

Generalizing Importance Weighting to A Universal Solver for Distribution Shift Problems

Tongtong Fang¹



Nan Lu^{2,3}



Gang Niu³



Masashi Sugiyama^{3,1}



¹ Univ. of Tokyo, Japan,

² Univ. of Tübingen, Germany,

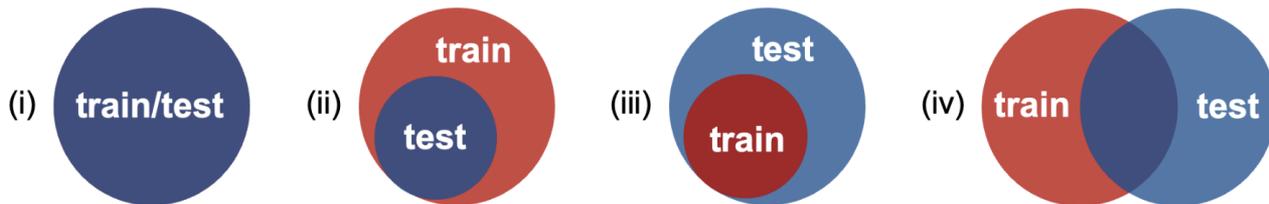
³ RIKEN, Japan

Two-levels of distribution shift

- Two levels of distribution shift (DS) $p_{\text{tr}}(\mathbf{x}, y) \neq p_{\text{te}}(\mathbf{x}, y)$:
 - The data distribution itself changes, e.g., covariate shift
 - The underlying *support** of data distribution changes.

* The set where the probability density is non-zero.

- The relationship between the training and test support



- Existing methods are good at cases (i) & (ii)
- Cases (iii) & (iv) are common due to data-collection biases

A real-world example of case (iii)

Family Felidae (猫科, ネコ科)



Training data:

Tiger, ocelot, cat, leopard,
puma, caracal, lion...

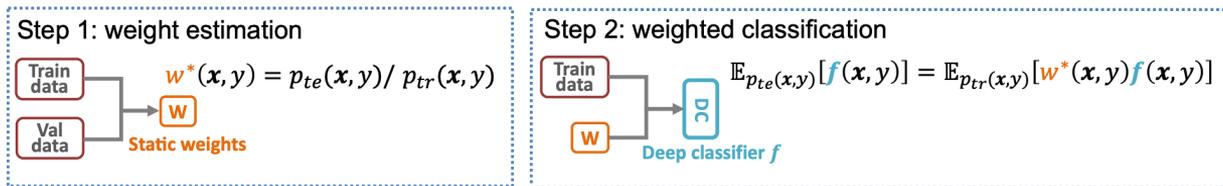


Test data:

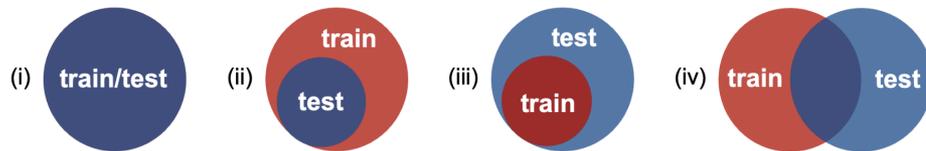
Tiger, ocelot, cat, leopard,
puma, caracal, lion...

Importance weighting (IW) & dynamic IW

- Problem setting:
 - A training set $\{(x_i^{\text{tr}}, y_i^{\text{tr}})\}_{i=1}^{n_{\text{tr}}} \sim p_{\text{tr}}(x, y)$ & a validation set $\{(x_i^{\text{v}}, y_i^{\text{v}})\}_{i=1}^{n_{\text{v}}} \sim p_{\text{te}}(x, y)$, $n_{\text{tr}} \gg n_{\text{v}}$
 - Estimate the risk: $R(\mathbf{f}) = \mathbb{E}_{p_{\text{te}}(x, y)}[\ell(\mathbf{f}(x), y)]$.
- IW is a golden solver for DS in cases (i) and (ii)^[1]



- Dynamic IW (DIW) makes IW work well for deep learning^[2]
- However, IW methods are problematic in cases (iii) & (iv)



[1] M. Sugiyama et al., Machine learning in non-stationary environments: Introduction to covariate shift adaptation. The MIT Press, 2012.

[2] T. Fang et al., Rethinking importance weighting for deep learning under distribution shift. In NeurIPS, 2020.

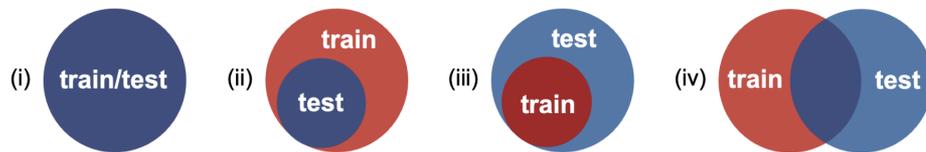
A deeper understanding of IW

- Recall that the risk is: $R(\mathbf{f}) = \mathbb{E}_{p_{\text{te}}(\mathbf{x}, y)}[\ell(\mathbf{f}(\mathbf{x}), y)]$
- The objective of IW is: $J(\mathbf{f}) = \mathbb{E}_{p_{\text{tr}}(\mathbf{x}, y)}[w^*(\mathbf{x}, y)\ell(\mathbf{f}(\mathbf{x}), y)]$

Definition: Given an (expected) objective $J(\mathbf{f})$, it is **risk-consistent** if $J(\mathbf{f}) = R(\mathbf{f})$ for any \mathbf{f} , i.e., the objective is equal to the original risk for any classifier; it is **classifier-consistent** if $\arg \min_{\mathbf{f}} J(\mathbf{f}) = \arg \min_{\mathbf{f}} R(\mathbf{f})$ where the minimization is taken over all measurable functions, i.e., the objective shares the optimal classifier with the original risk.

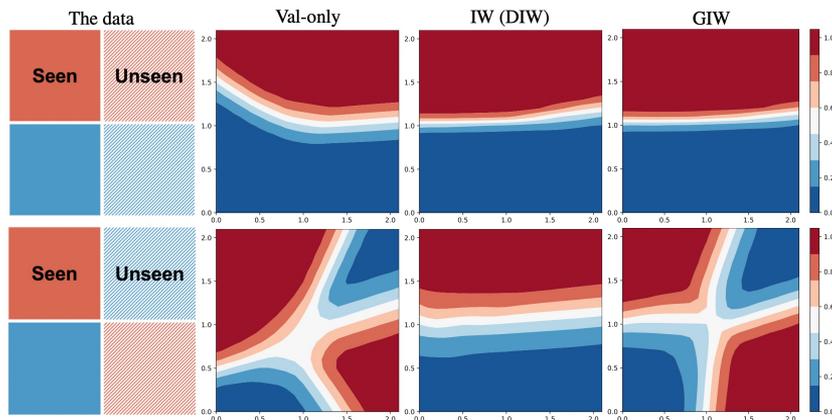
Theorem 1. *In cases (i) and (ii), IW is risk-consistent.*

Theorem 2. *In cases (iii) and (iv), IW is risk-inconsistent, and it holds that $J(\mathbf{f}) < R(\mathbf{f})$ for any \mathbf{f} .*



Two concrete examples

- **Task:** classify red/blue data
- **Data:**
 - Train: left two squares
 - Test: all four squares
 - Validation: 1 data per square
- **Methods:**
 - Val-only: use only validation data to train the model
 - IW: importance weighting
 - GIW: the proposed method



IW is classifier-consistent

IW is classifier-inconsistent

Generalized importance weighting (GIW)

- We split the test support \mathcal{S}_{te} into **in-training** $\mathcal{S}_{te} \cap \mathcal{S}_{tr}$ and **out-of-training** $\mathcal{S}_{te} \setminus \mathcal{S}_{tr}$ parts.
- A support-splitting variable $s \in \{0,1\}$, then

$$p(\mathbf{x}, y, s) = \begin{cases} p_{te}(\mathbf{x}, y) & \text{if } (\mathbf{x}, y) \in \mathcal{S}_{tr} \text{ and } s = 1, \text{ or } (\mathbf{x}, y) \in \mathcal{S}_{te} \setminus \mathcal{S}_{tr} \text{ and } s = 0, \\ 0 & \text{if } (\mathbf{x}, y) \in \mathcal{S}_{tr} \text{ and } s = 0, \text{ or } (\mathbf{x}, y) \in \mathcal{S}_{te} \setminus \mathcal{S}_{tr} \text{ and } s = 1. \end{cases}$$

- Expected objective of GIW is:

$$J_G(\mathbf{f}) = \underbrace{\alpha \mathbb{E}_{p_{tr}(\mathbf{x}, y)} [w^*(\mathbf{x}, y) \ell(\mathbf{f}(\mathbf{x}), y)]}_{\text{in-training (IT)}} + \underbrace{(1 - \alpha) \mathbb{E}_p(\mathbf{x}, y | s = 0) [\ell(\mathbf{f}(\mathbf{x}), y)]}_{\text{out-of-training (OOT)}}$$

where $\alpha = p(s = 1)$.

- In cases (i) and (ii), GIW is reduced to IW since $\alpha = 1$.

Theorem 4. *GIW is always risk-consistent for distribution shift problems.*

Implementation of GIW

- Split validation data and estimate α
 - Pretrain on training data to obtain a feature extractor
 - Train a one-class SVM^[3] on latent representation of training data
 - Split validation data into **IT** $\{(\mathbf{x}_i^{v1}, y_i^{v1})\}_{i=1}^{n_{v1}}$ & **OOT** $\{(\mathbf{x}_i^{v2}, y_i^{v2})\}_{i=1}^{n_{v2}}$ parts by the one-class SVM
 - Estimate α as $\hat{\alpha} = \frac{n_{v1}}{n_v}$
- Empirical objective of GIW is:

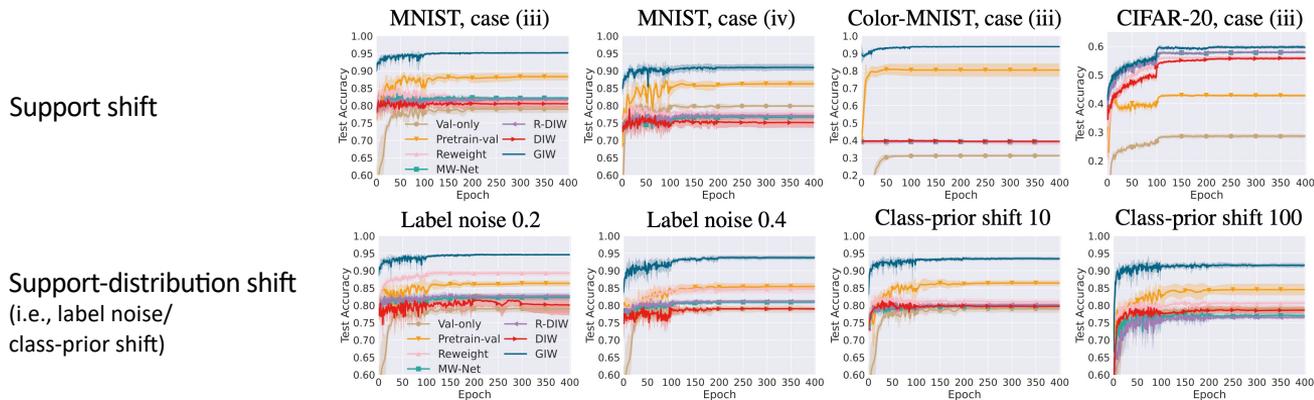
$$\hat{J}_G(\mathbf{f}) = \frac{n_{v1}}{n_v n_{tr}} \sum_{i=1}^{n_{tr}} \hat{w}(\mathbf{x}_i^{tr}, y_i^{tr}) \ell(\mathbf{f}(\mathbf{x}_i^{tr}), y_i^{tr}) + \frac{1}{n_v} \sum_{j=1}^{n_{v2}} \ell(\mathbf{f}(\mathbf{x}_j^{v2}), y_j^{v2}).$$

[3] B. Schölkopf et al., Support vector method for novelty detection. In NeurIPS, 1999.

Experiments

- Two distribution shift (DS) patterns
 - Support shift: DS solely comes from the support mismatch
 - Support-distribution shift: add DS on top of the support shift
- Setups & results*

Dataset	Task (classification for)	Training data	Test data	Model
MNIST	odd and even digits	4 digits (0-3)	10 digits (0-9)*	LeNet-5
Color-MNIST	10 digits	digits all in red	digits in red/blue/green	LeNet-5
CIFAR-20	objects in 20 superclasses	2 classes/superclass	5 classes/superclass	ResNet-18



*Results of support-distribution shift here are on MNIST; results on more datasets/baselines can be found in the paper.

Thanks for your attention!