IBIS2017 11/08

Positive-Unlabeled Learning with Non-Negative Risk Estimator

Ryuichi Gang Martins Masashi Kiryo^{1,2} Niu^{1,2} Christoffel du Sugiyama^{2,1} Plessis



Binary classification

•Classify input data $X \in \mathbb{R}^d$ to class $Y \in \{+1 \text{ (positive)}, -1 \text{ (negative)}\}$

•Supervised learning: Classifier $g: \mathbb{R}^d \to \mathbb{R}$ is learned from positive data and negative data



Supervised binary classification (PN learning)

•Goal: minimize expected risk

$$R(g) = \mathbb{E}_{(X,Y)\sim p(x,y)}[l(g(X),Y)] = \frac{\pi_{p}\mathbb{E}_{p}[l(g(X),+1)]}{\pi_{p}\mathbb{E}_{p}[l(g(X),+1)]} + \frac{\pi_{n}\mathbb{E}_{n}[l(g(X),-1)]}{\pi_{p}\mathbb{E}_{p}[l(g(X),+1)]}$$

risk for positive class risk for negative class

•Minimize empirical risk: approximation by data in hand

$$\widehat{R}_{pn}(g) = \frac{\pi_p}{n_p} \sum_{x \in \mathcal{X}_p} l(g(x), +1) + \frac{\pi_n}{n_n} \sum_{x \in \mathcal{X}_n} l(g(x), -1)$$

input data $X \in \mathbb{R}^d$ $\mathbb{E}_p[\cdot] \coloneqq \mathbb{E}_{X \sim p(x|Y=+1)}[\cdot]$ $\pi_p \coloneqq p(Y=+1)$ $\mathcal{X}_p = \{x_i^p\}_{i=1}^{n_p \ i.i.d.} p(x|Y=+1)$ class label $Y \in \{\pm 1\}$ $\mathbb{E}_n[\cdot] \coloneqq \mathbb{E}_{X \sim p(x|Y=-1)}[\cdot]$ $\pi_n \coloneqq p(Y=-1)$ $\mathcal{X}_n = \{x_i^n\}_{i=1}^{n_n \ i.i.d.} p(x|Y=-1)$ loss function $l: \mathbb{R} \times \{\pm 1\} \rightarrow \mathbb{R}$ \mathbb{R} \mathbb{R} \mathbb{R} \mathbb{R}

Classification when negative data is unavailable

- Example: click advertisement
 - Clicked : positive
 - Non-clicked : unlabeled (not interesting or unseen)

Learn a PN binary classifier from positive and unlabeled data (PU learning)



•Goal: minimize the same expected risk as PN learning

 $R(g) = \pi_{p}\mathbb{E}_{p}[l(g(X), +1)] + \frac{\pi_{n}\mathbb{E}_{n}[l(g(X), -1)]}{2}$

negative data unavailable

• Idea: unlabeled data = positive data + negative data

$$\mathbb{E}_{u}[l(g(X), -1)] = \pi_{p}\mathbb{E}_{p}[l(g(X), -1)] + \pi_{n}\mathbb{E}_{n}[l(g(X), -1)]$$

Risk can be expressed by only positive and unlabeled data

 $R_{pu}(g) = \pi_{p} \mathbb{E}_{p}[l(g(X), +1)] + \mathbb{E}_{u}[l(g(X), -1)] - \pi_{p} \mathbb{E}_{p}[l(g(X), -1)]$

 $\mathbb{E}_{u}[\cdot] = \mathbb{E}_{X \sim p(x)}[\cdot], \qquad \mathcal{X}_{u} = \{x_{i}^{u}\}_{i=1}^{n_{u}} \stackrel{i.i.d.}{\sim} p(x) \coloneqq \pi_{p}p(x|Y = +1) + \pi_{n}p(x|Y = -1)$

Theoretical properties of unbiased PU learning

[du Plessis+, NIPS 2014, ICML 2015; Niu+, NIPS 2016]

Risk estimator is unbiased

$$\mathbb{E}[\hat{R}_{pu}(g)] = R_{pu}(g) = R(g)$$

•For linear-in-parameter models, estimation error vanishes in the **optimal** parametric rate



cf. PN learning



• PU learning can be **better** than PN learning if

$$\frac{\pi_{\rm p}}{\sqrt{n_{\rm p}}} + \frac{1}{\sqrt{n_{\rm u}}} < \frac{\pi_{\rm n}}{\sqrt{n_{\rm n}}}$$

$$\pi_{\mathrm{p}} \coloneqq p(Y = +1), \pi_{\mathrm{n}} \coloneqq p(Y = -1)$$

Unbiased PU learning works well in linear-in-parameter models experimentally. [du Plessis+, NIPS 2014, ICML 2015; Niu+, NIPS 2016]



How about flexible models, like *deep nets*?



Unbiased PU learning with flexible model



Overfitting and negative risk of unbiased PU

• If classifier can perfectly separate P and U, training error w.r.t. 0-1 loss is:

$$\frac{\pi_{p}}{n_{p}} \sum_{x} \pi_{p} l(g(0), +1) + \frac{1}{n_{u}} \sum_{x \in \mathcal{X}_{u}} (\log(x), -1) - \frac{\pi_{p}}{n_{p}} \sum_{x} \pi_{p} l(g(1), -1) < 0$$

risk for positive class

risk for negative class



PU learning with non-negative risk estimator (non-negative PU learning)

10

·Idea:

Round-up risk for negative class to zero

•We propose the new risk estimator which is always non-negative

$$\tilde{R}_{pu(g)} = \frac{\pi_p}{n_p} \sum_{x \in \mathcal{X}_p} l(g(x), +1) + \max\left\{0, \frac{1}{n_u} \sum_{x \in \mathcal{X}_u} l(g(x), -1) - \frac{\pi_p}{n_p} \sum_{x \in \mathcal{X}_p} l(g(x), -1)\right\}$$

$$\frac{\pi_p := p(Y = +1)}{positive dataset \mathcal{X}_p = \{x_i^p\}_{i \ge 0}^{n_p} \frac{i.i.d.}{2} p(x|Y = +1)}$$

 $\begin{array}{ll} \text{classifier } g \colon \mathbb{R}^{n} \to \mathbb{R} \\ \text{loss function } l \colon \mathbb{R} \times \{\pm 1\} \to \mathbb{R} \\ \mathbb{E}_{n}[\cdot] \coloneqq \mathbb{E}_{X \sim p(x|Y=-1)}[\cdot] \end{array} \end{array} \qquad \begin{array}{ll} \text{positive dataset} \\ \text{unlabeled dataset} \\ \text{unlabeled dataset} \end{array}$

positive dataset $\mathcal{X}_{p} = \{x_{i}^{p}\}_{i=1}^{n_{p}} \stackrel{i.i.d.}{\sim} p(x|Y = +1)$ unlabeled dataset $\mathcal{X}_{u} = \{x_{i}^{u}\}_{i=1}^{n_{u}} \stackrel{i.i.d.}{\sim} p(x)$

Theoretical analysis of non-negative PU learning

- Risk estimator is consistent and its bias decreases exponentially
 - Bias is negligible in practice

$$\mathcal{O}_p\left(\exp\left(-\frac{1}{\left(\frac{\pi_p^2}{n_p}+\frac{1}{n_u}\right)}\right)\right)$$

Risk estimator may reduce mean squared error

Non-negative risk estimator is more stable

$$\mathbb{E}_{\chi_{\mathrm{p}},\chi_{\mathrm{u}}}\left[\left(\tilde{R}_{\mathrm{pu}}(g)-R(g)\right)^{2}\right] \leq \mathbb{E}_{\chi_{\mathrm{p}},\chi_{\mathrm{u}}}\left[\left(\hat{R}_{\mathrm{pu}}(g)-R(g)\right)^{2}\right]$$

non-negative PU

unbiased PU

•For linear-in-parameter models, estimation error vanishes in the **optimal** parametric rate



Large-scale algorithm

•Want to use mini-batch SGD

 $(\mathcal{X}_{\mathrm{p}}^{i}, \mathcal{X}_{\mathrm{u}}^{i})$: *i*-th mini-batch (*i* = 1, ..., *N*)

•Our objective function

$$\frac{\pi_{p}}{n_{p}}\sum_{i=1}^{N}\sum_{x\in\mathcal{X}_{p}^{i}}l(g(x),+1) + \max\left\{0,\sum_{i}^{N}\left(\frac{1}{n_{u}}\sum_{x\in\mathcal{X}_{u}^{i}}l(g(x),-1)-\frac{\pi_{p}}{n_{p}}\sum_{x\in\mathcal{X}_{p}}l(g(x),-1)\right)\right\}$$
•Sum of risks in mini-batches
$$\sum_{i=1}^{N}\left[\frac{\pi_{p}}{n_{p}}\sum_{x\in\mathcal{X}_{p}^{i}}l(g(x),+1) + \max\left\{0,\frac{1}{n_{u}}\sum_{x\in\mathcal{X}_{u}^{i}}l(g(x),-1)-\frac{\pi_{p}}{n_{p}}\sum_{x\in\mathcal{X}_{p}}l(g(x),-1)\right\}\right]$$

12



Conclusions

- •We proposed a non-negative risk estimator for PU learning which improves on the state-of-the-art unbiased risk estimators.
- •The new risk estimator is more robust against overfitting, and training very flexible model given limited P data becomes possible.
- •A large-scale PU learning algorithm was also developed.
- Extensive theoretical analyses were presented.
- •Intensive experiments were carried out as well.

- •Natarajan, Nagarajan, et al. "Learning with noisy labels." Advances in neural information processing systems. 2013.
- •M. C. du Plessis, G. Niu, and M. Sugiyama. Analysis of learning from positive and unlabeled data. In *NIPS*, 2014.
- •M. C. du Plessis, G. Niu, and M. Sugiyama. Convex formulation for learning from positive and unlabeled data. In *ICML*, 2015.
- •G. Niu, M. C. du Plessis, T. Sakai, Y. Ma, and M. Sugiyama. Theoretical comparisons of positive unlabeled learning against positive-negative learning. In *NIPS*, 2016.
- •J. T. Springenberg, A. Dosovitskiy, T. Brox, and M. Riedmiller. Striving for simplicity: The all convolutional net. 2015.