

Positive-Unlabeled Classification under Class Prior Shift and Asymmetric Error

Nontawat Charoenphakdee^{1,2} and Masashi Sugiyama^{2,1}

The University of Tokyo¹

RIKEN AIP²



東京大学
THE UNIVERSITY OF TOKYO



Supervised binary classification (PN classification)

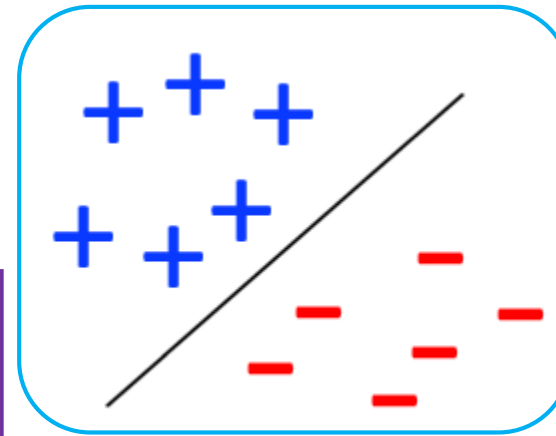
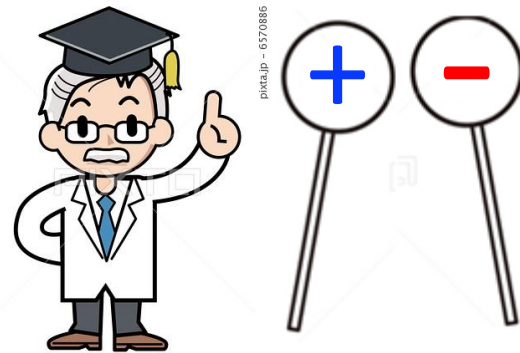
Positive and Negative data are given.

Data collection

Features (input)

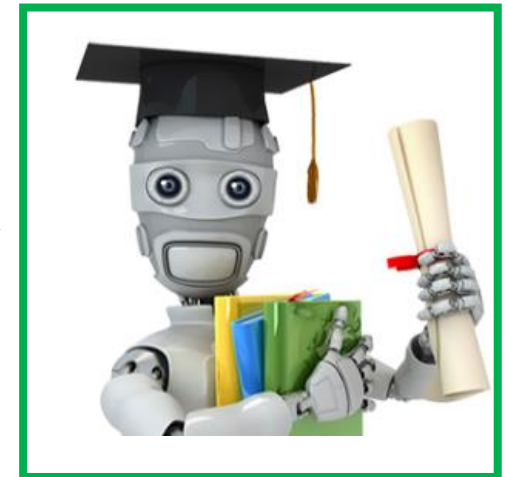


Labels (output)



Machine learning

Binary Classifier



Positive-unlabeled classification (PU classification)

Positive and **Unlabeled** data are given.

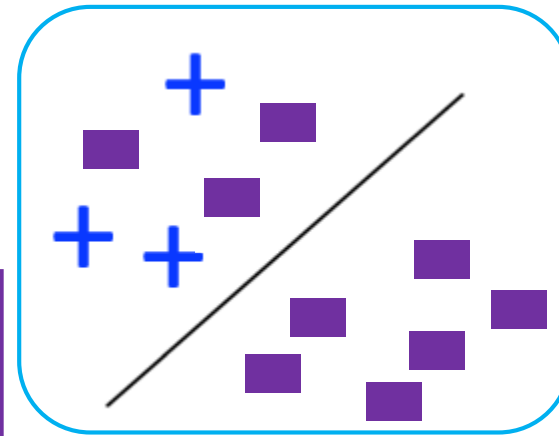
Data collection

Features (input)

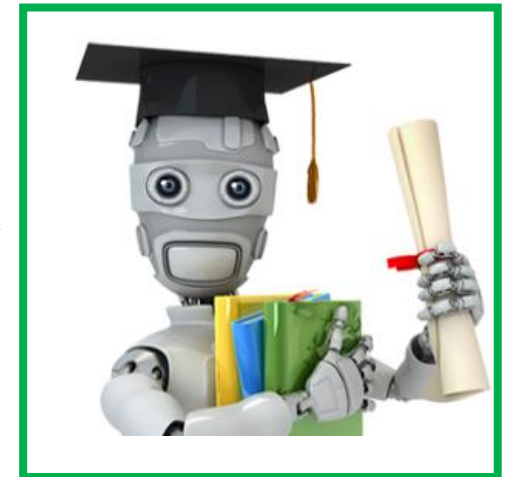
Labels (output)



pixtalip - 6570886



Binary Classifier



Why PU classification?

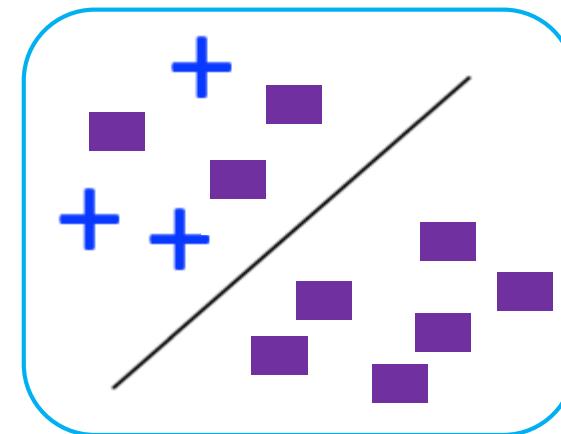
Unlabeled data are **cheaper** to obtain.

Sometimes, **negative** data are **hard to describe**.

In some real-world applications, **collecting negative data is impossible**.

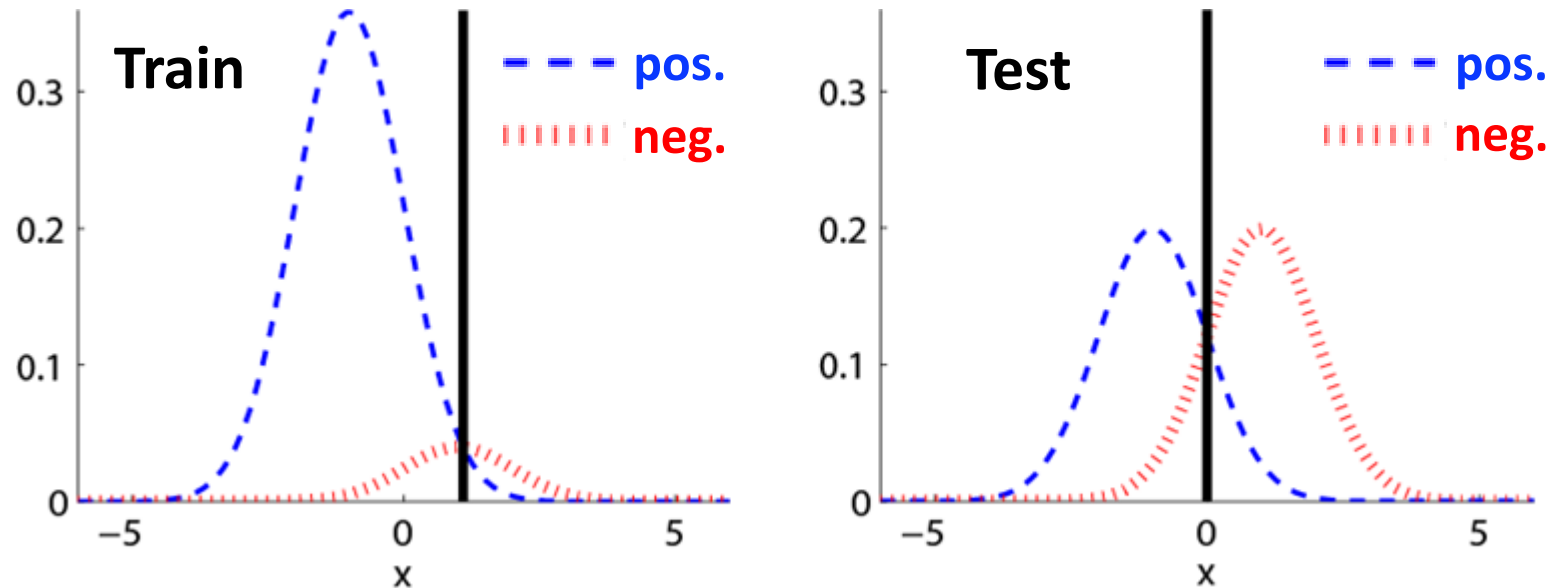
Applications:

- Bioinformatics (Yang+, 2012, Singh-Blom+ 2013, Ren+, 2015)
- Text classification (Li+, 2003)
- Time series classification (Nguyen+, 2011)
- Medical diagnosis (Zuluaga+, 2011)
- Remote-sensing classification (Li+, 2011)



Class prior shift

The ratio of **positive-negative** in the **training** and **test** data are different.



Decision boundary is also shifted



Lead to low accuracy!

Examples:

- Collect unlabeled data from **the internet**.
- Collect unlabeled data from **all users/patients/etc.** for **personalized application**.

Class prior shift (cont.)

Existing **PU classification** work assumes class prior of **training** and **test** data are the same (du Plessis+, 2014 2015, Kiryo+, 2017).

Existing class prior shift work is not applicable since they require **positive-negative** data (Saerens, 2002, du Plessis+, 2012).

PU classification under class prior shift

Observed

Given: Two sets of data

$$\begin{aligned} \pi &: p(y = 1) \\ \text{pos}(\mathbf{x}) &: p(\mathbf{x} | y = 1) \\ \text{neg}(\mathbf{x}) &: p(\mathbf{x} | y = -1) \end{aligned}$$

Positive $X_P := \{\mathbf{x}_i^P\}_{i=1}^{n_P} \stackrel{\text{i.i.d.}}{\sim} \text{pos}(\mathbf{x})$

Unlabeled $X_U := \{\mathbf{x}_j^U\}_{j=1}^{n_U} \stackrel{\text{i.i.d.}}{\sim} \pi_{\text{tr}} \text{pos}(\mathbf{x}) + (1 - \pi_{\text{tr}}) \text{neg}(\mathbf{x})$

Unobserved

$$\pi_{\text{tr}} \neq \pi_{\text{te}} : \text{Class prior shift!}$$

Test $X_{\text{te}} := \{\mathbf{x}_k^{\text{te}}\}_{k=1}^{n_{\text{te}}} \stackrel{\text{i.i.d.}}{\sim} \pi_{\text{te}} \text{pos}(\mathbf{x}) + (1 - \pi_{\text{te}}) \text{neg}(\mathbf{x})$

Q: Does class prior shift heavily degrade the performance?

Classifier may fail miserably under class prior shift...

Accuracy reported in mean and std. error of 10 trials with density ratio method.

Accuracy drops heavily!!

Our method

Dataset	Accuracy (no shift)
<i>banana</i>	90.1 (0.6)
<i>ijcnn1</i>	72.9 (0.4)
<i>MNIST</i>	86.0 (0.4)
<i>susy</i>	79.5 (0.5)
<i>cod-rna</i>	87.4 (0.6)
<i>magic</i>	76.7 (0.5)

Accuracy (shifted)
82.3 (0.5)
37.8 (0.7)
69.8 (0.7)
57.5 (0.9)
78.5 (0.6)
60.6 (1.4)

Accuracy (shifted)
87.9 (0.3)
71.7 (0.3)
82.5 (0.6)
75.9 (0.5)
84.7 (0.4)
79.0 (0.5)

No shift: $\pi_{tr} = \pi_{te} = 0.3$

Shift! $\pi_{tr} = 0.7, \pi_{te} = 0.3$

Problem setting

$$\begin{array}{ll} \pi : p(y = 1) & \mathbb{E}_{\mathcal{P}}[\cdot] : \mathbb{E}_{\mathbf{x} \sim \text{pos}(\mathbf{x})}[\cdot] \\ \text{pos}(\mathbf{x}) : p(\mathbf{x} | y = 1) & \mathbb{E}_{\mathcal{N}}[\cdot] : \mathbb{E}_{\mathbf{x} \sim \text{neg}(\mathbf{x})}[\cdot] \\ \text{neg}(\mathbf{x}) : p(\mathbf{x} | y = -1) & \end{array}$$

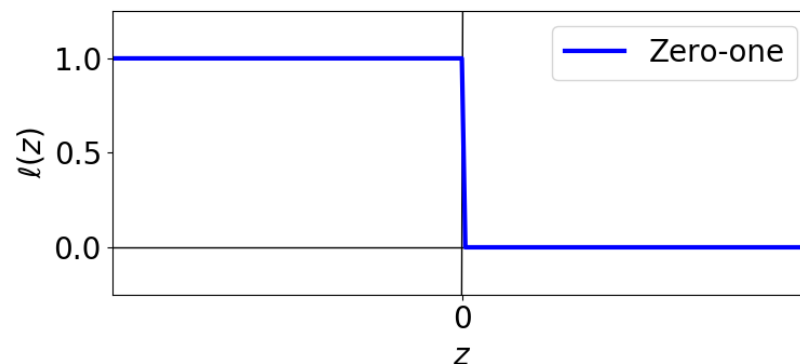
- **Given:** Two sets of data and test class prior π_{te}

Positive $X_{\text{P}} := \{\mathbf{x}_i^{\text{P}}\}_{i=1}^{n_{\text{P}}} \stackrel{\text{i.i.d.}}{\sim} \text{pos}(\mathbf{x})$

Unlabeled $X_{\text{U}} := \{\mathbf{x}_j^{\text{U}}\}_{j=1}^{n_{\text{U}}} \stackrel{\text{i.i.d.}}{\sim} \pi_{\text{tr}} \text{pos}(\mathbf{x}) + (1 - \pi_{\text{tr}}) \text{neg}(\mathbf{x})$

- **Goal:** Find a prediction function g that minimizes

$$R_{\text{Shift}}^{\ell_{0-1}}(g) = \pi_{\text{te}} \mathbb{E}_{\mathcal{P}}[\ell_{0-1}(g(\mathbf{x}))] + (1 - \pi_{\text{te}}) \mathbb{E}_{\mathcal{N}}[\ell_{0-1}(-g(\mathbf{x}))]$$



Proposed methods

We proposed two approaches for **PU classification** under **class prior shift**:

- **Risk minimization approach:**

Learn a classifier based on **empirical risk minimization** principle (Vapnik, 1998).

- **Density ratio approach:**

1. Estimate a **density ratio** of **positive** and **unlabeled** densities.
2. Use an appropriate threshold to classify.

Later, we will show that our methods are also applicable for

PU classification with **asymmetric error**.

Risk minimization approach

Consider the following classification risk:

$$R_{\text{Shift}}^{\ell_{0-1}}(g) = \pi_{\text{te}} \mathbb{E}_{\mathbf{P}} [\ell_{0-1}(g(\mathbf{x}))] + (1 - \pi_{\text{te}}) \mathbb{E}_{\mathbf{N}} [\ell_{0-1}(-g(\mathbf{x}))]$$

With $\mathbb{E}_{\mathbf{u}} [\cdot] = \pi_{\text{tr}} \mathbb{E}_{\mathbf{P}} [\cdot] + (1 - \pi_{\text{tr}}) \mathbb{E}_{\mathbf{N}} [\cdot]$, we can rewrite $R_{\text{Shift}}^{\ell_{0-1}}(g)$ as

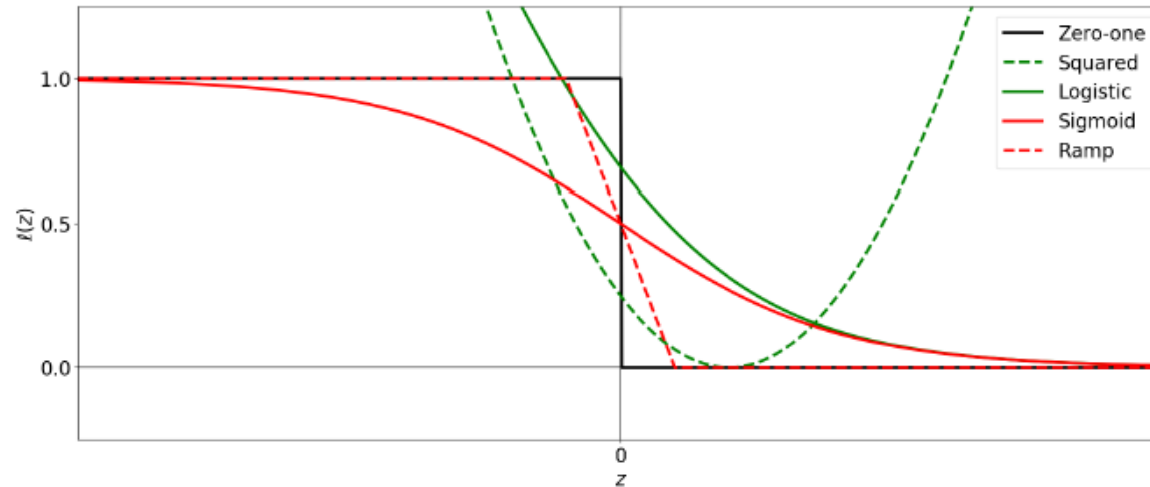
$$R_{\text{Shift}}^{\ell_{0-1}}(g) = \mathbb{E}_{\mathbf{P}} \left[\pi_{\text{te}} \ell_{0-1}(g(\mathbf{x})) - \frac{\pi_{\text{tr}}(1 - \pi_{\text{te}})}{1 - \pi_{\text{tr}}} \ell_{0-1}(-g(\mathbf{x})) \right] + \frac{1 - \pi_{\text{te}}}{1 - \pi_{\text{tr}}} \mathbb{E}_{\mathbf{u}} [\ell_{0-1}(-g(\mathbf{x}))]$$

Equivalent to existing methods (du Plessis+, 2015) if $\pi_{\text{tr}} = \pi_{\text{te}}$.

No access to distribution: we minimize **empirical error** (Vapnik, 1998):

$$\hat{R}_{\text{PU-shift}}^{\ell_{0-1}}(g) = \frac{1}{n_{\mathbf{P}}} \sum_{i=1}^{n_{\mathbf{P}}} \left[\pi_{\text{te}} \ell_{0-1}(g(\mathbf{x}_i^{\mathbf{P}})) - \frac{\pi_{\text{tr}}(1 - \pi_{\text{te}})}{1 - \pi_{\text{tr}}} \ell_{0-1}(-g(\mathbf{x}_i^{\mathbf{P}})) \right] + \frac{1}{n_{\mathbf{U}}} \frac{1 - \pi_{\text{te}}}{1 - \pi_{\text{tr}}} \sum_{j=1}^{n_{\mathbf{U}}} \ell_{0-1}(-g(\mathbf{x}_j^{\mathbf{U}}))$$

Surrogate losses for binary classification



Directly minimize 0-1 loss is difficult.

- NP-Hard, discontinuous, not differentiable (Ben-david+, 2003, Feldman+, 2012)

In practice, minimize a **surrogate loss** (regularization can also be added):

$$\widehat{R}_{\text{PU-shift}}^{\ell}(g) = \frac{1}{n_{\text{P}}} \sum_{i=1}^{n_{\text{P}}} \left[\pi_{\text{te}} \ell(g(\mathbf{x}_i^{\text{p}})) - \frac{\pi_{\text{tr}}(1 - \pi_{\text{te}})}{1 - \pi_{\text{tr}}} \ell(-g(\mathbf{x}_i^{\text{p}})) \right] + \frac{1}{n_{\text{U}}} \frac{1 - \pi_{\text{te}}}{1 - \pi_{\text{tr}}} \sum_{j=1}^{n_{\text{U}}} \ell(-g(\mathbf{x}_j^{\text{u}}))$$

Density ratio estimation

Goal: Estimate the density ratio:

$$r(\mathbf{x}) = \frac{p_{\text{nu}}(\mathbf{x})}{p_{\text{de}}(\mathbf{x})}$$

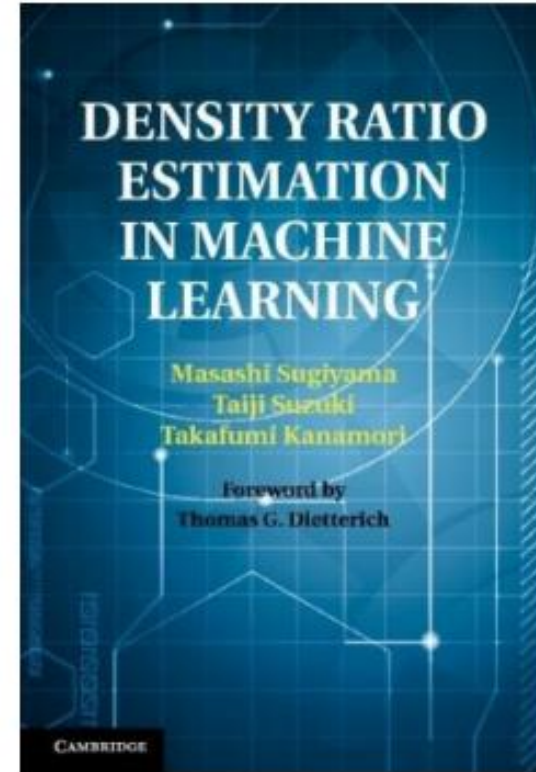
from two sets of data

$$X_{\text{nu}} := \{\mathbf{x}_i^{\text{nu}}\}_{i=1}^{n_{\text{nu}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{nu}}(\mathbf{x})$$

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

Applications: outlier detection (Hido+, 2011),
change-point detection (Liu+, 2013), robot control (Hachiya+, 2009)
event detection in images/movies/text (Yamanaka, 2011, Matsugu, 2011, Liu, 2012), etc.

Naïve approach: estimate $\hat{p}_{\text{nu}}(\mathbf{x})$, $\hat{p}_{\text{de}}(\mathbf{x})$ separately then perform division $\frac{\hat{p}_{\text{nu}}(\mathbf{x})}{\hat{p}_{\text{de}}(\mathbf{x})}$.
Does not work well (estimation error is amplified from division operation).



Please check this book to learn more about density ratio estimation (Sugiyama+, 2012)

Unconstrained least-squares important fitting (uLSIF)

Goal: Estimate the density ratio: $r(\mathbf{x}) = \frac{p_{\text{nu}}(\mathbf{x})}{p_{\text{de}}(\mathbf{x})}$ (Kanamori+, 2012)

How: estimate \hat{r} by minimizing squared loss objective:

$$\text{SQ}(\hat{r}) = \int \left(\hat{r}(\mathbf{x}) - r(\mathbf{x}) \right)^2 p_{\text{de}}(\mathbf{x}) d\mathbf{x}$$

Squared loss decomposition:

$$\text{SQ}(\hat{r}) = \int \left(\hat{r}(\mathbf{x}) \right)^2 p_{\text{de}}(\mathbf{x}) d\mathbf{x} - 2 \int \hat{r}(\mathbf{x}) p_{\text{nu}}(\mathbf{x}) d\mathbf{x} + \text{Constant}$$

Empirical minimization (constant can be safely ignored):

$$\widehat{\text{SQ}}(\hat{r}) = \frac{1}{n_{\text{de}}} \sum_{j=1}^{n_{\text{de}}} \left(\hat{r}(\mathbf{x}_j^{\text{de}}) \right)^2 - \frac{2}{n_{\text{nu}}} \sum_{i=1}^{n_{\text{nu}}} \hat{r}(\mathbf{x}_i^{\text{nu}})$$

Unconstrained least-squares important fitting (cont.)

Model: linear-in parameter model

(Kanamori+, 2012)

$$\hat{r}(\mathbf{x}) = \sum_b \theta_b \phi_b(\mathbf{x}) = \boldsymbol{\theta}^\top \boldsymbol{\phi}(\mathbf{x})$$

$\phi_b(\mathbf{x})$: basis function
(e.g., Gaussian kernel)

$$\widehat{\mathbf{H}} = \frac{1}{n_{\text{de}}} \sum_{j=1}^{n_{\text{de}}} \phi(\mathbf{x}_j^{\text{de}}) \phi(\mathbf{x}_j^{\text{de}})^\top$$

$$\hat{\mathbf{h}} = \frac{1}{n_{\text{nu}}} \sum_{i=1}^{n_{\text{nu}}} \phi(\mathbf{x}_i^{\text{nu}})$$

Objective:

$$\min_{\boldsymbol{\theta}} \left[\frac{1}{2} \boldsymbol{\theta}^\top \widehat{\mathbf{H}} \boldsymbol{\theta} - \hat{\mathbf{h}}^\top \boldsymbol{\theta} + \frac{\lambda}{2} \boldsymbol{\theta}^\top \boldsymbol{\theta} \right]$$

λ : regularization parameter
 \mathbf{I} : identity matrix

Global solution can be computed **analytically**: $\hat{\boldsymbol{\theta}} = (\widehat{\mathbf{H}} + \lambda \mathbf{I})^{-1} \hat{\mathbf{h}}$

Parameter tuning (regularization, basis) can be done by **cross-validation**.

Density ratio approach

$$\text{pos}(\mathbf{x}) : p(\mathbf{x} | y = 1)$$

$$\text{neg}(\mathbf{x}) : p(\mathbf{x} | y = -1)$$

$$\text{unl}(\mathbf{x}) = \pi_{\text{tr}} \text{pos}(\mathbf{x}) + (1 - \pi_{\text{tr}}) \text{neg}(\mathbf{x})$$

Consider Bayes-optimal classifier of binary classification (no prior shift)

$$f_{\text{Bayes}}^*(\mathbf{x}) = \text{sign} \left[p(y = +1 | \mathbf{x}) - \frac{1}{2} \right]$$

We can rewrite it as

$$f_{\text{Bayes}}^*(\mathbf{x}) = \text{sign} \left[\pi_{\text{tr}} \frac{\text{pos}(\mathbf{x})}{\text{unl}(\mathbf{x})} - \frac{1}{2} \right]$$

Density ratio!

Another formulation is

$$f_{\text{Bayes}}^*(\mathbf{x}) = \text{sign} \left[\pi_{\text{tr}} - \frac{1}{2} \frac{\text{unl}(\mathbf{x})}{\text{pos}(\mathbf{x})} \right]$$

Q1: How to modify when class prior shift occurs?

Q2: Which formulation is preferable?

Q1: Density ratio approach (shift)

Consider Bayes-optimal classifier of binary classification

$$f_{\text{Bayes}}^*(\mathbf{x}) = \text{sign} \left[p(y = +1|\mathbf{x}) - \frac{1}{2} \right]$$

$$\begin{aligned} \text{pos}(\mathbf{x}) &: p(\mathbf{x}|y = 1) \\ \text{neg}(\mathbf{x}) &: p(\mathbf{x}|y = -1) \\ \text{unl}(\mathbf{x}) &= \pi_{\text{tr}}\text{pos}(\mathbf{x}) + (1 - \pi_{\text{tr}})\text{neg}(\mathbf{x}) \end{aligned}$$

We can rewrite it as

$$f_{\text{Bayes}}^*(\mathbf{x}) = \text{sign} \left[\pi_{\text{tr}} \frac{\text{pos}(\mathbf{x})}{\text{unl}(\mathbf{x})} - \frac{\pi_{\text{tr}}(1 - \pi_{\text{te}})}{\pi_{\text{te}} + \pi_{\text{tr}} - 2\pi_{\text{tr}}\pi_{\text{te}}} \right]$$

Density ratio!

Another formulation is

$$f_{\text{Bayes}}^*(\mathbf{x}) = \text{sign} \left[\frac{\pi_{\text{te}} + \pi_{\text{tr}} - 2\pi_{\text{tr}}\pi_{\text{te}}}{(1 - \pi_{\text{te}})} - \frac{\text{unl}(\mathbf{x})}{\text{pos}(\mathbf{x})} \right]$$

Simply modifying the threshold can solve this problem!

Q2: Difficulty of density ratio estimation

In general, density ratio is **unbounded**. 😞

$$r(\mathbf{x}) = \frac{p_{\text{nu}}(\mathbf{x})}{p_{\text{de}}(\mathbf{x})}$$

$r(\mathbf{x})$ is unbounded when $p_{\text{de}}(\mathbf{x}) = 0$.

This raises issues of robustness and stability.

We show that the density ratio $\frac{p_{\text{pos}}(\mathbf{x})}{p_{\text{unl}}(\mathbf{x})}$ is bounded in PU classification. 😊

Q2: Density ratio in PU

$$\text{pos}(\mathbf{x}) : p(\mathbf{x} | y = 1)$$

$$\text{neg}(\mathbf{x}) : p(\mathbf{x} | y = -1)$$

$$\text{unl}(\mathbf{x}) = \pi_{\text{tr}} \text{pos}(\mathbf{x}) + (1 - \pi_{\text{tr}}) \text{neg}(\mathbf{x})$$

In **PU classification**, density ratio $\frac{\text{pos}(\mathbf{x})}{\text{unl}(\mathbf{x})}$ is bounded.

$$0 \leq \frac{\text{pos}(\mathbf{x})}{\text{unl}(\mathbf{x})} \leq \frac{1}{\pi_{\text{tr}}}$$

Lower and upper bounded 

$$\pi_{\text{tr}} \leq \frac{\text{unl}(\mathbf{x})}{\text{pos}(\mathbf{x})}$$

Unbounded from above 

Insight: estimate $\frac{\text{pos}(\mathbf{x})}{\text{unl}(\mathbf{x})}$ is preferable.

Our experimental results agree with this observation.

Experiments: class prior shift train 0.7 -> test 0.3

Datasets: *banana, ijcnn1, MNIST, susy, cod-rna, magic*

Methods:

- Density ratio $\frac{p_{\text{pos}}(\mathbf{x})}{p_{\text{unl}}(\mathbf{x})}$ (**$\frac{p}{u}$ uLSIF**)
- Density ratio $\frac{p_{\text{unl}}(\mathbf{x})}{p_{\text{pos}}(\mathbf{x})}$ (**$\frac{u}{p}$ uLSIF**)
- Linear-in input model (Lin): Double hinge loss (**DH-Lin**), squared loss (**Sq-Lin**)
- Kernel model (Ker): Double hinge loss (**DH-Ker**), squared loss (**Sq-Ker**)

Parameter selection: (regularization, kernel width) 5-fold cross-validation.

We also investigated when wrong test class prior is given.

Results reported in mean and std. error of accuracy of 10 trials.

Outperforming methods are bolded based on one-sided t-test with significance level 5%.

Dataset information and more experiments and can be found in the paper.

Results: class prior shift $\pi_{tr} = 0.7, \pi_{te} = 0.3$

Dataset	π^g	$\frac{u}{p}$ uLSIF	$\frac{p}{u}$ uLSIF	DH-Lin	DH-Ker	Sq-Lin	Sq-Ker
banana	π' 0.3	83.0(1.0)	86.4 (0.5)	70.2(0.5)	78.3(1.0)	70.0(0.0)	83.4(0.4)
ijcnn1		70.8(0.6)	74.2 (0.7)	70.0(0.1)	69.8(0.2)	71.5(0.3)	69.2(0.5)
MNIST		79.3(0.5)	81.7 (0.5)	74.0(1.1)	82.4 (1.0)	52.3(1.4)	83.4 (0.9)
susy		74.3(0.5)	76.0 (0.3)	72.7(0.6)	70.0(0.0)	75.5 (1.4)	74.7(0.7)
cod-rna		82.1(1.0)	82.8(0.8)	87.3 (0.7)	77.3(0.8)	85.2 (1.1)	80.2(1.0)
magic		71.5(0.7)	75.8 (0.6)	72.7(1.1)	70.8(0.4)	75.0 (1.0)	72.9(0.7)
banana	0.5	84.7(1.1)	88.7 (0.7)	54.9(1.4)	81.7(1.6)	53.6(1.2)	83.8(1.3)
ijcnn1		64.9 (1.4)	66.6 (1.0)	60.4(1.4)	51.6(3.0)	62.2(1.2)	48.2(2.8)
MNIST		81.9(0.4)	84.1 (0.6)	72.5(1.0)	82.5(0.7)	52.9(1.1)	81.9(0.9)
susy		75.9 (1.1)	77.0 (0.6)	67.5(1.4)	75.5(0.6)	71.6(1.0)	72.8(1.1)
cod-rna		85.3 (0.7)	85.4 (0.5)	86.2 (0.7)	80.1(1.1)	86.5 (0.9)	81.2(1.2)
magic		67.6(0.8)	73.6 (0.9)	72.6 (0.7)	62.4(1.9)	71.8 (0.7)	68.9(0.8)
banana	π 0.7	80.6 (1.3)	82.1 (1.1)	31.8(0.9)	48.9(1.5)	30.0(0.0)	69.9(1.1)
ijcnn1		35.2(1.4)	42.4 (0.9)	30.0(0.0)	30.0(0.0)	32.4(0.5)	30.9(0.4)
MNIST		79.9 (0.7)	72.6(0.6)	71.1(1.1)	64.8(1.1)	64.0(0.6)	74.2(1.0)
susy		35.6(3.1)	44.2 (2.9)	30.0(0.0)	30.0(0.0)	42.0 (1.5)	36.8(1.3)
cod-rna		77.7 (2.2)	77.8 (2.1)	79.6 (0.7)	67.8(0.8)	78.2(0.5)	68.3(1.0)
magic		51.6(0.3)	60.3 (1.5)	56.2 (2.7)	32.8(0.7)	58.7 (1.4)	50.1(1.6)

Correct test prior is given

Wrong test prior is given

Traditional PU

Preferable method in our experiments

(density ratio $\frac{p}{u}$ uLSIF)

PU classification with asymmetric error

- **Given:** Given two sets of sample:

Positive $X_P := \{\mathbf{x}_i^P\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} \text{pos}(\mathbf{x})$

Unlabeled $X_U := \{\mathbf{x}_i^U\}_{i=1}^{n'} \stackrel{\text{i.i.d.}}{\sim} \pi_{\text{tr}} \text{pos}(\mathbf{x}) + (1 - \pi_{\text{tr}}) \text{neg}(\mathbf{x})$

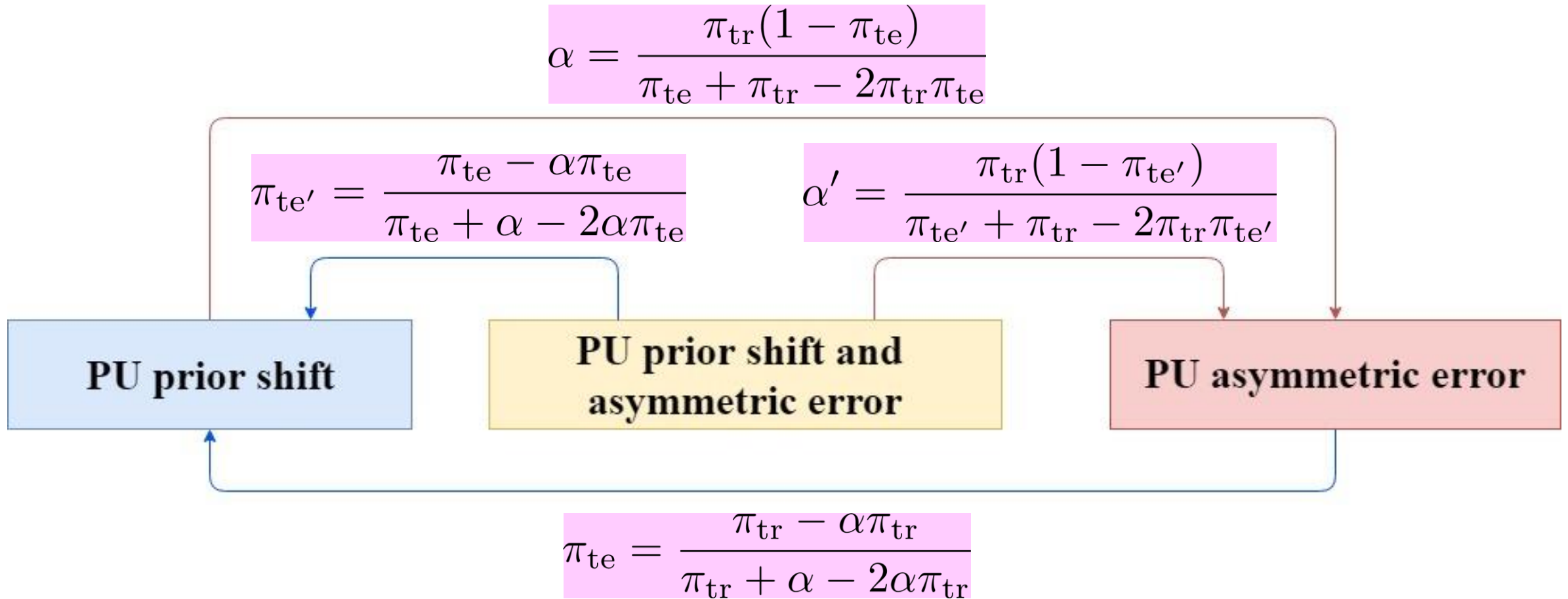
$$\begin{aligned} \pi_{\text{tr}} &: p(y = 1) & \mathbb{E}_P[\cdot] &: \mathbb{E}_{\mathbf{x} \sim \text{pos}(\mathbf{x})}[\cdot] \\ \text{pos}(\mathbf{x}) &: p(\mathbf{x} | y = 1) & \mathbb{E}_N[\cdot] &: \mathbb{E}_{\mathbf{x} \sim \text{neg}(\mathbf{x})}[\cdot] \\ \text{neg}(\mathbf{x}) &: p(\mathbf{x} | y = -1) & & \\ \alpha \in (0, 1) &: \text{false negative error} & & \end{aligned}$$

- **Goal:** Find a prediction function g that minimizes

$$R_{\text{Asym}}^{\ell}(g) = (1 - \alpha) \pi_{\text{tr}} \mathbb{E}_P[\ell(g(\mathbf{x}_P))] + \alpha (1 - \pi_{\text{tr}}) \mathbb{E}_N[\ell(-g(\mathbf{x}_N))]$$

Reduce to symmetric error when $\alpha = 0.5$

The equivalence of prior shift and asymmetric error



We can relate these problems based on the analysis of Bayes-optimal classifier.

Conclusion

Class prior shift may **heavily degrade** the performance of positive-unlabeled classification (**PU classification**). 😞

- Proposed **two approaches** for handling this problem **effectively**:
 - Risk minimization approach
 - Density ratio approach
- Showed the **equivalence** of **class prior shift** and **asymmetric error** problems in **PU classification**.
 - Our methods are **applicable** for **both problems**. 😊
 - Also **applicable** when considering **both problems simultaneously**. 😊
- Poster: #31: May 2nd from 7:00-9:00PM