

Unsupervised Domain Adaptation Based on Source-guided Discrepancy



Seiichi Kuroki^{1,2} Nontawat Charoenphakdee^{1,2} Han Bao^{1,2}
Junya Honda^{1,2} Issei Sato^{1,2} Masashi Sugiyama^{2,1}
1: The University of Tokyo 2: RIKEN AIP



Abstract

Proposed a new discrepancy measure for unsupervised domain adaptation.

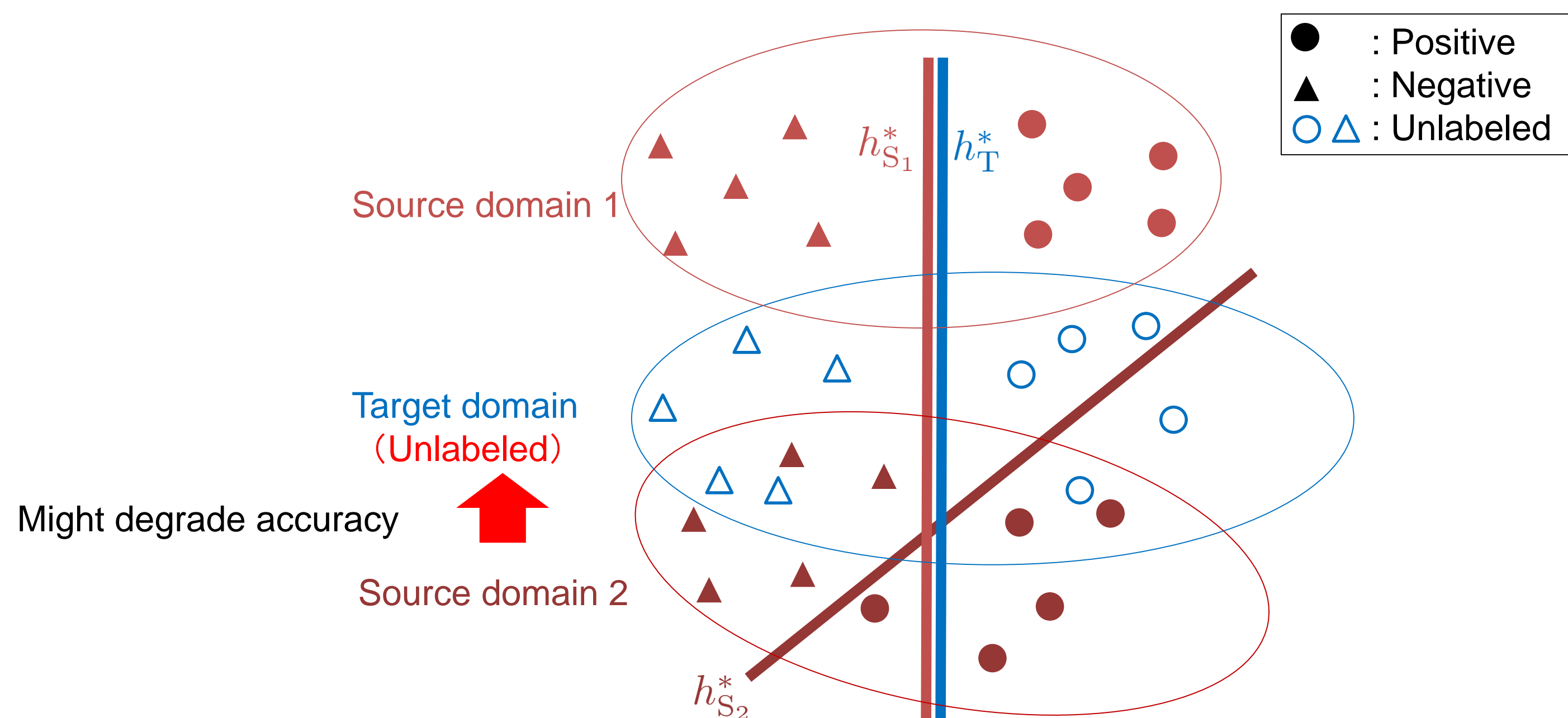
- Exploits all available information including labels in the source domain unlike existing discrepancy measures.
- Has a tighter generalization error bound and is computationally efficient.

Unsupervised Domain Adaptation

In unsupervised domain adaptation, we consider the following setting.

- Distributions in the source and target domains are related but not the same.
 - **Labeled** data in the source domain and **unlabeled** data in the target domain.
- Ex. Natural language processing, speech recognition, computer vision

How can we utilize labeled data in the source domain?



To select a good source, labels in the source domain might be useful.

Related Work

\mathcal{X} -disc [Mansour et al., 2009]

The difference between expected losses of the two domains for the worst pair of hypotheses:

$$\text{disc}_{\mathcal{H}}^{\ell}(P_T, P_S) = \sup_{h, h' \in \mathcal{H}} |R_T^{\ell}(h, h') - R_S^{\ell}(h, h')|$$

- The computation of \mathcal{X} -disc is **intractable**.

$d_{\mathcal{H}}$ [Ben-David et al., 2010]

A **computationally efficient** proxy of \mathcal{X} -disc:

$$d_{\mathcal{H}}(P_T, P_S) = \sup_{h \in \mathcal{H}} |R_T^{\ell_{01}}(h, 1) - R_S^{\ell_{01}}(h, 1)|$$

	Using source labels	Generalization error bound	Computation complexity
\mathcal{X} -disc	No	Loose	High
$d_{\mathcal{H}}$	No	N/A	Low
S-disc [Proposed]	Yes	Tight	Low

Proposed Measure: Source-guided Discrepancy (S-disc)

Explicitly use the best hypothesis h_S^* in the source domain:

$$\varsigma_{\mathcal{H}}^{\ell}(P_T, P_S) = \sup_{h \in \mathcal{H}} |R_T^{\ell}(h, h_S^*) - R_S^{\ell}(h, h_S^*)|$$

1. Estimation of h_S^* requires only labeled data in the source domain.
→ Can be estimated from samples.
2. No need to consider a pair of hypotheses.
→ **Computationally efficient**.
3. The following inequality holds:
 $|R_T^{\ell}(h, h_S^*) - R_S^{\ell}(h, h_S^*)| \leq \varsigma_{\mathcal{H}}^{\ell}(P_T, P_S) \leq \text{disc}_{\mathcal{H}}^{\ell}(P_T, P_S)$
→ Can give a **tighter bound** than \mathcal{X} -disc.

Generalization Error Bound

Theorem 1 If ℓ obeys the triangular inequality,

$$R_T^{\ell}(h, f_T) - R_T^{\ell}(h_S^*, f_T) \leq \varsigma_{\mathcal{H}}^{\ell}(P_T, P_S) + R_S^{\ell}(h, h_S^*) + R_T^{\ell}(h_S^*, h_T^*).$$

Regret arising from using h instead of h_S^*

Estimation error in the source domain

Gap of the best classifiers: Uncontrollable in this problem

→ The Lower **S-disc**, the better generalization.

Estimation Algorithm

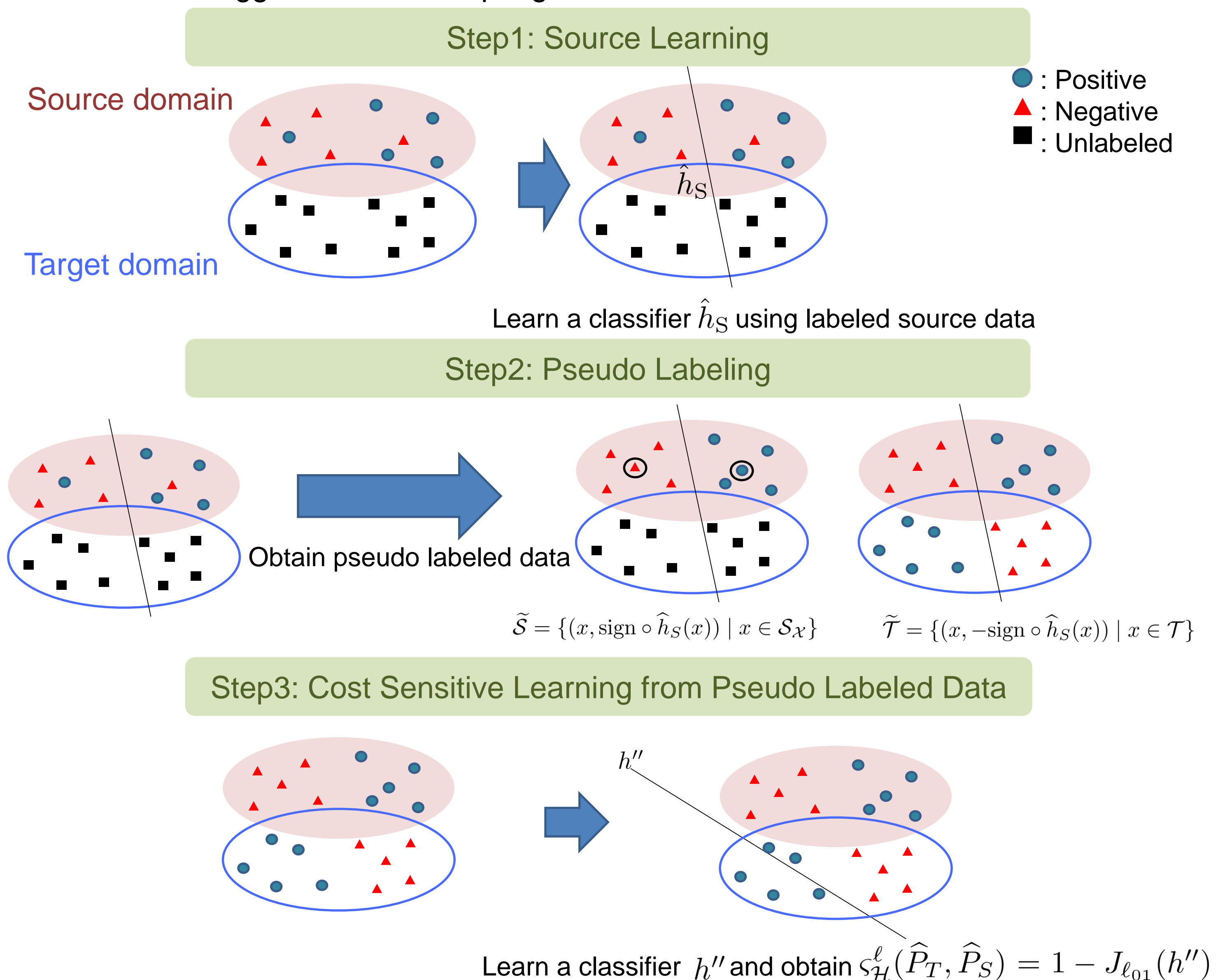
S-disc estimation can be reduced to a cost-sensitive classification for the 0-1 loss:

Theorem 2 For a symmetric hypothesis space \mathcal{H} ($h \in \mathcal{H}$ implies $-h \in \mathcal{H}$),

$$\varsigma_{\mathcal{H}}^{\ell_{01}}(\hat{P}_T, \hat{P}_S) = 1 - \min_{h \in \mathcal{H}} \left[\frac{1}{n_S} \sum_{j=1}^{n_S} \ell(h(x_j^S), h_S^*(x_j^S)) + \frac{1}{n_T} \sum_{i=1}^{n_T} \ell(h(x_i^T), -h_S^*(x_i^T)) \right].$$

Can be expressed as the binary classification loss.

This theorem suggests a three-step algorithm for S-disc estimation described as follows:



Source Selection Experiments

Toy Dataset

Method: SVM with linear kernel

Data: 200 data points per class for each of two sources S_1, S_2 , and target T

We obtain the following results:

$$\varsigma_{\mathcal{H}}^{\ell}(\hat{P}_T, \hat{P}_{S_1}) = 0.27, \quad \varsigma_{\mathcal{H}}^{\ell}(\hat{P}_T, \hat{P}_{S_2}) = 0.49$$

$$d_{\mathcal{H}}(\hat{P}_T, \hat{P}_{S_1}) = 0.69, \quad d_{\mathcal{H}}(\hat{P}_T, \hat{P}_{S_2}) = 0.49.$$

→ S-disc regards S_1 is better while $d_{\mathcal{H}}$ regards S_2 is better.

S-disc is the better discrepancy to measure the quality of sources.

Benchmark Dataset

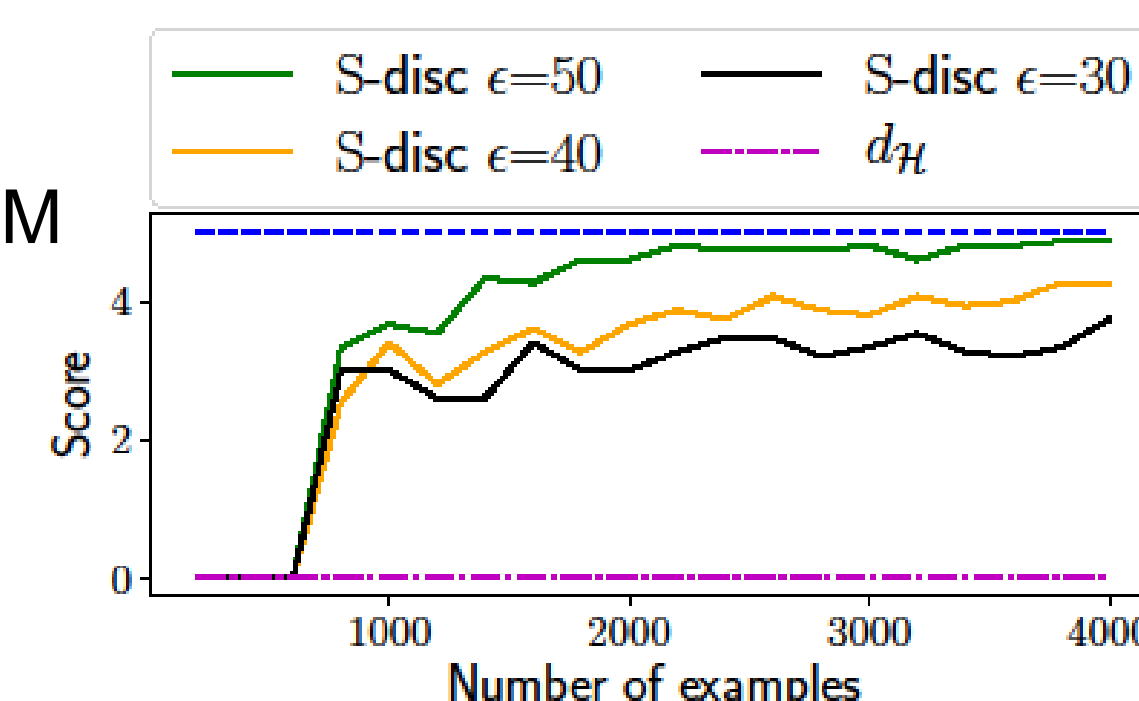
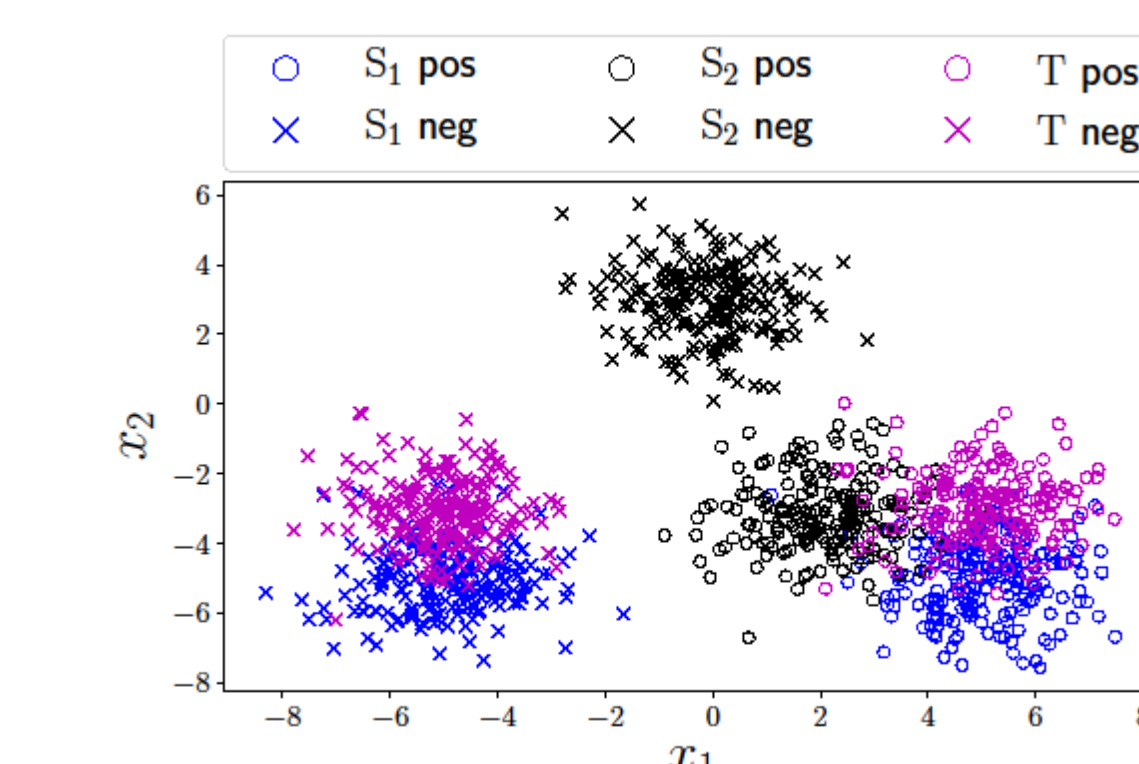
Method : Logistic Regression

Target : MNIST

Sources: Five Clean MNIST-M and Five Noisy MNIST-M

Task: Classify between even and odd digits

Score = # of clean sources from top 5 sources chosen by each discrepancy measure



1. S-disc achieved a better performance as the number of examples increases.
2. $d_{\mathcal{H}}$ cannot distinguish between noisy and clean sources.

References

- [1] S. Ben-David, J. Blitzer, K. Crammer, A. Kulesza, F. Pereira, and J. W. Vaughan. A theory of learning from different domains. Machine Learning, 79(1-2):151-175, 2010.
- [2] Y. Mansour, M. Mohri, and A. Rostamizadeh. Domain adaptation: Learning bounds and algorithms. In COLT, 2009.
- [3] Mohri, M., and Medina, A. M. 2012. New analysis and algorithm for learning with drifting distributions. In ALT, 124-138.