

### Provably Consistent Partial-Label Learning

Lei Feng<sup>1\*</sup> Jiaqi Lv<sup>2</sup> Bo Han<sup>3</sup> Miao Xu<sup>4,5</sup> Gang Niu<sup>5</sup> Xin Geng<sup>2</sup> Bo An<sup>1†</sup> Masashi Sugiyama<sup>5,6</sup> <sup>1</sup>School of Computer Science and Engineering, Nanyang Technological University, Singapore <sup>2</sup>School of Computer Science and Engineering, Southeast University, Nanjing, China <sup>3</sup>Department of Computer Science, Hong Kong Baptist University, China <sup>4</sup>The University of Queensland, Australia <sup>5</sup>Center for Advanced Intelligence Project, RIKEN, Japan <sup>6</sup>The University of Tokyo, Japan

NeurIPS 2020

# Ordinary Supervised Learning



#### Input space:

represented by a single instance (feature vector) characterizing its properties

#### Target space:

associated with a single label characterizing its semantics



## Basic Assumption: Strong Supervision

Successful Learning needs: supervised information + regularities

Strong Supervision Assumption

Sufficient labeling

abundant labeled training data are available

Accurate labeling

object labeling is accurate

Explicit labeling

object labeling is unique and unambiguous

However, supervision is usually weak in practice

approximation error

hypothesis space

complexity

estimation error

training error

# Partial Label Learning (PLL)



 Each object is associated with multiple candidate labels
Only one candidate label is the

unknown ground-truth label



#### PLL Methods

- Identification-based disambiguation: treat the ground-truth label as latent variable identified via iterative refining procedure
- Averaging-based disambiguation: treat all candidate labels in an equal manner and make final prediction by averaging
- Transformation strategy: binary decomposition, dictionary learning, graph matching, regression

However, previous works never consider the generation process of the candidate sets



Learning with ordinary labels

$$R(f) = \mathbb{E}_{p(\boldsymbol{x},y)}[\mathcal{L}(f(\boldsymbol{x}),y)] \qquad \mathcal{Y} = [k]$$

#### Learning with partial labels

$$\widetilde{\mathcal{D}} = \{(\boldsymbol{x}_i, Y_i)\}_{i=1}^n$$
  $Y_i \in \mathcal{C}$   $\mathcal{C} = \{2^{\mathcal{Y}} \setminus \emptyset \setminus \mathcal{Y}\}$   
 $|\mathcal{C}| = 2^k - 2$ 

 $p(y_i \in Y_i \mid \boldsymbol{x}_i, Y_i) = 1, \ \forall (\boldsymbol{x}_i, y_i) \in \mathcal{X} \times \mathcal{Y}, \ \forall Y_i \in \mathcal{C}.$ 

### PLL: Data Generation Model

Partially labeled data distribution

$$\widetilde{p}(\boldsymbol{x}, Y) = \sum_{i=1}^{k} p(Y \mid y = i) p(\boldsymbol{x}, y = i) \qquad p(Y \mid y = i) = \begin{cases} \frac{1}{2^{k-1}-1} & \text{if } i \in Y, \\ 0 & \text{if } i \notin Y. \end{cases}$$



**Lemma 1.** Given any instance x with its correct label y, for any unknown label set Y that is uniformly sampled from C, the equality  $p(y \in Y | x) = 1/2$  holds.

### Preliminaries: Importance Reweighting



#### **Risk-Consistent Method**



### **Risk-Consistent Method**

$$\begin{aligned} \text{Risk:} \quad R_{\text{rc}}(f) &= \frac{1}{2} \mathbb{E}_{\widetilde{p}(\boldsymbol{x},Y)} \left[ \sum_{i=1}^{k} \frac{p(y=i|\boldsymbol{x})}{\sum_{j\in Y} p(y=j|\boldsymbol{x})} \mathcal{L}\big(f(\boldsymbol{x}),i\big) \right] \qquad \{\boldsymbol{x}_{o},Y_{o}\}_{o=1}^{n} \\ \text{Empirical Risk:} \quad \widehat{R}_{\text{rc}}(f) &= \frac{1}{2n} \sum_{o=1}^{n} \left( \sum_{i=1}^{k} \frac{p(y_{o}=i|\boldsymbol{x}_{o})}{\sum_{j\in Y_{o}} p(y_{o}=j|\boldsymbol{x}_{o})} \mathcal{L}\big(f(\boldsymbol{x}_{o}),i\big) \right) \end{aligned}$$

#### Algorithm 1 RC Algorithm

**Input:** Model f, epoch  $T_{\max}$ , iteration  $I_{\max}$ , partially labeled training set  $\widetilde{\mathcal{D}} = \{(\boldsymbol{x}_i, Y_i)\}_{i=1}^n$ . 1: Initialize  $p(y_i = j | \boldsymbol{x}_i) = 1, \forall j \in Y_i$ , otherwise  $p(y_i = j | \boldsymbol{x}_i) = 0$ ; 2: for  $t = 1, 2, ..., T_{\max}$  do Shuffle  $\mathcal{D} = \{(x_i, Y_i)\}_{i=1}^n;$ 3: for  $j = 1, \ldots, I_{\text{max}}$  do 4:**Fetch** mini-batch  $\mathcal{D}_i$  from  $\mathcal{D}$ ; 5:Update model f by  $\widehat{R}_{rc}$  in Eq. (9);  $\searrow M$  step: update parameters of the model 6: **Update**  $p(y_i \mid \boldsymbol{x}_i)$  by Eq. (10); E step: estimate the label distribution 7: end for 8:  $p(y = i \mid \boldsymbol{x}) = g_i(\boldsymbol{x})$  if  $i \in Y$ , otherwise  $p(y = i \mid \boldsymbol{x}) = 0, \ \forall (\boldsymbol{x}, Y) \sim \widetilde{p}(\boldsymbol{x}, Y)$ . 9: end for **Output:** f.

### Classifier-Consistent Method

#### Algorithm 2 CC Algorithm

Input: Model f, epoch  $T_{\max}$ , iteration  $I_{\max}$ , partially labeled training set  $\widetilde{\mathcal{D}} = \{(\boldsymbol{x}_i, Y_i)\}_{i=1}^n$ ; 1: for  $t = 1, 2, \ldots, T_{\max}$  do

2: Shuffle the partially labeled training set  $\widetilde{\mathcal{D}} = \{(\boldsymbol{x}_i, Y_i)\}_{i=1}^n;$ 

3: **for** 
$$j = 1, ..., I_{\max} \, \mathbf{do}_{j}$$

- 4: **Fetch** mini-batch  $\mathcal{D}_j$  from  $\mathcal{D}$ ;
- 5: **Update** model f by minimizing the empirical risk estimator  $\widehat{R}_{cc}$  in Eq. (12);
- 6: end for

7: end for

Output: f.

	1				
	Texture	Yeast	Dermatology	Har	20Newsgroups
$\begin{array}{c} \mathrm{RC} \\ \mathrm{CC} \end{array}$	$99.24{\pm}0.14\%$ $98.02{\pm}2.91\%$ •	$59.89{\pm}1.27\%$ <b>59.97<math>{\pm}1.57\%</math></b>	$\begin{array}{c} 99.41{\pm}1.00\% \\ \textbf{99.73}{\pm}\textbf{0.85\%} \end{array}$	$\begin{array}{c} 98.03{\pm}0.09\%\\ \textbf{98.10}{\pm}\textbf{0.18\%}\end{array}$	$\begin{array}{c} {\bf 75.99 {\pm} 0.53\%} \\ {\rm 75.97 {\pm} 0.54\%} \end{array}$
$\begin{array}{c} \text{SURE} \\ \text{CLPL} \\ \text{PLECOC} \\ \text{PLSVM} \\ \text{PL}k\text{NN} \\ \text{IPAL} \end{array}$	$95.38 \pm 0.28\% \bullet$ $91.93 \pm 0.97\% \bullet$ $69.69 \pm 4.82\% \bullet$ $49.38 \pm 9.99\% \bullet$ $96.78 \pm 0.31\% \bullet$ $99.45 \pm 0.23\%$	$54.39\pm1.32\%$ • $54.58\pm2.11\%$ • $37.37\pm9.73\%$ • $45.70\pm8.01\%$ • $47.79\pm2.41\%$ • $48.99\pm3.84\%$ •	$97.48 \pm 0.32\% \bullet$ $99.62 \pm 0.85\%$ $87.84 \pm 5.30\% \bullet$ $80.00 \pm 7.53\% \bullet$ $80.54 \pm 5.06\% \bullet$ $98.65 \pm 2.27\% \bullet$	$97.43 \pm 0.24\% \bullet$ $97.48 \pm 0.18\% \bullet$ $96.97 \pm 0.29\% \bullet$ $91.64 \pm 1.43\% \bullet$ $94.17 \pm 0.59\% \bullet$ $96.55 \pm 0.40\% \bullet$	$69.82 \pm 0.26\% \bullet$ $71.44 \pm 0.55\% \bullet$ $15.32 \pm 7.86\% \bullet$ $32.59 \pm 8.91\% \bullet$ $27.18 \pm 0.65\% \bullet$ $48.36 \pm 0.85\% \bullet$

Table 3: Test performance (mean $\pm$ std) of each method using linear model on UCI datasets.

Table 4: Test performance (mean±std) of each method using linear model on real-world datasets.

	Lost	MSRCv2	BirdSong	Soccer Player	Yahoo! News
RC	$\textbf{79.43}{\pm\textbf{3.26\%}}$	$46.56 \pm 2.71\%$	$71.94{\pm}1.72\%$	57.00±0.97%	$68.23{\pm}0.83\%$
CC	$79.29 \pm 3.19\%$	$47.22 \pm 3.02\%$	$72.22{\pm}1.71\%$	$56.32{\pm}0.64\%$	$68.14{\pm}0.81\%$
SURE	$71.33{\pm}3.57\%{ullet}$	$46.88{\pm}4.67\%$	$58.92{\pm}1.28\%{\bullet}$	$49.41{\pm}086\%{\bullet}$	$45.49{\pm}1.15\%{\bullet}$
$\operatorname{CLPL}$	$74.87 {\pm} 4.30 \% {\bullet}$	$36.53{\pm}4.59\%{\bullet}$	$63.56{\pm}1.40\%{\bullet}$	$36.82{\pm}1.04\%{\bullet}$	$46.21{\pm}0.90\%{\bullet}$
PLECOC	$49.03 \pm 8.36\% \bullet$	$41.53 \pm 3.25\% \bullet$	$71.58{\pm}1.81\%$	$53.70{\pm}2.02\%{\bullet}$	$66.22{\pm}1.01\%{\bullet}$
PLSVM	$75.31 \pm 3.81\% \bullet$	$35.85{\pm}4.41\%{\bullet}$	$49.90{\pm}2.07\%{\bullet}$	$46.29 \pm 0.96\% \bullet$	$56.85 {\pm} 0.91\% {\bullet}$
PLkNN	$36.73 \pm 2.99\% \bullet$	$41.36 \pm 2.89\% \bullet$	$64.94{\pm}1.42\%{\bullet}$	$49.62 \pm 0.67\% \bullet$	$41.07 \pm 1.02\% \bullet$
$\operatorname{IPAL}$	$72.12{\pm}4.48\%{\bullet}$	$50.80{\pm}4.46\%{\circ}$	$72.06{\pm}1.55\%$	$55.03 \pm 0.77\% \bullet$	$66.79 {\pm} 1.22 \% {\bullet}$

# Thanks