

Positive-Unlabeled Classification under Class Prior Shift and Asymmetric Error



Nontawat Charoenphakdee^{1,2} Masashi Sugiyama^{2,1}
1: The University of Tokyo 2: RIKEN AIP



Summary

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

We propose **two frameworks** for solving this problem:

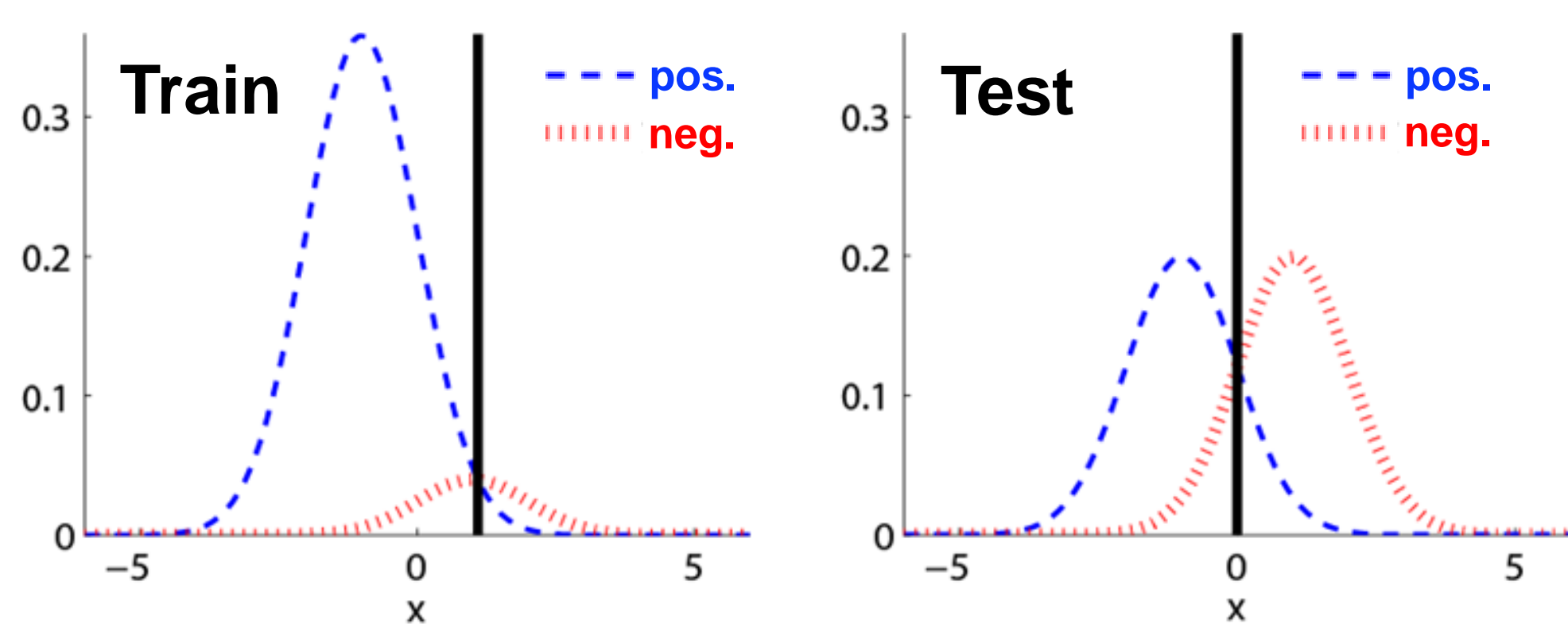
- Risk minimization framework
- Density ratio framework

We prove the **equivalence** of **class prior shift** and **asymmetric error** problems in **PU classification**.

Our methods are applicable for both problems!

Class prior shift

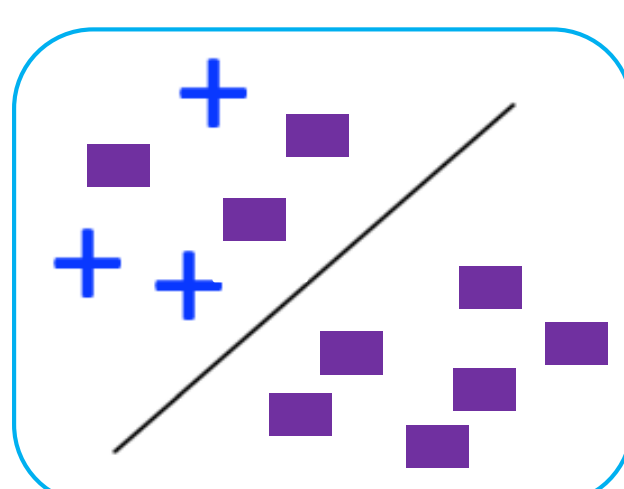
Positive-negative ratio in **training** and **test** data are different.



Decision boundary also shifts → **Lead to low accuracy!**

Practical examples: learn a classifier for **a specific user** from the internet and many users' information.

PU classification



Given: Two sets of data

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

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

$\pi : p(y=1)$
 $\text{pos}(\mathbf{x}) : p(\mathbf{x}|y=1)$
 $\text{neg}(\mathbf{x}) : p(\mathbf{x}|y=-1)$
 $\mathbb{E}_P[\cdot] : \mathbb{E}_{\mathbf{x} \sim \text{pos}(\mathbf{x})}[\cdot]$
 $\mathbb{E}_N[\cdot] : \mathbb{E}_{\mathbf{x} \sim \text{neg}(\mathbf{x})}[\cdot]$

Goal: Minimize either

Class prior shift classification risk:

$$R_{\text{Shift}}^{\ell_{0-1}}(g) = \pi_{\text{te}} \mathbb{E}_P[\ell_{0-1}(g(\mathbf{x}))] + (1 - \pi_{\text{te}}) \mathbb{E}_N[\ell_{0-1}(-g(\mathbf{x}))], \text{ or}$$

Asymmetric error classification risk:

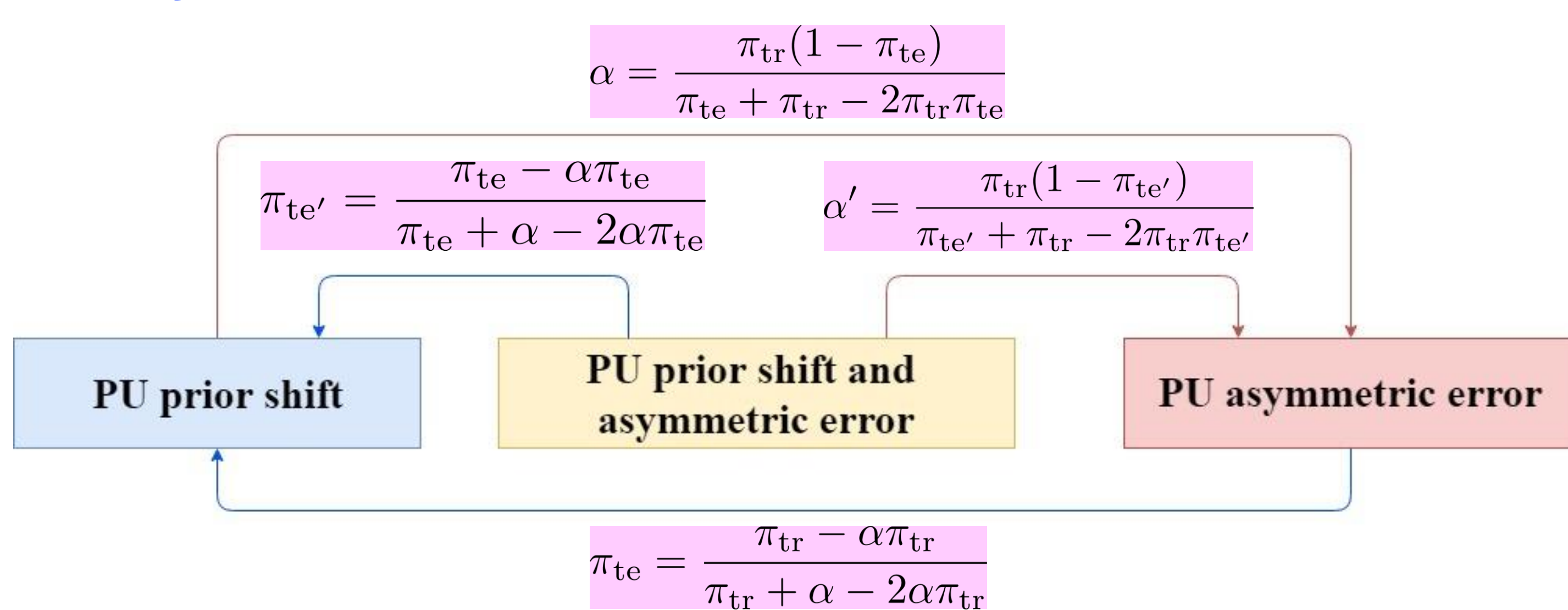
$$R_{\text{Asym}}^{\ell}(g) = (1 - \alpha) \pi_{\text{tr}} \mathbb{E}_P[\ell(g(\mathbf{x}))] + \alpha(1 - \pi_{\text{tr}}) \mathbb{E}_N[\ell(-g(\mathbf{x}))]$$

$\alpha \in (0, 1) : \text{false negative error}$

Existing PU classification work: **no class prior shift, no asymmetric error** (du Plessis+, 2015, Kiryo+, 2017).

Existing class prior shift / asymmetric error work: **require positive-negative data** (Saerens, 2002, Scott+, 2012, du Plessis+, 2012).

Equivalence of class prior shift and asymmetric error in PU classification



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

Risk minimization approach

$$\text{Risk: } R_{\text{Shift}}^{\ell_{0-1}}(g) = \pi_{\text{te}} \mathbb{E}_P[\ell_{0-1}(g(\mathbf{x}))] + (1 - \pi_{\text{te}}) \mathbb{E}_N[\ell_{0-1}(-g(\mathbf{x}))]$$

Using the following identity: $\mathbb{E}_U[\cdot] = \pi_{\text{tr}} \mathbb{E}_P[\cdot] + (1 - \pi_{\text{tr}}) \mathbb{E}_N[\cdot]$

We can equivalently express the risk as

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

Coincides with the existing method (du Plessis+, 2015) when $\pi_{\text{tr}} = \pi_{\text{te}}$.

Minimize empirical risk with surrogate loss (Bartlett+, 2006).

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

(regularization can also be added.)

Density ratio approach

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

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

can be equivalently expressed as

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

Another formulation:

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

Q: Which formulation is preferable?

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

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

$$0 \leq \frac{\text{pos}(\mathbf{x})}{\text{unl}(\mathbf{x})} \leq \frac{1}{\pi_{\text{tr}}} \quad \text{Lower and upper bounded} \quad \text{😊}$$

$$\pi_{\text{tr}} \leq \frac{\text{unl}(\mathbf{x})}{\text{pos}(\mathbf{x})} \quad \text{Unbounded from above} \quad \text{😞}$$

Naive approach: estimate $\widehat{\text{pos}}(\mathbf{x})$ and $\widehat{\text{unl}}(\mathbf{x})$ separately then calculate $\frac{\widehat{\text{pos}}(\mathbf{x})}{\widehat{\text{unl}}(\mathbf{x})}$.

Division operation amplifies the estimation error!

More effective direct approach:

unconstrained Least-squares Important Fitting (uLSIF) (Kanamori+, 2012).

Experiments $\pi_{\text{tr}} = 0.7, \pi_{\text{te}} = 0.3$

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

Methods:

Density ratio ($\frac{u}{u}$ LSIF, $\frac{u}{p}$ LSIF)

Linear-in input model: Double hinge (**DH-Lin**) and squared (**Sq-Lin**) losses. Kernel model (Ker): Double hinge (**DH-Ker**) and squared (**Sq-Ker**) losses.

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

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

Correct test prior is given

Wrong test prior is given

Traditional PU

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

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

References

- Saerens, Marco, Patrice Latinne, and Christine Decaestecker. "Adjusting the outputs of a classifier to new a priori probabilities: a simple procedure." *Neural computation* 14.1 (2002): 21-41.
- Scott, Clayton. "Calibrated asymmetric surrogate losses." *Electronic Journal of Statistics* 6 (2012): 958-992.
- Du Plessis, Marthinus Christoffel, and Masashi Sugiyama. "Semi-supervised learning of class balance under class-prior change by distribution matching." *Neural Networks* 50 (2014): 110-119.
- Du Plessis, Marthinus C., Gang Niu, and Masashi Sugiyama. "Analysis of learning from positive and unlabeled data." *NeurIPS*. 2014.
- Du Plessis, Marthinus, Gang Niu, and Masashi Sugiyama. "Convex formulation for learning from positive and unlabeled data." *ICML*. 2015.
- Kiryo, Ryuichi, et al. "Positive-unlabeled learning with non-negative risk estimator." *NeurIPS*. 2017.
- Bartlett, Peter L., Michael I. Jordan, and Jon D. McAuliffe. "Convexity, classification, and risk bounds." *JASA* 101.473 (2006): 138-156.
- Kanamori, Takafumi, Taiji Suzuki, and Masashi Sugiyama. "Statistical analysis of kernel-based least-squares density-ratio estimation." *Machine Learning* 86.3 (2012): 335-367.