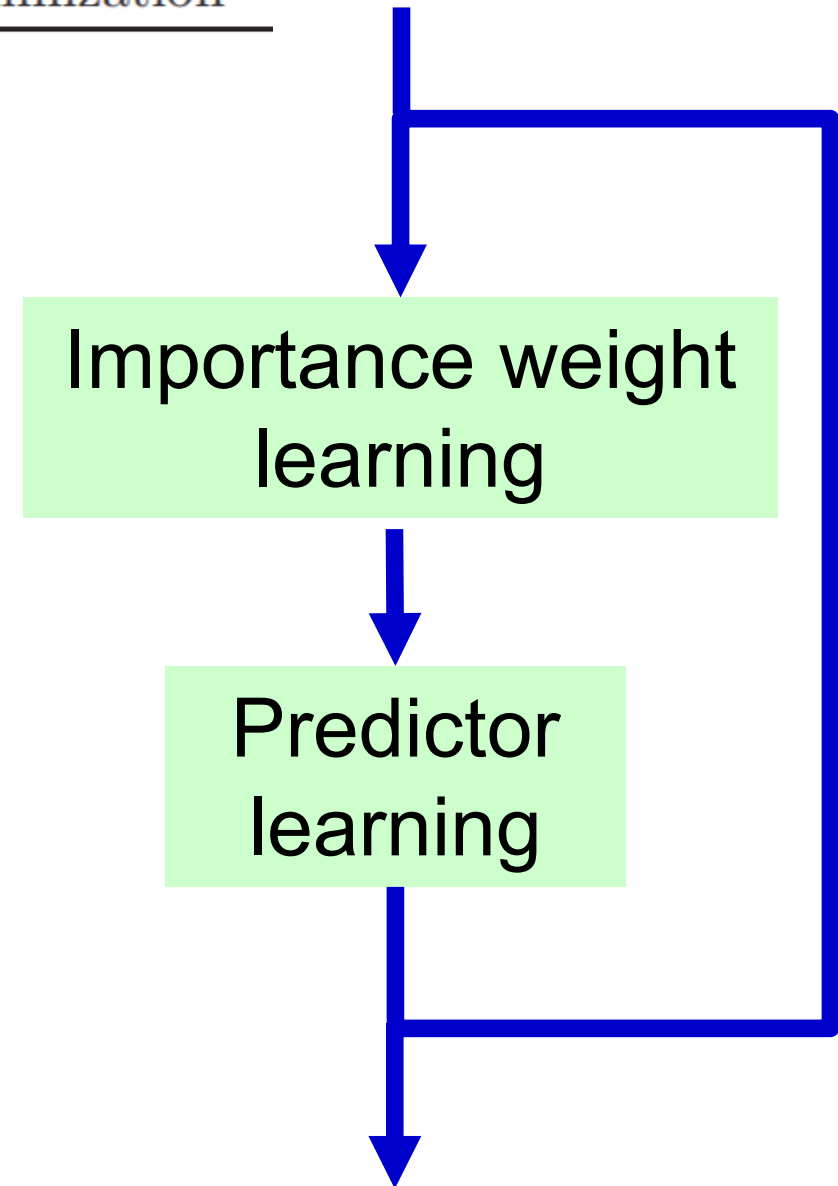
$$\{\mathbf{x}_j^{\text{te}}\}_{j=1}^{n_{\text{te}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{te}}(\mathbf{x})$$

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

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

Algorithm 2 Gradient-based Alternating Minimization

```
1:  $\mathcal{Z}^{\text{tr}}, \mathcal{X}^{\text{te}} \leftarrow \{(\mathbf{x}_i^{\text{tr}}, y_i^{\text{tr}})\}_{i=1}^{n_{\text{tr}}}, \{\mathbf{x}_i^{\text{te}}\}_{i=1}^{n_{\text{te}}}$ 
2:  $\mathcal{A} \leftarrow$  a gradient-based optimizer
3:  $f \leftarrow$  an arbitrary classifier
4: for round = 0, 1, ..., numOfRounds - 1 do
5:   for epoch = 0, 1, ..., numOfEpochsForG - 1 do
6:     for  $i = 0, 1, \dots, \text{numOfMiniBatches} - 1$  do
7:        $\mathcal{Z}_i^{\text{tr}}, \mathcal{X}_i^{\text{te}} \leftarrow \text{sampleMiniBatch}(\mathcal{Z}^{\text{tr}}, \mathcal{X}^{\text{te}})$ 
8:        $g \leftarrow \mathcal{A}(g, \nabla_g \hat{J}_{\text{UB}}(f, g; \mathcal{Z}_i^{\text{tr}} \cup \mathcal{X}_i^{\text{te}}))$ 
9:     end for
10:  end for
11:  for epoch = 0, 1, ..., numOfEpochsForF - 1 do
12:    for  $i = 0, 1, \dots, \text{numOfMiniBatches} - 1$  do
13:       $\mathcal{Z}_i^{\text{tr}} \leftarrow \text{sampleMiniBatch}(\mathcal{Z}^{\text{tr}})$ 
14:       $w_j \leftarrow \max(g(\mathbf{x}_j), 0), \forall (\mathbf{x}_j, \cdot) \in \mathcal{Z}_i^{\text{tr}}$ 
15:       $w_j \leftarrow w_j / \sum_j w_j, \forall j$ 
16:       $L_i \leftarrow \sum_{(\mathbf{x}_j, y_j) \in \mathcal{Z}_i^{\text{tr}}} w_j \ell_{\text{UB}}(f(\mathbf{x}_j), y_j)$ 
17:       $f \leftarrow \mathcal{A}(f, \nabla_f L_i)$ 
18:    end for
19:  end for
20: end for
```



Experimental Evaluation

Table 3 Mean test classification accuracy averaged over 5 trials on image datasets with neural networks. The numbers in the brackets are the standard deviations. For each dataset, the best method and comparable ones based on the *paired t-test* at the significance level 5% are described in bold face.

Dataset	Shift Level (a, b)	ERM	EIWERM	RIWERM	one-step
Fashion-MNIST	(2, 4)	81.71(0.17)	84.02(0.18)	84.12(0.06)	85.07(0.08)
	(2, 5)	72.52(0.54)	76.68(0.27)	77.43(0.29)	78.83(0.20)
	(2, 6)	60.10(0.34)	65.73(0.34)	66.73(0.55)	69.23(0.25)
Kuzushiji-MNIST	(2, 4)	77.09(0.18)	80.92(0.32)	81.17(0.24)	82.45(0.12)
	(2, 5)	65.06(0.26)	71.02(0.50)	72.16(0.19)	74.03(0.16)
	(2, 6)	51.24(0.30)	58.78(0.38)	60.14(0.93)	62.70(0.55)

Shimodaira (JSPI2000)

Yamada et al. (NIPS2011, NeCo2013)



Contents

1. Transfer learning
 - A) Joint upper-bound minimization
 - B) Dynamic importance weighting
2. Weakly supervised classification
3. Future outlook

Dynamic Importance Weighting

20

Fang et al. (NeurIPS2020)

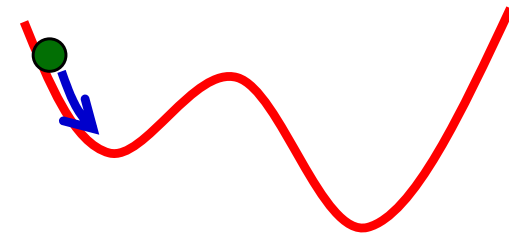
- Deep learning adopts **stochastic optimization**:

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

- Let's learn

- Importance weight r
- predictor f

dynamically in the **mini-batch-wise** manner.



- Suppose we are given

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

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

Huang, et al. (NeurIPS2007)

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

- **No covariate shift assumption is needed!**

Algorithm 1 Dynamic importance weighting (in a mini-batch).

Require: a training mini-batch \mathcal{S}^{tr} , a validation mini-batch \mathcal{S}^{v} , the current model f_{θ_t}

- 1: forward the input parts of \mathcal{S}^{tr} & \mathcal{S}^{v}
 - 2: compute the loss values as \mathcal{L}^{tr} & \mathcal{L}^{v}
 - 3: match \mathcal{L}^{tr} & \mathcal{L}^{v} to obtain \mathcal{W}
 - 4: weight the empirical risk $\widehat{R}(f_{\theta})$ by \mathcal{W}
 - 5: backward $\widehat{R}(f_{\theta})$ and update θ
-

Experimental Evaluation

Table 4: Mean accuracy (standard deviation) in percentage on Fashion-MNIST (F-MNIST for short), CIFAR-10/100 under label noise (5 trials). Best and comparable methods (paired t -test at significance level 5%) are highlighted in bold. p/s is short for pair/symmetric flip.

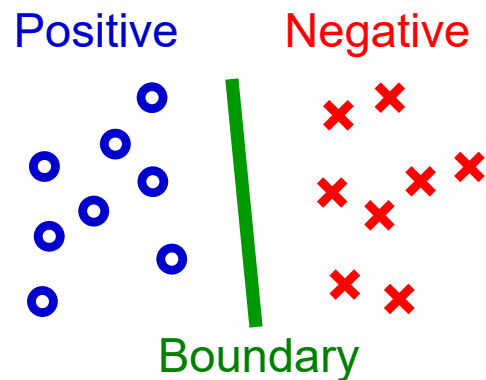
	Noise	Clean	Uniform	Random	IW	Reweight	DIW
F-MNIST	0.3 p	71.05 (1.03)	76.89 (1.06)	84.62 (0.68)	82.69 (0.38)	88.74 (0.19)	88.19 (0.43)
	0.4 s	73.55 (0.80)	77.13 (2.21)	84.58 (0.76)	80.54 (0.66)	85.94 (0.51)	88.29 (0.18)
	0.5 s	73.55 (0.80)	73.70 (1.83)	82.49 (1.29)	78.90 (0.97)	84.05 (0.51)	87.67 (0.57)
CIFAR-10	0.3 p	45.62 (1.66)	77.75 (3.27)	83.20 (0.62)	45.02 (2.25)	82.44 (1.00)	84.44 (0.70)
	0.4 s	45.61 (1.89)	69.59 (1.83)	76.90 (0.43)	44.31 (2.14)	76.69 (0.57)	80.40 (0.69)
	0.5 s	46.35 (1.24)	65.23 (1.11)	71.56 (1.31)	42.84 (2.35)	72.62 (0.74)	76.26 (0.73)
CIFAR-100	0.3 p	10.82 (0.44)	50.20 (0.53)	48.65 (1.16)	10.85 (0.59)	48.48 (1.52)	53.94 (0.29)
	0.4 s	10.82 (0.44)	46.34 (0.88)	42.17 (1.05)	10.61 (0.53)	42.15 (0.96)	53.66 (0.28)
	0.5 s	10.82 (0.44)	41.35 (0.59)	34.99 (1.19)	10.58 (0.17)	36.17 (1.74)	49.13 (0.98)



Contents

1. Transfer learning
2. **Weakly supervised classification**
3. Future outlook

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

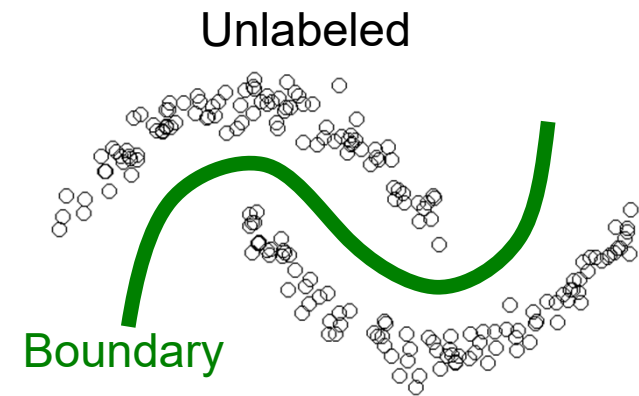


n : Number of labeled samples

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

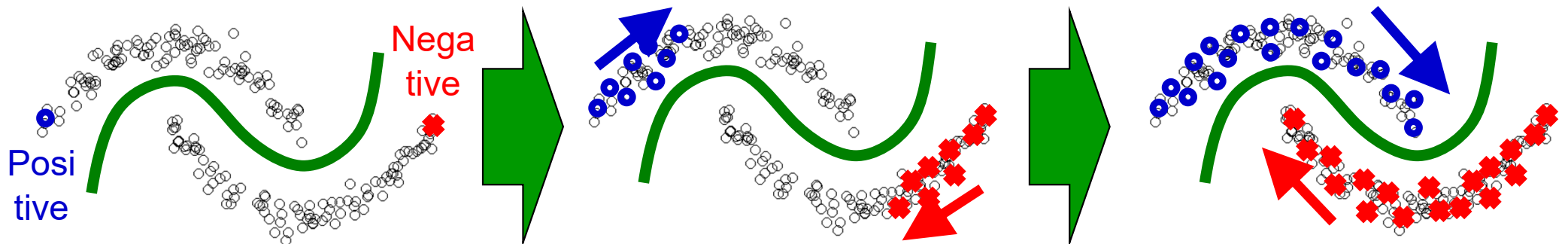
■ Unsupervised classification:

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



■ Semi-supervised classification:

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

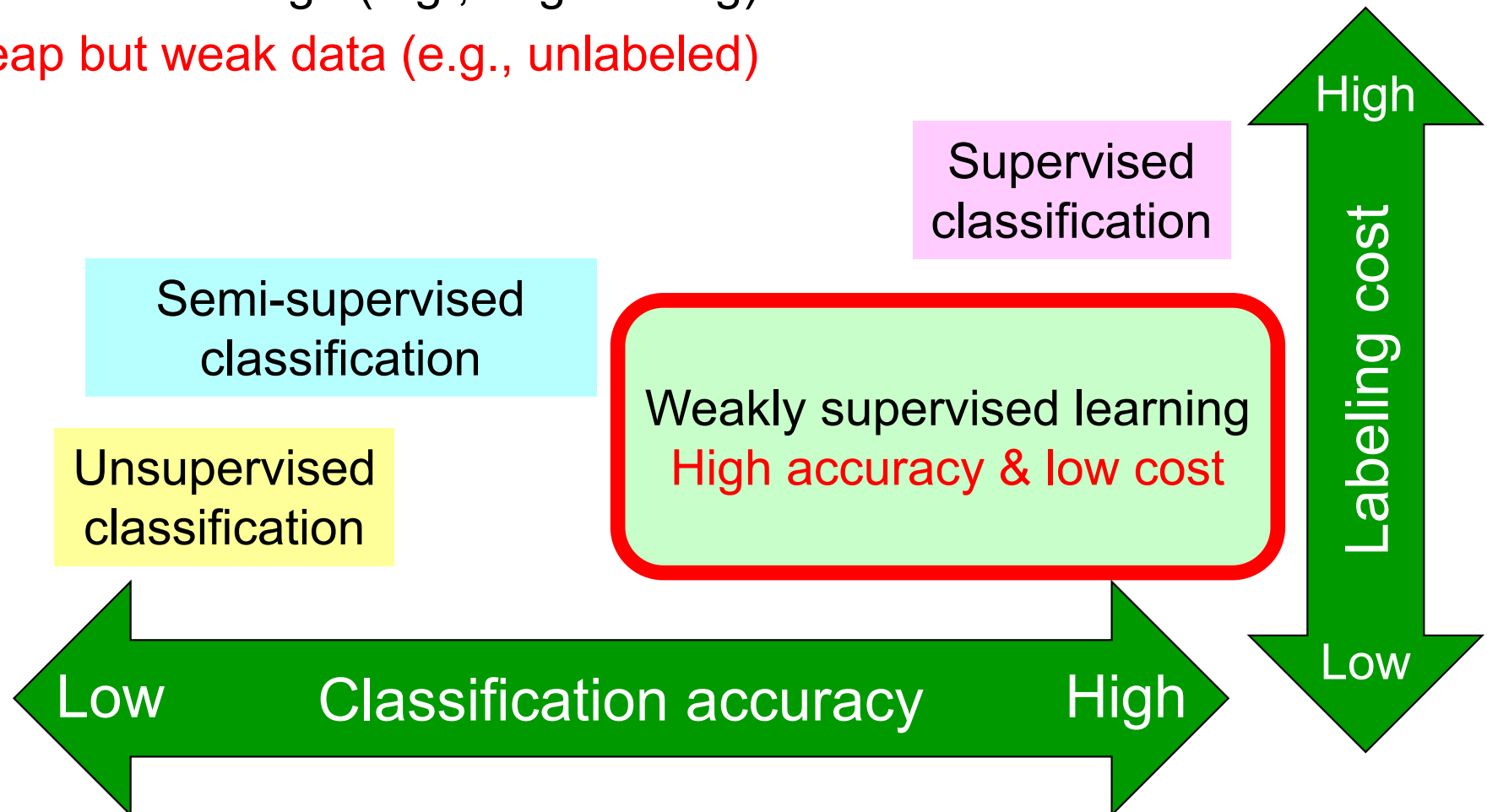


Weakly Supervised Learning

26

■ Coping with labeling cost:

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





Contents

27

1. Transfer learning
2. **Weakly supervised classification**
 - A) **Positive-unlabeled classification**
 - B) Extensions
3. Future outlook

- **Given:** Positive and unlabeled samples

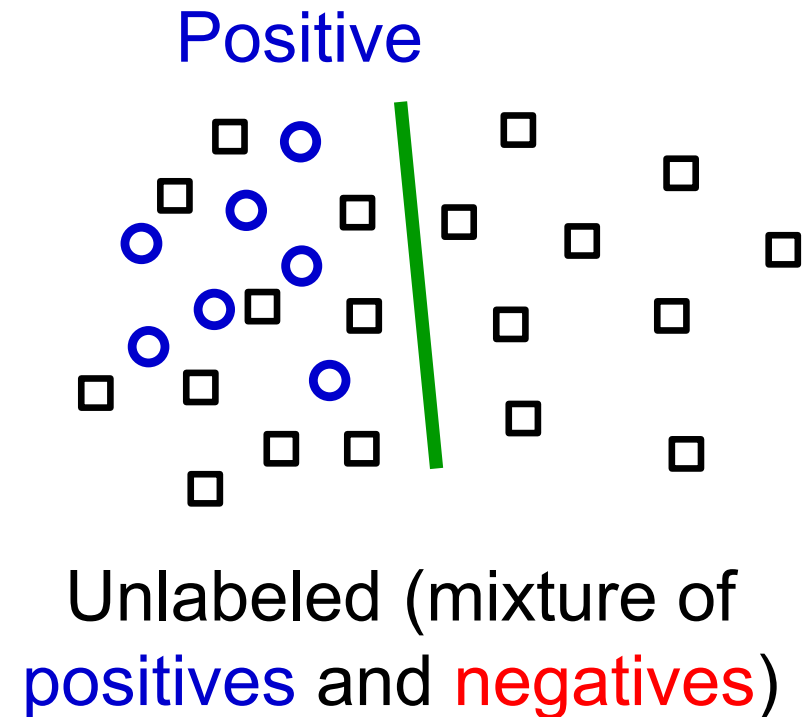
$$\{\mathbf{x}_i^P\}_{i=1}^{n_P} \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}|y = +1)$$

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

- **Goal:** Obtain a PN classifier

- **Example:** Ad-click prediction

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



- Risk of classifier f :

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

ℓ : loss function

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

$\pi = p(y = +1)$: Class-prior probability
(assumed known; **can be estimated**)

Scott & Blanchard (AISTATS2009)

Blanchard et al. (JMLR2010)

du Plessis et al. (IEICE2014, MLJ2017)

Ramaswamy et al. (ICML2016)

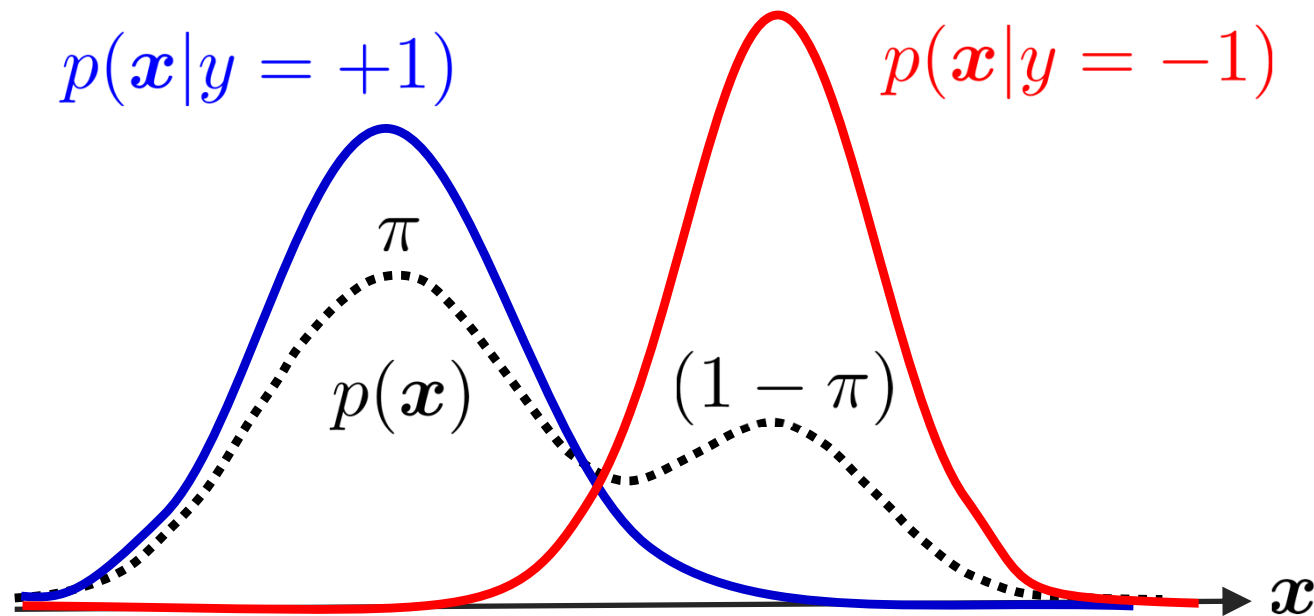
Yao et al. (arXiv2020)

- Since we do not have N data in the PU setting, the risk cannot be directly estimated.

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

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

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



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

- This allows us to eliminate the N-density:

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

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

- Unbiased risk estimation is possible from PU data, just by replacing expectations by sample averages!

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

- Replacing expectations by sample averages gives an empirical risk:

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

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

- Optimal convergence rate is attained:

Niu et al. (NIPS2016)

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

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

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

with probability $1 - \delta$

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

■ Estimation error bounds for PU and PN:

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

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

$$\hat{f}_{\text{PN}} = \underset{f}{\operatorname{argmin}} \hat{R}_{\text{PN}}(f)$$

with probability $1 - \delta$

$$\hat{R}_{\text{PN}}(f) = \frac{1}{n} \sum_{i=1}^n \ell(y_i f(\mathbf{x}_i))$$

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

■ Comparison: PU bound is smaller than PN if

$$\frac{\pi}{\sqrt{n_{\text{P}}}} + \frac{1}{\sqrt{n_{\text{U}}}} < \frac{1 - \pi}{\sqrt{n_{\text{N}}}}$$

- PU can be better than PN, provided many PU data!

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

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

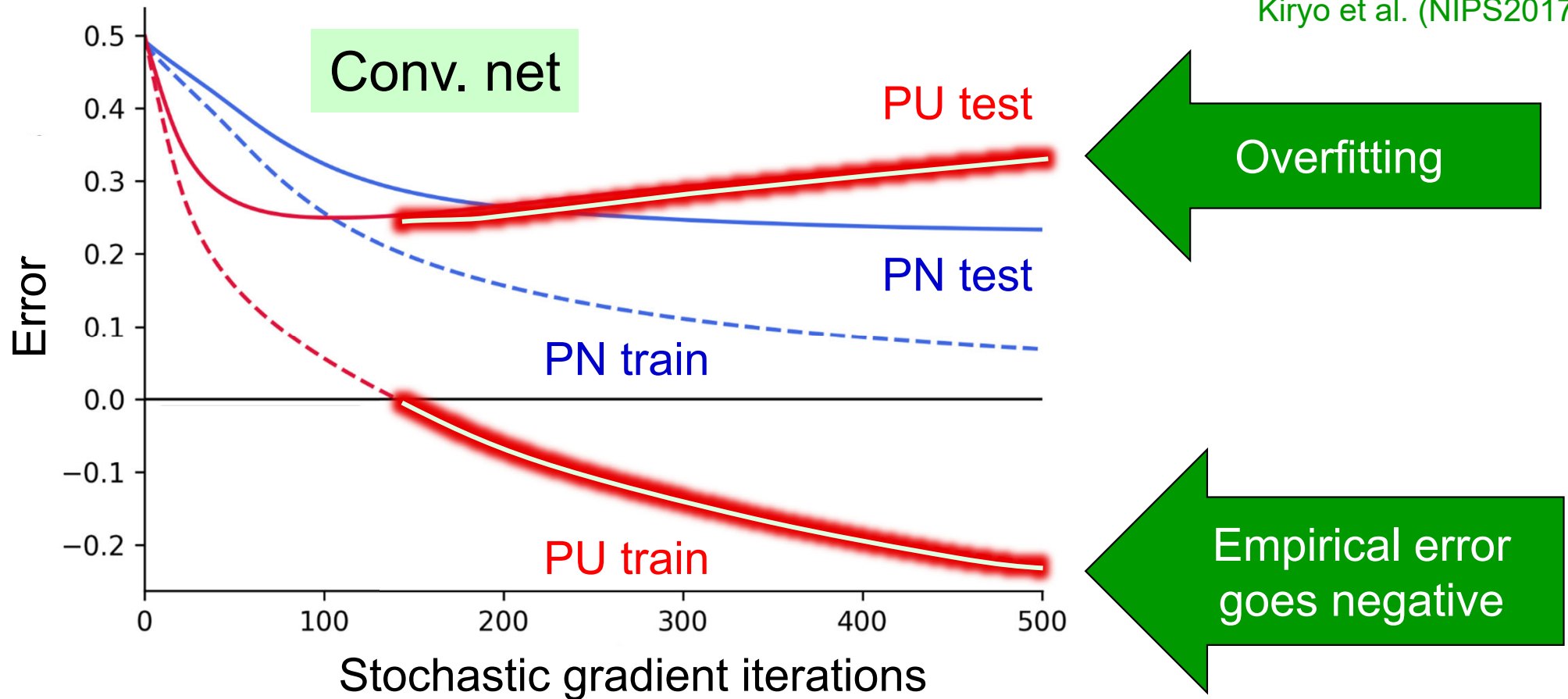
$$R^-(f) = \mathbb{E}_{p(\mathbf{x})} \left[\ell \left(-f(\mathbf{x}) \right) \right] - \pi \mathbb{E}_{p(\mathbf{x}|y=+1)} \left[\ell \left(-f(\mathbf{x}) \right) \right]$$

- If $\ell(m) \geq 0, \forall m$ $R^-(f) \geq 0$
- However, **its PU empirical approximation can be negative** due to “difference of approximations”.

$$\hat{R}_{\text{PU}}^-(f) = \frac{1}{n_U} \sum_{i=1}^{n_U} \ell \left(-f(\mathbf{x}_i^U) \right) - \frac{\pi}{n_P} \sum_{i=1}^{n_P} \ell \left(-f(\mathbf{x}_i^P) \right) \not\geq 0$$

- This problem is more critical for flexible models such as **deep neural networks**.

Kiryo et al. (NIPS2017)



- We constrain the sample approximation term **to be non-negative** through back-prop training:

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

- Now the risk estimator is biased. Is it really good?

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

- $\tilde{R}_{\text{PU}}(f)$ is still **consistent** and **its bias decreases exponentially**: $\mathcal{O}(e^{-n_{\text{P}}-n_{\text{U}}})$ $n_{\text{P}}, n_{\text{U}}$: # of P, U samples
 - In practice, we can ignore the bias of $\tilde{R}_{\text{PU}}(f)$!
- Mean-squared error of $\tilde{R}_{\text{PU}}(f)$ is not more than the original one.
 - In practice, $\tilde{R}_{\text{PU}}(f)$ is more reliable!
- Risk of $\operatorname{argmin}_f \tilde{R}_{\text{PU}}(f)$ for linear models attains the **optimal convergence rate**: $\mathcal{O}_p \left(\frac{1}{\sqrt{n_{\text{P}}}} + \frac{1}{\sqrt{n_{\text{U}}}} \right)$
 - Learned function is still optimal.

Practical Implementation for Deep Learning

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

- Use mini-batch stochastic gradient descent:
 - If $\hat{R}_{\text{PU}}^-(f) \geq 0$, perform gradient descent as usual.
 - If $\hat{R}_{\text{PU}}^-(f) < 0$, perform gradient **ascent**:
 - For bad data, step back the gradient
(to avoid converging to a poor local optimum)
and recompute the gradient with a new mini-batch.

Experiments

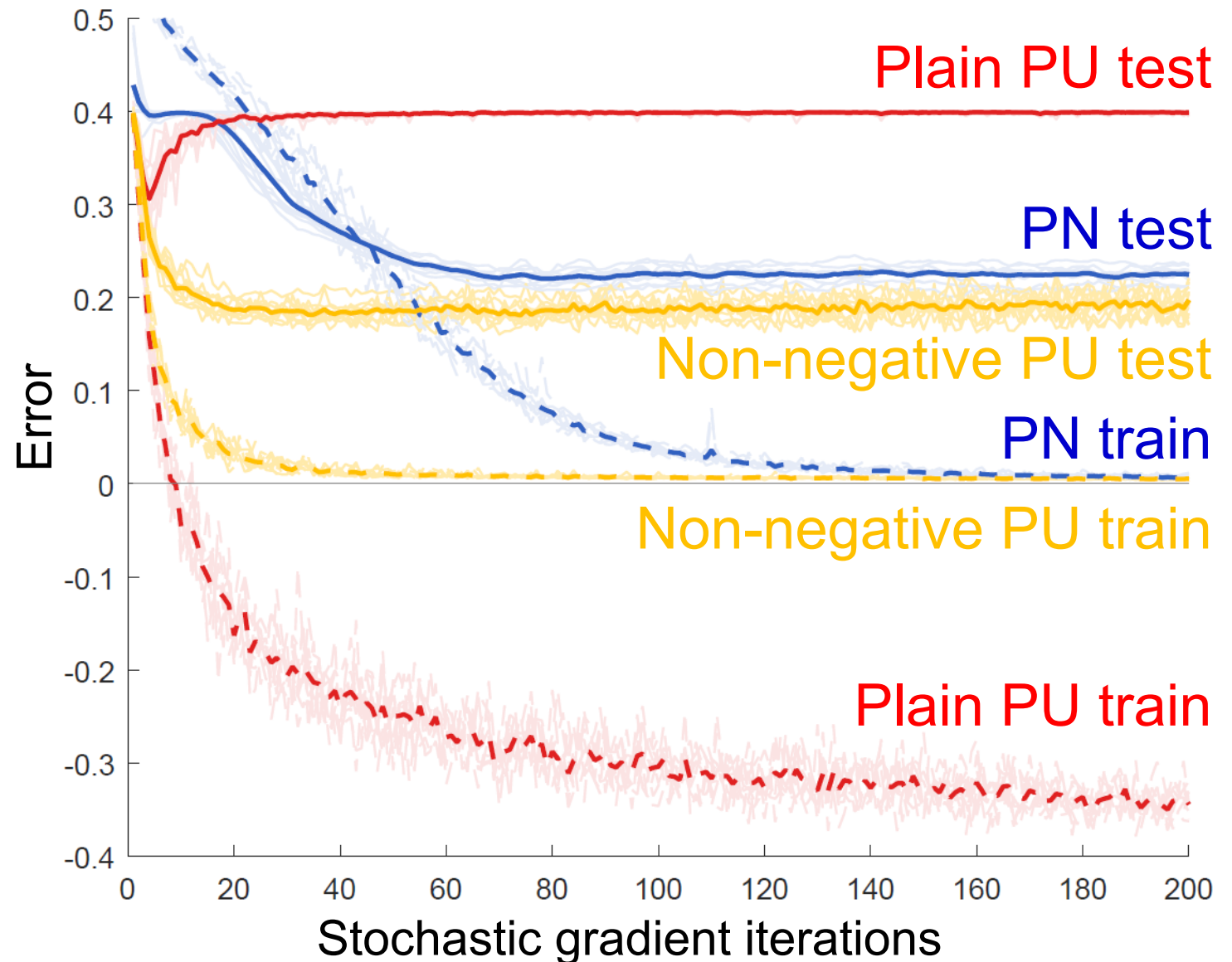
- With a large number of unlabeled data, non-negative PU can even outperform PN!

- Binary CIFAR-10:
 - Positive (airplane, automobile, ship, truck)
 - Negative (bird, cat, deer, dog, frog, horse)
- 13-layer CNN with ReLU

$$n_P = 1000$$

$$n_U = 50000$$

$$\pi = 0.4$$



- **Risk-rewriting**: Rewrite the classification risk only in terms of weak data.

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

- Standard empirical risk minimization.
 - Optimal convergence guarantee.
 - Compatible with any loss, regularization, model, and optimizer.
 - Applicable to various weak data (shown next).
- **Non-negative risk correction**: Utilize intrinsic non-negativity to mitigate overfitting.
 - Non-negativity of loss, convexity, etc.
 - Applicable to various weak data. Lu et al. (ICLR2019)
 - Applicable to noisy-label learning. Han et al. (ICML2020)



Contents

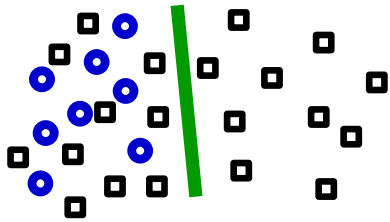
40

1. Transfer learning
2. **Weakly supervised classification**
 - A) Positive-unlabeled classification
 - B) **Extensions**
3. Future outlook

Various Binary Weak Labels

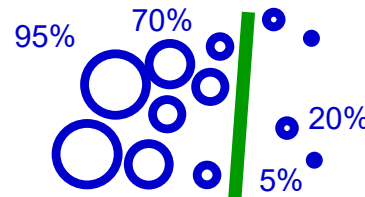
- Various weakly supervised classification problems can be solved by risk-rewriting **systematically!**

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



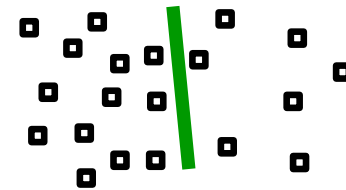
du Plessis et al.
(NIPS2014, ICML2015, MLJ2017)
Niu et al. (NIPS2016),
Kiryo et al. (NIPS2017)
Hsieh et al. (ICML2019)

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

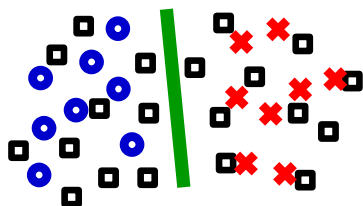


Ishida et al. (NeurIPS2018)
Shinoda et al. (IJCAI2021)

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

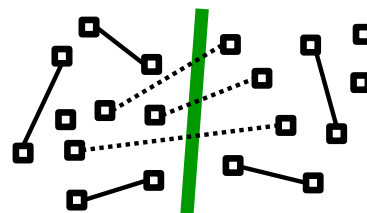


Semi-Supervised (PU+PN)
(first theoretically guaranteed method)

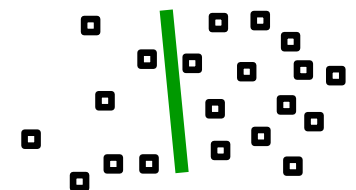


Sakai et al. (ICML2017, ML2018)

Similar-Dissimilar (SD)
(delicate information)



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



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

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

■ Risk rewriting is still possible in multi-class problems!

■ **Multi-class weak-labels:**

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

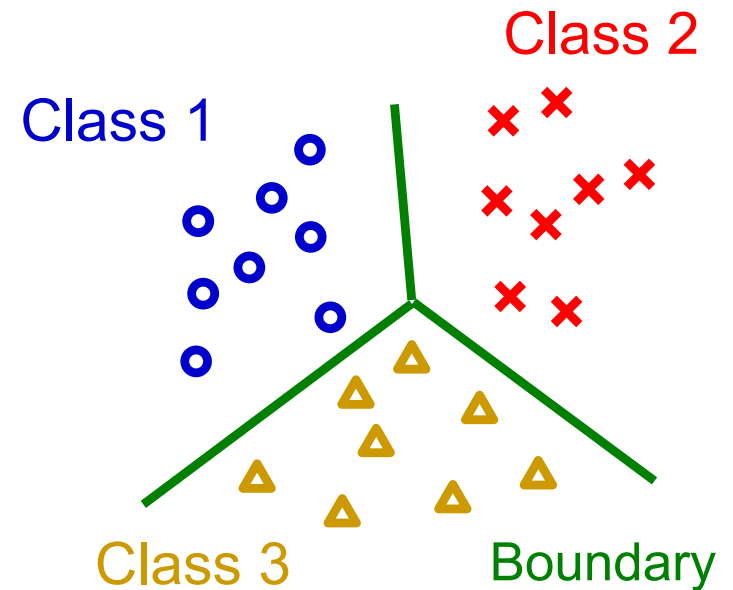
$$1/\sqrt{n}$$

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

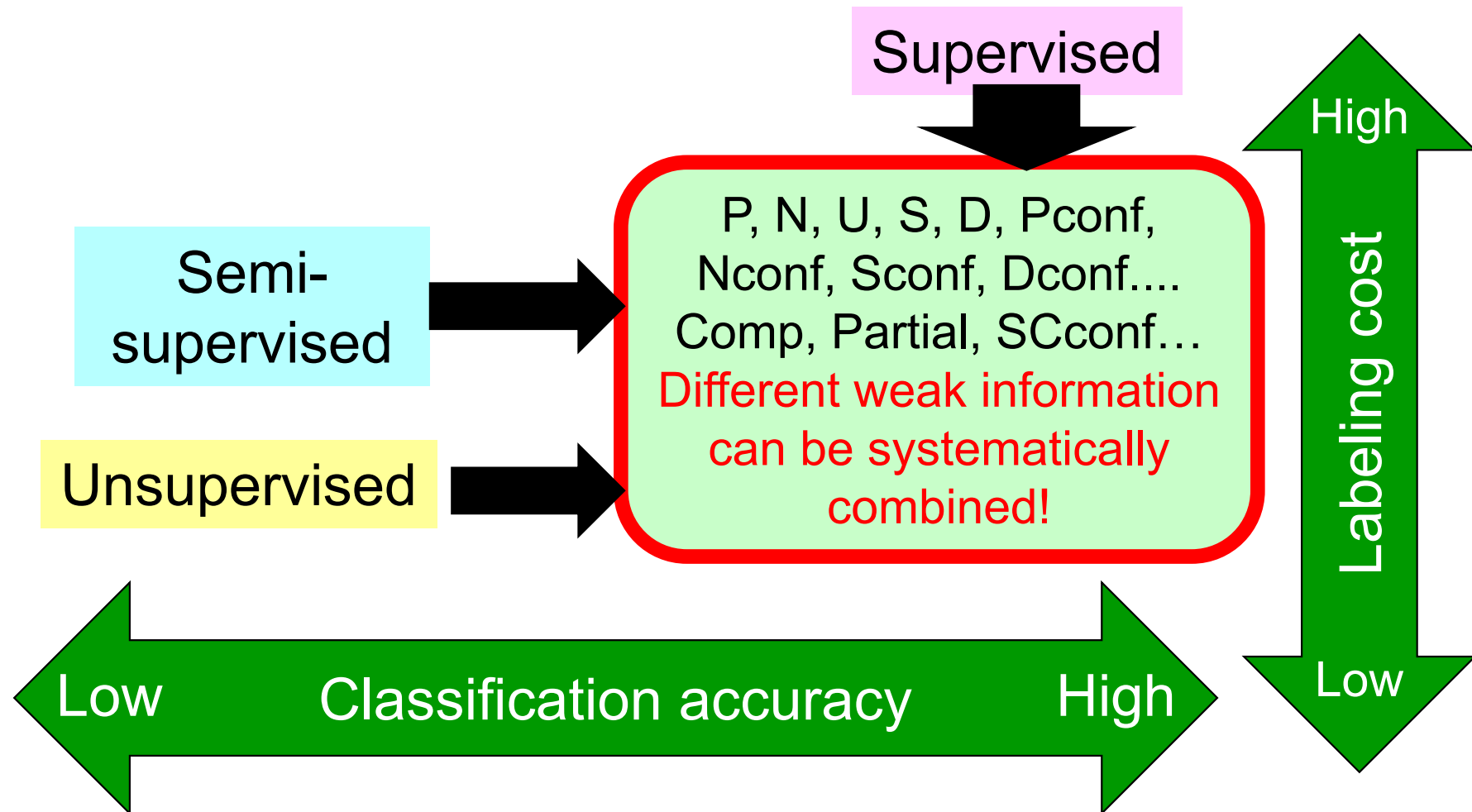
- **Partial labels:** Specify a subset of classes that contains the correct one (“1 or 2”).

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)



Summary: Empirical Risk Minimization Framework for Weakly Supervised Learning



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

Coming
soon



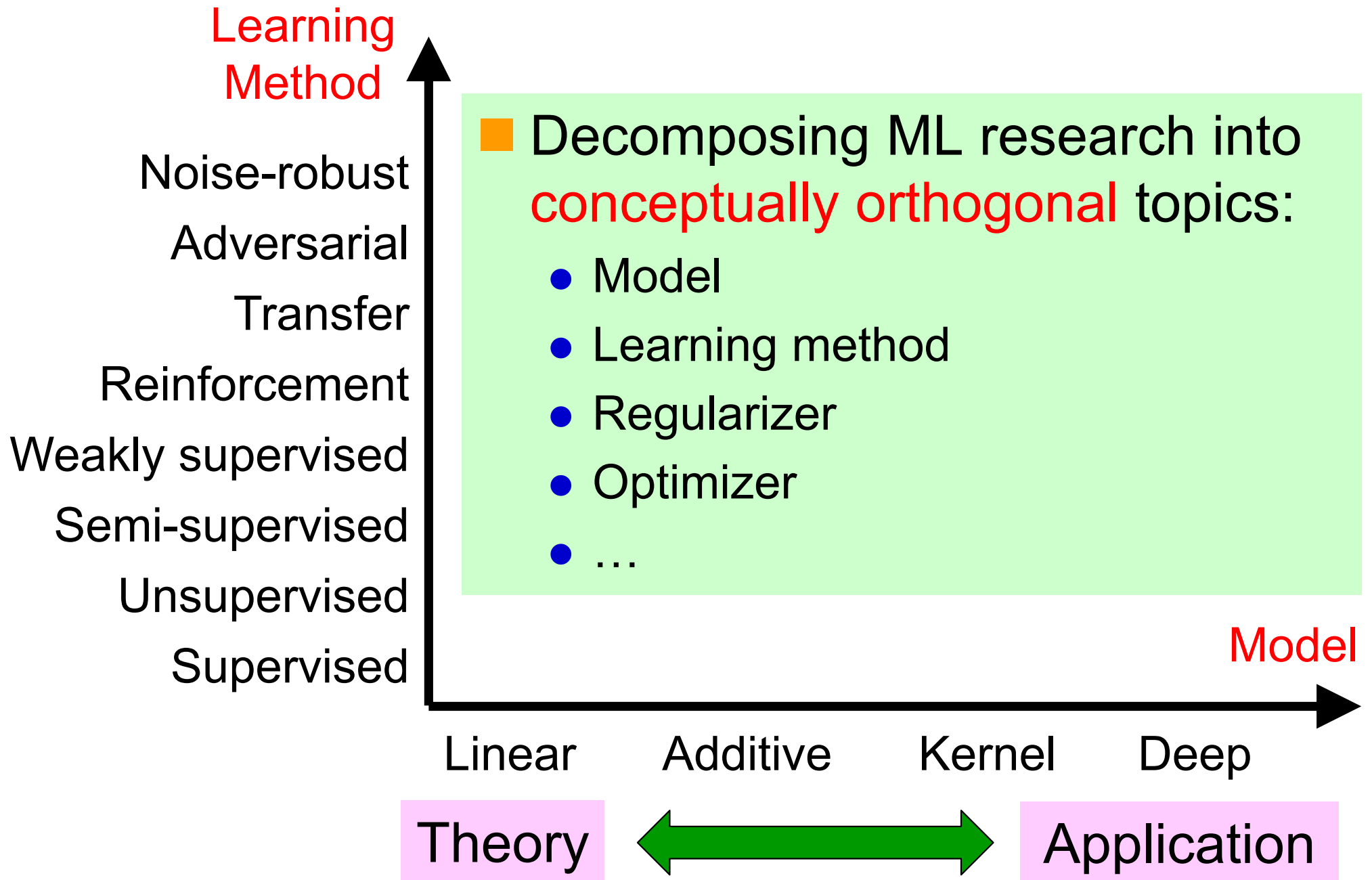
Contents

1. Transfer learning
2. Weakly supervised classification
3. Future outlook

Challenges in Reliable Machine Learning

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

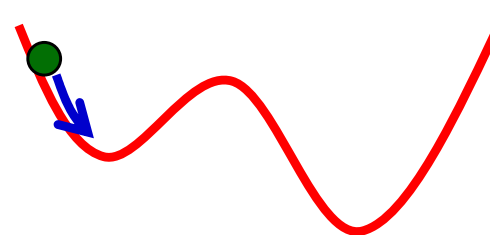
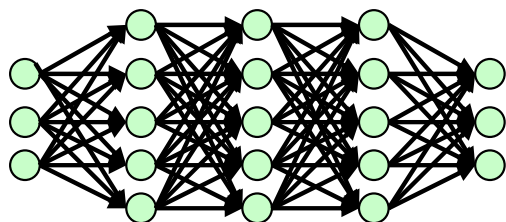
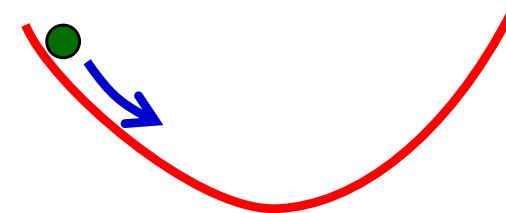
Axes of ML Research



Technological Breakthroughs

47

- Classical convex learning methods allow us to analyze the global solution.
- Since optimization in deep learning is complex, stochastic gradient descent is used.



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