

Disambiguation-Free Partial Label Learning

(非消歧偏标记学习)

Min-Ling Zhang

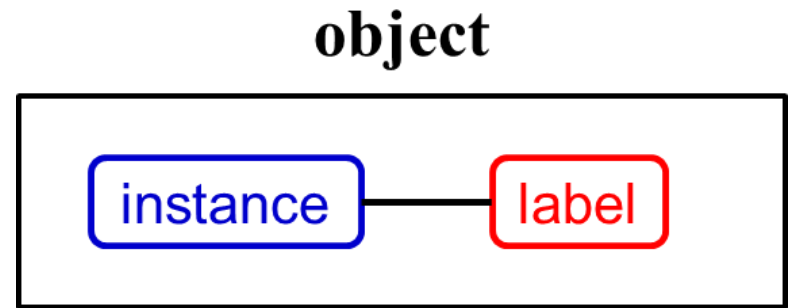
School of Computer Science and Engineering,
MOE Key Lab. of Computer Network & Information Integration
Southeast University, China



April 21, Xiamen

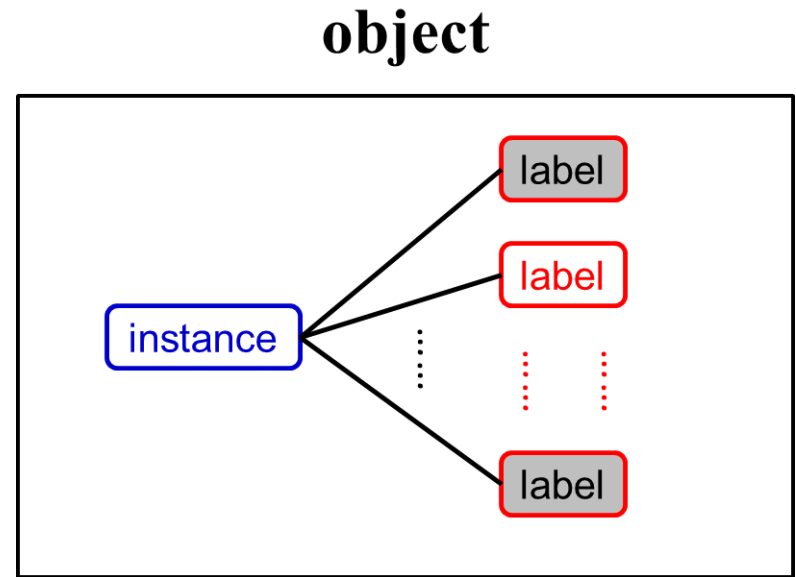
Partial Label (PL) Learning

Traditional Supervised Learning



Partial Label Learning (PLL)

- ❑ Multiple candidate labels
- ❑ Only one valid (but unknown)



The Problem

Difficulty: weak supervision

The ground-truth label of the PL training example is **concealed** in its candidate label set

Common strategy: disambiguation

Try to **disambiguate** the set of candidate labels

→ Prone to be misled by the *false positive* label(s)



Question: Are there other strategies of learning from PL examples without relying on disambiguating the candidate label set?

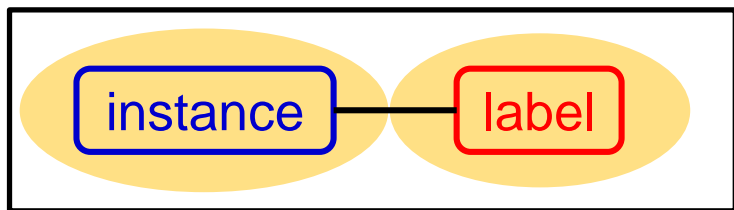
Outline

- Introduction
- The PL-ECOC Approach
- Experiments
 - Controlled UCI Data Sets
 - Real-world Data Sets
- Conclusion



Traditional Supervised Learning

object



Input Space

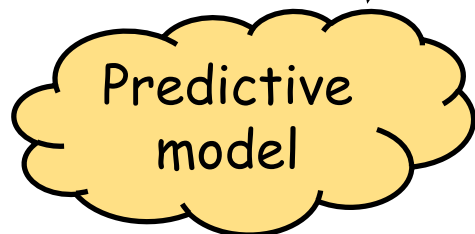
represented by a **single instance** (feature vector) characterizing its properties

Output Space

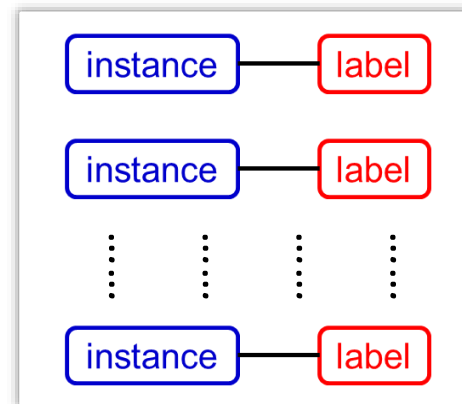
associated with a **single label** characterizing its semantics

instance

label



Supervised
Learning
Algorithm



Basic Assumption: Strong Supervision



Key factor for successful learning

(encoding *semantics* and *regularities* for the learning problem)

Strong supervision assumption

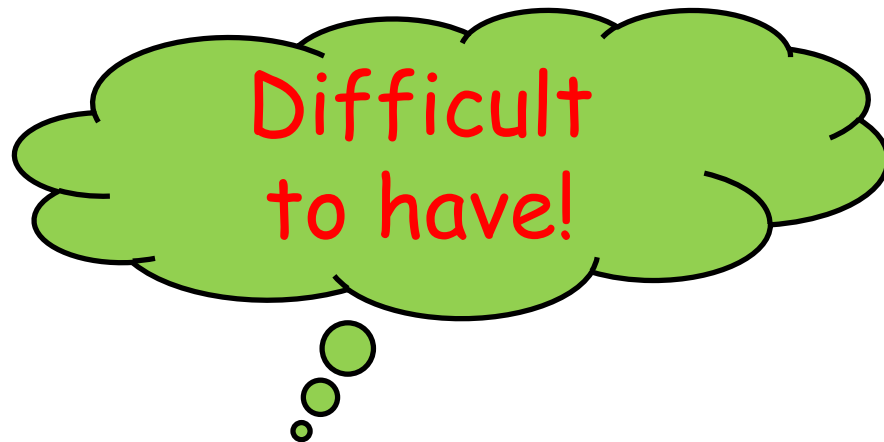
□ Sufficient labeling

abundant labeled training data are available

□ Explicit labeling

object labeling is unique and unambiguous

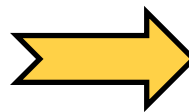
But, Supervision Is Usually Weak



Constrained by:

- ❑ Limited resources
- ❑ Physical environment
- ❑ Problem properties
- ❑

Strong supervision
(sufficient & explicit)

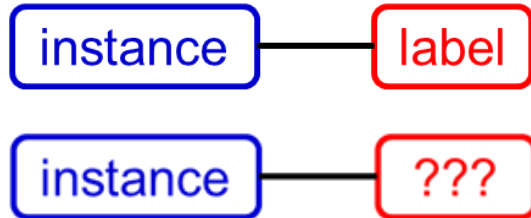


Strong
generalization ability

In practice, we usually have to learn with
weak supervision

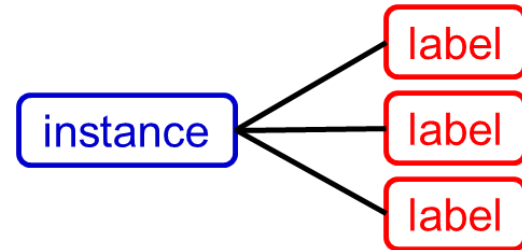
For Example...

semi-supervised learning



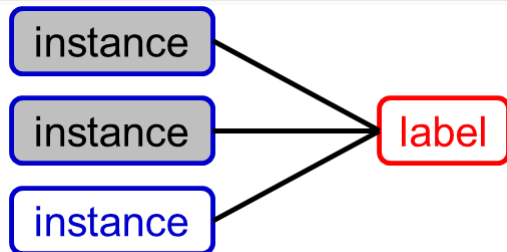
insufficient labeling

multi-label learning



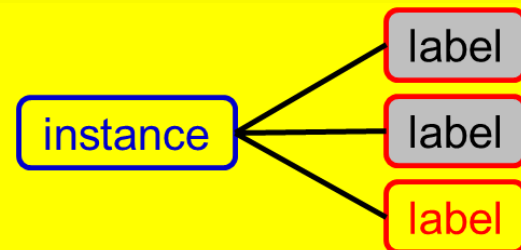
non-unique labeling

multi-instance learning



bag-level labeling

partial-label learning



ambiguous labeling

Partial Label

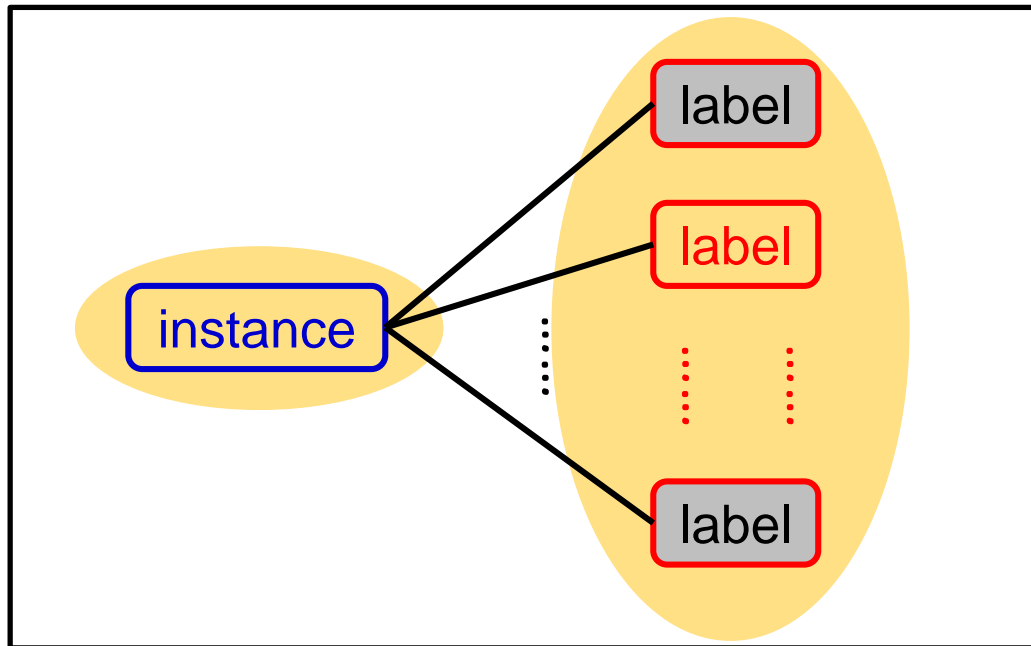


Widely exist in real-world applications

- ❑ Computer vision [Cour et al., JMLR11] [Tang & Zhang, AAAI'17]
- ❑ Image classification [Zeng et al., CVPR'13] [Chen et al., CVPR'13]
- ❑ Learning from crowds [Raykar et al., JMLR10]
- ❑ Ecoinformatics [Liu & Dietterich, NIPS'12] [Zhang & Yu, IJCAI'15]
- ❑

Partial-Label Learning (PLL)

object



□ Each object is associated with **multiple candidate labels**

□ Only one of the candidate label is the **unknown ground-truth label**

Partial-Label Learning (PLL)

Partial Label vs. Unlabel/Multi-Label

Partial-label vs. Unlabel

Commonness: ground-truth label is unknown

Difference: ground-truth label is confined

Partial-label vs. Multi-label

Commonness: multiple labels being assigned

Difference: only one assigned label being valid

Outline

- Introduction
- The PL-ECOC Approach
- Experiments
 - Controlled UCI Data Sets
 - Real-world Data Sets
- Conclusion



Formal Definition of PLL

Settings

\mathcal{X} : d -dimensional feature space \mathbb{R}^d

\mathcal{Y} : label space with q labels $\{y_1, y_2, \dots, y_q\}$

Inputs

\mathcal{D} : training set with m examples $\{(\mathbf{x}_i, Y_i) \mid 1 \leq i \leq m\}$

$\mathbf{x}_i \in \mathcal{X}$ is a d -dimensional feature vector $(\mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{id})^T$

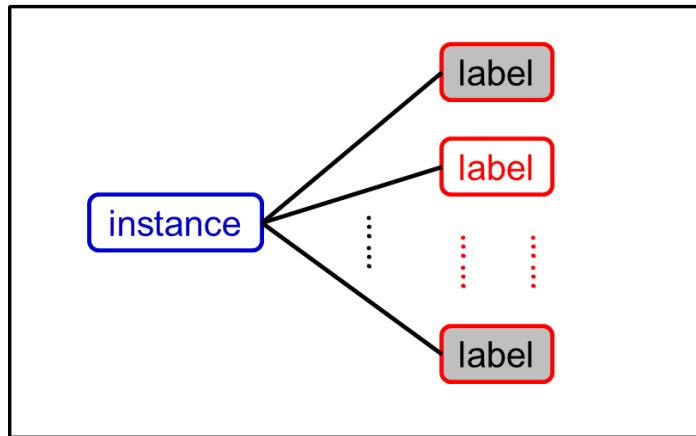
$Y_i \subseteq \mathcal{Y}$ is the candidate label set for \mathbf{x}_i , with its (unknown)
ground-truth label $y_i \in Y_i$

Outputs

h : multi-class predictor $\mathcal{X} \rightarrow \mathcal{Y}$

Key Challenge

object



Ambiguous labeling

ground-truth label not accessible by the learning algorithm

Common strategy: Disambiguation

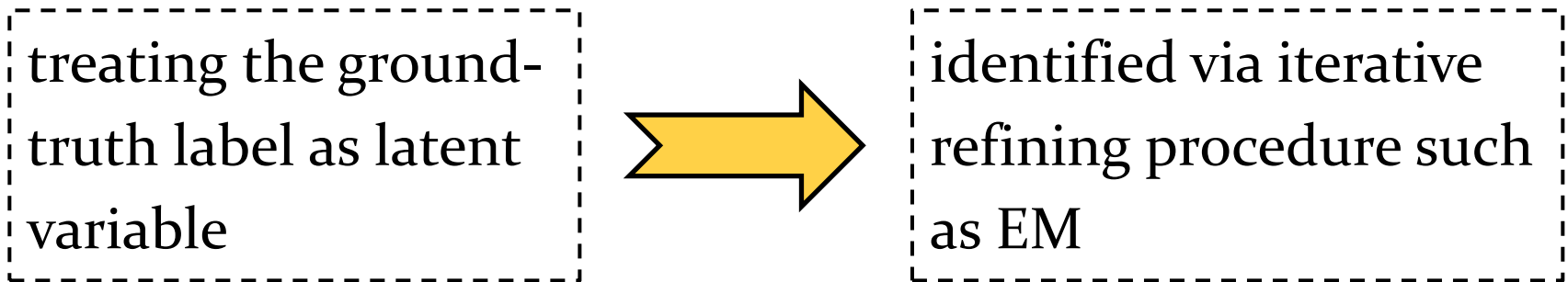
- ❑ Disambiguation by ground-truth label identification
- ❑ Disambiguation by candidate label averaging

Existing Approaches

Disambiguation by Identification

[Jin & Ghahramani, NIPS'03] [Nguyen & Caruana, KDD'08]

[Liu & Dietterich, NIPS'12] [Chen et al., CVPR'13] [Zhang et al., KDD'16]



Potential weakness:

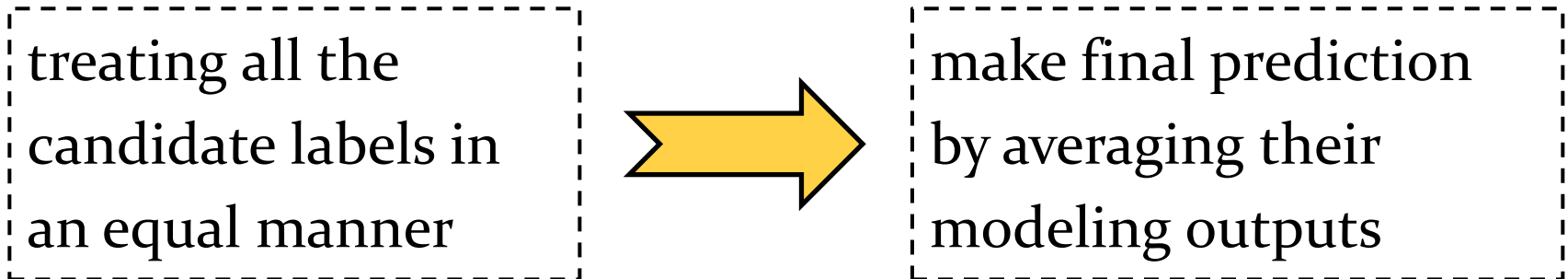
the identified label may turn out to be the false positive label

Existing Approaches

Disambiguation by Averaging

[Hullermeier & Beringer, IDA06] [Cour et al., CVPR'09]

[Cour et al., JMLR11] [Zhang & Yu, IJCAI'15]



Potential weakness:

ground-truth output
overwhelmed by false
positive outputs

The PL-ECOC Approach

Goal of PLL Induce a **multi-class predictor** $h : \mathcal{X} \rightarrow \mathcal{Y}$

Popular Binary
Decomposition

❑ One-vs-Rest (#classifiers: q)

❑ One-vs-One (#classifiers: $q(q-1)/2$)

Not applicable due to the unknown ground-truth label

PL-ECOC (Partial-label Learning with Error-Correcting
Output Codes)

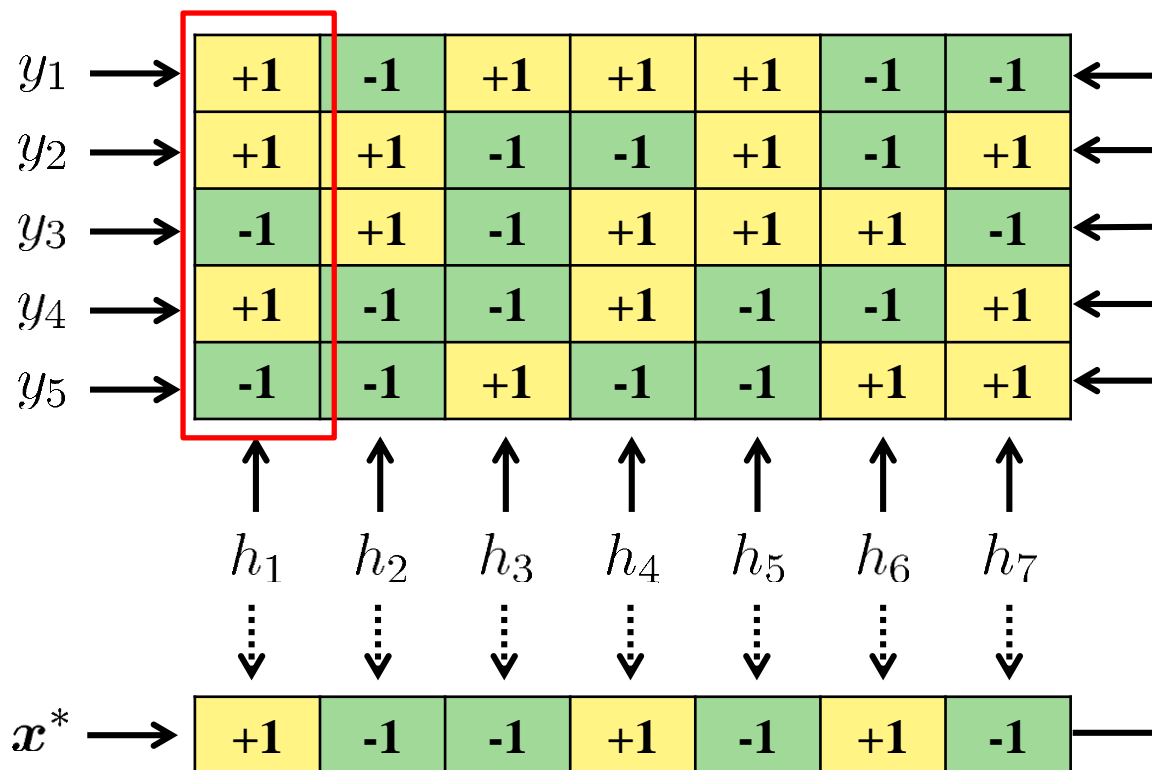
Two major
advantages

❑ Naturally enable binary decomposition

❑ Disambiguation-free

The PL-ECOC Approach (Cont.)

Illustrative procedure of ECOC



For each **multi-class** example (x_i, y_i)

$$\square h_1(x_i) = +1$$

if $y_i \in \{y_1, y_2, y_4\}$

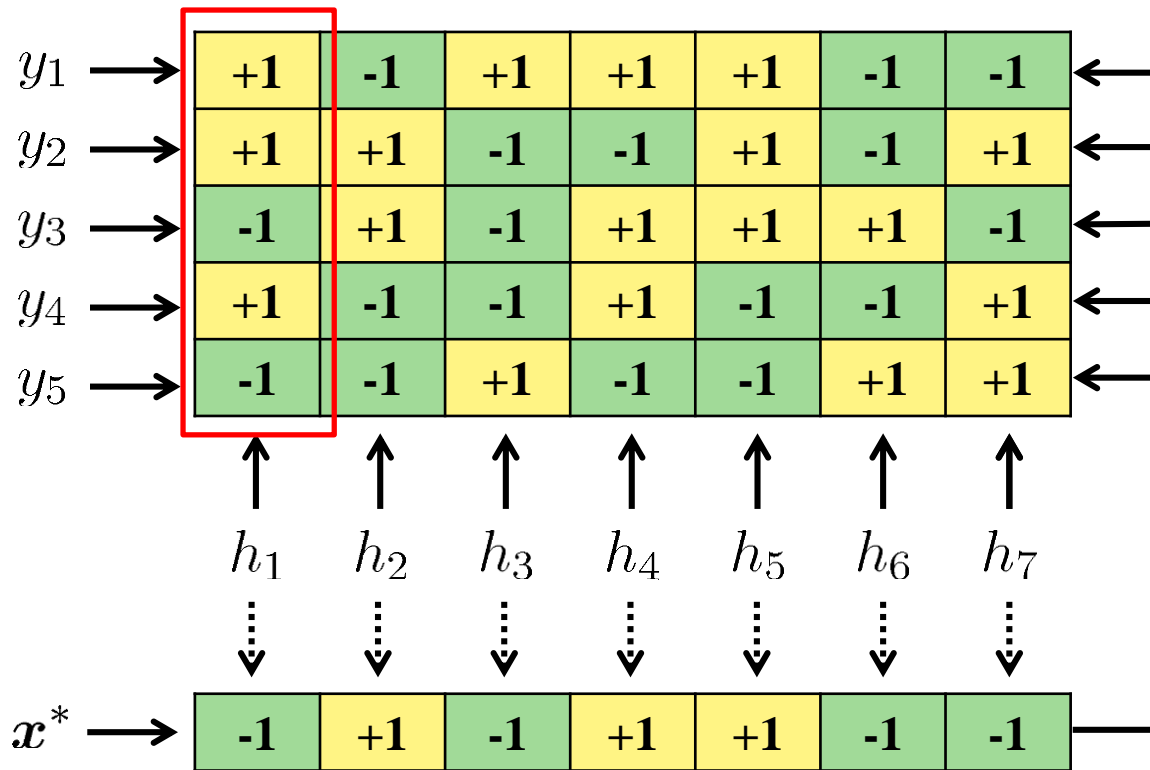
$$\square h_1(x_i) = -1$$

if $y_i \in \{y_3, y_5\}$

Identify the class with closest codeword to test instance x^*

The PL-ECOC Approach (Cont.)

Illustrative procedure of PL-ECOC



For each **partial-label** example (x_i, Y_i)

☐ $h_1(x_i) = +1$
if $Y_i \subseteq \{y_1, y_2, y_4\}$

☐ $h_1(x_i) = -1$
if $Y_i \subseteq \{y_3, y_5\}$

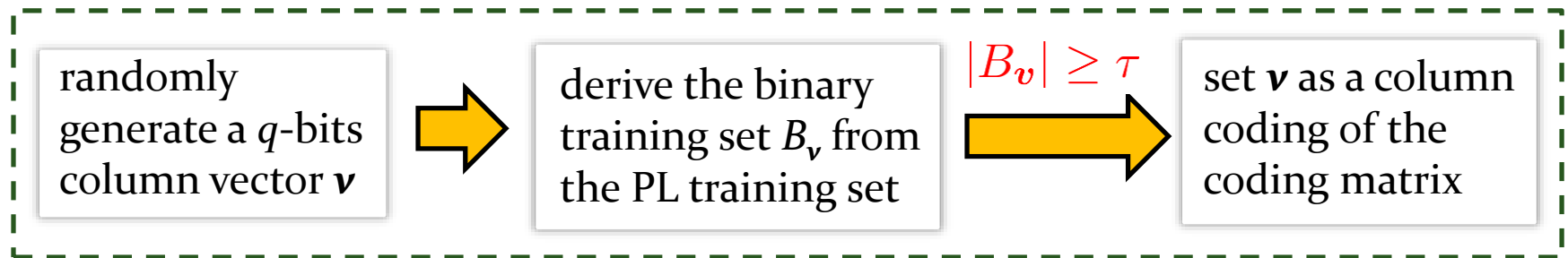
☐ ignored w.r.t. h_1
otherwise

make prediction in the same way as ECOC

The PL-ECOC Approach (Cont.)

Complete Pipeline of PL-ECOC

■ Coding matrix generation



Repeat until reaching the ECOC coding length L

■ Binary classifier induction

induce a total of L binary classifiers, one for each column coding

■ Make prediction for unseen instance

identify the class whose codeword is closest to the classifiers' outputs on unseen instance

Outline

- Introduction
- The PL-ECOC Approach
- Experiments
 - Experimental Setup
 - Controlled UCI Data Sets
 - Real-world Data Sets
- Conclusion



Experimental Setup

Comparing Algorithms

PL-ECOC : Coding length $= \lceil 10 \cdot \log_2(q) \rceil$; Base learner: Libsvm

averaging-based
disambiguation

CLPL: Base learner: SVM with squared hinge loss

PL-kNN: # nearest neighbors = 5

identification-based
disambiguation

PL-SVM: Regularization parameter pool $\{10^{-3}, \dots, 10^3\}$

LSB-CMM: # mixture components = q

Experimental Protocol

Ten-fold cross-validation + Pairwise t -test

Controlled UCI Data Sets

Controlled UCI Data Sets			
Data set	# Examples	# Features	# Class Labels
Ecoli	336	7	8
Dermatology	364	23	6
Vehicle	846	18	4
Segment	2,310	18	7
Abalone	4,177	7	29
Satimage	6,435	36	7
Usps	9,298	256	10
Pendigits	10,992	16	10
Letter	20,000	16	26

Generating an **artificial** PL data set from an UCI data set with three controlling parameters p, r, ϵ

Controlled UCI Data Sets

Controlled UCI Data Sets			
Data set	# Examples	# Features	# Class Labels
Ecoli	336	7	8
Dermatology	364	23	6
Vehicle	846	18	4
Segment	2,310	18	7
Abalone	4,177	7	29
Satimage	6,435	36	7
Usps	9,298	256	10
Pendigits	10,992	16	10
Letter	20,000	16	26

Generating an **artificial** PL data set from an UCI data set with three controlling parameters p, r, ϵ

p : Proportion of examples which are partially labeled ($|S_i| \neq 1$)

r : # false positive labels in candidate label set ($|S_i| = r + 1$)

ϵ : Co-occurring probability for one extra candidate label

Fix r ($=1, 2, 3$), varying $p \in \{0.1, \dots, 0.7\}$

Fix r ($=1$), p ($=1$), varying $\epsilon \in \{0.1, \dots, 0.7\}$

**28 configurations
per UCI data set**

Controlled UCI Data Sets (Cont.)

TABLE 3

Win/tie/loss counts (pairwise t -test at 0.05 significance level) on the classification performance of PL-ECOC against each comparing algorithm on the controlled UCI data sets.

PL-ECOC against		Data Sets (names in abbreviation)										In Total
		Eco.	Der.	Veh.	Seg.	Aba.	Sat.	Usp.	Pen.	Let.	Subtotal	
PL-KNN	[Figure 1]	0/7/0	1/6/0	7/0/0	3/4/0	7/0/0	0/7/0	7/0/0	5/2/0	7/0/0	37/26/0	156/96/0
	[Figure 2]	0/7/0	3/4/0	7/0/0	2/5/0	7/0/0	0/7/0	7/0/0	7/0/0	5/2/0	38/25/0	
	[Figure 3]	0/7/0	2/5/0	7/0/0	4/3/0	7/0/0	1/6/0	7/0/0	7/0/0	5/2/0	40/23/0	
	[Figure 4]	2/5/0	3/4/0	7/0/0	2/5/0	7/0/0	3/4/0	7/0/0	6/1/0	4/3/0	41/22/0	
CLPL	[Figure 1]	0/7/0	0/7/0	6/1/0	7/0/0	7/0/0	7/0/0	7/0/0	7/0/0	7/0/0	48/15/0	181/71/0
	[Figure 2]	0/7/0	0/7/0	3/4/0	7/0/0	7/0/0	7/0/0	7/0/0	7/0/0	7/0/0	45/18/0	
	[Figure 3]	0/7/0	0/7/0	3/4/0	7/0/0	7/0/0	7/0/0	7/0/0	7/0/0	7/0/0	45/18/0	
	[Figure 4]	0/7/0	1/6/0	0/7/0	7/0/0	7/0/0	7/0/0	7/0/0	7/0/0	7/0/0	43/20/0	
PL-SVM	[Figure 1]	0/7/0	0/7/0	7/0/0	7/0/0	7/0/0	7/0/0	7/0/0	7/0/0	7/0/0	49/14/0	195/57/0
	[Figure 2]	0/7/0	0/7/0	7/0/0	7/0/0	7/0/0	7/0/0	7/0/0	7/0/0	7/0/0	49/14/0	
	[Figure 3]	0/7/0	0/7/0	7/0/0	7/0/0	7/0/0	7/0/0	7/0/0	7/0/0	7/0/0	49/14/0	
	[Figure 4]	0/7/0	0/7/0	6/1/0	7/0/0	7/0/0	7/0/0	7/0/0	7/0/0	7/0/0	48/15/0	
LSB-CMM	[Figure 1]	7/0/0	0/7/0	1/6/0	7/0/0	0/7/0	7/0/0	7/0/0	7/0/0	7/0/0	43/20/0	179/73/0
	[Figure 2]	7/0/0	0/7/0	1/6/0	7/0/0	0/7/0	7/0/0	7/0/0	7/0/0	7/0/0	43/20/0	
	[Figure 3]	7/0/0	0/7/0	4/3/0	7/0/0	0/7/0	7/0/0	7/0/0	7/0/0	7/0/0	46/17/0	
	[Figure 4]	7/0/0	2/5/0	1/6/0	7/0/0	2/5/0	7/0/0	7/0/0	7/0/0	7/0/0	47/16/0	

Out of 252 statistical tests (28 configurations x 9 UCI data sets)

- None of the comparing algorithms significantly outperformed PL-ECOC
- PL-ECOC outperforms PL-KNN and CLPL in 61.9% and 71.8% cases respectively
- PL-ECOC outperforms PL-SVM and LSB-CMM in 77.3% and 71.0% cases respectively

Real-World Data Sets

Real-World Data Sets					
Data set	# Examples	# Features	# Class Labels	Avg. # CLs	Domain
Lost	1,122	108	16	2.23	<i>automatic face naming</i> [11]
MSRCv2	1,758	48	23	3.16	<i>object classification</i> [21]
BirdSong	4,998	38	13	2.18	<i>bird song classification</i> [4]
Soccer Player	17,472	279	171	2.09	<i>automatic face naming</i> [27]
LYN 10	18,313	163	11	2.02	<i>automatic face naming</i> [17]
LYN 20	19,027	163	21	2.01	
LYN 50	20,308	163	54	1.97	
LYN 100	21,390	163	101	1.94	
LYN 200	22,991	163	219	1.91	

automatic face naming

instance: face cropped from image/video

candidate labels: names extracted from associated captions/subtitles

object classification

instance: image segmentation

candidate labels: objects appearing within the same image

bird song classification

instance: singing syllable of the bird

candidate labels: bird species jointly singing within 10-seconds period

URL: <http://cse.seu.edu.cn/PersonalPage/zhangml/Resources.htm#partial data>

Real-World Data Sets (Cont.)

TABLE 4

Predictive accuracy (mean \pm std) of each comparing algorithm on the real-world PL data sets. In addition, \bullet/\circ indicates whether the performance of PL-ECOC is statistically superior/inferior to the comparing algorithm on each data set (pairwise t -test at 0.05 significance level).

	PL-ECOC	PL-KNN	CLPL	PL-SVM	LSB-CMM
Lost	0.703 \pm 0.052	0.424 \pm 0.041 \bullet	0.742 \pm 0.038 \circ	0.729 \pm 0.040	0.707 \pm 0.055
MSRCv2	0.505 \pm 0.027	0.448 \pm 0.037 \bullet	0.413 \pm 0.039 \bullet	0.482 \pm 0.043	0.456 \pm 0.031 \bullet
BirdSong	0.740 \pm 0.016	0.614 \pm 0.024 \bullet	0.632 \pm 0.017 \bullet	0.663 \pm 0.032 \bullet	0.717 \pm 0.024 \bullet
Soccer Player	0.537 \pm 0.020	0.497 \pm 0.014 \bullet	0.368 \pm 0.010 \bullet	0.443 \pm 0.014 \bullet	0.525 \pm 0.015
LYN 10	0.694 \pm 0.010	0.460 \pm 0.012 \bullet	0.605 \pm 0.013 \bullet	0.692 \pm 0.009	0.703 \pm 0.010 \circ
LYN 20	0.697 \pm 0.012	0.469 \pm 0.015 \bullet	0.585 \pm 0.010 \bullet	0.686 \pm 0.011 \bullet	0.702 \pm 0.011
LYN 50	0.694 \pm 0.008	0.472 \pm 0.014 \bullet	0.540 \pm 0.012 \bullet	0.666 \pm 0.002 \bullet	0.679 \pm 0.007 \bullet
LYN 100	0.680 \pm 0.012	0.459 \pm 0.010 \bullet	0.507 \pm 0.011 \bullet	0.655 \pm 0.010 \bullet	0.673 \pm 0.010
LYN 200	0.662 \pm 0.010	0.457 \pm 0.014 \bullet	0.462 \pm 0.009 \bullet	0.636 \pm 0.010 \bullet	0.648 \pm 0.007 \bullet

- On *BirdSong*, *LYN 50* and *LYN 200*, PL-ECOC is **superior** to all the comparing algorithms
- On *Soccer Player*, *LYN 20*, *LYN 100* and *MSRCv2*, PL-ECOC is **superior or at least comparable** to all the comparing algorithms
- On *Lost* and *LYN 10*, PL-ECOC is **inferior** to the comparing algorithms in only two cases (CLPL on *Lost*; LSB-CMM on *LYN 10*)

Sensitivity Analysis for Coding Length

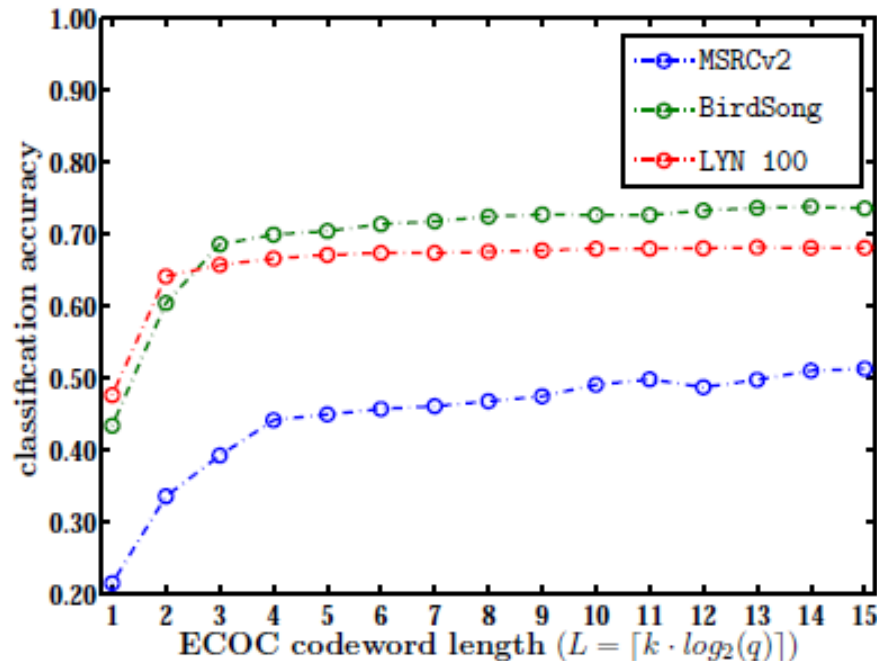


Fig. 5. Classification accuracy of PL-ECOC changes as the codeword length L increases from $\lceil \log_2(q) \rceil$ to $\lceil 15 \cdot \log_2(q) \rceil$ with step-size $\lceil \log_2(q) \rceil$.

- Accuracy improves as the coding length increases
- Becomes stable as coding length approaches $\lceil 10 \cdot \log_2(q) \rceil$

Outline

- Introduction
- The PL-ECOC Approach
- Experiments
 - Controlled UCI Data Sets
 - Real-world Data Sets
- Conclusion



Conclusion

Main Contribution

Propose a new strategy to learn from partial label data, which is free of disambiguation

Key Technique

Treat the candidate label set as an entirety, and then adapt the ECOC procedure

Future Work

Investigate variants of PL-ECOC, new strategies to learn from partial label data, etc.