

# Classification from Pairwise Similarity and Unlabeled Data

Han Bao<sup>1,2</sup>, Gang Niu<sup>2</sup>, Masashi Sugiyama<sup>2,1</sup>

<sup>1</sup>The University of Tokyo, Japan / <sup>2</sup>RIKEN, Japan

July 13<sup>th</sup>, 2018



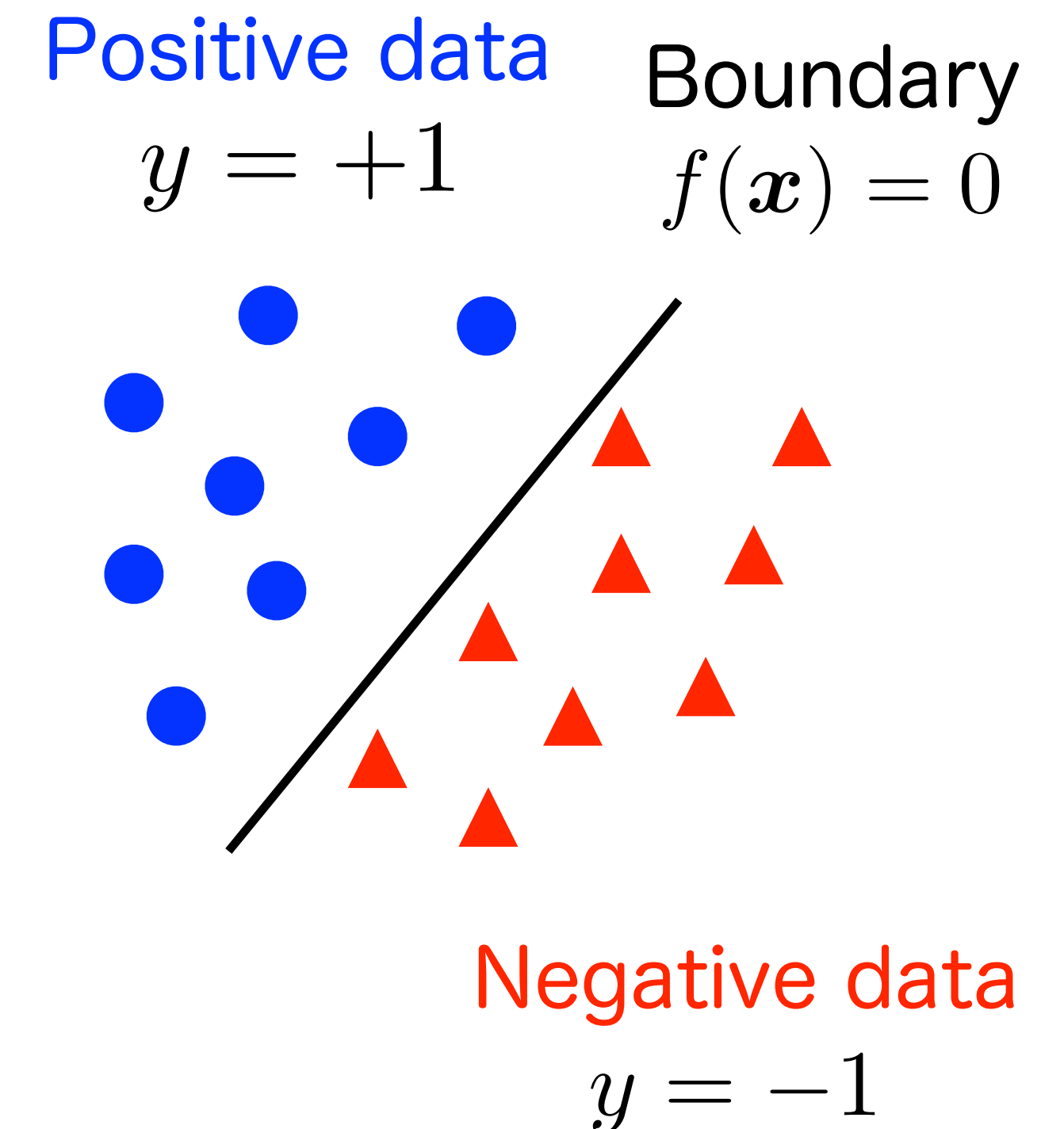
東京大学  
THE UNIVERSITY OF TOKYO



RIKEN

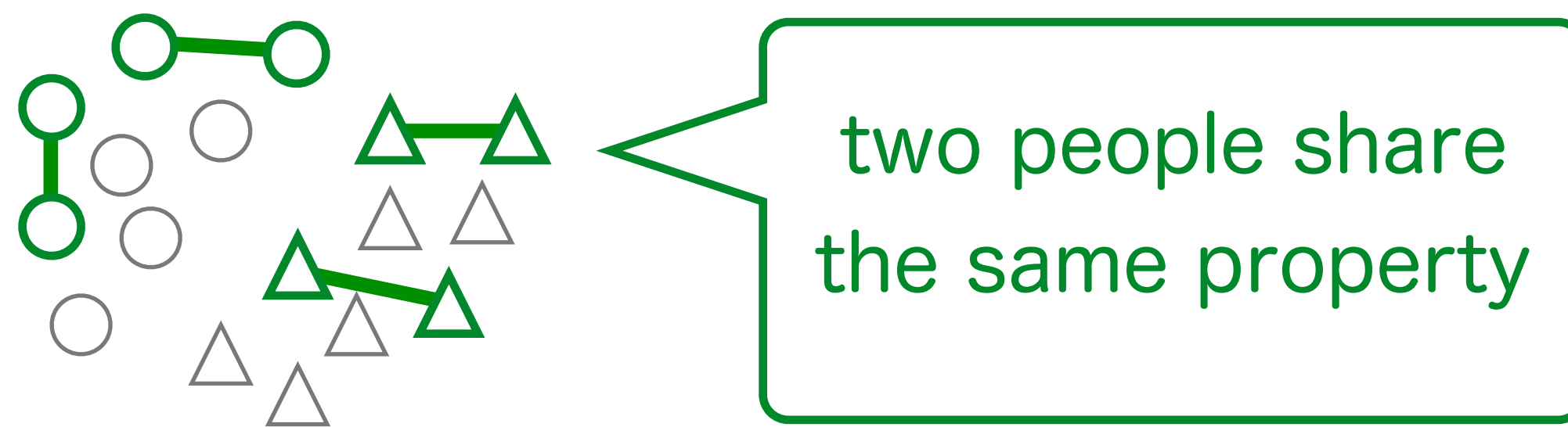
# Gentle Start: Binary Classification

- Training data  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}, y)$   
where data  $\mathbf{x}_i \in \mathbb{R}^d$  is labeled as  $y_i \in \{+1, -1\}$
- Goal: find a classifier  $f : \mathbb{R}^d \rightarrow \mathbb{R}$
- Method: minimize classification error
  - ▶ empirical risk minimization (ERM)



# Motivation: Pairwise Information in Classification

- Classification of sensitive matters
  - ▶ e.g., politics, religion, opinion on racial issue
  - ▶ hard to obtain explicit label
  - ▶ instead asking “Which person do you share the same belief as?”
  - ▶ cf. randomized response technique [Warner 1965]



<http://leanintokyo.org/wp-content/uploads/2017/12/MeToo.jpg>

# Related: Semi-supervised Clustering

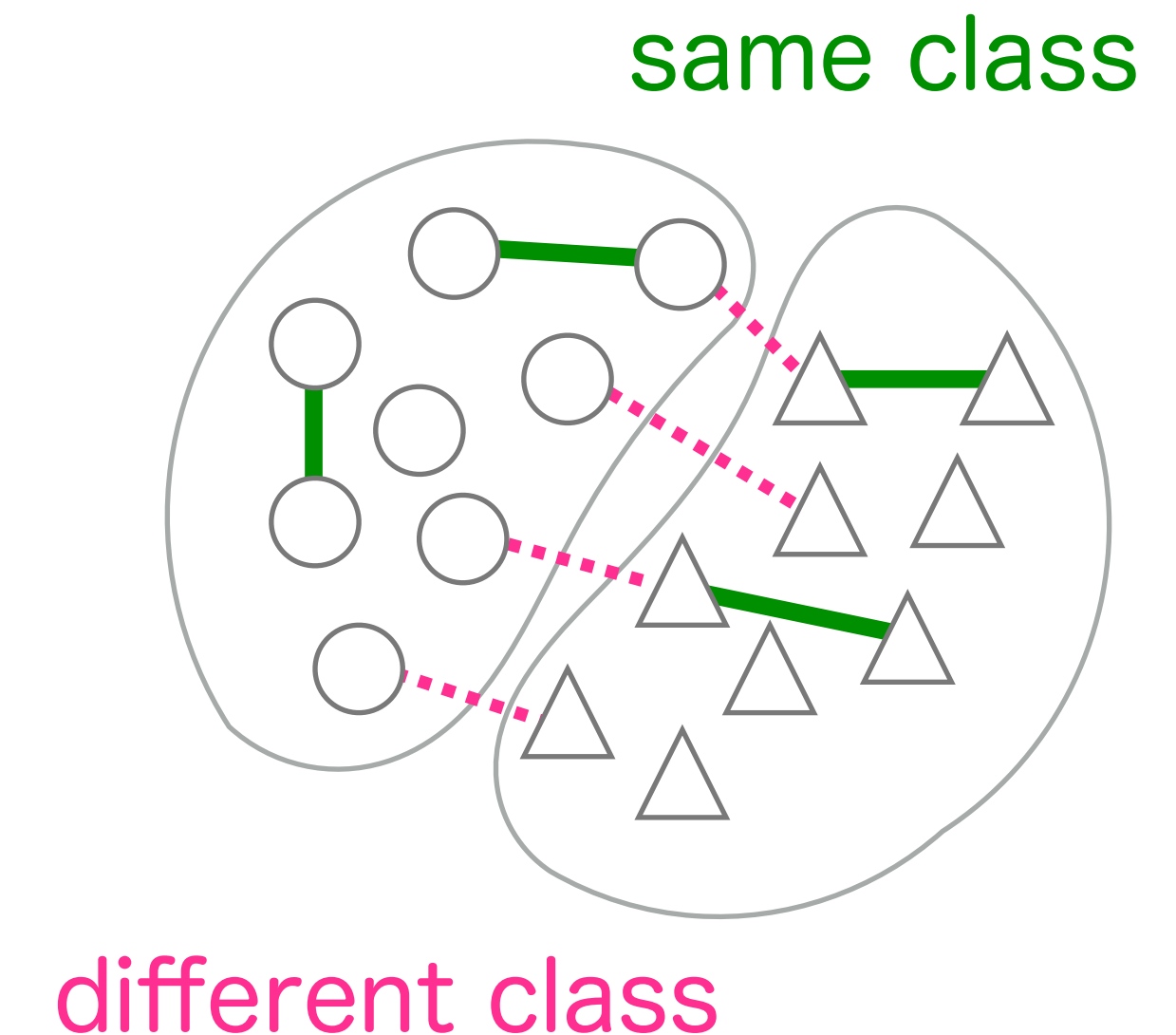
[Wagstaff+ ICML2001; many other papers]

## ■ Clustering from

- ▶ unlabeled  $\mathcal{U} = \{\mathbf{x}_i\}$
- ▶ similar  $\mathcal{S} = \{(\mathbf{x}_i, \mathbf{x}'_i)\}$
- ▶ dissimilar  $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{x}'_i)\}$

## ■ Offspring of unsupervised clustering

- ▶ **Problem: Cluster assumption**  
(manifold assumption, low-density separation)
- ▶ does not hold for many datasets



# Our work: SU Classification

## Binary classification from

▶ S(imilar) data

$$\mathcal{S} = \{(\underline{x_i}, x'_i)\}$$

▶ U(nlabeled) data

$$\mathcal{U} = \{x_i\}$$

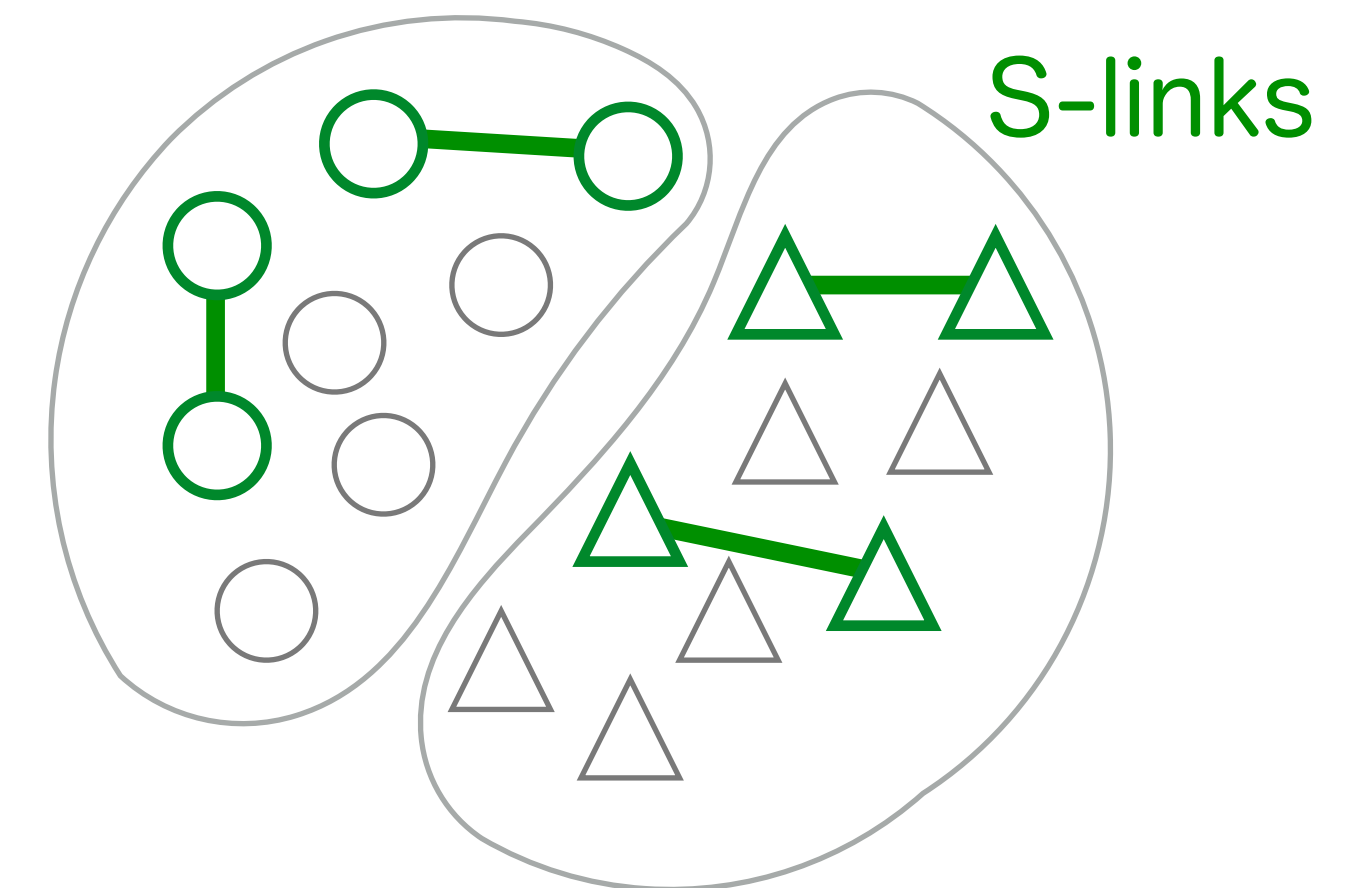
belong to the  
same class

## Formulation based on classification error is available

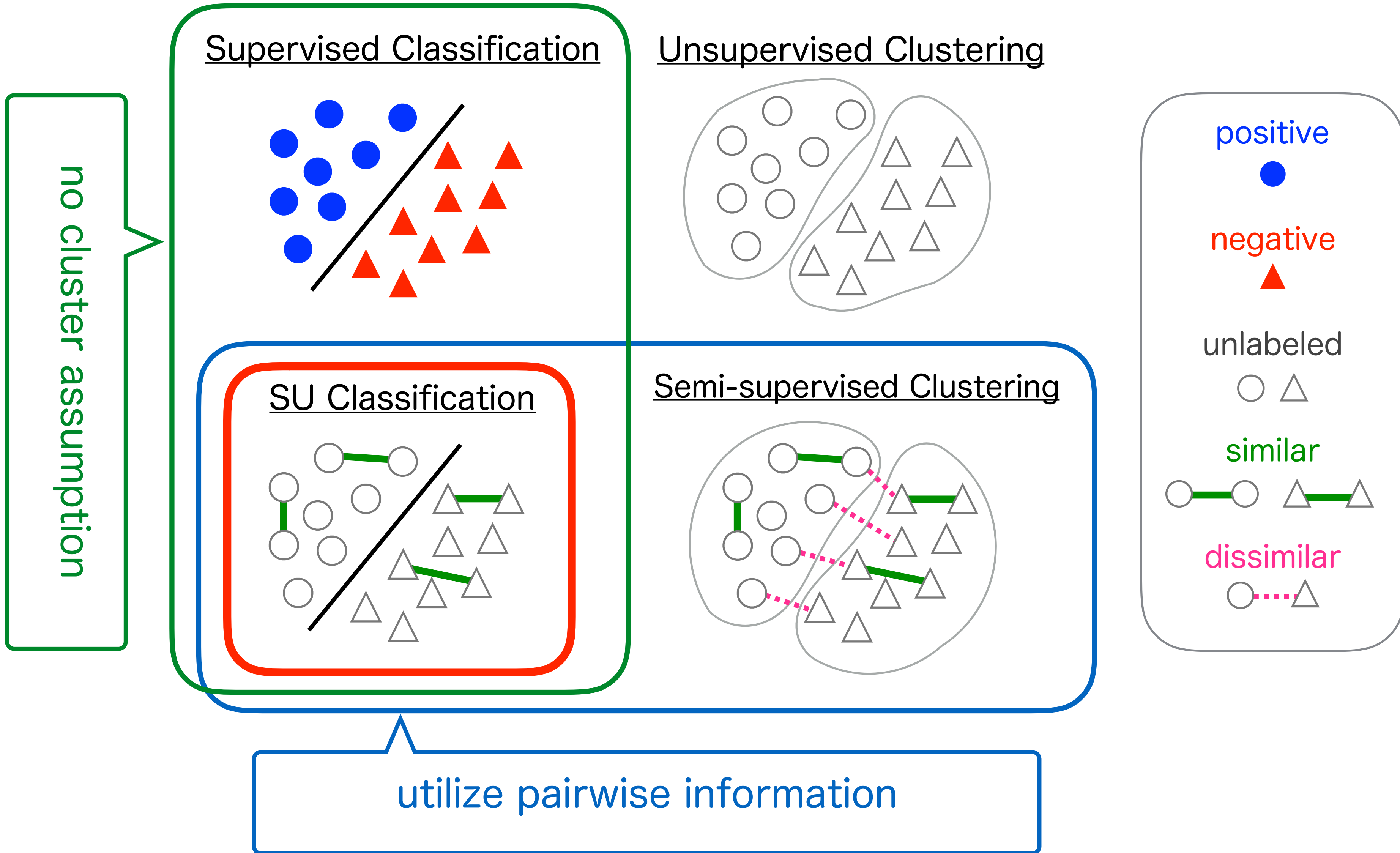
## Situation: prediction of sensitive matters

▶ e.g., politics, religion

▶ “Which person do you share the same belief as?”



# Summary



# Empirical Risk Minimization

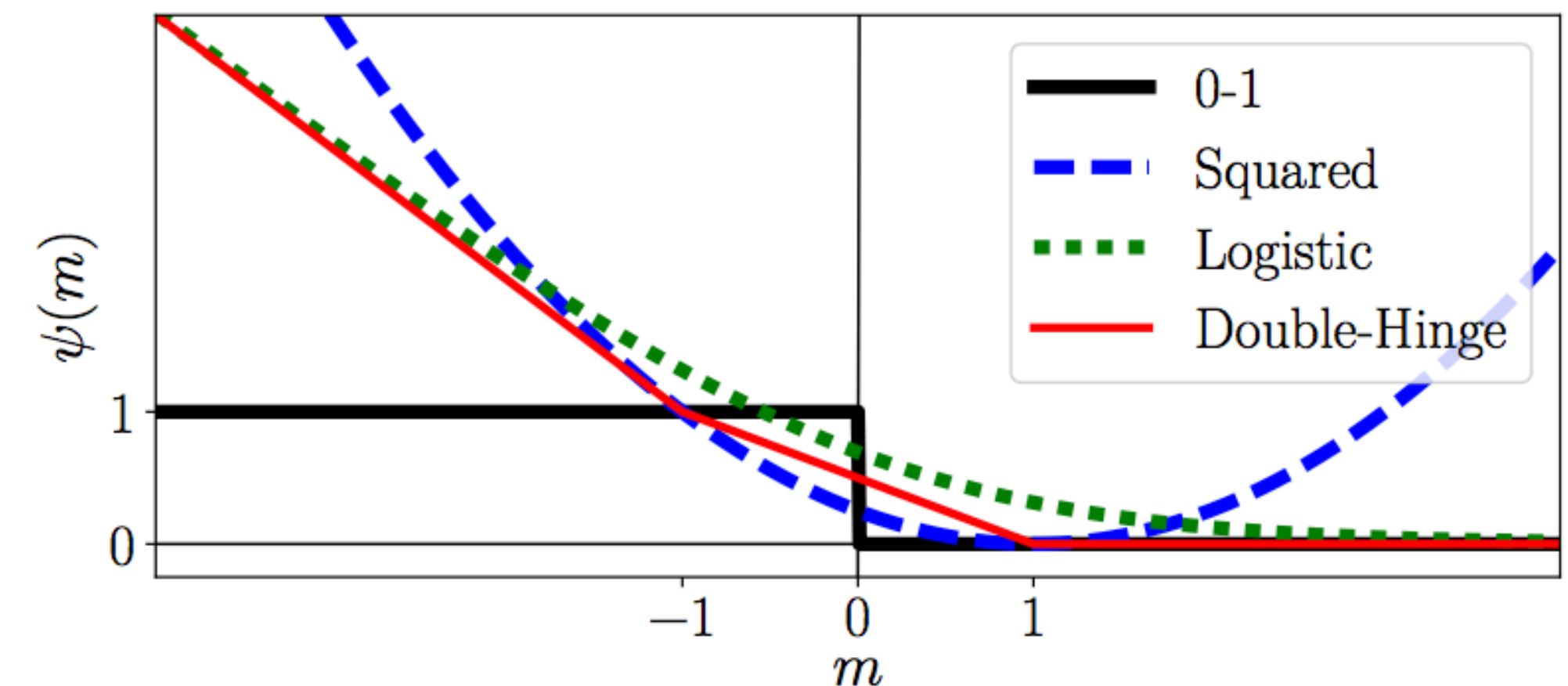
- Goal: Minimize classification risk (= error rate)

Classification risk  $R_{\text{PN}}(f) = \mathbb{E}[\ell(yf(\mathbf{x}))]$   
 expectation over  $p(\mathbf{x}, y)$

Empirical risk  $\hat{R}_{\text{PN}}(f) = \frac{1}{n} \sum_{i=1}^n \ell(y_i f(\mathbf{x}_i))$

- Margin loss  $\ell$ : give small/large penalty for positive/negative value of  $yf(\mathbf{x})$

- ▶  $\text{sign}(y) = \text{sign}(f(\mathbf{x})) \Leftrightarrow$  correct
- ▶  $\text{sign}(y) \neq \text{sign}(f(\mathbf{x})) \Leftrightarrow$  incorrect



example of loss functions

# Assumption: Data Generating Process

- Pairwise Similarity (S): belong to the same class

$$\begin{aligned} \{(\mathbf{x}_{S,i}, \mathbf{x}'_{S,i})\}_{i=1}^{n_S} &\sim p_S(\mathbf{x}, \mathbf{x}') = p(\mathbf{x}, \mathbf{x}' \mid y = y' = +1 \vee y = y' = -1) \\ &= \frac{\pi_+^2 p_+(\mathbf{x})p_+(\mathbf{x}') + \pi_-^2 p_-(\mathbf{x})p_-(\mathbf{x}')}{\pi_+^2 + \pi_-^2} \end{aligned}$$

- Unlabeled Data (U)

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

conditional density

$$p_+(\mathbf{x}) \triangleq p(\mathbf{x} \mid y = +1)$$

$$p_-(\mathbf{x}) \triangleq p(\mathbf{x} \mid y = -1)$$

class prior

$$\pi_+ \triangleq p(y = +1)$$

$$\pi_- \triangleq p(y = -1)$$



# Classification risk can be estimated from SU data

## ■ Theorem:

$$\hat{R}_{\text{SU},\ell}(f) = \underbrace{\frac{\pi_{\text{S}}}{n_{\text{S}}} \sum_{i=1}^{n_{\text{S}}} \frac{\mathcal{L}_{\text{S},\ell}(f(\mathbf{x}_{\text{S},i})) + \mathcal{L}_{\text{S},\ell}(f(\mathbf{x}'_{\text{S},i}))}{2}}_{\text{risk for S data}} + \underbrace{\frac{1}{n_{\text{U}}} \sum_{i=1}^{n_{\text{U}}} \mathcal{L}_{\text{U},\ell}(f(\mathbf{x}_{\text{U},i}))}_{\text{risk for U data}}$$

no explicit labels  
needed

is unbiased to  $R_{\text{PN},\ell}(f) = \mathbb{E}[\ell(yf(\mathbf{x}))]$   
(original classification risk)

- ▶ minimize  $\hat{R}_{\text{SU},\ell}(f) \Rightarrow$  minimize classification risk
- ▶  $\hat{R}_{\text{SU},\ell}(f)$  can be computed only from SU data

N.B.:  $\pi_{\text{S}}$  can be estimated

$$\begin{aligned} \pi_+ &\triangleq p(y = +1) & \ell &: \text{surrogate loss} \\ \pi_- &\triangleq p(y = -1) & \mathcal{L}_{\text{S},\ell}(z) &\triangleq \frac{\ell(z) - \ell(-z)}{2\pi_+ - 1} \\ \pi_{\text{S}} &\triangleq \pi_+^2 + \pi_-^2 & \mathcal{L}_{\text{U},\ell}(z) &\triangleq \frac{-\pi_- \ell(z) + \pi_+ \ell(-z)}{2\pi_+ - 1} \end{aligned}$$

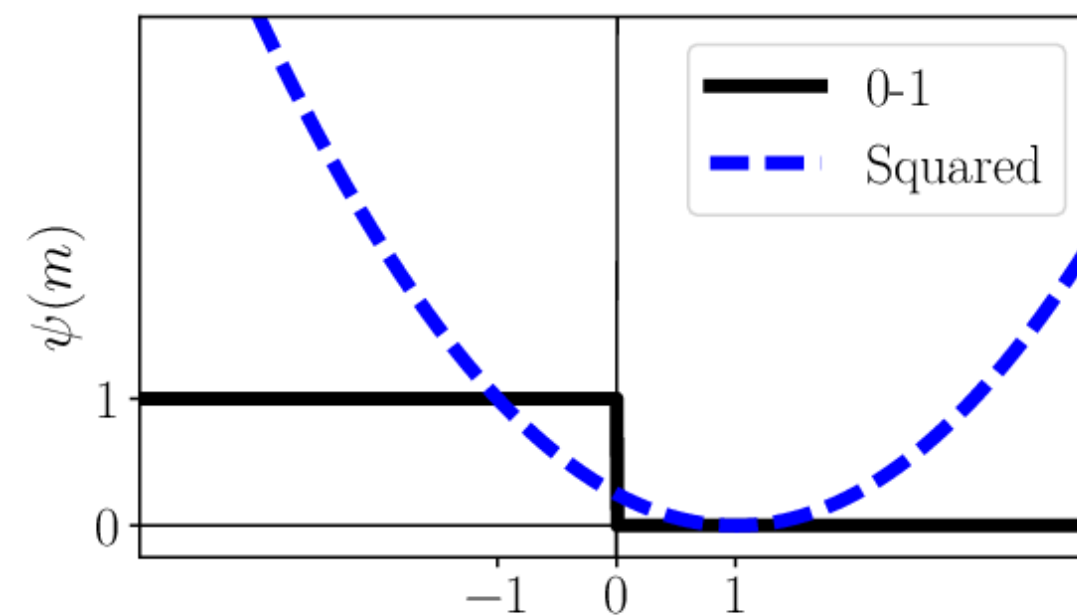
# Good loss function gives convex objective

- Theorem: If the loss function  $\ell$  satisfies

$$\ell(z) - \ell(-z) = -z,$$

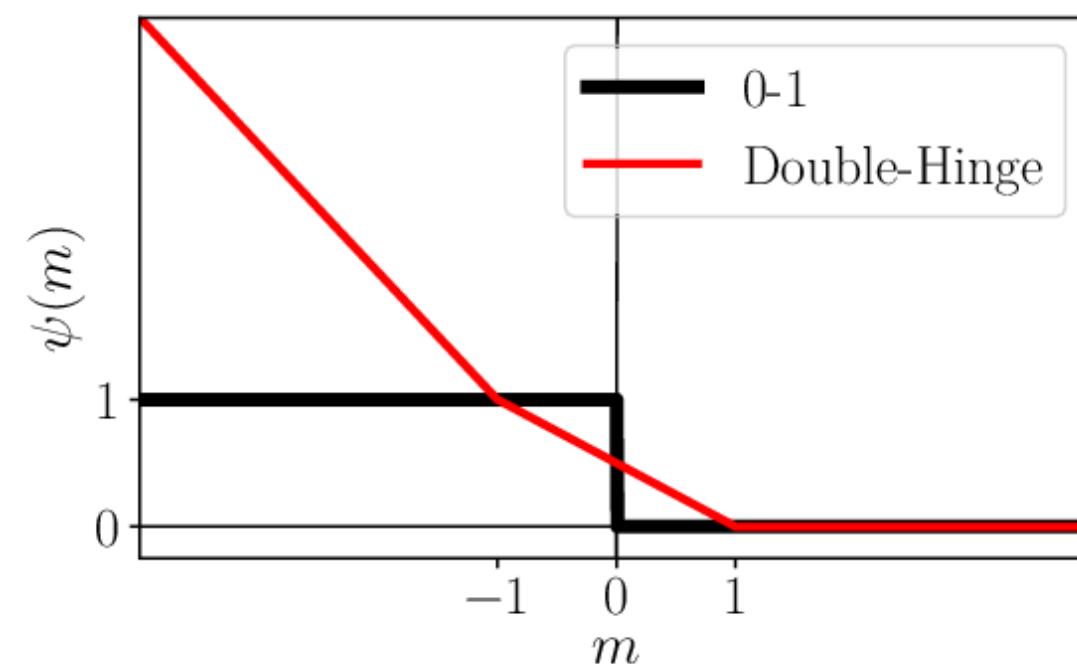
then the objective (w/ model  $f(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x})$  and  $l_2$ -reguralization) is convex.

- Example: squared loss, double-hinge loss, logistic loss



squared loss

⇒ analytical solution (linear system)



double-hinge loss

⇒ quadratic program

computationally efficient  
to obtain solution

# Estimation Error Bound

$$\frac{R_{\text{PN}}(\hat{f}) - R_{\text{PN}}(f^*)}{\text{estimation error of risk of empirical minimizer}} = \mathcal{O}_p \left( \frac{1}{\sqrt{2n_S}} + \frac{1}{\sqrt{n_U}} \right)$$

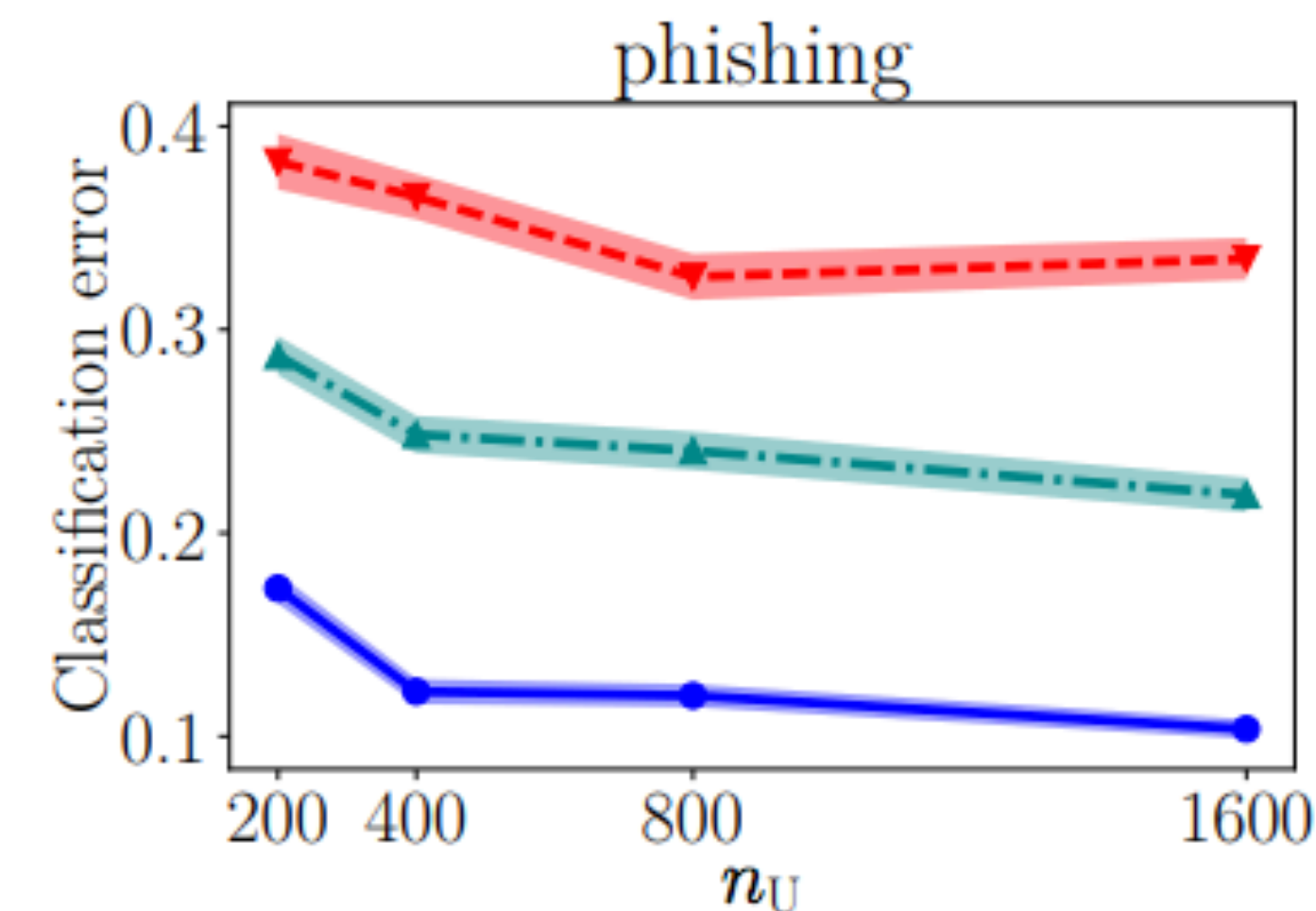
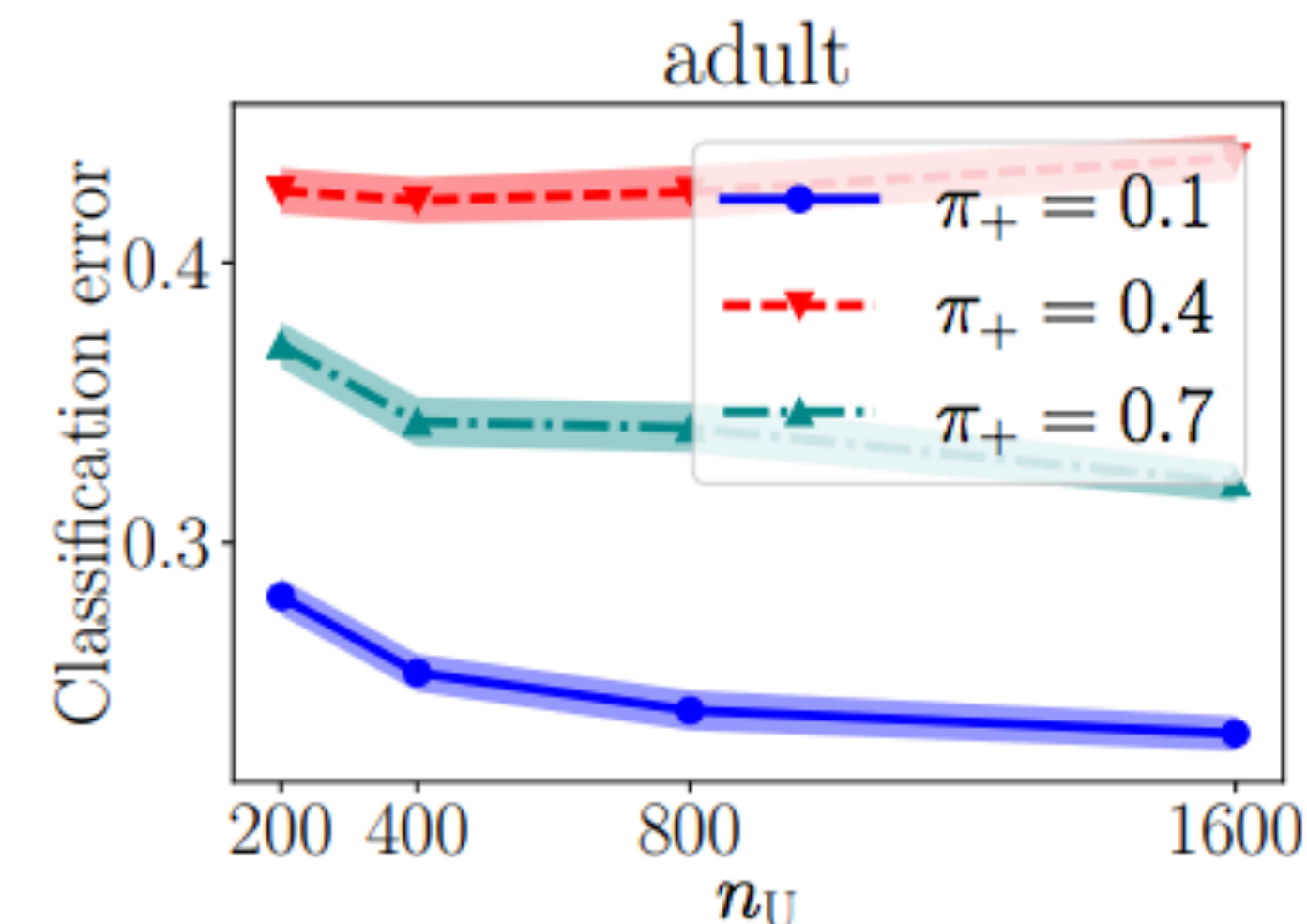
# S data
# U data

- consistency: estimation error  $\rightarrow 0$  as  $n \rightarrow \infty$
- optimal parametric convergence rate [Mendelson 2008]

$$R_{\text{PN}}(f) = \mathbb{E}[\ell(yf(\mathbf{x}))]$$

$$f^* = \operatorname{argmin}_f R_{\text{PN}}(f) \quad \text{true minimizer}$$

$$\hat{f} = \operatorname{argmin}_f \hat{R}_{\text{SU}}(f) \quad \text{empirical minimizer}$$



# Experiments (Benchmark Datasets)

classification accuracies are shown

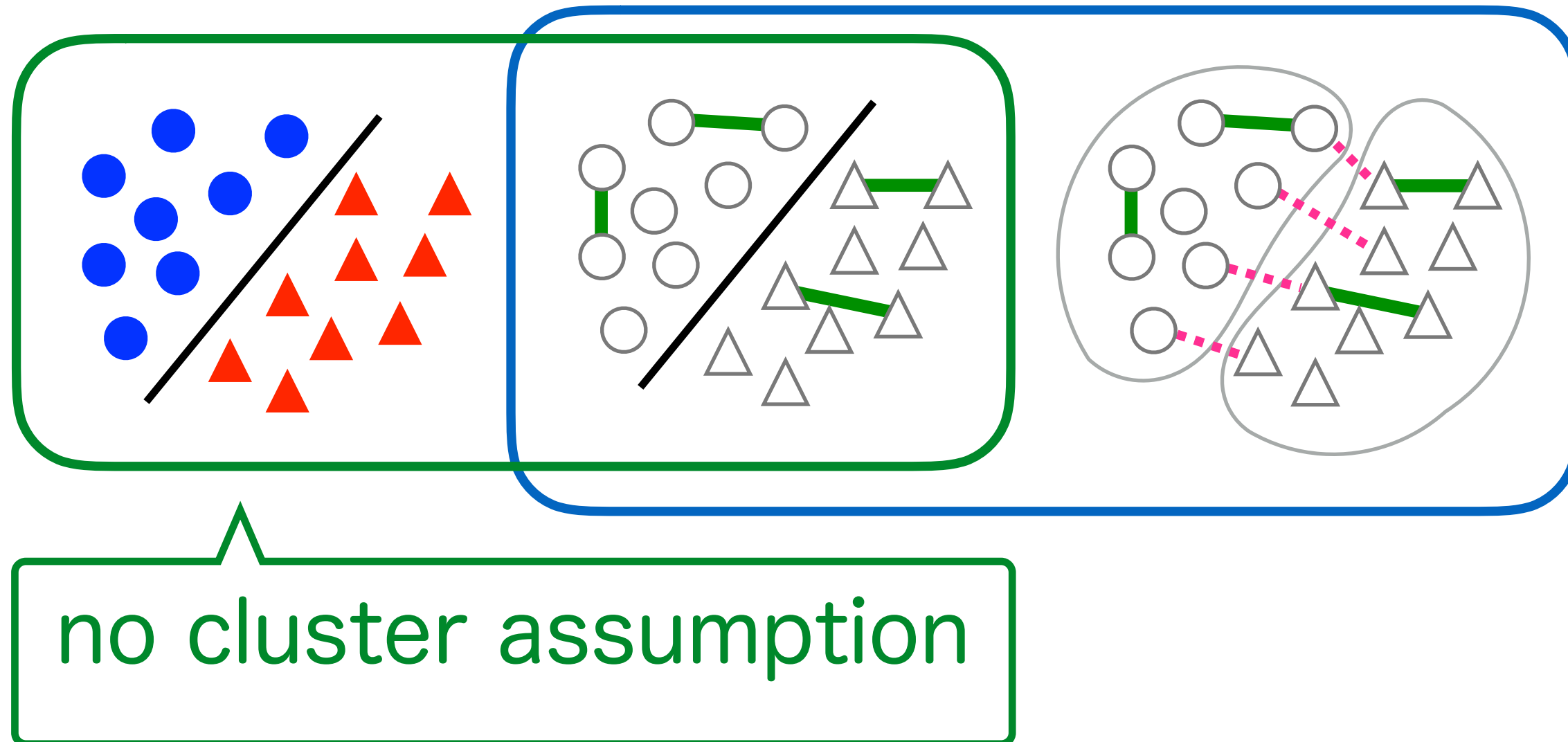
Dataset	Dim	SU(proposed)		Baselines					
		Squared	Double-hinge	KM	ITML	SERAPH	3SMIC	DIMC	IMSAT(linear)
adult	123	64.5 (1.2)	<b>84.5 (0.8)</b>	58.1 (1.1)	57.9 (1.1)	66.5 (1.7)	58.5 (1.3)	63.7 (1.2)	69.8 (0.9)
banana	2	<b>67.5 (1.2)</b>	<b>68.2 (1.2)</b>	54.3 (0.7)	54.8 (0.7)	55.0 (1.1)	61.9 (1.2)	64.3 (1.0)	<b>69.8 (0.9)</b>
cod-rna	8	<b>82.8 (1.3)</b>	71.0 (0.9)	63.1 (1.1)	62.8 (1.0)	62.5 (1.4)	56.6 (1.2)	63.8 (1.1)	69.1 (0.9)
higgs	28	55.1 (1.1)	<b>69.3 (0.9)</b>	<b>66.4 (1.6)</b>	<b>66.6 (1.3)</b>	63.4 (1.1)	57.0 (0.9)	65.0 (1.1)	<b>69.7 (1.4)</b>
ijcnn1	22	65.5 (1.3)	<b>73.6 (0.9)</b>	54.6 (0.9)	55.8 (0.7)	59.8 (1.2)	58.9 (1.3)	66.2 (2.2)	68.5 (1.1)
magic	10	66.0 (2.0)	<b>69.0 (1.3)</b>	53.9 (0.6)	54.5 (0.7)	55.0 (0.9)	59.1 (1.4)	63.1 (1.1)	<b>70.0 (1.1)</b>
phishing	68	75.0 (1.4)	<b>91.3 (0.6)</b>	64.4 (1.0)	61.9 (1.1)	62.4 (1.1)	60.1 (1.3)	64.8 (1.4)	69.4 (0.8)
phoneme	5	<b>67.8 (1.5)</b>	<b>70.8 (1.0)</b>	65.2 (0.9)	66.7 (1.4)	<b>69.1 (1.4)</b>	61.3 (1.1)	64.5 (1.2)	<b>69.2 (1.1)</b>
spambase	57	69.7 (1.4)	<b>85.5 (0.5)</b>	60.1 (1.8)	54.4 (1.1)	65.4 (1.8)	61.5 (1.3)	63.6 (1.3)	70.5 (1.1)
susy	18	59.8 (1.3)	<b>74.8 (1.2)</b>	55.6 (0.7)	55.4 (0.9)	58.0 (1.0)	57.1 (1.2)	65.2 (1.0)	70.4 (1.2)
w8a	300	62.1 (1.5)	<b>86.5 (0.6)</b>	71.0 (0.8)	69.5 (1.5)	0.0 (0.0)	60.5 (1.5)	65.0 (2.0)	70.2 (1.2)
waveform	21	77.8 (1.3)	<b>87.0 (0.5)</b>	56.1 (0.8)	54.8 (0.7)	56.5 (0.9)	56.5 (0.9)	65.0 (0.9)	69.7 (1.1)

- ▶ linear-in-input model /  $n_U = n_S = 500$  /  $l_2$ -reguralization
- ▶ baseline: unsupervised / semi-supervised clustering

# Summary

## ■ Motivation

utilize pairwise information



no cluster assumption

- SU classification does not need explicit labels
- Properties: convexity, estimation error bound

## Poster: #67

arXiv URL

