

Instance Selection in Semi-Supervised Learning

Yuanyuan Guo, Harry Zhang, and Xiaobo Liu

data

Semi-Supervised Learning (SSL)

In many real-world applications:

- Labeled data (L): scarce; expensive/difficult to collect.
- Unlabeled data (U): abundant; relatively easy to obtain.

□ SSL uses labeled data and unlabeled data to learn hypotheses.

Usage	Supervised S learning	iemi-supervised learning	Unsupervised learning	
{(x,y)} labeled data	Yes	Yes	No	
{x} unlabeled data	No	Yes	Yes	

Motivation

• Selecting the most confident unlabeled instances ("CF") is a common instance selection method in two standard SSL methods: self-training and co-training.



t = t + 1



□ Shortcomings of "*CF*" :

- Label noise may be added to the training set.
- It is not necessarily superior to that of randomly selecting unlabeled instances.
- Our pervious work showed that the original labeled instances are more reliable than the self-labeled instances that are labeled by the classifier.

Results on 26 benchmark UCI datasets

	Self	-training	Co-training						
Accuracy	CF	ISBOLD	CF	ISBOLD					
Mean	71.80	74.61	70.34	73.71					
Paired t-test		6/20/0		7/19/0					
	Self	-training	Co-	training					
AUC	Self <i>CF</i>	-training ISBOLD	Co- CF	training ISBOLD					
AUC Mean	Self <i>CF</i> 80.57	-training <i>ISBOLD</i> 83.12	Co- <i>CF</i> 78.56	training <i>ISBOLD</i> 81.74					

ble 1: Accuracy of CF vs ISBOLD in self-training and co-training (a) self-training (b) co-training						
taset	CF	ISBOLD		Dataset	CF	ISBOLD
ance-scale	59.52	66.21		balance-scale	59.10	67.17
ast-cancer	65.09	65.61		breast-cancer	70.41	71.00
ast-w	96.67	96.34		breast-w	96.85	96.47
e	74.54	75.38		colic	76.60	75.76
c.ORIG	55.05	60.57		colic.ORIG	55.19	62.04
lit-a	80.68	80.78		credit-a	81.36	79.67
lit-g	60.62	66.03 v		credit-g	63.04	67.72 v
oetes	70.55	70.53		diabetes	67.51	69.58
rt-c	81.55	81.15		heart-c	82.77	80.13
rt-h	83.06	82.41		heart-h	81.46	78.60
rt-statlog	81.37	80.74		heart-statlog	82.03	80.30
atitis	79.70	78.34		hepatitis	81.04	80.21
osphere	80.97	79.86		ionosphere	81.50	83.08
	90.31	90.05		iris	80.79	78.98
/s-kp	67.26	80.07 v		kr-vs-kp	59.22	77.36 v
or	88.26	87.92		labor	77.21	78.43
er	40.38	57.39 v		letter	36.67	56.05 v
shroom	91.90	92.57 v		mushroom	91.74	92.38 v
ment	63.49	72.88 v		segment	61.49	71.64 v
	91.54	94.15		sick	93.40	93.56
ar	55.72	57.93		sonar	55.43	58.08
ce	82.05	85.48 v		splice	73.91	82.63 v
icle	41.79	48.35		vehicle	41.57	47.86
е	87.89	88.53		vote	88.21	88.60
vel	18.75	21.78		vowel	18.83	23.36
eform-5000	77.98	78.87		waveform-5000	71.61	75.91 v
an	71.80	74.61		mean	70.34	73.71
/1		6/20/0		w/t/l		7/19/0

- > Steps: In each iteration, after the selection of the most confident unlabeled instances,
 - a. Compute the accuracy of the current classifier on the original labeled;
 - b. Check whether the accuracy is lower than that in last iteration;
 - c. If so, discard the selected instances; Otherwise, add them to the training set in the next iteration.
- Reasons:
 - Performance on the original labeled data reflects the performance on the future testing set.
 - It prevents adding unlabeled instances that will possibly degrade the performance.

\succ Evaluation:

- 4-fold stratified cross-validation (10 runs), 26 UCI datasets
- |L|: |U| = 5% : 95%
- Measurements: Average accuracy and AUC; Paired t-test
- Implemented in Weka

Table 2: AUC of CF vs ISBOLD in self-training and co-training (b) co-trainin (a) self-training CF ISBOLD CF ISBOL Dataset Dataset 60.44 65.34 alance-scale alance-scale 63.51 64.37 63.98 63.48reast-cancer breast-cance 99.07 99.08 breast-w 99.22 99.19 breast-w 79.24 78.43 78.99 79.08 51.62 58.49 colic.ORIC 49.62 55.82 colic.ORI0 88.05 86.35 86.81 86.79 credit-a credit-a 55.33 61.62 56.56 65.24 credit-g credit-g 78.03 76.36 72.61 74.95 diabetes liabetes 83.97 83.92 84.02 83.80 heart-c heart-c 83.77 83.50 83.74 83.74 heart-h heart-h 88.93 88.64 90.03 88.03 heart-statlog heart-statl 83.02 80.99 hepatitis 78.38 73.19 hepatitis 86.86 86.68 ionosphere 87.89 88.92 ionosphere 98.33 98.29 93.21 92.27 74.65 89.03 66.86 86.39 kr-vs-kp kr-vs-kp 96.59|96.7287.76 85.18 86.09 93.08 82.98 92.57 letter letter 98.04 98.81 97.89 98.75 mushroon mushroon 90.86 95.24 87.93 94.82 v segment segmen 87.74 93.83 91.51 93.96 sick 59.59 62.93 58.64 62.21 sonat sonat 94.40 96.23 88.65 94.87 splice 59.56 67.09 59.63 66.95 rehicle vehicl 96.31 96.52 96.31 96.46 57.65 64.49 1 57.97 66.44 vowel vowel waveform-5000 88.85 90.96 v waveform-5000 84.22 89.54 v 78.56 81.74 80.57 83.12 9/17/08/18/0 w/t/l



Potential real-word applications

- Text mining (e.g., web page classification), Natural language processing (e.g., information extraction)
- Online data stream learning
- Bioinformatics, such as gene expression data analysis
- Medical application (improving computer-aided diagnosis using undiagnosed samples)
- Computer Vision (such as face detection, object detection and object tracking)
- Image classification, traffic classification, intrusion detection, spam detection, and more