Revisiting Sample Selection Approach to Positive-Unlabeled Learning: Turning Unlabeled Data into Positive rather than Negative

> Miao Xu et.al. RIKEN **arXiv:1901.10155**

1 Introduction

censoring PU learning



U data is collected first and then some of P data in the U data is labeled case-control PU learning:



P and U data are collected separately

e.g. anomaly detection

e.g. recommendation

1 Introduction

PU Learning Methods

sample selection



Select N data from U data and performs ordinary PN learning afterwards

importance reweighting



Take U data as N data with weights

e.g. uPU, nnPU

2 ERM-based PU Learning Review



3 Proposed Method

Motivation: promote performance without extra annotation cost – Self training

What to select?

How to select?

How to use the selected data?

3.1 What to Select



How to select P from U?

The effect of adding positive data for PU learning on test loss of the CIFAR-10 dataset.

Affect of noise



0: no noise 0.2: 20% random label noise 0.4: 40% random label noise

* how to select P data as *pure* as possible?

- 1. feed forward all U data into the neural network and calculate their sigmoid loss when treating them as negative
- 2. divide the range [0, 1] equally into 100 bins and calculate the number of instances falling into each bin



How about change a PU learning method?



The uPU algorithm's sigmoid loss histogram of U data on CIFAR-10.

Different Loss Function



 $\ell_{\text{sigmoid}}(ty) = 1/(1 + \exp(ty))$ $\ell_{\text{logistic}}(ty) = \ln(1 + \exp(-ty))$



Logistic loss

3.3 How to Use the Selected Data



Experimental results showing the data with largest loss vs. decision boundary in epoch 200. (a) The top 200 largest loss data. (b) The top 300 largest loss data. (c) The top 400 largest loss data.

3.4 aaPU-adaptively augmented PU learning



 $R_{\mathbf{p}}^+ = \mathbb{E}_{X \sim p(x|Y=1)}[\ell(g(X;\theta))],$

$$R_{\mathbf{p}}^{-} = \mathbb{E}_{X \sim p(x|Y=1)}[\ell(-g(X;\theta))],$$

$$R_{n} = \mathbb{E}_{X \sim p(x|Y=-1)}[\ell(-g(X;\theta))],$$

$$R_{\mathrm{u}} = \mathbb{E}_{X \sim p(x)}[\ell(-g(X;\theta))].$$

$$\widehat{R}_{aaPU} = \frac{\pi}{n_{p}} \sum_{x_{i} \in \mathcal{X}_{p} \cup \mathcal{S}} \ell_{\text{logistic}}(g(x_{i};\theta)) + \max\left(\frac{1}{n_{u}} \sum_{x_{i} \in \mathcal{X}_{u}} \ell_{\text{logistic}}(-g(x_{i};\theta)) - \frac{\pi}{n_{p}} \sum_{x_{i} \in \mathcal{X}_{p}} \ell_{\text{logistic}}(-g(x_{i};\theta)), 0\right)$$

Algorithm

\mathbf{Input}		
\mathcal{X}_{p}	positive training data	
\mathcal{X}_{u}	unlabeled training data	
$n_{\mathbf{u}}$	number of unlabeled training data	
T	maximum number of epochs	
η	learning rate of stochastic gradient descent	
${\mathcal S}$	selected data	
Outpu	t	
heta	model parameter θ for $g(x; \theta)$	
1: Initia	lize θ and $S = \emptyset$	
2: for t =	$= 1, 2, \ldots, T \mathbf{do}$	
3: Sh	uffle $\mathcal{X}_{\mathrm{p}}, \mathcal{X}_{\mathrm{u}}, \mathcal{S}$ into N_t mini-batches	
4: De	enoted by $(\mathcal{X}_{\mathbf{p}}^{i}, \mathcal{X}_{\mathbf{n}}^{i}, \mathcal{S}^{i})$ the <i>i</i> -th mini-batch	
5: fo	$\mathbf{r} \ i = 1, 2, \dots, N_t \ \mathbf{do}$	
6:	Update θ by $\theta = \theta - \eta \Delta_{\theta} \widehat{R}_{aaPU}(g; \mathcal{X}_{p}^{i}; \mathcal{X}_{u}^{i}; \mathcal{S}^{i})$	
7: Ca	lculate L_u using Eq. (11)	$L_{u} = [\ell(-q(x_{1};\theta)), \ldots, \ell(-q(x_{n_{u}};\theta))]$
8: Sel	lect the first n_s largest L_{ui} keeping $x_i \notin \mathcal{X}_p$	
9: Un	date \mathcal{S} using corresponding x_i -s	

4 Experiments - Synthetic Data



Synthetic Data



Results on synthetic data comparing aaPU and nnPU

4 Experiments - Real Data



N: non animal

When to add: begin to add P data when the validation loss becomes stable How much to add: select the number of positive data selected from [40; 80; 160; 320]

5 Conclusion

- 1) Propose aaPU, a new sample selection method which can boost performance of PU
- 2) Automatically selects P data from U data during training, and uses these selected data in future epochs
- 3) Data with large loss are selected to estimate the positive risk
- 4) Experiments validate that the proposed aaPU can work better than nnPU