



Dual Set Multi-Label Learning

Chong Liu ¹, Peng Zhao ¹, Sheng-Jun Huang ²,
Yuan Jiang ¹, and Zhi-Hua Zhou ¹

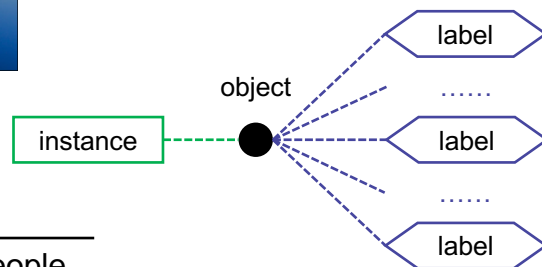
¹ LAMDA Group, Nanjing University, China

² Nanjing University of Aeronautics and Astronautics, China

- Introduction
- Potential Solutions and Deficiencies
- Our Approach
- Theoretical Results
- Experiments
- Conclusion

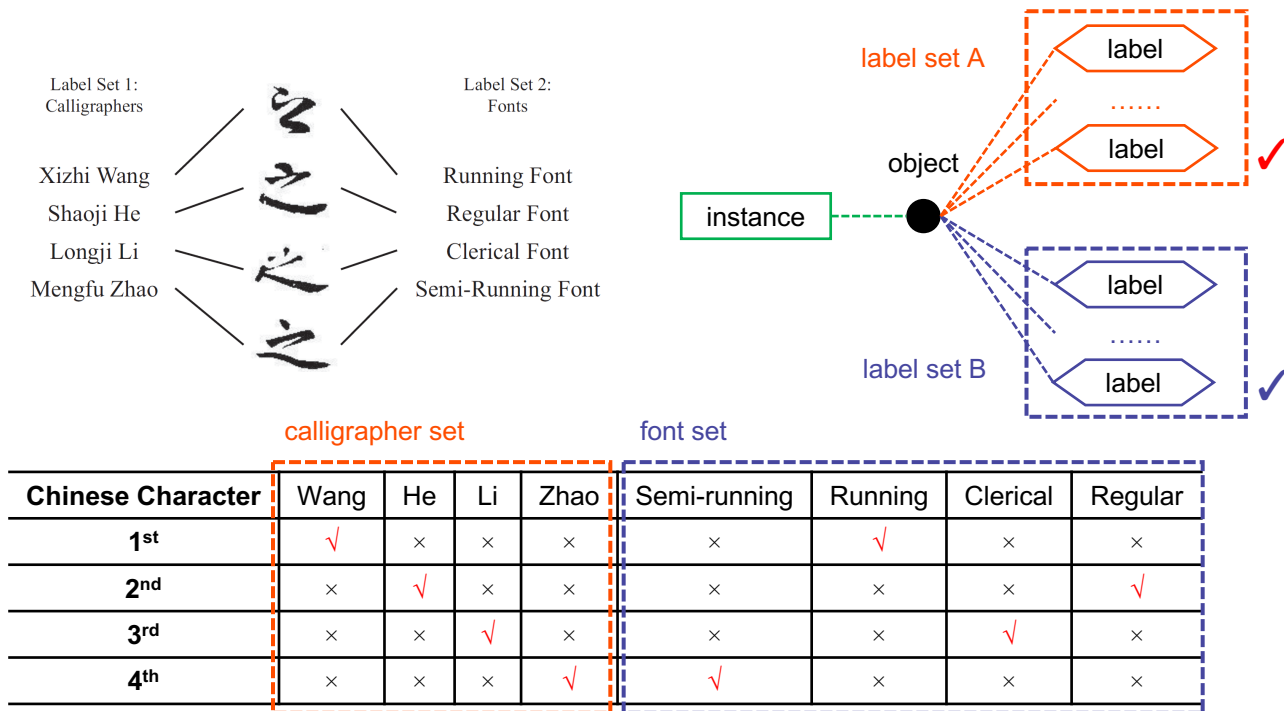
Introduction

- An example of traditional multi-label learning



Sky	Lake	Road	Mountain	Bird	People
✓	✓	×	✓	×	×

- An example different from traditional multi-label learning



Introduction

- Similar cases are popular among our lives, such as

movie classification

Production Company Set	Genre Set
20 th Century Fox	Action
Warner Bros. Pictures	Adventure
Columbia Pictures	Comedy
Paramount Pictures	Horror
Universal Pictures	Science Fiction
Walt Disney Pictures	War



car classification

Production Company Set	Type Set
Audi	Economy
BMW	Family
Mercedes-Benz	Sedan
Opel	Luxury vehicle
Porsche	Sports
Volkswagen	Commercial



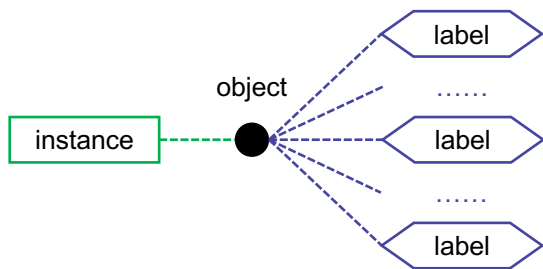
Problem Formulation

- Definition

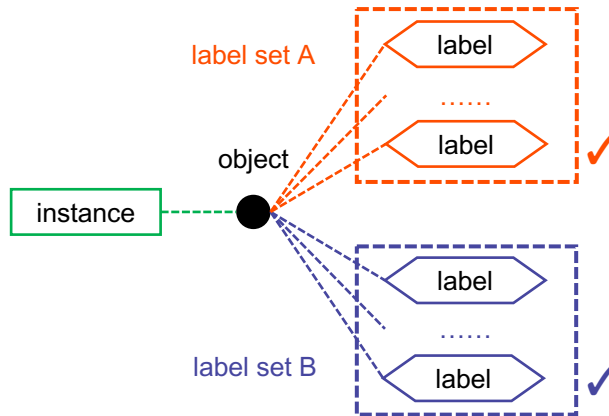
Definition 1. (*Dual Set Multi-Label Learning*) Given the training set \mathcal{D} , the task is to learn a mapping function from the input space to the output space,

$$h : \mathcal{X} \rightarrow \mathcal{Y}^a \times \mathcal{Y}^b.$$

For an unseen instance $\mathbf{x} \in \mathcal{X}$, the mapping function $h(\cdot)$ predicts $h(\mathbf{x}) \subseteq \mathcal{Y}^a \times \mathcal{Y}^b$ as the dual labels for \mathbf{x} .



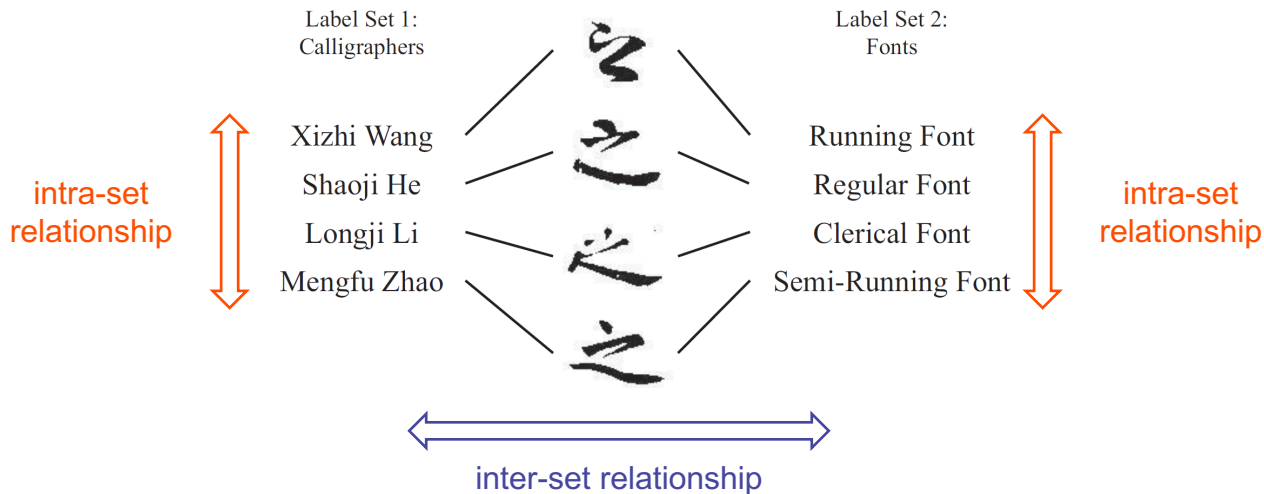
traditional multi-label learning



dual set multi-label learning

Problem Formulation

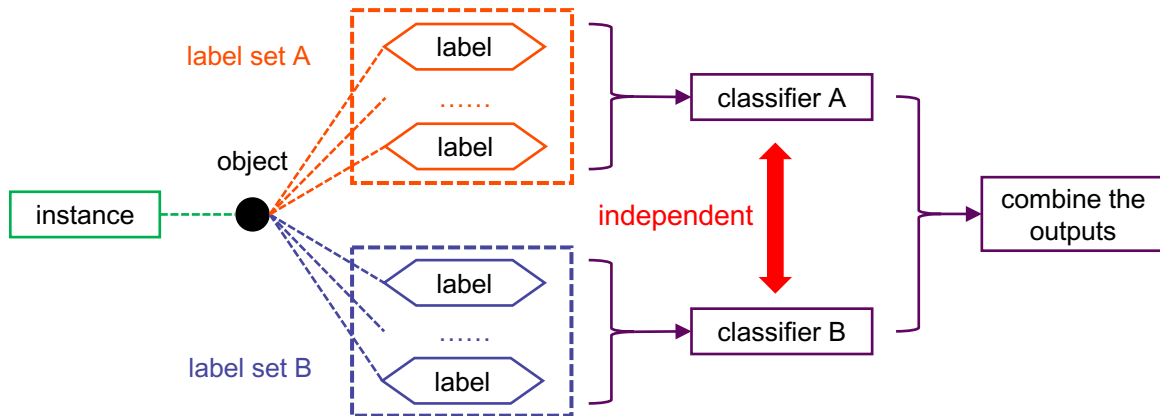
- Key challenge: exploiting label relationships
 - **Intra-set**: the exclusive relationship within the same set
 - **Inter-set**: the pairwise label set relationship



- Introduction
- Potential Solutions and Deficiencies
- Our Approach
- Theoretical Results
- Experiments
- Conclusion

- Independent Decomposition
 - Decomposing the original problem into two classification problems
- Co-occurrence Based Decomposition
 - Decomposing the original problem into a new multi-class problem
- Label Stacking
 - Transforming the original problem into sequential problems

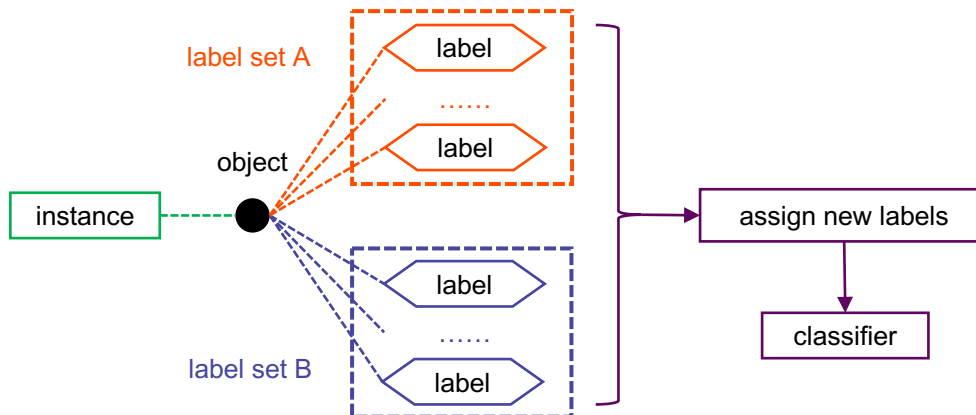
- Independent Decomposition
 - Decomposing the original problem into two classification problems



Deficiency:

Inter-set relationship is neglected.

- Co-occurrence Based Decomposition
 - Decomposing the original problem into a new multi-class problem



- How do we assign new labels by label co-occurrence?

Potential Solutions

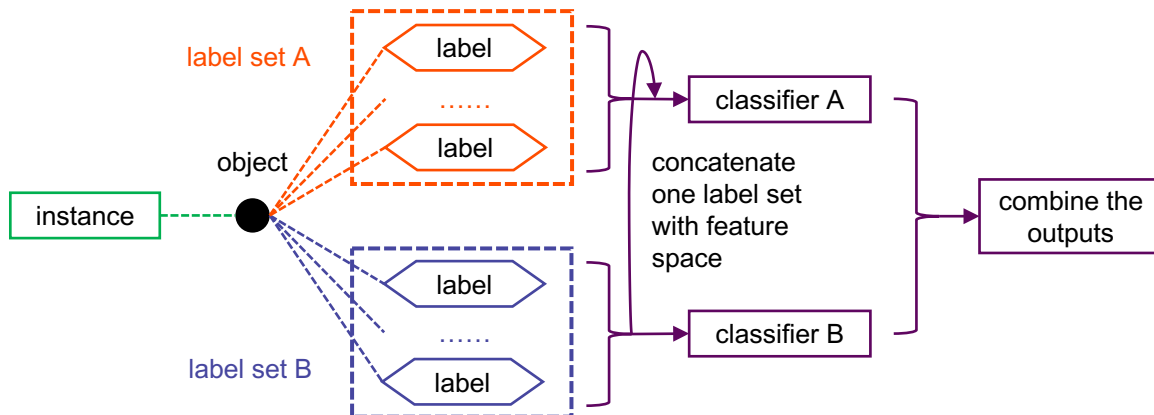
- Co-occurrence Based Decomposition
 - An example showing how to assign new labels

Instance	label set A				label set B				New multi-class label
	A-1	A-2	A-3	A-4	B-1	B-2	B-3	B-4	
1 st	✓	×	×	×	×	✓	×	×	1
2 nd	×	✓	×	×	×	×	×	✓	2
3 rd	×	×	✓	×	×	×	✓	×	3
4 th	×	×	×	✓	✓	×	×	×	4
5 th	×	×	✓	×	×	×	✓	×	3
6 th	×	✓	×	×	×	×	×	✓	2
7 th	×	×	✓	×	×	×	✓	×	3
8 th	✓	×	×	×	×	✓	×	×	1

Deficiency:

It is unable to handle new label co-occurrence cases.

- Label Stacking
 - Transforming the original problem into two sequential problems

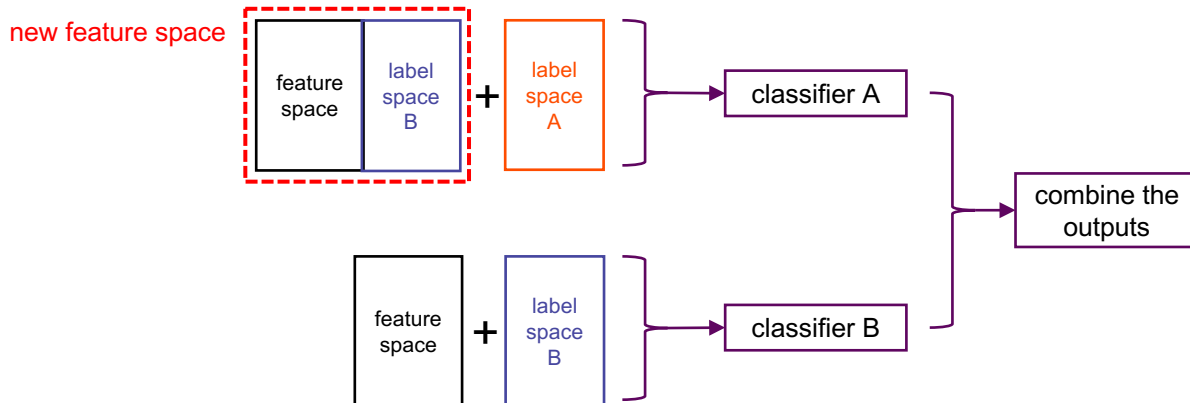


- How do we train classifier A and B?

Potential Solutions

- Label Stacking

- An example showing how to train classifier A and B



Deficiency:

Only one label set helps the other one.

- Introduction
- Potential Solutions and Deficiencies
- Our Approach
- Theoretical Results
- Experiments
- Conclusion

- Key Problem
 - How to find a better way to exploit intra-set and inter-set label relationship simultaneously?
- Key ideas
 - Multi-class classifiers are used to exploit intra-set label relationship.
 - Model-reuse mechanism and distribution adjusting mechanism are used to make label sets help each other, all of which exploit inter-set label relationship.
- Boosting framework is used to carry out these ideas.

- The DSML algorithm
 - How does it work?

Algorithm 1 The DSML algorithm

Input: Training set $\mathcal{D} = \{(\mathbf{x}_i, y_i^a, y_i^b) | 1 \leq i \leq m\}$, base learning algorithm \mathcal{A} , number of rounds T , weight tuning parameter B

Training process:

```

1:  $w_{1,i}^a = w_{1,i}^b = 1/m$ ;
2: for  $t = 1$  to  $T$  do
3:    $(X_t^a, y_t^a) \leftarrow \text{Sample}(\mathcal{D}, w_t^a)$ 
4:    $(X_t^b, y_t^b) \leftarrow \text{Sample}(\mathcal{D}, w_t^b)$ 
5:   Training three models  $h_t^{raw}$ ,  $h_t^a$  and  $h_t^b$  with model-reuse mechanism by Eq. (1), (2) and (3)
6:   Calculating error rate  $\epsilon_t^a$  and  $\epsilon_t^b$  by Eq. (4) and (5)
7:   if  $\epsilon_t^a > (L_1 - 1)/L_1$  or  $\epsilon_t^b > (L_2 - 1)/L_2$  then
8:     Break
9:   end if
10:  Updating model weight  $\alpha_t^a$  and  $\alpha_t^b$  by Eq. (6) and (7)
11:  Updating sample distribution  $w_{t+1}^a$  and  $w_{t+1}^b$  by  $\alpha_t^a$ ,  $\alpha_t^b$  and  $B$  with distribution adjusting mechanism according to Eq. (8) and (9)
12:  Performing normalization to  $w_{t+1}^a$  and  $w_{t+1}^b$ 
13: end for
    
```

Output: Predict labels for dual set: $f^a(\mathbf{x})$ and $f^b(\mathbf{x})$ by Eq. (10) and (11)

Step	Task
1	Initialize sample weights
3-4	Resample
5	Train base learners with model-reuse mechanism
6	Calculate error rates
7-9	Check error rates
10-12	Update sample weights with distribution adjusting mechanism

Our Approach

- The DSML algorithm
 - Training base learners with **model-reuse mechanism**

Algorithm 1 The DSML algorithm

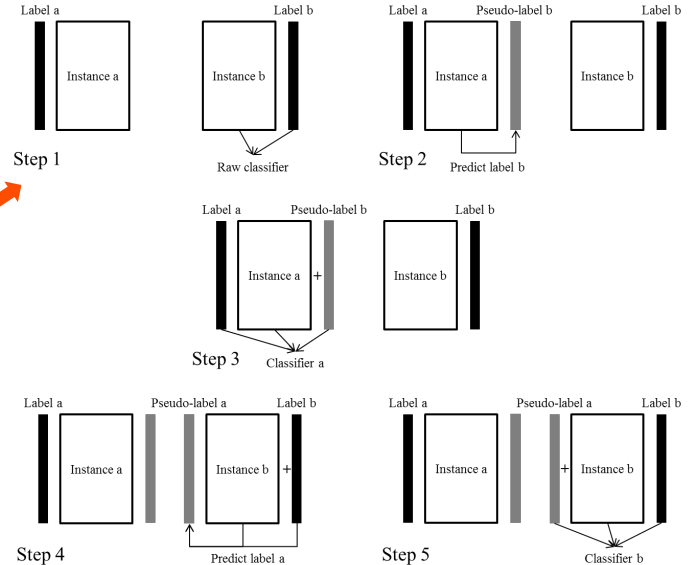
Input: Training set $\mathcal{D} = \{(x_i, y_i^a, y_i^b) | 1 \leq i \leq m\}$, base learning algorithm \mathcal{A} , number of rounds T , weight tuning parameter B

Training process:

```

1:  $w_{1,i}^a = w_{1,i}^b = 1/m$ ;
2: for  $t = 1$  to  $T$  do
3:    $(X_t^a, y_t^a) \leftarrow \text{Sample}(\mathcal{D}, w_t^a)$ 
4:    $(X_t^b, y_t^b) \leftarrow \text{Sample}(\mathcal{D}, w_t^b)$ 
5:   Training three models  $h_t^{aw}, h_t^a$  and  $h_t^b$  with model-reuse mechanism by Eq. (1), (2) and (3)
6:   Calculating error rate  $\epsilon_t^a$  and  $\epsilon_t^b$  by Eq. (4) and (5)
7:   if  $\epsilon_t^a > (L_1 - 1)/L_1$  or  $\epsilon_t^b > (L_2 - 1)/L_2$  then
8:     Break
9:   end if
10:  Updating model weight  $\alpha_t^a$  and  $\alpha_t^b$  by Eq. (6) and (7)
11:  Updating sample distribution  $w_{t+1}^a$  and  $w_{t+1}^b$  by  $\alpha_t^a$ ,  $\alpha_t^b$  and  $B$  with distribution adjusting mechanism according to Eq. (8) and (9)
12:  Performing normalization to  $w_{t+1}^a$  and  $w_{t+1}^b$ 
13: end for
    
```

Output: Predict labels for dual set: $f^a(x)$ and $f^b(x)$ by Eq. (10) and (11)



- The DSML algorithm
 - Calculating error rate and updating model weight

Algorithm 1 The DSML algorithm

Input: Training set $\mathcal{D} = \{(\mathbf{x}_i, y_i^a, y_i^b) | 1 \leq i \leq m\}$, base learning algorithm \mathcal{A} , number of rounds T , weight tuning parameter B

Training process:

- 1: $w_{1,i}^a = w_{1,i}^b = 1/m$;
- 2: **for** $t = 1$ to T **do**
- 3: $(X_s^a, y_s^a) \leftarrow \text{Sample}(\mathcal{D}, w_t^a)$
- 4: $(X_s^b, y_s^b) \leftarrow \text{Sample}(\mathcal{D}, w_t^b)$
- 5: Training three models h_t^{raw} , h_t^a and h_t^b with model-reuse mechanism by Eq. (1), (2) and (3)
- 6: Calculating error rate ϵ_t^a and ϵ_t^b by Eq. (4) and (5)
- 7: **if** $\epsilon_t^a > (L_1 - 1)/L_1$ or $\epsilon_t^b > (L_2 - 1)/L_2$ **then**
- 8: Break
- 9: **end if**
- 10: Updating model weight α_t^a and α_t^b by Eq. (6) and (7)
- 11: Updating sample distribution w_{t+1}^a and w_{t+1}^b by α_t^a , α_t^b and B with distribution adjusting mechanism according to Eq. (8) and (9)
- 12: Performing normalization to w_{t+1}^a and w_{t+1}^b
- 13: **end for**

Output: Predict labels for dual set: $f^a(x)$ and $f^b(x)$ by Eq. (10) and (11)

$$\epsilon_t^a = \sum_{i=1}^m \mathbb{I}[h_t^a([X_s^a, \hat{Y}^b]_i) \neq (y_s^a)_i]$$

$$\epsilon_t^b = \sum_{i=1}^m \mathbb{I}[h_t^b([X_s^b, \hat{Y}^a]_i) \neq (y_s^b)_i]$$

$$\alpha_t^a = \frac{1}{L_1} \left[\log \frac{1 - \epsilon_t^a}{\epsilon_t^a} + \log(L_1 - 1) \right]$$

$$\alpha_t^b = \frac{1}{L_2} \left[\log \frac{1 - \epsilon_t^b}{\epsilon_t^b} + \log(L_2 - 1) \right]$$

Our Approach

- The DSML algorithm
 - Updating sample weight with **distribution adjusting mechanism**

Algorithm 1 The DSML algorithm

Input: Training set $\mathcal{D} = \{(\mathbf{x}_i, y_i^a, y_i^b) | 1 \leq i \leq m\}$, base learning algorithm \mathcal{A} , number of rounds T , weight tuning parameter B

Training process:

- 1: $w_{1,i}^a = w_{1,i}^b = 1/m$;
- 2: **for** $t = 1$ to T **do**
- 3: $(X_t^a, y_t^a) \leftarrow \text{Sample}(\mathcal{D}, w_t^a)$
- 4: $(X_t^b, y_t^b) \leftarrow \text{Sample}(\mathcal{D}, w_t^b)$
- 5: Training three models h_t^{raw} , h_t^a and h_t^b with model-reuse mechanism by Eq. (1), (2) and (3)
- 6: Calculating error rate ϵ_t^a and ϵ_t^b by Eq. (4) and (5)
- 7: **if** $\epsilon_t^a > (L_1 - 1)/L_1$ or $\epsilon_t^b > (L_2 - 1)/L_2$ **then**
- 8: Break
- 9: **end if**
- 10: Updating model weight α_t^a and α_t^b by Eq. (6) and (7)
- 11: Updating sample distribution w_{t+1}^a and w_{t+1}^b by α_t^a , α_t^b and B with distribution adjusting mechanism according to Eq. (8) and (9)
- 12: Performing normalization to w_{t+1}^a and w_{t+1}^b
- 13: **end for**

Output: Predict labels for dual set: $f^a(\mathbf{x})$ and $f^b(\mathbf{x})$ by Eq. (10) and (11)

B is the distribution adjusting parameter

$$w_{t+1,i}^a = w_{t,i}^a \exp(\alpha_t^a \cdot \mathbb{I}[y_i^a \neq \hat{y}_i^a]) B^{\mathbb{I}[y_i^b \neq \hat{y}_i^b]}$$

$$w_{t+1,i}^b = w_{t,i}^b \exp(\alpha_t^b \cdot \mathbb{I}[y_i^b \neq \hat{y}_i^b]) B^{\mathbb{I}[y_i^a \neq \hat{y}_i^a]}$$

- Introduction
- Potential Solutions and Deficiencies
- Our Approach
- Theoretical Results
- Experiments
- Conclusion

- Superiority of learning by splitting the label set

Theorem 1. *For dual-set multi-label learning problems, h^a and h^b are classifiers trained on the instance space \mathcal{X} and label space $\mathcal{Y}^a, \mathcal{Y}^b$ respectively. h is a classifier trained directly from $\mathcal{X} \times [\mathcal{Y}^a \times \mathcal{Y}^b]$, namely,*

$$h : \mathbf{x} \rightarrow \arg \max_{y^a, y^b \in [\mathcal{Y}^a \times \mathcal{Y}^b]} h(\mathbf{x}, y),$$

where $y = [y^a, y^b]$, then margin of learning from dual label set is larger than that of directly learning from all labels:

$$\min\{\bar{\rho}_{h^a}(\mathbf{x}, y^a), \bar{\rho}_{h^b}(\mathbf{x}, y^b)\} \geq \bar{\rho}_h(\mathbf{x}, y).$$

margin of multi-margin of learning
class learningdirectly from all labels

Remark:

It shows the effectiveness of learning by splitting the label set into two disjoint label sets, which implies that we should explicitly considering the dual label sets.

- Generalization bound of learning by splitting the label set

Theorem 2. Let $H = \{(x, y^a, y^b) \in \mathcal{X} \times [\mathcal{Y}^a \times \mathcal{Y}^b] \rightarrow \mathbf{w}^T \phi(x) | \sum_{\ell=1}^{L_1+L_2} \|\mathbf{w}\|_{\mathbb{H}}^2 \leq \Lambda^2\}$ be a hypothesis set with $y^a = 1, \dots, L_1, y^b = 1, \dots, L_2$, where $\phi : \mathcal{X} \rightarrow \mathbb{H}$ is a feature mapping induced by some positive definite kernel κ . Assume that $S \subset \{x : \kappa(x, x) \leq r^2\}$, and fix $\rho > 0$, then for any $\delta > 0$, with probability at least $1 - \delta$, the following generalization bound holds for all $h^{spl} = [h^a, h^b] \in H$:

$$\underbrace{R(h^{spl})}_{\text{generalization error}} \leq \underbrace{\hat{R}_\rho(h^{spl})}_{\text{empirical error}} + \underbrace{\frac{2r\Lambda}{\rho} \sqrt{\frac{\overbrace{\max\{L_1, L_2\}}^{O(\sqrt{L})}}{m}}}_{O(1/\sqrt{m})} + 3\sqrt{\frac{\log(2/\delta)}{m}}.$$

Remark:

The convergence rate of the generalization error is standard as $O(1/\sqrt{m})$. And the error bound exhibits a radical dependence on the maximal number of labels in dual label sets.

- Introduction
- Potential Solutions and Deficiencies
- Our Approach
- Theoretical Results
- **Experiments**
- Conclusion

- Datasets

- We collected or adapted three real-world dataset. Now take Calligrapher-Font dataset for example
 - We collected 23195 calligraphic images
 - We transformed each of them into 512-dimensional feature vector
 - There are 14 calligraphers and 5 kinds of fonts

- Statistics of three datasets

Dataset	No. of instances	No. of dimensions	Size of label set A	Size of label set B
Calligrapher-Font	23195	512	14	5
Brand-Type	2247	4096	7	3
Frequency-Gender	3157	19	5	2

- Evaluation Measures
 - Accuracy of the label set A
 - Accuracy of the label set B
 - Overall accuracy

Definition 4. Let $\mathcal{Z} = \{z_i, y_i^a, y_i^b | 1 \leq i \leq n\}$ denote the testing set where n is the total number of testing instances and let h^a, h^b be the underlying classifiers learned from the training process associated with two label sets respectively. Three accuracies are defined to evaluate the performance,

$$Accuracy_a = \frac{1}{n} \sum_{i=1}^n \mathbb{I}[h^a(z_i) = y_i^a],$$

$$Accuracy_b = \frac{1}{n} \sum_{i=1}^n \mathbb{I}[h^b(z_i) = y_i^b],$$

$$Accuracy_{all} = \frac{1}{n} \sum_{i=1}^n \mathbb{I}[h^a(z_i) = y_i^a] \cdot \mathbb{I}[h^b(z_i) = y_i^b].$$

Experiments

- Comparing DSML with other algorithms
 - Multi-class RBF neural networks are used as base learner for **DSML** and **potential solutions**.
 - The outputs of **classical multi-label learning approaches** are modified to fit dual set multi-label learning.
 - 5-fold cross-validation performance of these algorithms (mean \pm std.)

Dataset	Measure	Algorithms							
		DSML	Ind. Dec.	Co-Occ. Dec.	Label Stacking	ML-KNN	ML-RBF	BP-MLL	RankSVM
<i>Cal.-Font</i>	<i>Accy.-a</i>	.6562 \pm .0059	.5967 \pm .0082	N/A	.6019 \pm .0088	.6337 \pm .0075	.6372 \pm .0045	.1493 \pm .0051	N/A
	<i>Accy.-b</i>	.7223 \pm .0079	.6751 \pm .0040	N/A	.6801 \pm .0078	.7101 \pm .0030	.7100 \pm .0087	.4104 \pm .0670	N/A
	<i>Accy.-all</i>	.5672 \pm .0087	.4836 \pm .0099	.5609 \pm .0050	.4889 \pm .0094	.5570 \pm .0048	.5396 \pm .0066	.0764 \pm .0077	N/A
<i>Brand-Type</i>	<i>Accy.-a</i>	.5723 \pm .0226	.5661 \pm .0129	N/A	.5968 \pm .0254	.4722 \pm .0160	.5207 \pm .0223	.1206 \pm .0182	.5238 \pm .0352
	<i>Accy.-b</i>	.7730 \pm .0249	.7677 \pm .0092	N/A	.7637 \pm .0225	.7245 \pm .0115	.7405 \pm .0126	.3000 \pm .0509	.7517 \pm .0137
	<i>Accy.-all</i>	.4949 \pm .0227	.4744 \pm .0105	.4784 \pm .0294	.4735 \pm .0302	.3912 \pm .0078	.4201 \pm .0160	.0538 \pm .0053	.4183 \pm .0345
<i>Freq.-Gndr</i>	<i>Accy.-a</i>	.8521 \pm .0091	.8321 \pm .0212	N/A	.8375 \pm .0170	.5879 \pm .0091	.7570 \pm .0144	.4004 \pm .1464	.0326 \pm .0135
	<i>Accy.-b</i>	.9547 \pm .0061	.9579 \pm .0067	N/A	.9550 \pm .0051	.6953 \pm .0196	.9661 \pm .0047	.5014 \pm .0271	.5382 \pm .0643
	<i>Accy.-all</i>	.8220 \pm .0082	.8039 \pm .0214	.8068 \pm .0187	.8096 \pm .0183	.4587 \pm .0161	.7387 \pm .0134	.1704 \pm .0847	.0127 \pm .0116

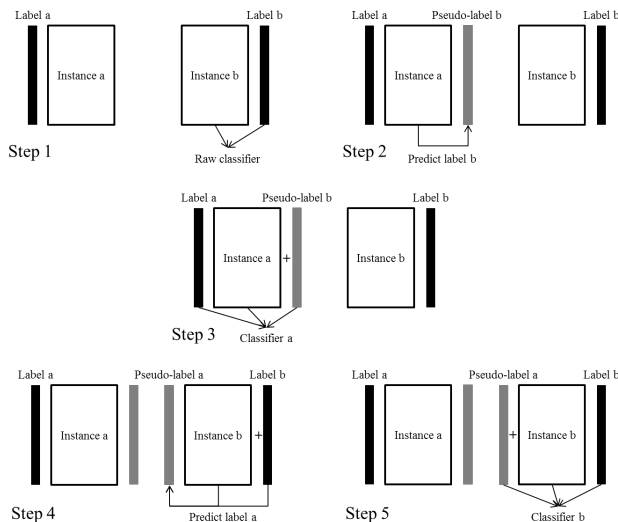
our approach

potential solutions

classical multi-label approaches

- **DSML is better than others**

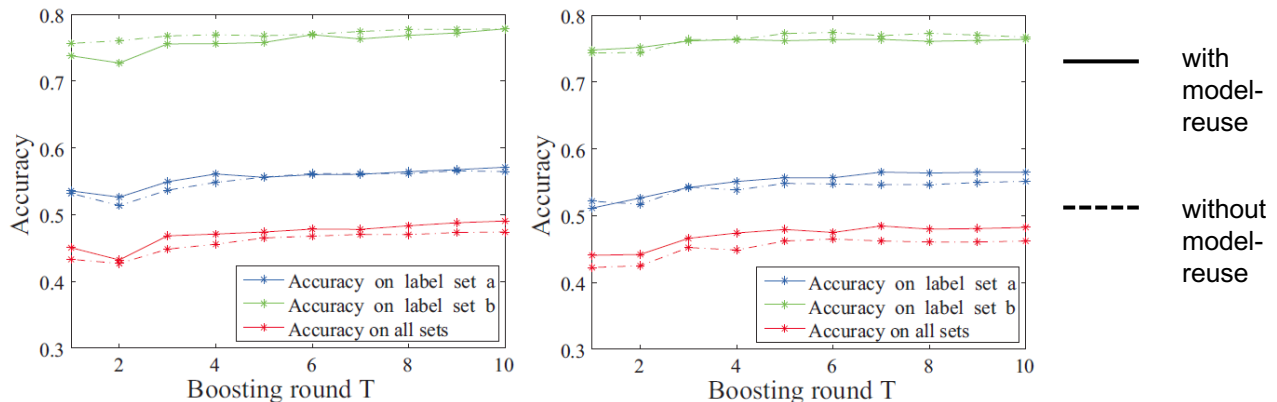
- Study on **model-reuse mechanism**



- 5-fold cross-validation performance of DSML on the *Brand-Type* dataset (mean)
 - Boosting round increases to 10
 - Distribution adjusting parameter is set to 1.00 and 1.10

Experiments

- Study on **model-reuse mechanism**
 - 5-fold cross-validation Performance of DSML on the *Brand-Type* dataset (mean)
 - Boosting round increases to 10
 - Distribution adjusting parameter is set to 1.00 and 1.10



- It validates the effectiveness of model-reuse mechanism
 - Similar phenomena can be observed in other datasets

- Study on **distribution adjusting mechanism**

- B is the distribution adjusting parameter

$$w_{t+1,i}^a = w_{t,i}^a \exp(\alpha_t^a \cdot \mathbb{I}[y_i^a \neq \hat{y}_i^a]) B^{\mathbb{I}[y_i^b \neq \hat{y}_i^b]}$$

$$w_{t+1,i}^b = w_{t,i}^b \exp(\alpha_t^b \cdot \mathbb{I}[y_i^b \neq \hat{y}_i^b]) B^{\mathbb{I}[y_i^a \neq \hat{y}_i^a]}$$

- When $B = 1.00$, algorithms perform without distribution adjusting mechanism
- 5-fold cross-validation performance of DSML algorithm (mean±std.)

Dataset	Measure	Distribution Adjusting Parameter B							
		1.00	1.01	1.02	1.03	1.05	1.10	1.15	1.20
<i>Cal.-Font</i>	<i>Accy-a</i>	.6536 ± .0054	.6576 ± .0064	.6567 ± .0051	.6557 ± .0067	.6562 ± .0059	.6541 ± .0033	.6546 ± .0076	.6528 ± .0060
	<i>Accy-b</i>	.7225 ± .0060	.7244 ± .0062	.7249 ± .0043	.7263 ± .0046	.7223 ± .0079	.7246 ± .0041	.7210 ± .0037	.7230 ± .0054
	<i>Accy-all</i>	.5656 ± .0078	.5697 ± .0062	.5674 ± .0043	.5690 ± .0058	.5672 ± .0087	.5698 ± .0043	.5659 ± .0078	.5660 ± .0045
<i>Brand-Type</i>	<i>Accy-a</i>	.5710 ± .0296	.5657 ± .0259	.5706 ± .0303	.5706 ± .0206	.5723 ± .0226	.5648 ± .0185	.5710 ± .0201	.5603 ± .0343
	<i>Accy-b</i>	.7784 ± .0142	.7668 ± .0185	.7659 ± .0193	.7650 ± .0212	.7730 ± .0249	.7641 ± .0107	.7788 ± .0182	.7699 ± .0182
	<i>Accy-all</i>	.4905 ± .0324	.4847 ± .0227	.4856 ± .0257	.4882 ± .0231	.4949 ± .0227	.4824 ± .0073	.4922 ± .0228	.4833 ± .0340
<i>Freq.-Gndr.</i>	<i>Accy-a</i>	.8413 ± .0110	.8432 ± .0107	.8432 ± .0177	.8413 ± .0140	.8521 ± .0091	.8435 ± .0137	.8473 ± .0162	.8476 ± .0119
	<i>Accy-b</i>	.9541 ± .0071	.9531 ± .0041	.9512 ± .0073	.9554 ± .0074	.9547 ± .0061	.9515 ± .0040	.9557 ± .0054	.9560 ± .0038
	<i>Accy-all</i>	.8131 ± .0060	.8134 ± .0118	.8134 ± .0158	.8119 ± .0166	.8220 ± .0082	.8128 ± .0151	.8172 ± .0153	.8175 ± .0155

- Best result appears when $B = 1.05$

- Introduction
- Potential Solutions and Deficiencies
- Our Approach
- Theoretical Results
- Experiments
- Conclusion

- **Dual Set Multi-Label Learning** is proposed as a novel learning framework.
- A boosting-like **DSML approach** is designed to address this kind of problem which outperforms other compared algorithms.
- Theoretical and empirical analyses are presented to show **it is better to learn with dual label sets** than to learn directly from all labels.



Thank you for listening.

Q & A