

# Classification from Pairwise Similarity and Unlabeled Data

Han Bao<sup>†‡</sup> Gang Niu<sup>‡</sup> Masashi Sugiyama<sup>†‡</sup>

<sup>‡</sup>The University of Tokyo  
RIKEN AIP



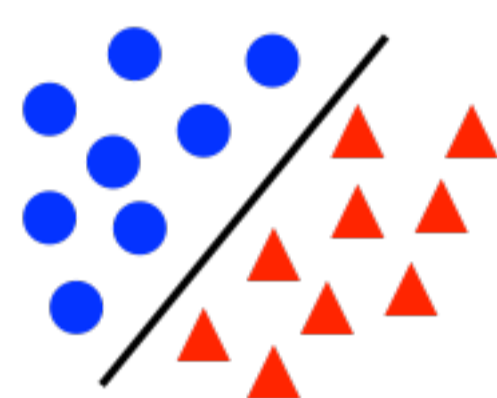
## Abstract

Supervised classification requires vast amount of labeled data, which is costly. In order to mitigate this bottleneck, we proposed a classification problem where only pairwise similarity (two examples belong to the same class) and unlabeled data points are needed.

## Introduction

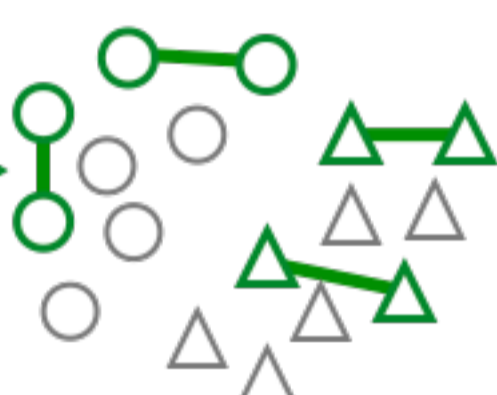
### Pairwise Data in Classification

Supervised classification:  
Explicit labels might be difficult to obtain...



Instead,

two people share the same property



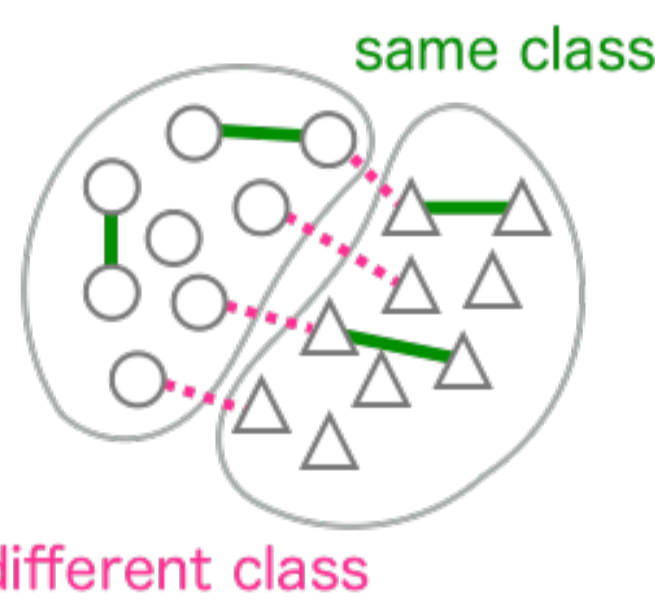
Example: classification of sensitive matters (e.g., politics, religion, opinion on racial issue)  
⇒ "Which person do you share the same belief as?"

### Related Work:

#### Semi-supervised Clustering<sup>[1]</sup>

Clustering with

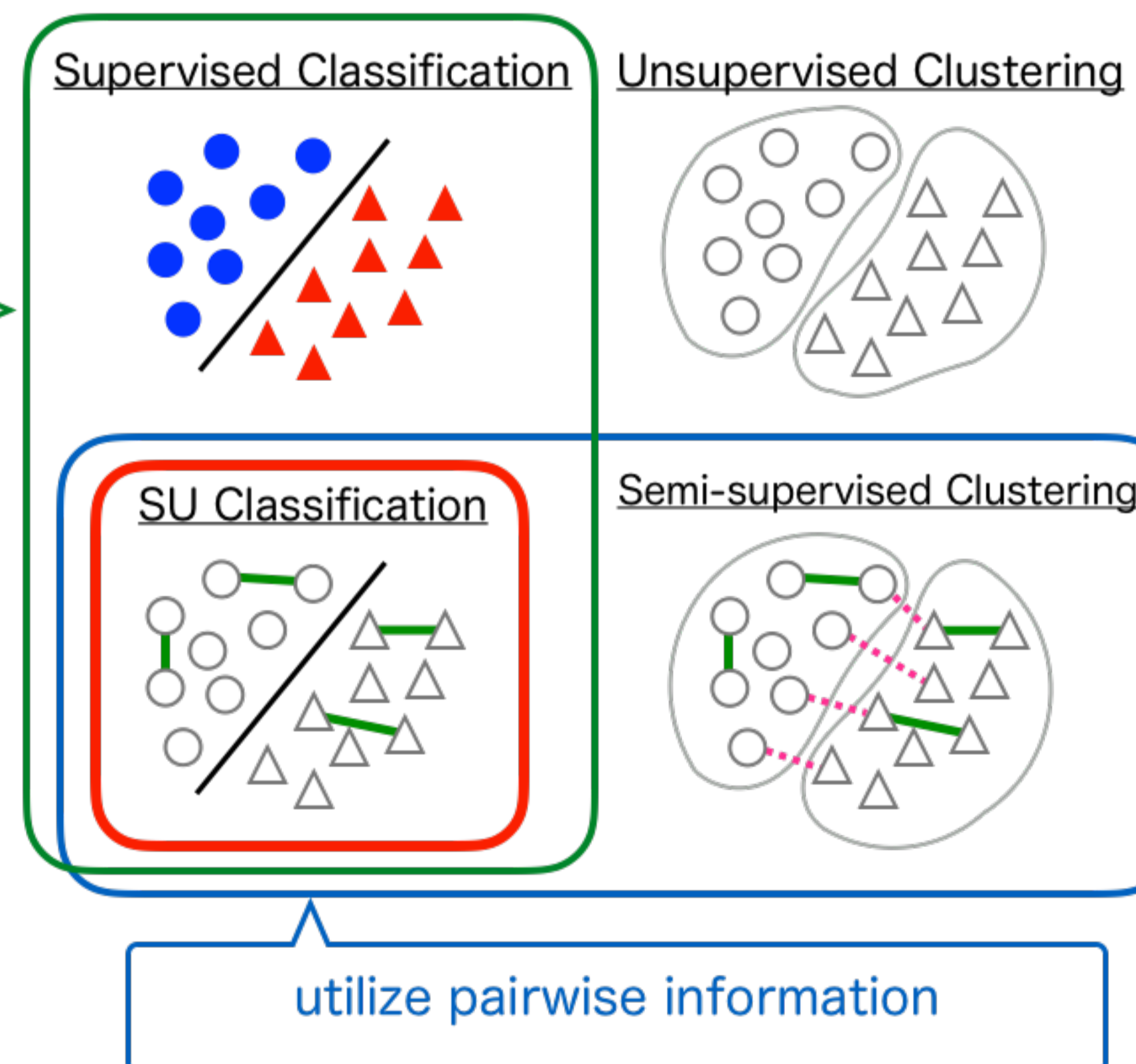
- Similar pair (must-link)
- Dissimilar pair (cannot-link)



### Problem

Strong assumption on datasets  
e.g., cluster assumption, manifold assumption

### Goal of This Work



★ Classify data using pairwise information w/o strong assumption

### Empirical Risk Minimization

Supervised classification = Minimize misclassification rate

Classification Risk  $R_{PN}(f) = \mathbb{E}[\ell_{0-1}(yf(\mathbf{x}))]$

↑ unbiased

Empirical Risk  $\hat{R}_{PN}(f) = \frac{1}{n} \sum_{i=1}^n \ell_{0-1}(y_i f(\mathbf{x}_i))$

Risk of labeled data

### Similar and Unlabeled Data

#### S(imilar) Data Pairs

$$\{(\mathbf{x}_{S,i}, \mathbf{x}'_{S,i})\}_{i=1}^{n_S} \sim p(\mathbf{x}, \mathbf{x}' | y = y' = +1 \vee y = y' = -1)$$

#### U(nlabeled) Data Points

$$\{\mathbf{x}_{U,i}\}_{i=1}^{n_U} \sim p(\mathbf{x})$$

## Unbiased Risk Estimator

$R_{PN,\ell}(f) = \mathbb{E}[\ell(yf(\mathbf{x}))]$  misclassification rate of  $f$

↑ unbiased

$$\hat{R}_{SU,\ell}(f) = \frac{\pi_S \sum_{i=1}^{n_S} \mathcal{L}_{S,\ell}(f(\mathbf{x}_{S,i})) + \mathcal{L}_{S,\ell}(f(\mathbf{x}'_{S,i}))}{2} + \frac{1}{n_U} \sum_{i=1}^{n_U} \mathcal{L}_{U,\ell}(f(\mathbf{x}_{U,i}))$$

Risk of similar data  
Risk of unlabeled data

$\pi_+ \triangleq p(y = +1)$   $\ell$ : loss function

$\pi_- \triangleq p(y = -1)$   $\pi_S \triangleq \pi_+^2 + \pi_-^2$

$$\mathcal{L}_{S,\ell}(z) \triangleq \frac{\ell(z) - \ell(-z)}{2\pi_+ - 1} \quad \mathcal{L}_{U,\ell}(z) \triangleq \frac{-\pi_- \ell(z) + \pi_+ \ell(-z)}{2\pi_+ - 1}$$

★ Empirical risk can be minimized without explicitly labeled data

Prop. 1: Error  $\rightarrow 0$  asymptotically

### Estimation Error Bound

$$R(\hat{f}) - R(f^*) = \mathcal{O}_p\left(\frac{1}{\sqrt{2n_S}} + \frac{1}{\sqrt{n_U}}\right)$$

#S data      #U data

Estimation error of risk of empirical minimizer  $\hat{f}$

Def. (Rademacher complexity): For function class  $\mathcal{H}$ ,

$$\mathfrak{R}(\mathcal{H}; n, \mu) \triangleq \mathbb{E}_{Z_1, \dots, Z_n, \sim \mu} \mathbb{E}_\sigma \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i h(Z_i) \right]$$

where  $\sigma = (\sigma_1, \dots, \sigma_n)$  are Rademacher variables

Assumption:

$\mathcal{F}$  satisfies the following assumption:

$$\exists R \in \mathbb{R}_{>0} \text{ s.t. } \|f\|_\infty \leq R \quad (\forall f \in \mathcal{F})$$

$$\mathfrak{R}(\mathcal{F}; n, \mu) \leq \frac{C_{\mathcal{F}}}{\sqrt{n}} \quad (\forall \mu) \quad (C_{\mathcal{F}}: \text{constant})$$

where  $\rho$  is Lipschitz constant of  $\ell$

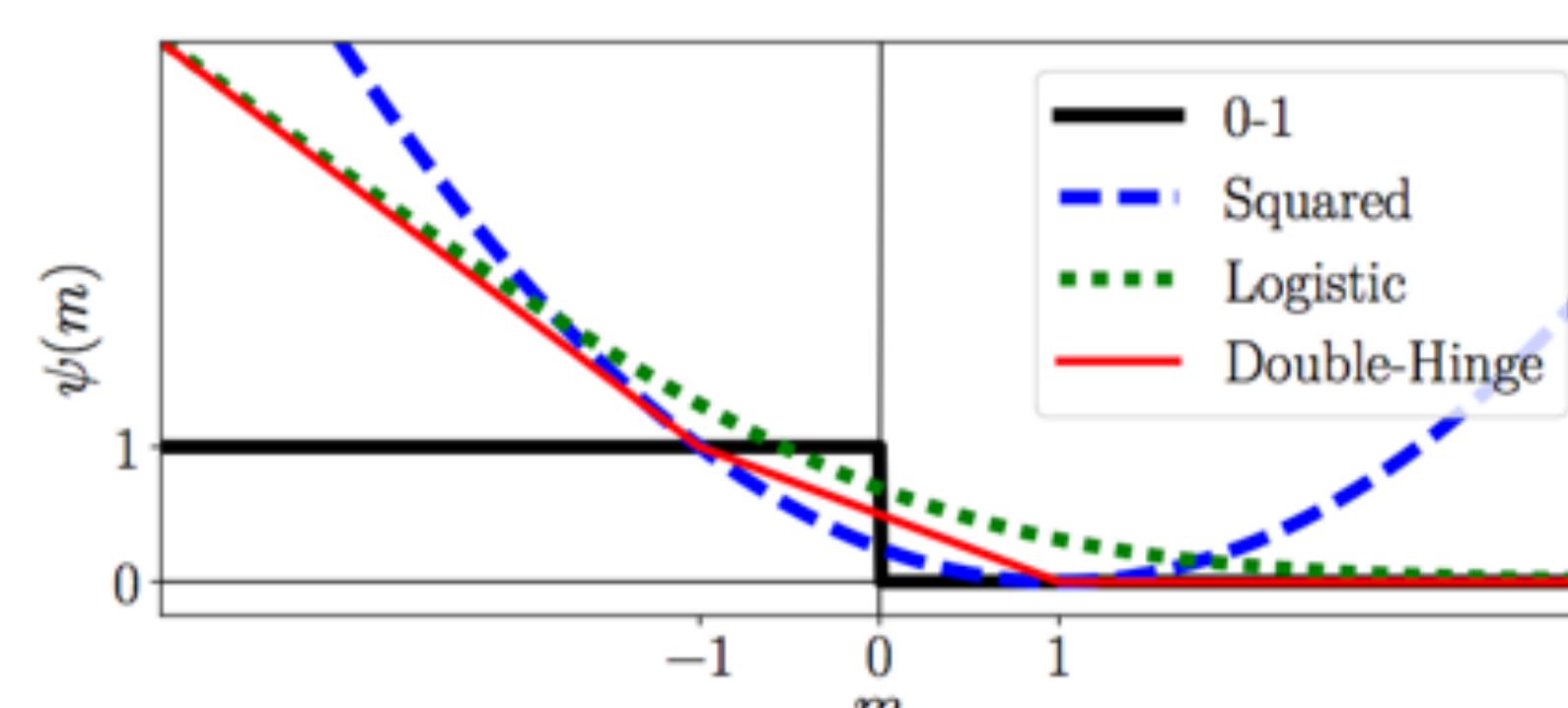
★ Estimation error converges to 0 in the optimal rate<sup>[2]</sup> w/o strong assumption

Prop. 2: Reduced to convex optimization

**Theorem** If  $\ell$  satisfies  $\ell(z) - \ell(-z) = -z$  then the optimization of  $\hat{R}_{SU,\ell}(f)$  with the classifier  $f(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) + b$

and the L2 regularization becomes convex.

**Example** Squared loss, Double-hinge loss



★ Computationally efficient to optimize / Unique global optimum

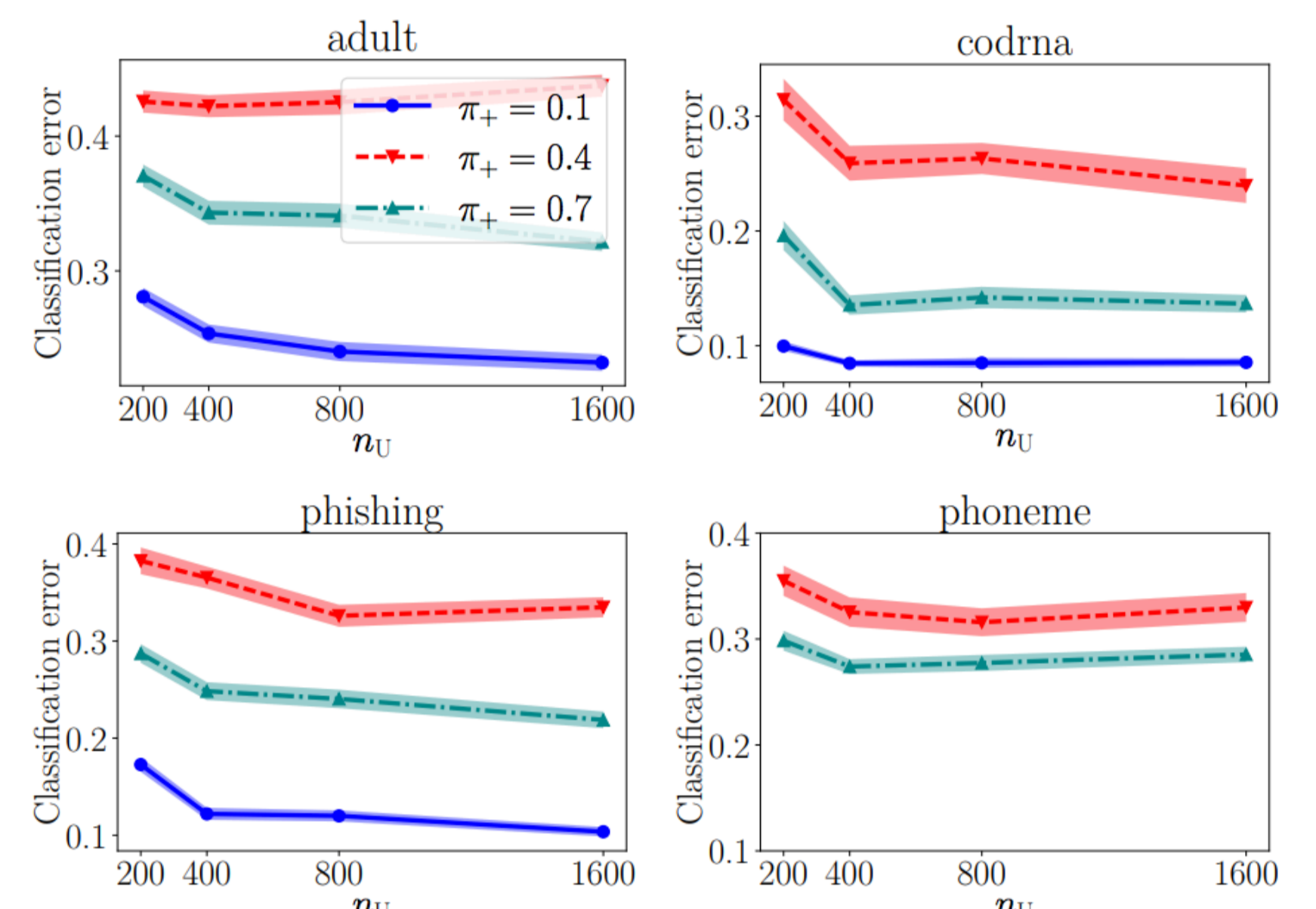
## Experiments

### Relationship between # of data and error

#### Setting

- linear-in-input model with squared loss
- regularization parameter  $\in \{10^{-1}, 10^{-4}, 10^{-7}\}$  (chosen by cross-validation)
- # of S data = 200
- # of U data  $\in \{200, 400, 800, 1600\}$
- true class-prior  $\in \{0.1, 0.4, 0.7\}$

#### Result



### Comparison with Baselines

#### Setting

- loss: squared / double-hinge loss
- # of S data = # of U data = 500
- true class-prior = 0.7
- (model, regularization are the same with the above)

#### Baselines

- SERAPH<sup>[3]</sup> clustering with semi-supervised metric learning
- 3SMIC<sup>[4]</sup> semi-supervised clustering with information maximization
- DirtyIMC<sup>[5]</sup> clustering by matrix completion

#### Results

Dataset	SU(proposed)		Baselines		
	Squared	Double-hinge	SERAPH	3SMIC	DMC
adult	66.3 (1.2)	<b>84.5 (0.8)</b>	66.5 (1.7)	58.5 (1.3)	63.7 (1.2)
banana	64.1 (1.7)	<b>68.2 (1.2)</b>	55.0 (1.1)	61.9 (1.2)	64.3 (1.0)
cod-rna	<b>82.5 (1.1)</b>	71.0 (0.9)	62.5 (1.4)	56.6 (1.2)	63.8 (1.1)
higgs	54.9 (1.6)	<b>69.3 (0.9)</b>	63.4 (1.1)	57.0 (0.9)	65.0 (1.1)
ijcnn1	68.2 (1.3)	<b>73.6 (0.9)</b>	59.8 (1.2)	58.9 (1.3)	66.2 (2.2)
magic	<b>65.9 (1.5)</b>	<b>69.0 (1.3)</b>	55.0 (0.9)	59.1 (1.4)	63.1 (1.1)
phishing	75.2 (1.3)	<b>91.3 (0.6)</b>	62.4 (1.1)	60.1 (1.3)	64.8 (1.4)
phoneme	<b>68.0 (1.4)</b>	<b>70.8 (1.0)</b>	<b>69.1 (1.4)</b>	61.3 (1.1)	64.5 (1.2)
spambase	69.5 (1.3)	<b>85.5 (0.5)</b>	65.4 (1.8)	61.5 (1.3)	63.6 (1.3)
susy	60.7 (1.0)	<b>74.8 (1.2)</b>	58.0 (1.0)	57.1 (1.2)	65.2 (1.0)
w8a	60.5 (1.2)	<b>86.5 (0.6)</b>	N/A	60.5 (1.5)	65.0 (2.0)
waveform	78.6 (1.6)	<b>87.0 (0.5)</b>	56.5 (0.9)	56.5 (0.9)	65.0 (0.9)

Classification accuracies for each dataset (with standard errors). Bold-faces are outperforming methods chosen by 5% t-test.

## References

- [1] Klein, D., Kamvar, S. D., and Manning, C. D. From instance-level constraints to space-level constraints: Making the most of prior knowledge in data clustering. In *ICML*, 2002.
- [2] Mendelson, S. Lower bounds for the empirical minimization algorithm. *IEEE Transactions on Information Theory*, 2008.
- [3] Niu, G., Dai, B., Yamada, M., and Sugiyama, M. Information-theoretic semi-supervised metric learning via entropy regularization. In *ICML*, 2012.
- [4] Calandriello, D., Niu, G., and Sugiyama, M. Semi-supervised information-maximization clustering. *Neural Networks*, 2014.
- [5] Chiang, K.-Y., Hsieh, C.-J., and Dhillon, I. S. Matrix completion with noisy side information. In *NIPS*, 2015.