



# Learning from Semi-Supervised Weak-Label Data

Hao-Chen Dong, Yu-Feng Li, Zhi-Hua Zhou  
[donghc@lamda.nju.edu.cn](mailto:donghc@lamda.nju.edu.cn)

LAMDA Group  
National Key Lab for Novel Software Technology  
Nanjing University, China

# What is the work about?

---

- In many multi-label learning tasks, it's difficult to get full relevant label set and label incompleteness significantly influences the performance of multi-label learning.
- Previous studies usually consider only one scenario in which labels are incomplete. In fact, there are various scenarios. The combination of these conditions have not been fully studied.
- We propose SSWL method to address semi-supervised weak-label problem which labels can be partially known and completely unknown.

# Outline

---

- Background
- Proposed Method
- Experiments
- Conclusion

# Background

- In many real-world tasks, one instance usually has more than one label. Traditional supervised learning based on one label per instance is out of its capability to cope with this problem.



*car,  
road,  
sky,  
grass*



*tower,  
sky*

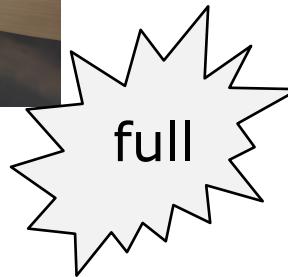
- **Multi-label learning** deals with instances associated with a set of labels and it has attracted much attention.

# Background

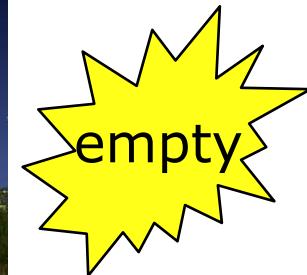
- But, it's difficult to get full label set.



*car,  
road,  
sky,  
grass*



*tower,  
sky*



- We can use some labeled instances and many unlabeled instances for training. This kind of multi-label learning problem is **semi-supervised multi-label learning problem**.

# Background

- But, it's difficult to get full label set.



*car,  
road,  
sky,  
grass*

partial



*tower,  
sky*

partial

- In many cases, we only use labeled instances for training, but we can't make sure if we get full relevant label set, it may be partial. This is **weak-label learning problem**.

# Background

---

- Semi-supervised multi-label learning tries to handle the issue of empty relevant label set.
- Transductive multi-label learning methods ([Liu, Jin and Yang 2006](#); [Chen et al. 2008](#); [Guo and Schuurmans 2012](#); [Kong, Ng and Zhou 2013](#)) that assume testing instances are from unlabeled instances.
- Pure semi-supervised multi-label learning methods ([Zhao and Guo 2015](#); [Zhan and Zhang 2017](#)) try to make multi-label prediction for any unseen instance.

# Background

---

- Weak-label learning focuses on the issue of partial relevant label set.
- WELL ([Sun, Zhang and Zhou 2010](#)) is based on the assumption that instance similarities are determined by a group of low-rank similarity matrixes.
- MLR-GL ([Bucak, Jin and Jain 2011](#)) uses group lasso to regularize the training errors.

# Background

- More general, we usually have instances with partial or empty label sets for training.



*car,  
road,  
sky,  
grass*



*tower,  
sky*

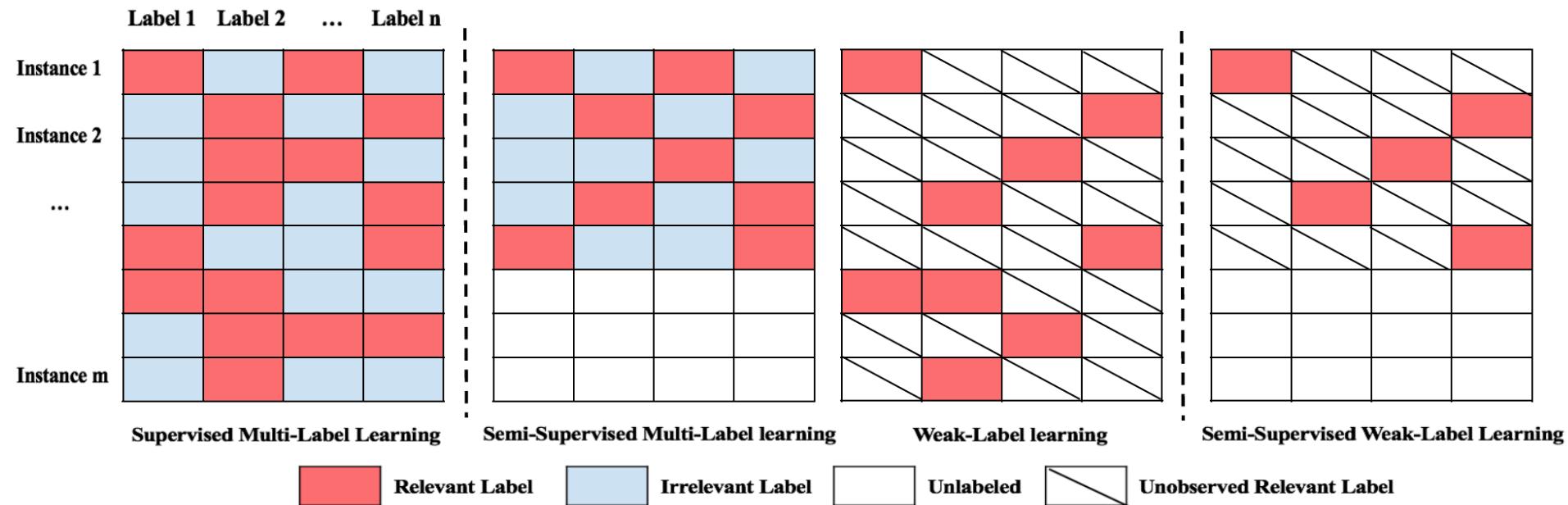
**partial**

**empty**

- We call this problem as **semi-supervised weak-label learning problem**. This problem considers that labels are partially known and completely unknown.

# Proposed Method

- We've summarized four problem settings.



- The number of known labels are not more than the existing problems'. To learn a promising predictive label, we need to complement missing relevant labels at first.

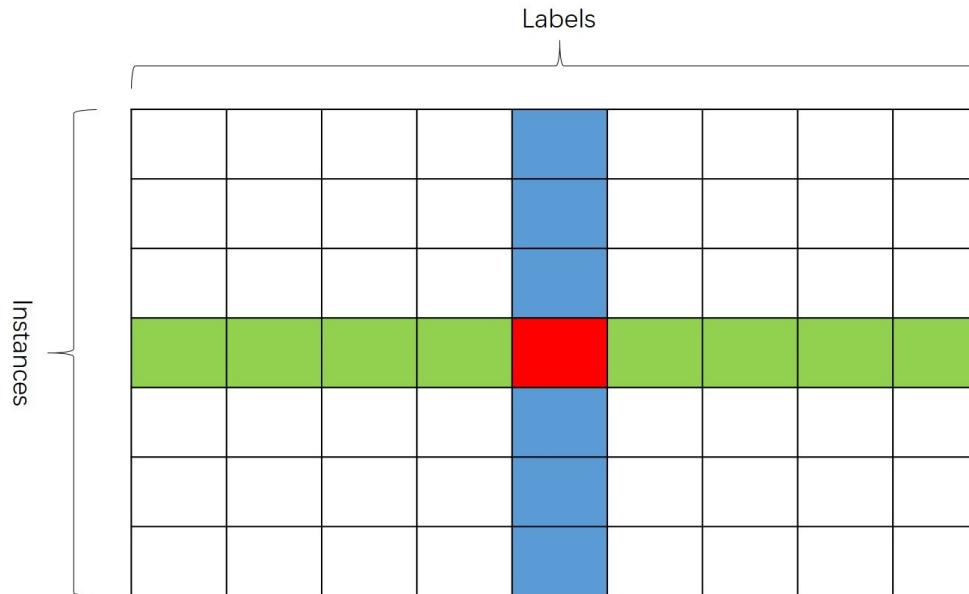
# Proposed Method

---

- Traditional label propagation method use instance similarity. The relevant label set of an instance can be derived by that of its nearest neighbors.
- In multi label learning, the assignment of one certain label on training instances can be derived by the assignments of its adjacent labels.
- In order to take label correlation into account, we propose to use both the instance and label similarity for the complementation of missing relevant labels.

# Proposed Method

- We have this picture in our mind.



- Then we can get a new regularization term to complement missing relevant labels.

$$\begin{aligned}
 \Omega(\hat{\mathbf{Y}}, \mathbf{S}, \mathbf{L}) &= \sum_{ij} (\hat{Y}_{ij} - \sum_{p \in \mathcal{N}_i} \sum_{q \in \hat{\mathcal{N}}_j} S_{ip} \hat{Y}_{pq} L_{qj})^2 \\
 &= \|\hat{\mathbf{Y}} - \mathbf{S}\hat{\mathbf{Y}}\mathbf{L}\|_F^2
 \end{aligned}$$

# Proposed Method

---

- In semi-supervised weak-label learning problem, we don't have the irrelevant labels, we only have some relevant labels.
- Basic ideas
  - We can find the relevant labels in the uncertainty labels, and the rest of them are irrelevant labels.
  - We can employ ensemble learning which is known to be more robust than a single model.

# Proposed Method

---

- We first build two models with the new regularization term for labeled and unlabeled instances respectively. Formally, let  $\mathbf{XW}$  and  $\mathbf{X}\bar{\mathbf{W}}$  denote two linear multi-label models.
- The first model  $\mathbf{XW}$  is initialized to predict the observed relevant labels, whose objective is formulated as  $\|(\mathbf{XW}) \circ \mathbf{C} - \mathbf{C}\|_F^2$ .
- The second model  $\mathbf{X}\bar{\mathbf{W}}$  is initialized to predict the uncertain elements in the label occurrence matrix  $\mathbf{C} \in \{0, 1\}^{m \times n}$ , whose objective is formulated as  $\|(\mathbf{X}\bar{\mathbf{W}}) \circ (\mathbf{E} - \mathbf{C}) + (\mathbf{E} - \mathbf{C})\|_F^2$ .
- We then leverage them via the promising co-regularization framework to derive a robust predictive result, whose objective is cast as  $\|(\mathbf{X}(\mathbf{W} - \bar{\mathbf{W}})) \circ (\mathbf{E} - \mathbf{C})\|_F^2$ .

# Proposed Method

---

- So we combine these and get the final objective function.

$$\begin{aligned}
 \min_{\mathbf{W}, \bar{\mathbf{W}}, \mathbf{L}} \quad & \|(\mathbf{X}\mathbf{W}) \circ \mathbf{C} - \mathbf{C}\|_F^2 + \alpha \Omega(\mathbf{U}, \mathbf{S}, \mathbf{L}) + \\
 & \beta \|(\mathbf{X}(\mathbf{W} - \bar{\mathbf{W}})) \circ (\mathbf{E} - \mathbf{C})\|_F^2 + \\
 & \zeta \|(\mathbf{X}\bar{\mathbf{W}}) \circ (\mathbf{E} - \mathbf{C}) + (\mathbf{E} - \mathbf{C})\|_F^2 \\
 \text{s.t.} \quad & \mathbf{U} = (\mathbf{X}\mathbf{W}) \circ \mathbf{C} + (\mathbf{X}\bar{\mathbf{W}}) \circ (\mathbf{E} - \mathbf{C})
 \end{aligned}$$

- Where  $\alpha, \beta, \zeta$  are the parameters.  $\mathbf{U} = (\mathbf{X}\mathbf{W}) \circ \mathbf{C} + (\mathbf{X}\bar{\mathbf{W}}) \circ (\mathbf{E} - \mathbf{C})$  is the integrated prediction of two models and we can smooth two models at the same time.
- It is worth noting that classical label propagation techniques can be realized as a special case of our proposal and our method can make prediction for any unseen instance.

# Proposed Method

---

- We observe the objective function is a bi-convex function.
- We extend an efficient block coordinate descend algorithm to optimize the objective function.

---

## **Algorithm 1:** SSWL Method

---

**Input :**  $\mathbf{X}$ :  $m \times d$  instance matrix

$\mathbf{C}$ :  $m \times n$  label occurrence matrix

$\mathbf{S}$ :  $m \times m$  similarity matrix of instances

**Output:**  $\mathbf{W}$  and  $\bar{\mathbf{W}}$ :  $d \times n$  coefficient matrixes

1 Initialize  $\mathbf{W}$ ,  $\bar{\mathbf{W}}$ ,  $\mathbf{L}$ ;

2 **while** *not converged* **do**

3 Fix  $\bar{\mathbf{W}}$  and  $\mathbf{L}$ , update  $\mathbf{W}$  by Eq.5;

4 Fix  $\mathbf{W}$  and  $\mathbf{L}$ , update  $\bar{\mathbf{W}}$  by Eq.7;

5 Fix  $\mathbf{W}$  and  $\bar{\mathbf{W}}$ , update  $\mathbf{L}$  by Eq.9;

6 **end**

---

# Proposed Method

---

- If we directly use the derivative of the objective function, we will get the complex matrix equation.
- We can use this theorem to transform the matrix equation to the normal linear equation.

**Theorem 1.** (Horn and Johnson 1991) Suppose a matrix  $\hat{\mathbf{X}}$  satisfies an equation,  $\sum_{i=1}^b \mathbf{A}_i \hat{\mathbf{X}} \mathbf{B}_i = \mathbf{V}$ , where  $\{\mathbf{A}_i\}_{i=1}^b$ ,  $\{\mathbf{B}_i\}_{i=1}^b$  and  $\mathbf{V}$  are known. To obtain the solution  $\hat{\mathbf{X}}$ , one could solve the following equivalent problem instead,  $(\sum_{i=1}^b \mathbf{B}'_i \otimes \mathbf{A}_i) \text{vec}(\hat{\mathbf{X}}) = \text{vec}(\mathbf{V})$ , which is a normal linear equation.

- In this way, we can use some efficient algorithms, such as the conjugate gradient algorithm, to solve the equation.

# Proposed Method

---

- We use the optimized solution of the objective function as our final  $\mathbf{L}$ .
- After getting the coefficient matrix  $\mathbf{W}$ , we use 0 as the threshold to discretize the predictive result. If the predictive result is greater than 0, we set the label as 1, otherwise we set the label as -1.

# Experiments

---

- Datasets
  - *TMC* dataset
  - *Yeast* dataset
  - *SceneImage* dataset
  - *msrc* dataset
- Comparative methods
  - state-of-the-art weak-label learning methods  
**MLR-GL** (Bucak, Jin, and Jain 2011) **SSML** (Zhao and Guo 2015)
  - state-of-the-art supervised multi-label learning method  
**ML-KNN** (Zhang and Zhou 2007)
  - three naive methods that directly decompose the task into multiple binary classification problems via treating labels independently.  
**Well-SVM** (Li et al. 2013) **S4VM** (Li and Zhou 2015) **BSVM**
  - our method without using unlabeled data  
**SSWL-wo**

# Experiments

---

- Incomplete label ratio (I.L. Ratio)
  - we randomly drop {0%, 20%, 40%, 60%} of the observed labels on the labeled training data.
- Evaluation criteria
  - **Micro-F1** and **Macro-F1**  
take both precision and recall into account  
the larger the value, the better the performance
  - **Hamming Loss (H.L.)**  
evaluates the fraction of misclassified instance-label pairs  
the smaller the value, the better the performance

# Experiments

---

- Results on *TMC* dataset are shown in Table 2

Table 2: Experimental results (mean $\pm$ std) on *TMC*.  $\uparrow$  ( $\downarrow$ ) indicates the larger (smaller) the better. The best performance and its comparable performances are bolded (pairwise t-tests at 95% significance level).

	I.L. Ratio	SSWL	SSWL-wo	Well-SVM	MLR-GL	SSML	ML-kNN	S4VM	BSVM
Micro-F1( $\uparrow$ )	0%	<b>.640 <math>\pm</math> .001</b>	.639 $\pm$ .001	.612 $\pm$ .003	.615 $\pm$ .001	.638 $\pm$ .001	.501 $\pm$ .001	.578 $\pm$ .002	.487 $\pm$ .002
	20%	<b>.602 <math>\pm</math> .003</b>	.578 $\pm$ .001	.556 $\pm$ .002	.596 $\pm$ .002	.580 $\pm$ .002	.213 $\pm$ .002	.506 $\pm$ .001	.292 $\pm$ .001
	40%	<b>.582 <math>\pm</math> .001</b>	.455 $\pm$ .004	.356 $\pm$ .002	.461 $\pm$ .003	.423 $\pm$ .001	.032 $\pm$ .001	.365 $\pm$ .003	.023 $\pm$ .002
	60%	<b>.570 <math>\pm</math> .002</b>	.505 $\pm$ .001	.113 $\pm$ .002	.563 $\pm$ .002	.160 $\pm$ .022	.012 $\pm$ .001	.215 $\pm$ .002	.007 $\pm$ .003
Macro-F1( $\uparrow$ )	0%	.618 $\pm$ .002	<b>.620 <math>\pm</math> .003</b>	.586 $\pm$ .002	.588 $\pm$ .002	.613 $\pm$ .002	.467 $\pm$ .001	.545 $\pm$ .002	.464 $\pm$ .002
	20%	<b>.582 <math>\pm</math> .001</b>	.568 $\pm$ .001	.519 $\pm$ .001	.567 $\pm$ .001	.543 $\pm$ .002	.175 $\pm$ .001	.457 $\pm$ .002	.244 $\pm$ .001
	40%	<b>.566 <math>\pm</math> .003</b>	.409 $\pm$ .005	.295 $\pm$ .003	.413 $\pm$ .002	.368 $\pm$ .002	.024 $\pm$ .001	.309 $\pm$ .002	.017 $\pm$ .001
	60%	<b>.553 <math>\pm</math> .002</b>	.494 $\pm$ .001	.089 $\pm$ .005	.537 $\pm$ .001	.125 $\pm$ .029	.008 $\pm$ .001	.279 $\pm$ .002	.004 $\pm$ .001
H.L.( $\downarrow$ )	0%	<b>.065 <math>\pm</math> .002</b>	.069 $\pm$ .001	.067 $\pm$ .002	.076 $\pm$ .001	.067 $\pm$ .001	.082 $\pm$ .002	.085 $\pm$ .002	.080 $\pm$ .001
	20%	.075 $\pm$ .002	.075 $\pm$ .001	.072 $\pm$ .001	.078 $\pm$ .002	<b>.071 <math>\pm</math> .002</b>	.092 $\pm$ .002	.087 $\pm$ .005	.088 $\pm$ .001
	40%	<b>.079 <math>\pm</math> .001</b>	.082 $\pm$ .002	.083 $\pm$ .002	.086 $\pm$ .002	.080 $\pm$ .001	.099 $\pm$ .001	.092 $\pm$ .002	.100 $\pm$ .001
	60%	.087 $\pm$ .003	.089 $\pm$ .002	.097 $\pm$ .003	<b>.084 <math>\pm</math> .002</b>	.095 $\pm$ .003	.101 $\pm$ .002	.111 $\pm$ .003	.101 $\pm$ .002

# Experiments

- Results on *TMC* dataset are shown in Table 2

Table 2: Experimental results (mean $\pm$ std) on *TMC*.  $\uparrow$  ( $\downarrow$ ) indicates the larger (smaller) the better. The best performance and its comparable performances are bolded (pairwise t-tests at 95% significance level).

	I.L. Ratio	SSWL	SSWL-wo	Well-SVM	MLR-GL	SSML	ML-kNN	S4VM	BSVM
Micro-F1( $\uparrow$ )	0%	<b>.640 <math>\pm</math> .001</b>	.639 $\pm$ .001	.612 $\pm$ .003	.615 $\pm$ .001	.638 $\pm$ .001	.501 $\pm$ .001	.578 $\pm$ .002	.487 $\pm$ .002
	20%	<b>.602 <math>\pm</math> .003</b>	.578 $\pm$ .001	.556 $\pm$ .002	.596 $\pm$ .002	.580 $\pm$ .002	.213 $\pm$ .002	.506 $\pm$ .001	.292 $\pm$ .001
	40%	<b>.582 <math>\pm</math> .001</b>	.455 $\pm$ .004	.356 $\pm$ .002	.461 $\pm$ .003	.423 $\pm$ .001	.032 $\pm$ .001	.365 $\pm$ .003	.023 $\pm$ .002
	60%	<b>.570 <math>\pm</math> .002</b>	.505 $\pm$ .001	.113 $\pm$ .002	.563 $\pm$ .002	.160 $\pm$ .022	.012 $\pm$ .001	.215 $\pm$ .002	.007 $\pm$ .003
Macro-F1( $\uparrow$ )	0%	.618 $\pm$ .002	<b>.620 <math>\pm</math> .003</b>	.586 $\pm$ .002	.588 $\pm$ .002	.613 $\pm$ .002	.467 $\pm$ .001	.545 $\pm$ .002	.464 $\pm$ .002
	20%	<b>.582 <math>\pm</math> .001</b>	.568 $\pm$ .001	.519 $\pm$ .001	.567 $\pm$ .001	.543 $\pm$ .002	.175 $\pm$ .001	.457 $\pm$ .002	.244 $\pm$ .001
	40%	<b>.566 <math>\pm</math> .003</b>	.409 $\pm$ .005	.295 $\pm$ .003	.413 $\pm$ .002	.368 $\pm$ .002	.024 $\pm$ .001	.309 $\pm$ .002	.017 $\pm$ .001
	60%	<b>.553 <math>\pm</math> .002</b>	.494 $\pm$ .001	.089 $\pm$ .005	.537 $\pm$ .001	.125 $\pm$ .029	.008 $\pm$ .001	.279 $\pm$ .002	.004 $\pm$ .001
H.L.( $\downarrow$ )	0%	<b>.065 <math>\pm</math> .002</b>	.069 $\pm$ .001	.067 $\pm$ .002	.076 $\pm$ .001	.067 $\pm$ .001	.082 $\pm$ .002	.085 $\pm$ .002	.080 $\pm$ .001
	20%	.075 $\pm$ .002	.075 $\pm$ .001	.072 $\pm$ .001	.078 $\pm$ .002	<b>.071 <math>\pm</math> .002</b>	.092 $\pm$ .002	.087 $\pm$ .005	.088 $\pm$ .001
	40%	<b>.079 <math>\pm</math> .001</b>	.082 $\pm$ .002	.083 $\pm$ .002	.086 $\pm$ .002	.080 $\pm$ .001	.099 $\pm$ .001	.092 $\pm$ .002	.100 $\pm$ .001
	60%	.087 $\pm$ .003	.089 $\pm$ .002	.097 $\pm$ .003	<b>.084 <math>\pm</math> .002</b>	.095 $\pm$ .003	.101 $\pm$ .002	.111 $\pm$ .003	.101 $\pm$ .002

- SSWL obtains quite promising performance and consistently performs robustly while the I.L. Ratio getting higher.

# Experiments

- Results on *TMC* dataset are shown in Table 2

Table 2: Experimental results (mean $\pm$ std) on *TMC*.  $\uparrow$  ( $\downarrow$ ) indicates the larger (smaller) the better. The best performance and its comparable performances are bolded (pairwise t-tests at 95% significance level).

	I.L. Ratio	SSWL	SSWL-wo	Well-SVM	MLR-GL	SSML	ML-kNN	S4VM	BSVM
Micro-F1( $\uparrow$ )	0%	<b>.640 <math>\pm</math> .001</b>	.639 $\pm$ .001	.612 $\pm$ .003	.615 $\pm$ .001	.638 $\pm$ .001	.501 $\pm$ .001	.578 $\pm$ .002	.487 $\pm$ .002
	20%	<b>.602 <math>\pm</math> .003</b>	.578 $\pm$ .001	.556 $\pm$ .002	.596 $\pm$ .002	.580 $\pm$ .002	.213 $\pm$ .002	.506 $\pm$ .001	.292 $\pm$ .001
	40%	<b>.582 <math>\pm</math> .001</b>	.455 $\pm$ .004	.356 $\pm$ .002	.461 $\pm$ .003	.423 $\pm$ .001	.032 $\pm$ .001	.365 $\pm$ .003	.023 $\pm$ .002
	60%	<b>.570 <math>\pm</math> .002</b>	.505 $\pm$ .001	.113 $\pm$ .002	.563 $\pm$ .002	.160 $\pm$ .022	.012 $\pm$ .001	.215 $\pm$ .002	.007 $\pm$ .003
Macro-F1( $\uparrow$ )	0%	.618 $\pm$ .002	<b>.620 <math>\pm</math> .003</b>	.586 $\pm$ .002	.588 $\pm$ .002	.613 $\pm$ .002	.467 $\pm$ .001	.545 $\pm$ .002	.464 $\pm$ .002
	20%	<b>.582 <math>\pm</math> .001</b>	.568 $\pm$ .001	.519 $\pm$ .001	.567 $\pm$ .001	.543 $\pm$ .002	.175 $\pm$ .001	.457 $\pm$ .002	.244 $\pm$ .001
	40%	<b>.566 <math>\pm</math> .003</b>	.409 $\pm$ .005	.295 $\pm$ .003	.413 $\pm$ .002	.368 $\pm$ .002	.024 $\pm$ .001	.309 $\pm$ .002	.017 $\pm$ .001
	60%	<b>.553 <math>\pm</math> .002</b>	.494 $\pm$ .001	.089 $\pm$ .005	.537 $\pm$ .001	.125 $\pm$ .029	.008 $\pm$ .001	.279 $\pm$ .002	.004 $\pm$ .001
H.L.( $\downarrow$ )	0%	<b>.065 <math>\pm</math> .002</b>	.069 $\pm$ .001	.067 $\pm$ .002	.076 $\pm$ .001	.067 $\pm$ .001	.082 $\pm$ .002	.085 $\pm$ .002	.080 $\pm$ .001
	20%	.075 $\pm$ .002	.075 $\pm$ .001	.072 $\pm$ .001	.078 $\pm$ .002	<b>.071 <math>\pm</math> .002</b>	.092 $\pm$ .002	.087 $\pm$ .005	.088 $\pm$ .001
	40%	<b>.079 <math>\pm</math> .001</b>	.082 $\pm$ .002	.083 $\pm$ .002	.086 $\pm$ .002	.080 $\pm$ .001	.099 $\pm$ .001	.092 $\pm$ .002	.100 $\pm$ .001
	60%	.087 $\pm$ .003	.089 $\pm$ .002	.097 $\pm$ .003	<b>.084 <math>\pm</math> .002</b>	.095 $\pm$ .003	.101 $\pm$ .002	.111 $\pm$ .003	.101 $\pm$ .002

- SSWL-wo also obtains good performance but is not as good as SSWL.

# Experiments

- Results on *TMC* dataset are shown in Table 2

Table 2: Experimental results (mean $\pm$ std) on *TMC*.  $\uparrow$  ( $\downarrow$ ) indicates the larger (smaller) the better. The best performance and its comparable performances are bolded (pairwise t-tests at 95% significance level).

	I.L. Ratio	SSWL	SSWL-wo	Well-SVM	MLR-GL	SSML	ML-kNN	S4VM	BSVM
Micro-F1( $\uparrow$ )	0%	<b>.640 <math>\pm</math> .001</b>	.639 $\pm$ .001	.612 $\pm$ .003	.615 $\pm$ .001	.638 $\pm$ .001	.501 $\pm$ .001	.578 $\pm$ .002	.487 $\pm$ .002
	20%	<b>.602 <math>\pm</math> .003</b>	.578 $\pm$ .001	.556 $\pm$ .002	.596 $\pm$ .002	.580 $\pm$ .002	.213 $\pm$ .002	.506 $\pm$ .001	.292 $\pm$ .001
	40%	<b>.582 <math>\pm</math> .001</b>	.455 $\pm$ .004	.356 $\pm$ .002	.461 $\pm$ .003	.423 $\pm$ .001	.032 $\pm$ .001	.365 $\pm$ .003	.023 $\pm$ .002
	60%	<b>.570 <math>\pm</math> .002</b>	.505 $\pm$ .001	.113 $\pm$ .002	.563 $\pm$ .002	.160 $\pm$ .022	.012 $\pm$ .001	.215 $\pm$ .002	.007 $\pm$ .003
Macro-F1( $\uparrow$ )	0%	.618 $\pm$ .002	<b>.620 <math>\pm</math> .003</b>	.586 $\pm$ .002	.588 $\pm$ .002	.613 $\pm$ .002	.467 $\pm$ .001	.545 $\pm$ .002	.464 $\pm$ .002
	20%	<b>.582 <math>\pm</math> .001</b>	.568 $\pm$ .001	.519 $\pm$ .001	.567 $\pm$ .001	.543 $\pm$ .002	.175 $\pm$ .001	.457 $\pm$ .002	.244 $\pm$ .001
	40%	<b>.566 <math>\pm</math> .003</b>	.409 $\pm$ .005	.295 $\pm$ .003	.413 $\pm$ .002	.368 $\pm$ .002	.024 $\pm$ .001	.309 $\pm$ .002	.017 $\pm$ .001
	60%	<b>.553 <math>\pm</math> .002</b>	.494 $\pm$ .001	.089 $\pm$ .005	.537 $\pm$ .001	.125 $\pm$ .029	.008 $\pm$ .001	.279 $\pm$ .002	.004 $\pm$ .001
H.L.( $\downarrow$ )	0%	<b>.065 <math>\pm</math> .002</b>	.069 $\pm$ .001	.067 $\pm$ .002	.076 $\pm$ .001	.067 $\pm$ .001	.082 $\pm$ .002	.085 $\pm$ .002	.080 $\pm$ .001
	20%	.075 $\pm$ .002	.075 $\pm$ .001	.072 $\pm$ .001	.078 $\pm$ .002	<b>.071 <math>\pm</math> .002</b>	.092 $\pm$ .002	.087 $\pm$ .005	.088 $\pm$ .001
	40%	<b>.079 <math>\pm</math> .001</b>	.082 $\pm$ .002	.083 $\pm$ .002	.086 $\pm$ .002	.080 $\pm$ .001	.099 $\pm$ .001	.092 $\pm$ .002	.100 $\pm$ .001
	60%	.087 $\pm$ .003	.089 $\pm$ .002	.097 $\pm$ .003	<b>.084 <math>\pm</math> .002</b>	.095 $\pm$ .003	.101 $\pm$ .002	.111 $\pm$ .003	.101 $\pm$ .002

- Our proposal works better than state-of-the-art weak-label learning and semi-supervised multi-label learning algorithms.

# Experiments

---

- Results on *yeast* dataset.

Table 3: Experimental results (mean $\pm$ std) on *yeast*.  $\uparrow$  ( $\downarrow$ ) indicates the larger (smaller) the better. The best performance and its comparable performances are bolded (pairwise t-tests at 95% significance level).

	I.L. Ratio	SSWL	SSWL-wo	Well-SVM	MLR-GL	SSML	ML-kNN	S4VM	BSVM
Micro-F1( $\uparrow$ )	0%	<b>.647 <math>\pm</math> .001</b>	.619 $\pm$ .002	.612 $\pm$ .002	.623 $\pm$ .003	.584 $\pm$ .001	.625 $\pm$ .002	.592 $\pm$ .001	.623 $\pm$ .002
	20%	<b>.638 <math>\pm</math> .001</b>	.626 $\pm$ .002	.534 $\pm$ .003	.618 $\pm$ .005	.510 $\pm$ .002	.506 $\pm$ .001	.511 $\pm$ .001	.509 $\pm$ .001
	40%	<b>.604 <math>\pm</math> .002</b>	.554 $\pm$ .003	.394 $\pm$ .003	.379 $\pm$ .004	.152 $\pm$ .001	.103 $\pm$ .002	.432 $\pm$ .001	.188 $\pm$ .003
	60%	<b>.616 <math>\pm</math> .002</b>	.568 $\pm$ .002	.241 $\pm$ .002	.209 $\pm$ .002	.046 $\pm$ .031	.002 $\pm$ .002	.320 $\pm$ .002	.019 $\pm$ .007
Macro-F1( $\uparrow$ )	0%	<b>.635 <math>\pm</math> .001</b>	.594 $\pm$ .001	.582 $\pm$ .003	.600 $\pm$ .001	.557 $\pm$ .002	.602 $\pm$ .001	.578 $\pm$ .002	.592 $\pm$ .002
	20%	<b>.618 <math>\pm</math> .001</b>	.613 $\pm$ .001	.494 $\pm$ .002	.593 $\pm$ .004	.476 $\pm$ .001	.470 $\pm$ .001	.476 $\pm$ .002	.478 $\pm$ .002
	40%	<b>.574 <math>\pm</math> .001</b>	.538 $\pm$ .005	.359 $\pm$ .004	.340 $\pm$ .002	.126 $\pm$ .001	.083 $\pm$ .001	.397 $\pm$ .001	.145 $\pm$ .002
	60%	<b>.595 <math>\pm</math> .002</b>	.554 $\pm$ .001	.194 $\pm$ .003	.177 $\pm$ .002	.039 $\pm$ .025	.001 $\pm$ .001	.280 $\pm$ .005	.016 $\pm$ .003
H.L.( $\downarrow$ )	0%	<b>.207 <math>\pm</math> .001</b>	.209 $\pm$ .001	.211 $\pm$ .001	.213 $\pm$ .002	.215 $\pm$ .002	.208 $\pm$ .001	.214 $\pm$ .001	.209 $\pm$ .001
	20%	<b>.210 <math>\pm</math> .001</b>	.216 $\pm$ .001	.221 $\pm$ .002	.211 $\pm$ .002	.224 $\pm$ .002	.225 $\pm$ .001	.253 $\pm$ .002	.224 $\pm$ .001
	40%	<b>.225 <math>\pm</math> .001</b>	.252 $\pm$ .004	.246 $\pm$ .003	.251 $\pm$ .001	.286 $\pm$ .002	.294 $\pm$ .001	.257 $\pm$ .001	.279 $\pm$ .002
	60%	<b>.231 <math>\pm</math> .003</b>	.268 $\pm$ .003	.278 $\pm$ .002	.275 $\pm$ .003	.299 $\pm$ .010	.305 $\pm$ .001	.286 $\pm$ .002	.302 $\pm$ .002

# Experiments

---

- Results on *SceneImage* dataset.

Table 4: Experimental results (mean $\pm$ std) on *SceneImage*.  $\uparrow$  ( $\downarrow$ ) indicates the larger (smaller) the better. The best performance and its comparable performances are bolded (pairwise t-tests at 95% significance level).

	I.L. Ratio	SSWL	SSWL-wo	Well-SVM	MLR-GL	SSML	ML-kNN	S4VM	BSVM
Micro-F1( $\uparrow$ )	0%	.589 $\pm$ .003	.583 $\pm$ .002	<b>.613 <math>\pm</math> .002</b>	.528 $\pm$ .002	.499 $\pm$ .002	.456 $\pm$ .001	.489 $\pm$ .001	.538 $\pm$ .002
	20%	<b>.572 <math>\pm</math> .002</b>	.545 $\pm$ .001	.558 $\pm$ .003	.419 $\pm$ .001	.394 $\pm$ .002	.389 $\pm$ .002	.458 $\pm$ .003	.392 $\pm$ .001
	40%	<b>.540 <math>\pm</math> .002</b>	.534 $\pm$ .002	.395 $\pm$ .001	.220 $\pm$ .001	.174 $\pm$ .001	.094 $\pm$ .001	.289 $\pm$ .001	.205 $\pm$ .001
	60%	<b>.521 <math>\pm</math> .003</b>	.517 $\pm$ .002	.251 $\pm$ .002	.000 $\pm$ .001	.019 $\pm$ .010	.000 $\pm$ .000	.300 $\pm$ .001	.010 $\pm$ .002
Macro-F1( $\uparrow$ )	0%	<b>.576 <math>\pm</math> .002</b>	.553 $\pm$ .001	.567 $\pm$ .001	.454 $\pm$ .003	.407 $\pm$ .001	.362 $\pm$ .001	.466 $\pm$ .001	.437 $\pm$ .001
	20%	<b>.550 <math>\pm</math> .003</b>	.505 $\pm$ .001	.494 $\pm$ .002	.320 $\pm$ .001	.295 $\pm$ .002	.295 $\pm$ .002	.417 $\pm$ .001	.284 $\pm$ .001
	40%	<b>.523 <math>\pm</math> .002</b>	.499 $\pm$ .001	.306 $\pm$ .001	.140 $\pm$ .001	.113 $\pm$ .001	.056 $\pm$ .001	.207 $\pm$ .002	.134 $\pm$ .002
	60%	<b>.510 <math>\pm</math> .004</b>	.495 $\pm$ .001	.165 $\pm$ .002	.000 $\pm$ .001	.011 $\pm$ .005	.000 $\pm$ .000	.230 $\pm$ .001	.005 $\pm$ .003
H.L.( $\downarrow$ )	0%	.184 $\pm$ .001	.186 $\pm$ .001	<b>.167 <math>\pm</math> .002</b>	.192 $\pm$ .001	.193 $\pm$ .001	.192 $\pm$ .002	.245 $\pm$ .002	.167 $\pm$ .001
	20%	.199 $\pm$ .001	.201 $\pm$ .001	<b>.182 <math>\pm</math> .002</b>	.199 $\pm$ .001	.203 $\pm$ .001	.204 $\pm$ .002	.236 $\pm$ .002	.194 $\pm$ .001
	40%	<b>.206 <math>\pm</math> .001</b>	.214 $\pm$ .001	<b>.206 <math>\pm</math> .001</b>	.219 $\pm$ .001	.227 $\pm$ .001	.240 $\pm$ .002	.226 $\pm$ .001	.221 $\pm$ .001
	60%	<b>.208 <math>\pm</math> .002</b>	.222 $\pm$ .001	.220 $\pm$ .002	.249 $\pm$ .001	.247 $\pm$ .003	.250 $\pm$ .002	.229 $\pm$ .002	.248 $\pm$ .002

# Experiments

---

- Results on *msrc* dataset.

Table 5: Experimental results (mean $\pm$ std) on *msrc*.  $\uparrow$  ( $\downarrow$ ) indicates the larger (smaller) the better. The best performance and its comparable performances are bolded (pairwise t-tests at 95% significance level).

	I.L. Ratio	SSWL	SSWL-wo	Well-SVM	MLR-GL	SSML	ML-kNN	S4VM	BSVM
Micro-F1( $\uparrow$ )	0%	.590 $\pm$ .003	.550 $\pm$ .003	.505 $\pm$ .003	.533 $\pm$ .003	<b>.604 <math>\pm</math> .002</b>	.355 $\pm$ .001	.351 $\pm$ .001	.487 $\pm$ .001
	20%	<b>.571 <math>\pm</math> .001</b>	.513 $\pm$ .003	.444 $\pm$ .001	.472 $\pm$ .002	.569 $\pm$ .003	.307 $\pm$ .001	.331 $\pm$ .002	.292 $\pm$ .002
	40%	.507 $\pm$ .001	.449 $\pm$ .005	.216 $\pm$ .002	.500 $\pm$ .001	<b>.533 <math>\pm</math> .002</b>	.122 $\pm$ .002	.316 $\pm$ .001	.047 $\pm$ .003
	60%	<b>.465 <math>\pm</math> .001</b>	.367 $\pm$ .002	.118 $\pm$ .003	.244 $\pm$ .001	.462 $\pm$ .001	.034 $\pm$ .001	.294 $\pm$ .001	.007 $\pm$ .001
Macro-F1( $\uparrow$ )	0%	.562 $\pm$ .003	.505 $\pm$ .004	.412 $\pm$ .002	.465 $\pm$ .004	<b>.612 <math>\pm</math> .001</b>	.291 $\pm$ .002	.378 $\pm$ .001	.464 $\pm$ .001
	20%	<b>.522 <math>\pm</math> .001</b>	.419 $\pm$ .004	.302 $\pm$ .002	.390 $\pm$ .002	.518 $\pm$ .002	.216 $\pm$ .001	.369 $\pm$ .002	.244 $\pm$ .002
	40%	<b>.491 <math>\pm</math> .002</b>	.396 $\pm$ .004	.151 $\pm$ .003	.424 $\pm$ .002	.472 $\pm$ .001	.072 $\pm$ .002	.317 $\pm$ .002	.032 $\pm$ .002
	60%	<b>.430 <math>\pm</math> .001</b>	.298 $\pm$ .001	.082 $\pm$ .002	.166 $\pm$ .001	.391 $\pm$ .002	.037 $\pm$ .001	.282 $\pm$ .001	.004 $\pm$ .001
H.L.( $\downarrow$ )	0%	<b>.083 <math>\pm</math> .001</b>	.085 $\pm$ .001	<b>.083 <math>\pm</math> .002</b>	.093 $\pm$ .002	.108 $\pm$ .002	.094 $\pm$ .001	.154 $\pm$ .001	.092 $\pm$ .001
	20%	<b>.085 <math>\pm</math> .001</b>	.091 $\pm$ .001	.086 $\pm$ .001	.101 $\pm$ .001	.105 $\pm$ .002	.095 $\pm$ .001	.163 $\pm$ .002	.088 $\pm$ .001
	40%	.090 $\pm$ .001	.098 $\pm$ .002	.098 $\pm$ .003	<b>.083 <math>\pm</math> .002</b>	.100 $\pm$ .001	.101 $\pm$ .001	.182 $\pm$ .002	.100 $\pm$ .002
	60%	<b>.092 <math>\pm</math> .001</b>	.099 $\pm$ .001	.103 $\pm$ .002	.100 $\pm$ .002	.094 $\pm$ .001	.106 $\pm$ .001	.201 $\pm$ .001	.101 $\pm$ .001

# Conclusion

---

- We consider semi-supervised weak-label learning problem which is a new kind of multi-label learning problem.
- We propose the SSWL method to address this problem.
  - both instance similarity and label similarity are considered for the complement of missing labels
  - ensemble of multiple models is employed which is more robust than a single model when the label information is insufficient
  - experiments on a number of real tasks validate the effectiveness of SSWL in handling the semi-supervised weak-label learning problem
- Code is on <http://lamda.nju.edu.cn/donghc/>

Thanks!

# Review opinions

---

**Masked Reviewer ID:**Assigned\_Reviewer\_1

**Review:**

Question	
[Summary] Please summarize the main claims/contributions of the paper in your own words.	This manuscript essentially elaborates the problem for semi-supervised weak-label learning and proposes a bi-convex optimization problem with block coordinate descend algorithm. The authors have given a clear comparison between other works and theirs. In its present form, the organization is clear and easy to follow.
[Relevance] Is this paper relevant to an AI audience?	Likely to be of interest to a large proportion of the community
[Significance] Are the results significant?	Significant
[Novelty] Are the problems or approaches novel?	Novel
[Soundness] Is the paper technically sound?	Technically sound
[Evaluation] Are claims well-supported by theoretical analysis or experimental results?	Very convincing
[Clarity] Is the paper well-organized and clearly written?	Good
[Detailed Comments] Please elaborate on your assessments and provide constructive feedback.	Comments with the order of appearance are as follows: 1. In this paper, the authors propose the semi-supervised problem, but the contribution of unlabeled data is not clear. More details should be provided. 2. In this paper, the authors propose ensemble learning model. While in this way, the relevant labels have less effect on irrelevant labels, which seems unreasonable. 3. In this paper, the authors only provide the optimization of label similarity $L$ , but how to determine the final $L$ , only select top K elements? More details should be provided. 4. For the experiments, the compared algorithms should have more state-of-the-art algorithms.
[QUESTIONS FOR THE AUTHORS] Please provide questions for authors to address during the author feedback period.	1. The contribution of unlabeled data is not clear. More details should be provided. 2. Explain why the ensemble learning model used. 3. More details should be provided for $L$ . 4. Add more state-of-the-art algorithms, e.g. Semi-supervised Multi-label Learning with Incomplete Labels in IJCAI 2015.
[OVERALL SCORE]	Accept (Top 50% accepted papers (est.))
Please acknowledge that you have read the author rebuttal. If your opinion has changed, please summarize the main reasons below.	After reading the author's rebuttal and reviews from fellow reviewers, I would like to maintain my previous decision.

# Review opinions

---

Masked Reviewer ID:Assigned\_Reviewer\_2

Review:

Question	
[Summary] Please summarize the main claims/contributions of the paper in your own words.	This paper focuses on multi-label learning. Specifically, the authors notice an important observation which is ignored in most previous multi-label works that the labels could be partially and completely missing simultaneously, called the semi-supervised weak-label learning. The authors demonstrate that treating them independently may harm the overall accuracy. To address this issue, this paper proposes a novel ensemble algorithm based on both label similarity and instance similarity and then building SSWL models based on the co-regularization strategy. Experiments on several benchmark datasets show that the SSWL accuracy is significantly improved based on the ensemble algorithm.
[Relevance] Is this paper relevant to an AI audience?	Likely to be of interest to a large proportion of the community
[Significance] Are the results significant?	Highly significant
[Novelty] Are the problems or approaches novel?	Very novel
[Soundness] Is the paper technically sound?	Technically sound
[Evaluation] Are claims well-supported by theoretical analysis or experimental results?	Sufficient
[Clarity] Is the paper well-organized and clearly written?	Good
[Detailed Comments] Please elaborate on your assessments and provide constructive feedback.	Generally, I am satisfied with this paper. It is an interesting paper whose observation is meaningful and solution is simple yet effective. Most multi-label papers focus on supervised learning or one type of missing labels. But this paper pays attention to two kinds of missing labels simultaneously, partially known or completely unknown. The demonstration in Figure 1 is convincing and impressive which clearly points out the difference between previous setting and the proposed new setup, and a novel method SSWL is proposed to ensemble two classifiers via the regularization from both instance and label similarities. This paper notices this problem and improves SSWL accuracy from the ensemble-learning perspective, which is very good. The experiments also show that ensemble learning is helpful in many application data sets for many multi-label models. I believe this work can inspire many other multi-label works focusing on other parts in an algorithm to cope with complicated missing label issue in multi-label learning.
[QUESTIONS FOR THE AUTHORS] Please provide questions for authors to address during the author feedback period.	I have some questions which I wish the authors could address in the rebuttal.  1. Based on Eq. (2), the authors proposed a new regularization for both instance and label similarities. Since label similarity seems a little bit complicated when combined with subsequent SSWL model training, I am more interested that can the authors provide the results if only instance similarity is used? How much performance gain is achieved with only instance similarity?  2. Ensemble learning is an effective way to help improve the learning performance. To demonstrate the necessity of ensemble learning, a comprehensive evaluation should be given. What is the average performance by simply combining two initial learners?  3. I notice that macro and micro F1s have more gain than hamming loss for the proposed method. Could the authors give some explanation?  4. The authors use an iterative algorithm to solve Eq. (3). How about the efficiency of the algorithm?
[OVERALL SCORE]	Accept (Top 15%)
Please acknowledge that you have read the author rebuttal. If your opinion has changed, please summarize the main reasons below.	Thanks for the feedback.

# Review opinions

---

**Masked Reviewer ID:**Assigned\_Reviewer\_3

**Review:**

Question	
[Summary] Please summarize the main claims/contributions of the paper in your own words.	The paper presents a optimization problem which solves the problem of multi-label classification in the semi-supervised learning setting. Novelty of the proposed approach is that it works with weak-labels.
[Relevance] Is this paper relevant to an AI audience?	Likely to be of interest to a large proportion of the community
[Significance] Are the results significant?	Highly significant
[Novelty] Are the problems or approaches novel?	Novel
[Soundness] Is the paper technically sound?	Technically sound
[Evaluation] Are claims well-supported by theoretical analysis or experimental results?	Very convincing
[Clarity] Is the paper well-organized and clearly written?	Good
[Detailed Comments] Please elaborate on your assessments and provide constructive feedback.	<p>The paper presents a optimization problem which solves the problem of multi-label classification in the semi-supervised learning setting. Novelty of the proposed approach is that it works with weak-labels.</p> <p>This is a quite interesting task and, in some sense, extends the task of learning with incompletely annotated instances.</p> <p>Despite some typos, the paper is well and clearly written. The part on related works is quite condensed and misses some recent references which adopt complete different approaches (not for weak-labels).</p> <p>This is an example:          Jurica Levatic, Michelangelo Ceci, Dragi Kocev, Saso Dzeroski:          Self-training for multi-target regression with tree ensembles. Knowl.-Based Syst. 123: 41-60 (2017)</p> <p>The method is quite standard and the authors clearly presented it. Theorem 1 is at the basis of the optimization solution the authors present. Probably, more details on this theorem shoud be given.</p> <p>Results are obtained with a correct experimental setting on several datasets coming from different domains. The only negative aspect that can be raised is that the number of used datasets is relatively small. Another limitation is that results should be integrated with statistical tests. However, from the results I see I'm pretty sure that statistical tests will confirm the superiority of SSWL.</p> <p>In my opinion, the paper is relatively solid, even is not extraordinarily novel from the methodological point of view. However, the extension to weak labels is enough to justify the paper.</p> <p>Typos (non-exhaustive list):          We further compare with a variant of our proposal that does not usING the unlabeled instances".</p>
[QUESTIONS FOR THE AUTHORS] Please provide questions for authors to address during the author feedback period.	See detailed comments
[OVERALL SCORE]	Accept (Top 15%)
Please acknowledge that you have read the author rebuttal. If your opinion has changed, please summarize the main reasons below.	I acknowledge I have read the author rebuttal. It did not change my opinion about the paper: I'm still positive about acceptance.