

# Leveraging Ensemble Diversity for Robust Self-Training under Sample Selection Bias

Ambroise Odonnat

Huawei Noah's Ark Lab  
Université Rennes 2  
INRIA, CNRS, IRISA

Séminaire au vert Obélix

August 27, 2024





# Outline

- ① Introduction
- ② Failure of Self-Training
- ③ Learning with the  $\mathcal{T}$ -similarity
- ④ Numerical Experiments
- ⑤ Discussion



# Outline

① Introduction

② Failure of Self-Training

③ Learning with the  $\mathcal{T}$ -similarity

④ Numerical Experiments

⑤ Discussion



# Introduction

In some applications, data acquisition is cheaper than labeling ...



# Introduction

... and supervised learning is inefficient.

# Semi-Supervised Learning (SSL)



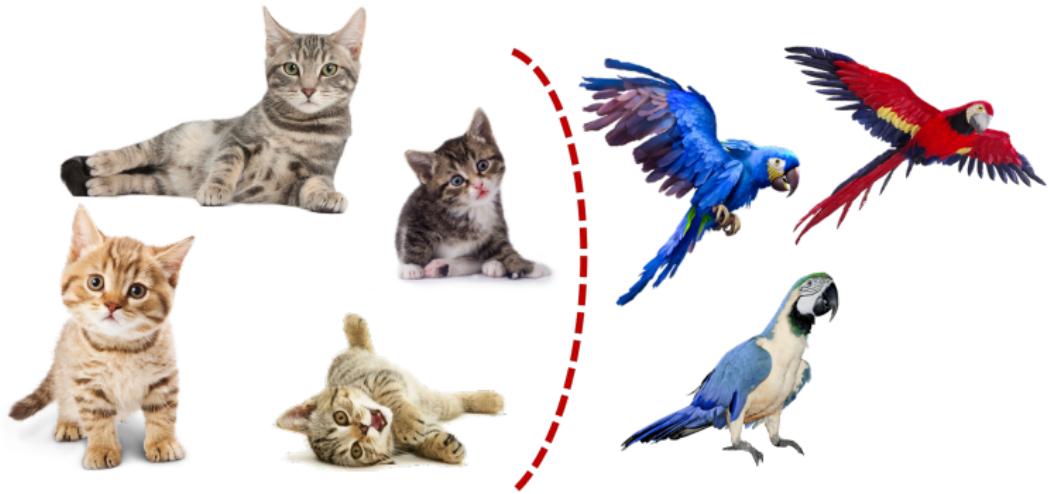
**SSL** → learn from a few labeled and many unlabeled examples.



# Semi-Supervised Learning (SSL)



**SSL** → learn from a few labeled and many unlabeled examples.





- Pseudo-labeling ([Amini et al., 2023](#)):
  - Unlabeled regularization ([Feofanov et al., 2023](#))
  - Self-training ([Feofanov et al., 2019](#))
- Graph-based algorithms ([van Engelen and Hoos, 2020](#)):
  - Label propagation
  - Label spreading
- Unsupervised preprocessing ([van Engelen and Hoos, 2020](#)):
  - Cluster-then-label
  - Feature extraction: auto-encoders, PCA
  - Pre-training: self-supervised learning, stacked auto-encoders

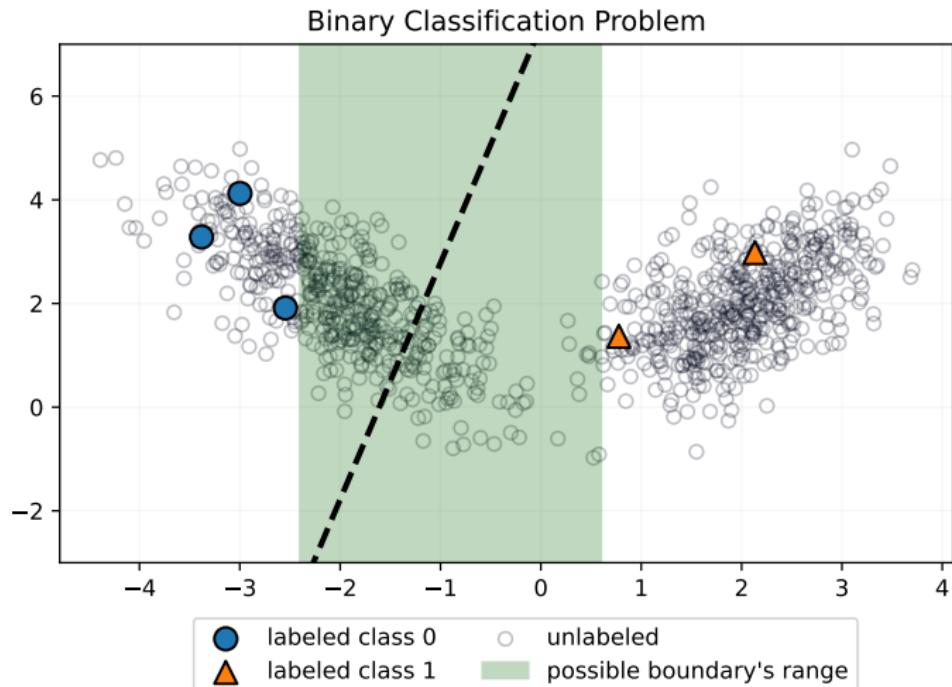


- **Pseudo-labeling** ([Amini et al., 2023](#)):
  - Unlabeled regularization ([Feofanov et al., 2023](#))
  - Self-training([Feofanov et al., 2019](#))
- Graph-based algorithms ([van Engelen and Hoos, 2020](#)):
  - Label propagation
  - Label spreading
- Unsupervised preprocessing ([van Engelen and Hoos, 2020](#)):
  - Cluster-then-label
  - Feature extraction: auto-encoders, PCA
  - Pre-training: self-supervised learning, stacked auto-encoders



# Low Density Separation

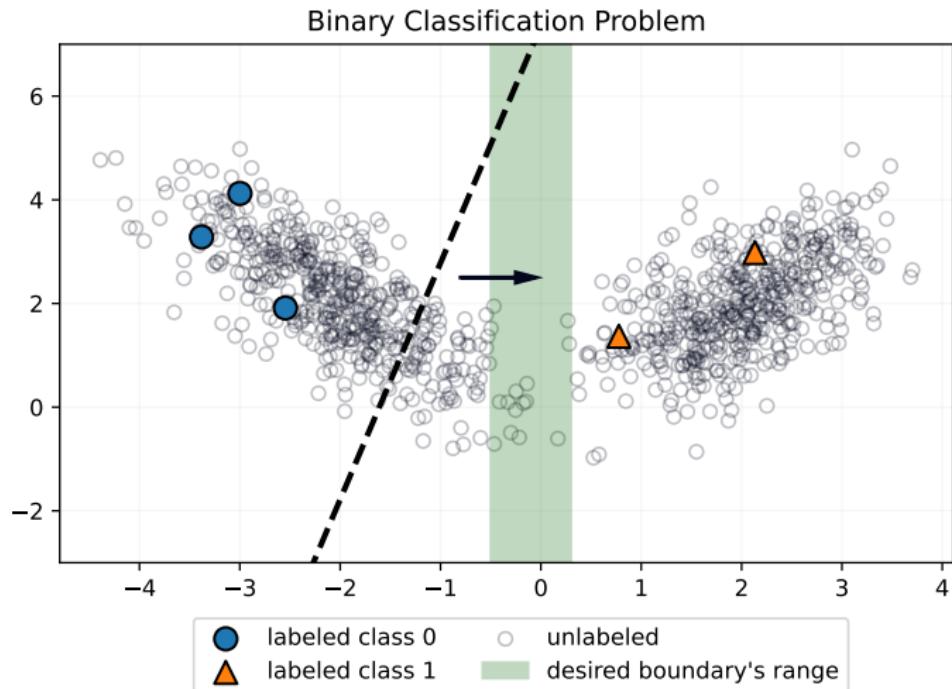
Range of possible supervised classifiers is vast: we need to make assumptions.





# Low Density Separation

Low Density Separation (LDS) assumption: push boundary away from regions of unlabeled data.



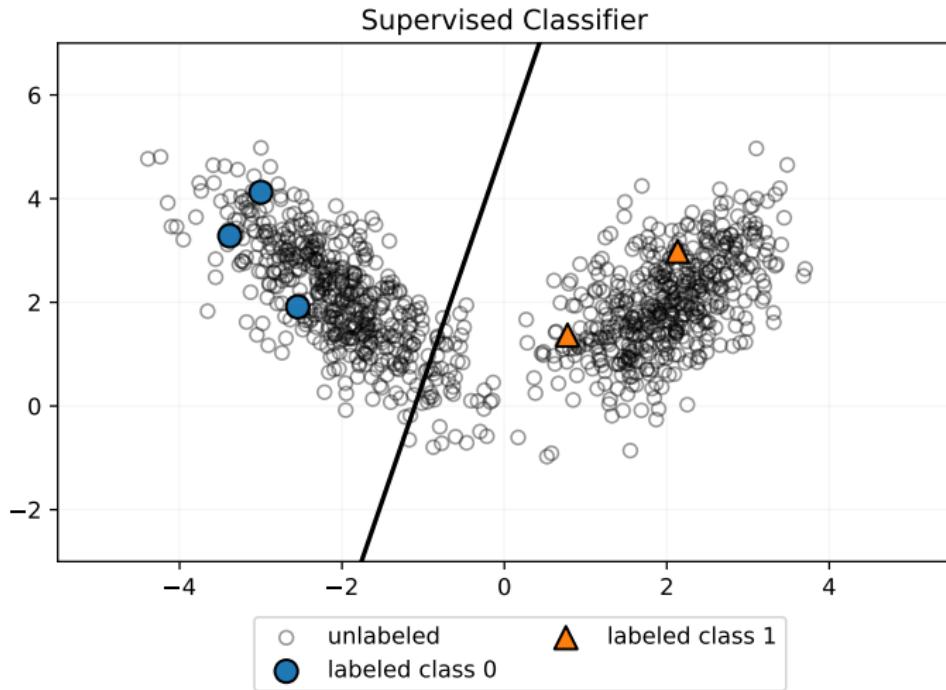


- **Pseudo-labeling** ([Amini et al., 2023](#)):
  - Unlabeled regularization ([Feofanov et al., 2023](#))
  - **Self-training** ([Feofanov et al., 2019](#))
- Graph-based algorithms ([van Engelen and Hoos, 2020](#)):
  - Label propagation
  - Label spreading
- Unsupervised preprocessing ([van Engelen and Hoos, 2020](#)):
  - Cluster-then-label
  - Feature extraction: auto-encoders, PCA
  - Pre-training: self-supervised learning, stacked auto-encoders



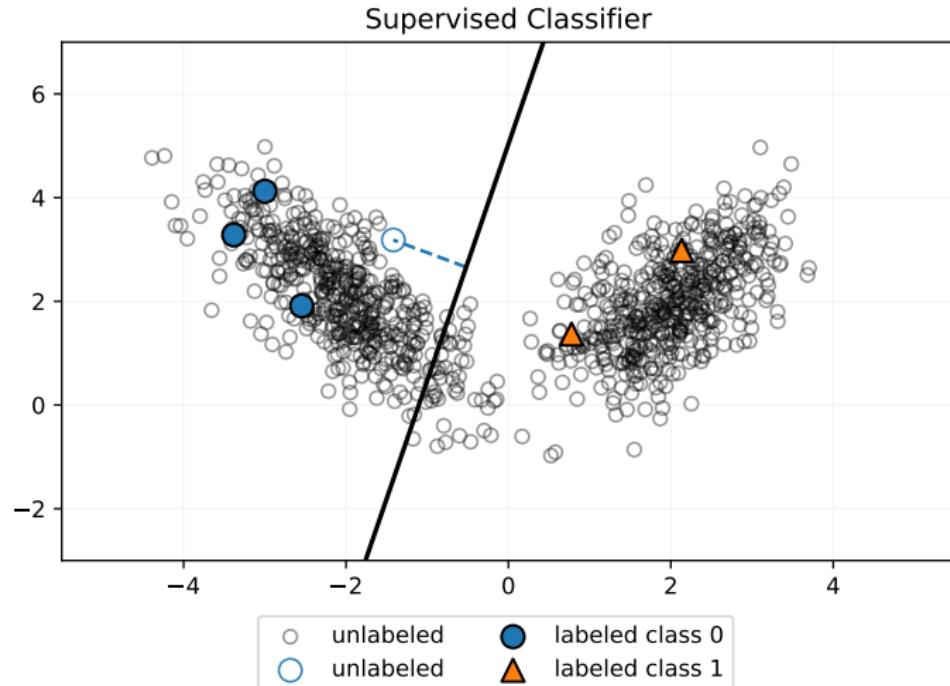
# Self-Training

Start from a supervised classifier trained on the labeled set.



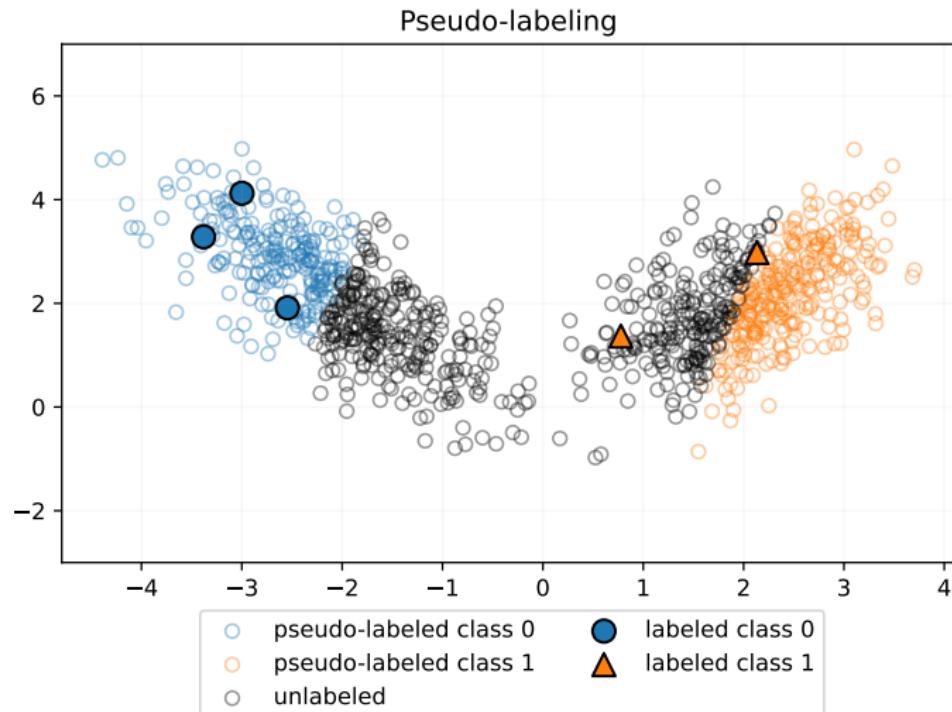


Predict labels and confidence scores for unlabeled data.





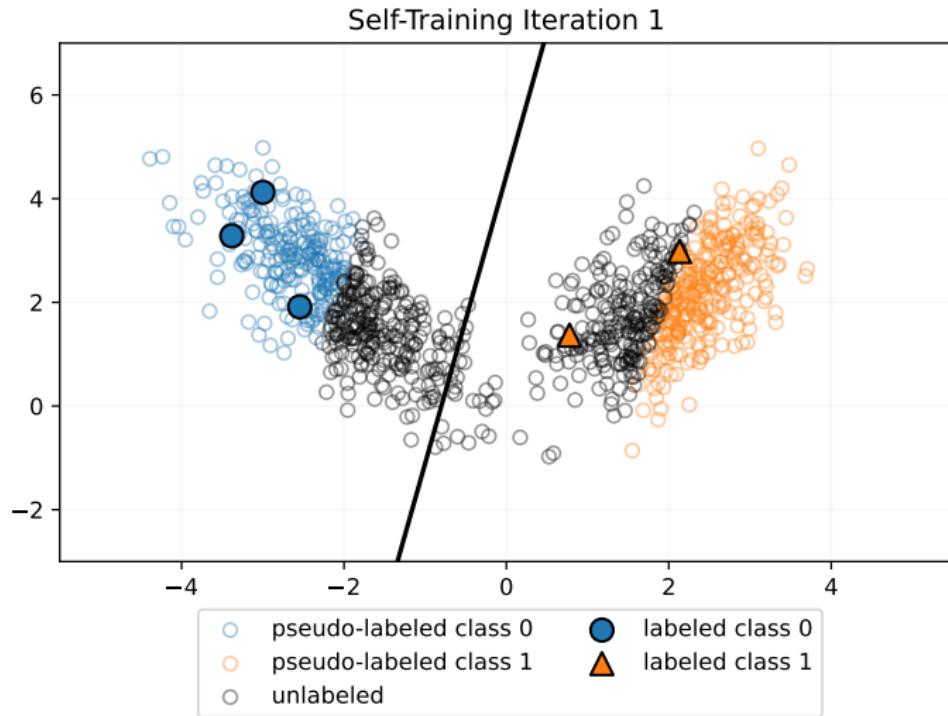
Pseudo-label most confident data and include in the labeled set.



# Self-Training

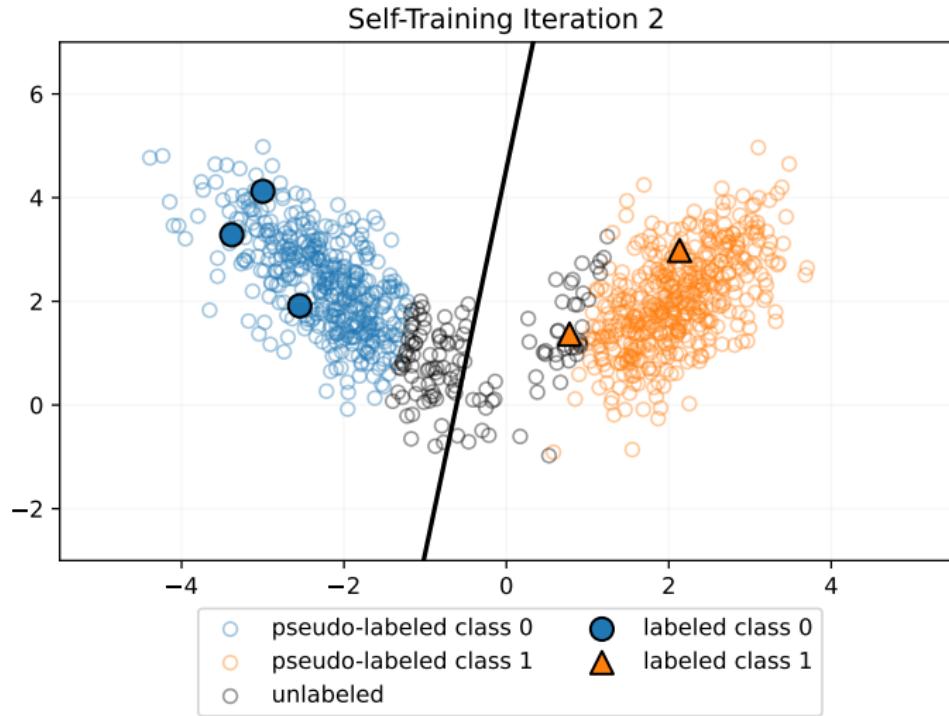


Retrain the model and repeat the same procedure again.





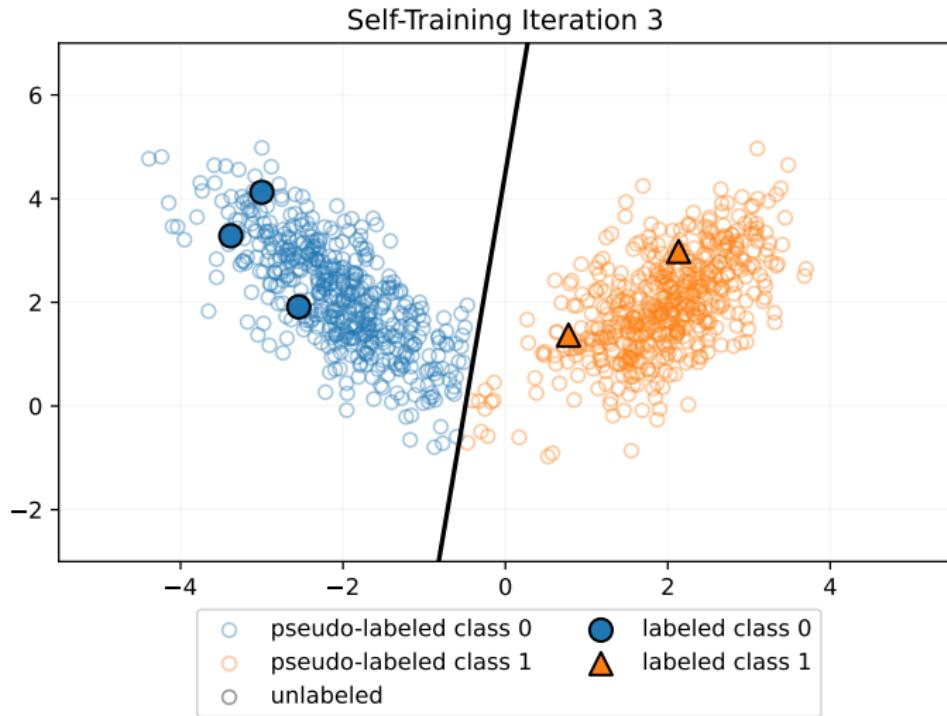
And again...





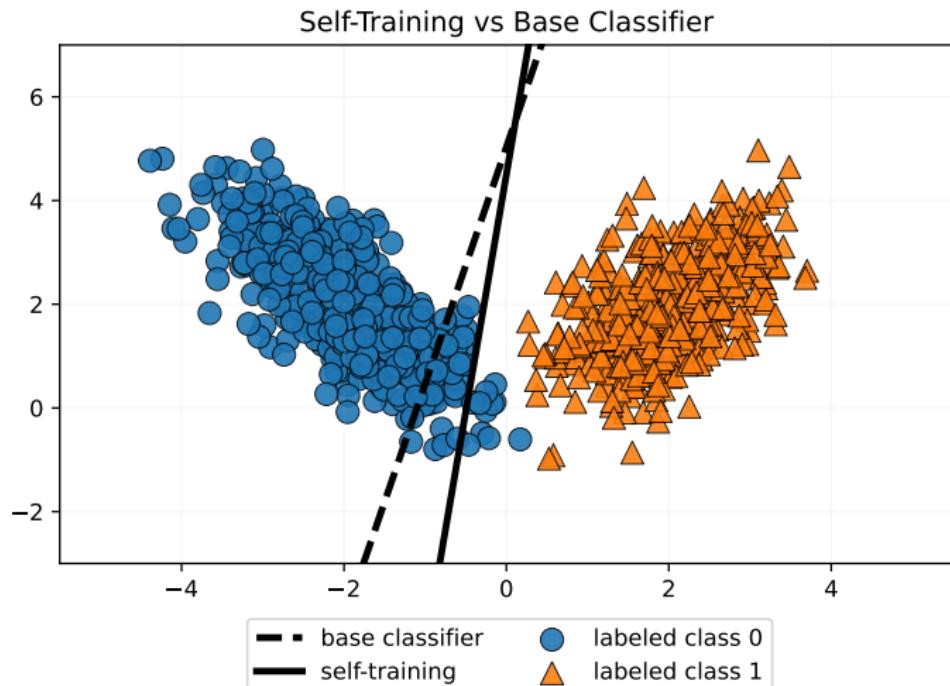
# Self-Training

Until there are no data to pseudo-label.





Self-training pushed the boundary away from the confident data.





# Outline

① Introduction

② Failure of Self-Training

③ Learning with the  $\mathcal{T}$ -similarity

④ Numerical Experiments

⑤ Discussion



# Two Fundamental Questions

- ① *Confidence Estimation* → How to rank unlabeled data?



# Two Fundamental Questions

- ① *Confidence Estimation* → How to rank unlabeled data?
- ② *Pseudo-Labeling Policy* → How to select unlabeled data for pseudo-labeling at each iteration?



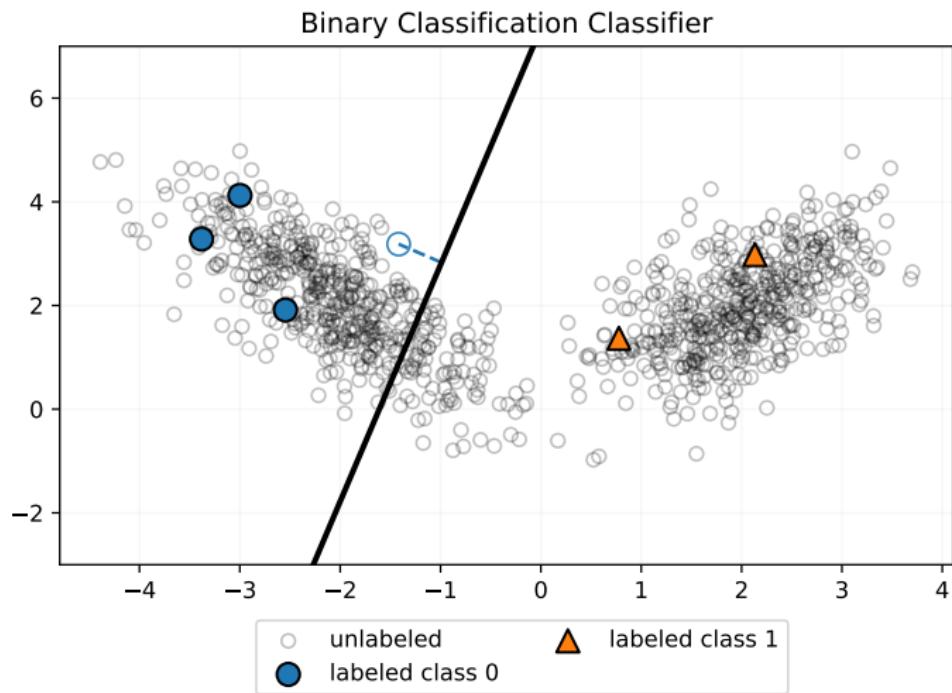
# Two Fundamental Questions

- ① *Confidence Estimation* → How to rank unlabeled data?
- ② *Pseudo-Labeling Policy* → How to select unlabeled data for pseudo-labeling at each iteration?

In this work, we focus on *Confidence Estimation*.



Requirements → trust the classifier's predictions.





# Failure Cases

Problem → not safe since the prediction can be wrong.



# Failure Cases

Biased prediction confidence  $\Rightarrow$  wrong direction can be chosen.  
→ This can occur when there is a distribution shift in the data.



- SSL assumption: labeled and unlabeled data are i.i.d.
- Confidence can be biased when this assumption does not hold



- SSL assumption: labeled and unlabeled data are i.i.d.
- Confidence can be biased when this assumption does not hold
- **Sample Selection Bias (SSB)**: data labeling subject to constraints



- SSL assumption: labeled and unlabeled data are i.i.d.
- Confidence can be biased when this assumption does not hold
- **Sample Selection Bias (SSB)**: data labeling subject to constraints
  - Creation of group study in clinical trials;
  - People with poor mobility less likely to be in street surveys;
  - Labeling can be constrained for privacy reasons.

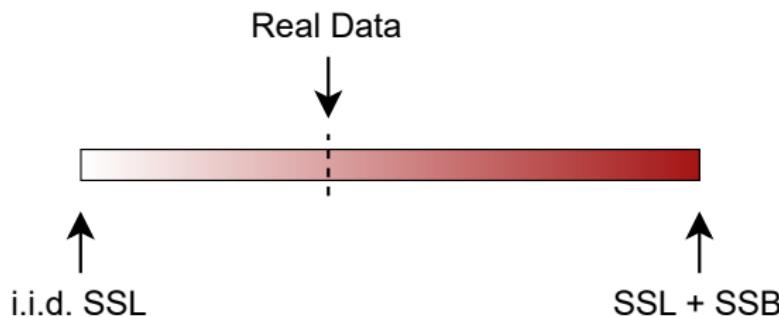


- SSL assumption: labeled and unlabeled data are i.i.d.
- Confidence can be biased when this assumption does not hold
- **Sample Selection Bias (SSB)**: data labeling subject to constraints
  - Creation of group study in clinical trials;
  - People with poor mobility less likely to be in street surveys;
  - Labeling can be constrained for privacy reasons.
- SSB has been studied but not in the case of SSL.



SSL + SSB combines SSL and Sample Selection Bias (SSB):

- ① Few labeled examples (SSL)
- ② Biased labeling procedure (SSB)



Goal → obtain a method good on **both** i.i.d. SSL and SSL + SSB.



Select the labeled set to violate the i.i.d. assumption.

- Binary selection variable  $s_i$  for each  $\mathbf{x}_i$ ;
- $s_i = 1$  if  $\mathbf{x}_i$  is labeled,  $s_i = 0$  otherwise;
- Model  $\mathbb{P}(s_i = 1|\mathbf{x}_i, y_i)$  to violate i.i.d. assumption.



Select the labeled set to violate the i.i.d. assumption.

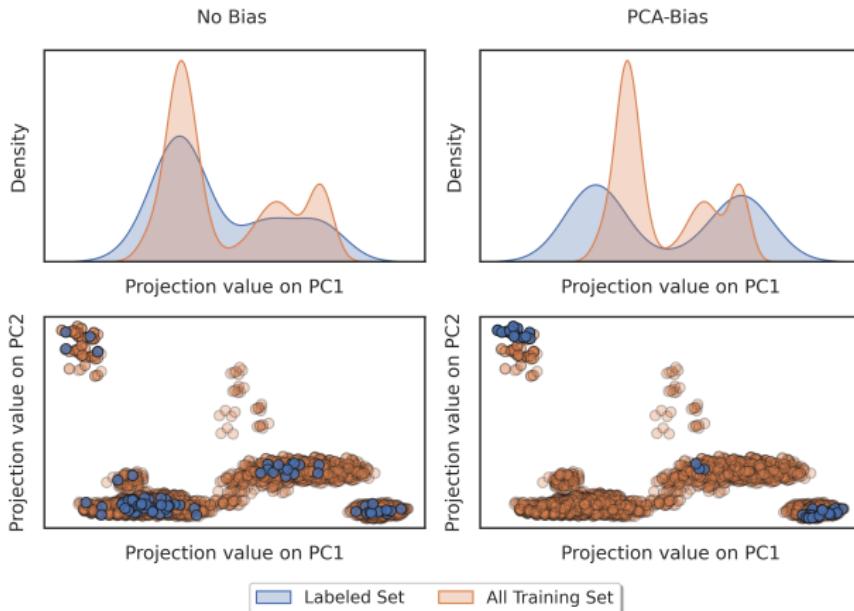
PCA-Bias → for each class  $c$ ,

- ① Apply PCA on training data of class  $c$ ;
- ② Compute  $\text{proj}_1(\mathbf{x}_i)$ , projection value on PC1;
- ③  $\mathbb{P}(s_i = 1 | \mathbf{x}_i, y_i = c) \propto \exp(r \cdot |\text{proj}_1(\mathbf{x}_i)|)$ ,  $r > 0$ .

# Implementation of SSL + SSB



Select the labeled set to violate the i.i.d. assumption.

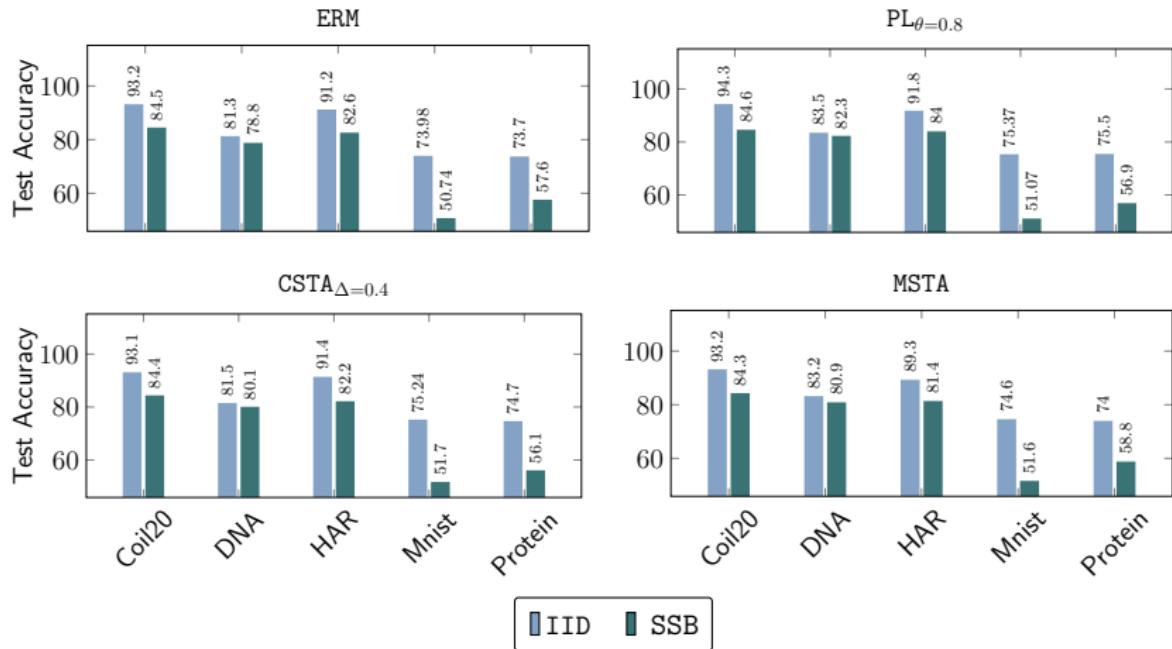




# Pseudo-Labeling Policies

- ERM corresponds to supervised learning on the labeled set
- PL $_{\theta=0.8}$  uses a fixed threshold  $\theta = 0.8$  ([Lee, 2013](#))
- CSTA $_{\Delta=0.4}$  takes  $\Delta\%$  most confident ([Cascante-Bonilla et al., 2020](#))
- MSTA optimizes the threshold to balance the error and the amount of data pseudo-labeled ([Feofanov et al., 2019](#))

# Failure of Self-Training under SSL + SSB



**Figure:** Test accuracies of the different baselines on 5 datasets. Full results to be found [here](#).



# Unreliable Model Selection

LOO over-optimistic w.r.t. generalization performance (Figure 1).

- Leave one labeled point out;
- Train on the remaining  $n_\ell - 1$ ;
- Test on the one left out;
- Repeat for each labeled point.

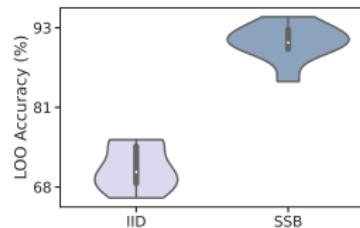


Figure: LOO on Mnist.



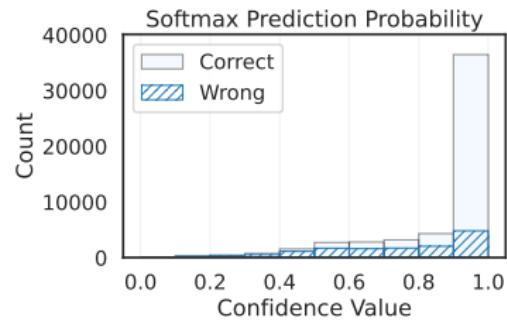
# Outline

- ➊ Introduction
- ➋ Failure of Self-Training
- ➌ Learning with the  $\mathcal{T}$ -similarity
- ➍ Numerical Experiments
- ➎ Discussion



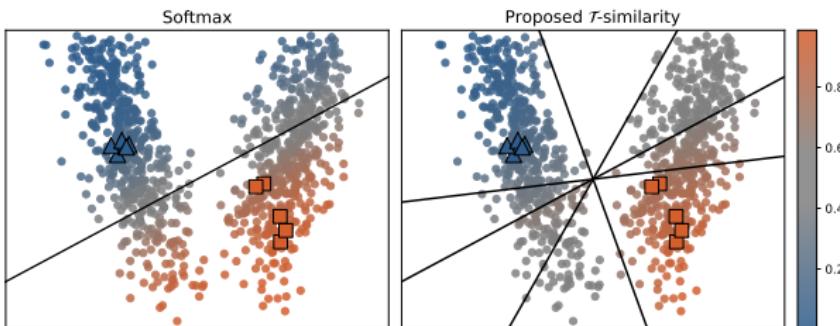
softmax-based confidence measure is unreliable in SSL + SSB.

- NNs are overconfident;
- softmax predictions biased towards the labeled set.



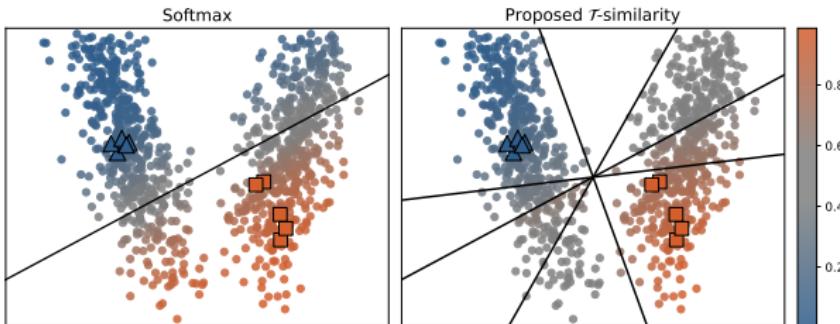
→ We propose a novel confidence measure for NNs.

# Leveraging Ensemble Diversity



$$\min_{\mathcal{T}} \frac{1}{M} \sum_{h \in \mathcal{T}} \underbrace{\left( \frac{1}{n_\ell} \sum_{(\mathbf{x}, y) \in \mathbf{X}_\ell \times \mathbf{y}_\ell} \ell(h(\mathbf{x}), y) \right)}_{\text{supervised loss}} + \underbrace{\frac{\gamma}{n_u M(M-1)} \sum_{h \neq \tilde{h} \in \mathcal{T}} \sum_{\mathbf{x} \in \mathbf{X}_u} h(\mathbf{x})^\top \tilde{h}(\mathbf{x})}_{\text{agreement loss}}$$

# Leveraging Ensemble Diversity



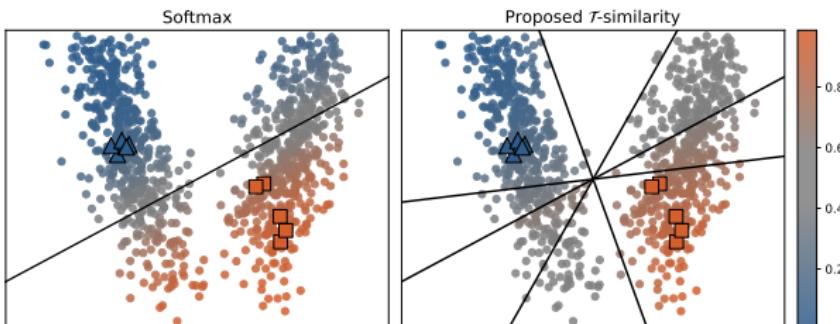
$$\min_{\mathcal{T}} \frac{1}{M} \sum_{h \in \mathcal{T}} \left( \underbrace{\frac{1}{n_\ell} \sum_{(\mathbf{x}, y) \in \mathbf{X}_\ell \times \mathbf{y}_\ell} \ell(h(\mathbf{x}), y)}_{\text{supervised loss}} \right) + \underbrace{\frac{\gamma}{n_u M(M-1)} \sum_{h \neq \tilde{h} \in \mathcal{T}} \sum_{\mathbf{x} \in \mathbf{X}_u} h(\mathbf{x})^\top \tilde{h}(\mathbf{x})}_{\text{agreement loss}}$$

We jointly train the ensemble to

- ① Fit very well the labeled data
- ② Disagree as much as possible on unlabeled data



# Leveraging Ensemble Diversity



$$\min_{\mathcal{T}} \frac{1}{M} \sum_{h \in \mathcal{T}} \left( \underbrace{\frac{1}{n_\ell} \sum_{(\mathbf{x}, y) \in \mathbf{X}_\ell \times \mathbf{y}_\ell} \ell(h(\mathbf{x}), y)}_{\text{supervised loss}} \right) + \underbrace{\frac{\gamma}{n_u M(M-1)} \sum_{h \neq \tilde{h} \in \mathcal{T}} \sum_{\mathbf{x} \in \mathbf{X}_u} h(\mathbf{x})^\top \tilde{h}(\mathbf{x})}_{\text{agreement loss}}$$

We jointly train the ensemble to

- ① Fit very well the labeled data
- ② Disagree as much as possible on unlabeled data

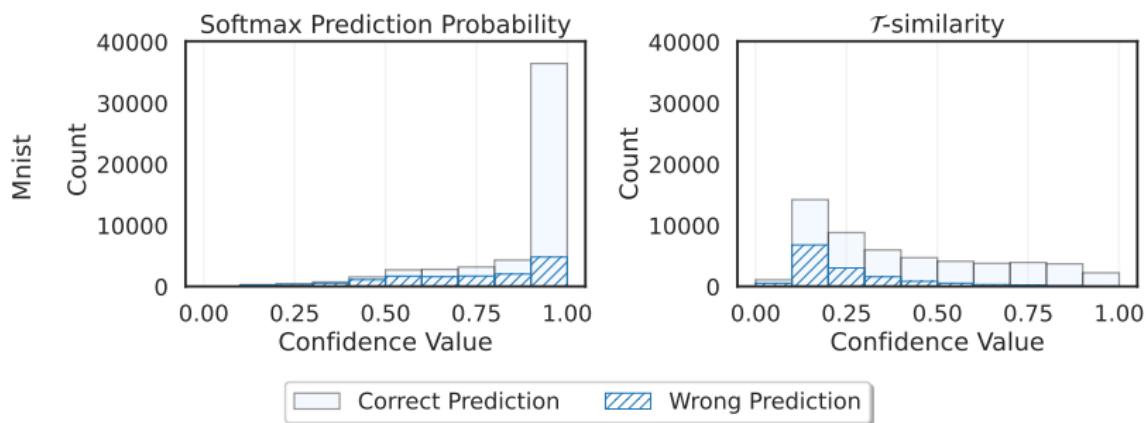


# $\mathcal{T}$ -similarity

- We define the  $\mathcal{T}$ -similarity as:

$$s_{\mathcal{T}}(\mathbf{x}) = \frac{1}{M(M-1)} \sum_{h \neq \tilde{h} \in \mathcal{T}} h(\mathbf{x})^\top \tilde{h}(\mathbf{x}).$$

- For any  $\mathbf{x}$ , we have  $0 \leq s_{\mathcal{T}}(\mathbf{x}) \leq 1$ .





# Practical Implementation

- ① Projection layers are learned through a classification head;
- ② Confidence estimator is ensemble of  $M=5$  linear heads that don't affect representation.

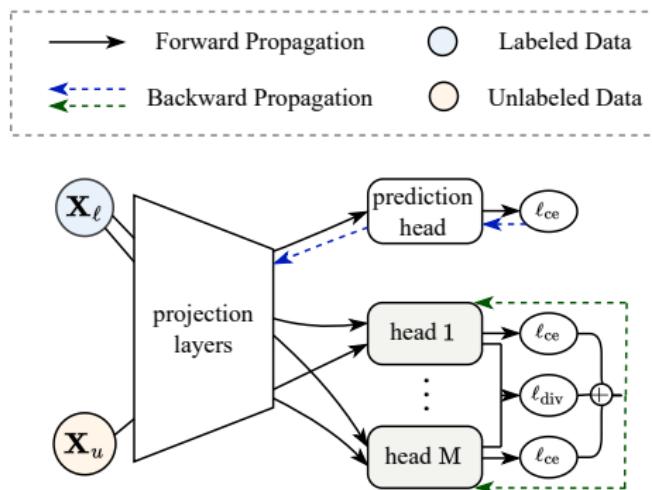


Figure: Architecture of the model.



# Theoretical Analysis

- Fixed representation of dimension  $d$ , binary linear classification
- Linear ensemble  $\mathbf{W} = \{\omega_m \in \mathbb{R}^d | 1 \leq m \leq M\}$
- Prediction of  $\omega_m$  on  $\mathbf{x}$  is  $\text{sign}(\omega_m^\top \mathbf{x})$

$$\begin{aligned}\mathcal{L}(\mathbf{W}) := & \underbrace{\frac{1}{Mn_\ell} \sum_{m=1}^M \sum_{i=1}^{n_\ell} (y_i - \omega_m^\top \mathbf{x}_i)^2}_{\text{label fidelity term}} + \underbrace{\frac{1}{M} \sum_{m=1}^M \lambda_m \|\omega_m\|^2}_{\text{regularization}} \\ & + \underbrace{\frac{\gamma}{n_u M(M-1)} \sum_{m \neq k} \sum_{i=n_\ell+1}^{n_\ell+n_u} w_m^\top \mathbf{x}_i w_k^\top \mathbf{x}_i}_{\text{agreement term}},\end{aligned}\tag{P}$$

where  $\gamma$  controls the influence of the diversity on the learning.



# Theoretical Analysis

- Fixed representation of dimension  $d$ , binary linear classification
- Linear ensemble  $\mathbf{W} = \{\mathbf{w}_m \in \mathbb{R}^d | 1 \leq m \leq M\}$
- Prediction of  $\omega_m$  on  $\mathbf{x}$  is  $\text{sign}(\omega_m^\top \mathbf{x})$

Theorem (O., Feofanov, Redko)

- ① *Convergence to a stationary point under mild assumption*
- ② *Lower-bound on the diversity of stationary points*
- ③ *Connection to contrastive learning*



# Convergence to a Stationary Point

- Labeled set  $(\mathbf{X}_\ell, \mathbf{y}_\ell) = (\mathbf{x}_i, y_i)_{i=1}^{n_\ell}$
- Unlabeled set  $\mathbf{X}_u = (\mathbf{x}_i)_{i=n_\ell+1}^{n_\ell+n_u}$
- Assumption **A**:  $\forall m \in \llbracket 1, M \rrbracket, \lambda_m > \frac{\gamma(M+1)}{n_u(M-1)} \lambda_{\max}(\mathbf{X}_u^\top \mathbf{X}_u)$ .

## Theorem (O., Feofanov, Redko)

*Under Assumption A,  $\mathcal{L}$  is strictly convex and coercive on  $\mathbb{R}^{d \times M}$ . Hence, the optimization problem (P) admits a unique solution  $\mathbf{W}^*$  that verifies*

$$\nabla \mathcal{L}(\mathbf{W}^*) = 0. \tag{1}$$



# Diversity of an Ensemble

$$\ell_{\text{div}}(\mathbf{W}, \mathbf{X}_u) = -\frac{1}{n_u M(M-1)} \sum_{m \neq k} \boldsymbol{\omega}_m^\top \mathbf{X}_u^\top \mathbf{X}_u \boldsymbol{\omega}_k.$$

Theorem (O., Feofanov, Redko)

$$\begin{aligned} \gamma \ell_{\text{div}}(\mathbf{W}^*, \mathbf{X}_u) &\geq \frac{1}{2n_\ell M} \sum_{m=1}^M \|\mathbf{y}_\ell - \mathbf{X}_\ell \boldsymbol{\omega}_m^*\|_2^2 \\ &\quad + \frac{1}{2M} \sum_{m=1}^M (\boldsymbol{\omega}_m^*)^\top \left( \lambda_m \mathbf{I}_d + \frac{\mathbf{X}_\ell^\top \mathbf{X}_\ell}{n_\ell} \right) \boldsymbol{\omega}_m^*. \end{aligned}$$

- ① Trade-off between supervised performance and margin term
- ② Assuming orthogonality, the predictors  $\boldsymbol{\omega}_m$  span the  $M$  directions of largest variance of the labeled data.



# Diversity of an Ensemble

$$\ell_{\text{div}}(\mathbf{W}, \mathbf{X}_u) = -\frac{1}{n_u M(M-1)} \sum_{m \neq k} \boldsymbol{\omega}_m^\top \mathbf{X}_u^\top \mathbf{X}_u \boldsymbol{\omega}_k.$$

Theorem (O., Feofanov, Redko)

$$\gamma \ell_{\text{div}}(\mathbf{W}^*, \mathbf{X}_u) \geq \frac{1}{2M} \left( \lambda + \frac{1}{n_\ell} \lambda_{\min} (\mathbf{X}_\ell^\top \mathbf{X}_\ell) \right) \|\mathbf{W}^*\|_F^2.$$

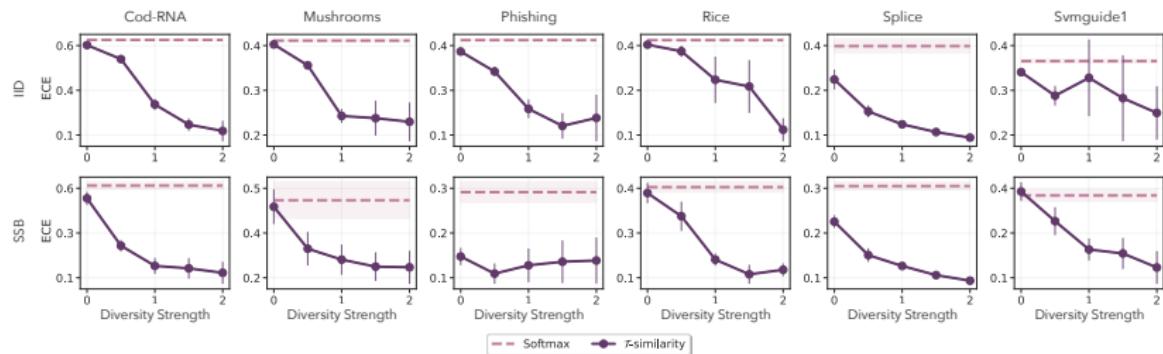
- ① Direction of smallest variance is also important for diversity
- ② Theorem shows the importance of representation learning



# Outline

- ① Introduction
- ② Failure of Self-Training
- ③ Learning with the  $\mathcal{T}$ -similarity
- ④ Numerical Experiments
- ⑤ Discussion

# Diversity provides Calibrated Confidence Measure



**Figure:** Increasing the diversity improves the classifier's calibration



# Pseudo-Labeling Policies

- ERM corresponds to supervised learning on the labeled set
- PL $_{\theta=0.8}$  uses a fixed threshold  $\theta = 0.8$  ([Lee, 2013](#))
- CSTA $_{\Delta=0.4}$  takes  $\Delta\%$  most confident ([Cascante-Bonilla et al., 2020](#))
- MSTA optimizes the threshold to balance the error and the amount of data pseudo-labeled ([Feofanov et al., 2019](#))

# Results in SSL+SSB



Dataset	ERM	$\text{PL}_{\theta=0.8}$		$\text{CSTA}_{\Delta=0.4}$		$\text{MSTA}$	
		softmax	$\mathcal{T}$ -similarity	softmax	$\mathcal{T}$ -similarity	softmax	$\mathcal{T}$ -similarity
Cod-RNA	74.51 ± 8.86	74.75 ± 8.14	<b>80.06 ± 3.55</b>	73.39 ± 7.36	<b>78.39 ± 4.66</b>	75.28 ± 8.79	<b>76.88 ± 7.67</b>
COIL-20	84.54 ± 2.19	<b>84.69 ± 3.56</b>	84.57 ± 2.85	84.38 ± 3.05	<b>84.57 ± 3.16</b>	<b>84.32 ± 2.34</b>	84.07 ± 2.85
Digits	75.68 ± 4.59	<b>80.47 ± 3.8</b>	78.2 ± 3.34	78.4 ± 3.28	<b>79.14 ± 3.5</b>	78.02 ± 5.15	<b>79.8 ± 5.92</b>
DNA	78.82 ± 2.31	<b>80.29 ± 2.24</b>	79.06 ± 2.31	80.12 ± 2.08	<b>80.76 ± 2.24</b>	80.89 ± 2.64	<b>84.09 ± 1.7</b>
DryBean	64.6 ± 3.89	<b>65.6 ± 4.18</b>	61.55 ± 4.91	<b>64.91 ± 3.72</b>	64.6 ± 3.53	66.24 ± 4.31	<b>67.0 ± 3.96</b>
HAR	82.57 ± 1.96	82.87 ± 3.02	<b>83.12 ± 2.27</b>	82.19 ± 2.61	<b>83.53 ± 3.77</b>	<b>81.35 ± 2.54</b>	81.16 ± 1.63
Mnist	50.74 ± 2.25	51.08 ± 2.55	<b>52.69 ± 2.42</b>	51.7 ± 3.52	<b>54.26 ± 1.82</b>	51.6 ± 2.58	<b>54.18 ± 2.34</b>
Mushrooms	69.45 ± 7.29	59.53 ± 10.46	<b>71.36 ± 6.63</b>	62.98 ± 7.25	<b>77.55 ± 7.65</b>	72.16 ± 7.59	<b>76.16 ± 13.04</b>
Phishing	67.42 ± 3.55	66.08 ± 5.66	<b>77.41 ± 3.93</b>	66.88 ± 5.64	<b>76.17 ± 8.58</b>	69.48 ± 4.37	<b>75.83 ± 7.52</b>
Protein	57.57 ± 6.33	57.45 ± 6.36	<b>57.61 ± 6.23</b>	56.09 ± 5.61	<b>57.74 ± 7.8</b>	58.81 ± 6.54	<b>59.88 ± 6.29</b>
Rice	79.19 ± 5.12	80.54 ± 4.31	<b>81.1 ± 4.28</b>	79.88 ± 4.48	<b>81.56 ± 3.61</b>	80.35 ± 4.89	<b>82.63 ± 5.63</b>
Splice	66.13 ± 4.47	67.14 ± 2.62	<b>67.45 ± 2.53</b>	67.28 ± 2.07	<b>68.05 ± 2.17</b>	66.08 ± 4.98	<b>66.32 ± 4.73</b>
Svmguide1	70.89 ± 10.98	70.35 ± 11.74	<b>81.07 ± 5.39</b>	69.84 ± 11.06	<b>74.46 ± 7.23</b>	71.04 ± 11.11	<b>73.13 ± 8.82</b>

- $\mathcal{T}$ -similarity is better overall;



# Results in SSL+SSB

Dataset	ERM	PL <sub><math>\theta=0.8</math></sub>		CSTA <sub><math>\Delta=0.4</math></sub>		MSTA	
		softmax	$\mathcal{T}$ -similarity	softmax	$\mathcal{T}$ -similarity	softmax	$\mathcal{T}$ -similarity
Cod-RNA	74.51 ± 8.86	74.75 ± 8.14	<b>80.06 ± 3.55</b>	73.39 ± 7.36	<b>78.39 ± 4.66</b>	75.28 ± 8.79	<b>76.88 ± 7.67</b>
COIL-20	84.54 ± 2.19	<b>84.69 ± 3.56</b>	84.57 ± 2.85	84.38 ± 3.05	<b>84.57 ± 3.16</b>	<b>84.32 ± 2.34</b>	84.07 ± 2.85
Digits	75.68 ± 4.59	<b>80.47 ± 3.8</b>	78.2 ± 3.34	78.4 ± 3.28	<b>79.14 ± 3.5</b>	78.02 ± 5.15	<b>79.8 ± 5.92</b>
DNA	78.82 ± 2.31	<b>80.29 ± 2.24</b>	79.06 ± 2.31	80.12 ± 2.08	<b>80.76 ± 2.24</b>	80.89 ± 2.64	<b>84.09 ± 1.7</b>
DryBean	64.6 ± 3.89	<b>65.6 ± 4.18</b>	61.55 ± 4.91	<b>64.91 ± 3.72</b>	64.6 ± 3.53	66.24 ± 4.31	<b>67.0 ± 3.96</b>
HAR	82.57 ± 1.96	82.87 ± 3.02	<b>83.12 ± 2.27</b>	82.19 ± 2.61	<b>83.53 ± 3.77</b>	<b>81.35 ± 2.54</b>	81.16 ± 1.63
Mnist	50.74 ± 2.25	51.08 ± 2.55	<b>52.69 ± 2.42</b>	51.7 ± 3.52	<b>54.26 ± 1.82</b>	51.6 ± 2.58	<b>54.18 ± 2.34</b>
Mushrooms	<b>69.45 ± 7.29</b>	<b>59.53 ± 10.46</b>	<b>71.36 ± 6.63</b>	<b>62.98 ± 7.25</b>	<b>77.55 ± 7.65</b>	<b>72.16 ± 7.59</b>	<b>76.16 ± 13.04</b>
Phishing	<b>67.42 ± 3.55</b>	<b>66.08 ± 5.66</b>	<b>77.41 ± 3.93</b>	<b>66.88 ± 5.64</b>	<b>76.17 ± 8.58</b>	<b>69.48 ± 4.37</b>	<b>75.83 ± 7.52</b>
Protein	57.57 ± 6.33	57.45 ± 6.36	<b>57.61 ± 6.23</b>	56.09 ± 5.61	<b>57.74 ± 7.8</b>	58.81 ± 6.54	<b>59.88 ± 6.29</b>
Rice	79.19 ± 5.12	80.54 ± 4.31	<b>81.1 ± 4.28</b>	79.88 ± 4.48	<b>81.56 ± 3.61</b>	80.35 ± 4.89	<b>82.63 ± 5.63</b>
Splice	66.13 ± 4.47	67.14 ± 2.62	<b>67.45 ± 2.53</b>	67.28 ± 2.07	<b>68.05 ± 2.17</b>	66.08 ± 4.98	<b>66.32 ± 4.73</b>
Svmguide1	70.89 ± 10.98	70.35 ± 11.74	<b>81.07 ± 5.39</b>	69.84 ± 11.06	<b>74.46 ± 7.23</b>	71.04 ± 11.11	<b>73.13 ± 8.82</b>

- $\mathcal{T}$ -similarity is better overall;
- Even go from degradation to improvement on **2 datasets**.



# Results in SSL+SSB

Dataset	ERM	PL $_{\theta=0.8}$		CSTA $_{\Delta=0.4}$		MSTA	
		softmax	T-similarity	softmax	T-similarity	softmax	T-similarity
Cod-RNA	74.51 ± 8.86	74.75 ± 8.14	<b>80.06 ± 3.55</b>	73.39 ± 7.36	<b>78.39 ± 4.66</b>	75.28 ± 8.79	<b>76.88 ± 7.67</b>
COIL-20	84.54 ± 2.19	<b>84.69 ± 3.56</b>	84.57 ± 2.85	84.38 ± 3.05	<b>84.57 ± 3.16</b>	<b>84.32 ± 2.34</b>	84.07 ± 2.85
Digits	75.68 ± 4.59	<b>80.47 ± 3.8</b>	78.2 ± 3.34	78.4 ± 3.28	<b>79.14 ± 3.5</b>	78.02 ± 5.15	<b>79.8 ± 5.92</b>
DNA	78.82 ± 2.31	<b>80.29 ± 2.24</b>	79.06 ± 2.31	80.12 ± 2.08	<b>80.76 ± 2.24</b>	80.89 ± 2.64	<b>84.09 ± 1.7</b>
DryBean	64.6 ± 3.89	<b>65.6 ± 4.18</b>	61.55 ± 4.91	<b>64.91 ± 3.72</b>	64.6 ± 3.53	66.24 ± 4.31	<b>67.0 ± 3.96</b>
HAR	82.57 ± 1.96	82.87 ± 3.02	<b>83.12 ± 2.27</b>	82.19 ± 2.61	<b>83.53 ± 3.77</b>	<b>81.35 ± 2.54</b>	81.16 ± 1.63
Mnist	50.74 ± 2.25	51.08 ± 2.55	<b>52.69 ± 2.42</b>	51.7 ± 3.52	<b>54.26 ± 1.82</b>	51.6 ± 2.58	<b>54.18 ± 2.34</b>
Mushrooms	<b>69.45 ± 7.29</b>	<b>59.53 ± 10.46</b>	<b>71.36 ± 6.63</b>	<b>62.98 ± 7.25</b>	<b>77.55 ± 7.65</b>	<b>72.16 ± 7.59</b>	<b>76.16 ± 13.04</b>
Phishing	<b>67.42 ± 3.55</b>	<b>66.08 ± 5.66</b>	<b>77.41 ± 3.93</b>	<b>66.88 ± 5.64</b>	<b>76.17 ± 8.58</b>	<b>69.48 ± 4.37</b>	<b>75.83 ± 7.52</b>
Protein	57.57 ± 6.33	57.45 ± 6.36	<b>57.61 ± 6.23</b>	56.09 ± 5.61	<b>57.74 ± 7.8</b>	58.81 ± 6.54	<b>59.88 ± 6.29</b>
Rice	79.19 ± 5.12	80.54 ± 4.31	<b>81.1 ± 4.28</b>	79.88 ± 4.48	<b>81.56 ± 3.61</b>	80.35 ± 4.89	<b>82.63 ± 5.63</b>
Splice	66.13 ± 4.47	67.14 ± 2.62	<b>67.45 ± 2.53</b>	67.28 ± 2.07	<b>68.05 ± 2.17</b>	66.08 ± 4.98	<b>66.32 ± 4.73</b>
Svmguide1	70.89 ± 10.98	70.35 ± 11.74	<b>81.07 ± 5.39</b>	69.84 ± 11.06	<b>74.46 ± 7.23</b>	71.04 ± 11.11	<b>73.13 ± 8.82</b>

- T-similarity is better overall;
- Even go from degradation to improvement on 2 datasets;
- Our approach remains similar to softmax in i.i.d. SSL.



# Outline

- ➊ Introduction
- ➋ Failure of Self-Training
- ➌ Learning with the  $\mathcal{T}$ -similarity
- ➍ Numerical Experiments
- ➎ Discussion



# Discussion

- ① Practical and principled framework to study SSL + SSB;
- ② Calibrated confidence measure;
- ③  $\mathcal{T}$ -similarity good both in i.i.d. SSL and SSL + SSB.

Future work → use  $\mathcal{T}$ -similarity for iterative self-training, domain adaptation, or uncertainty modeling.



This work has been accepted to AISTATS 2024, Valencia, Spain.  
You may find the links to the paper and the code below. to know  
more about my research, see my website: [ambroiseodt.github.io](https://ambroiseodt.github.io)  
and feel free to contact me.

- Paper: <https://arxiv.org/abs/2310.14814>
- Code: <https://github.com/ambroiseodt/tsim>

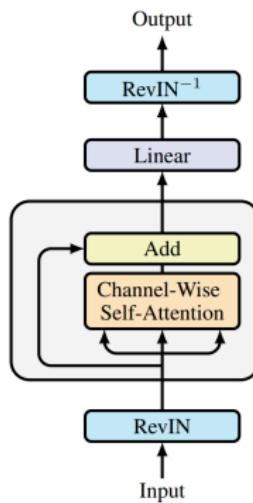
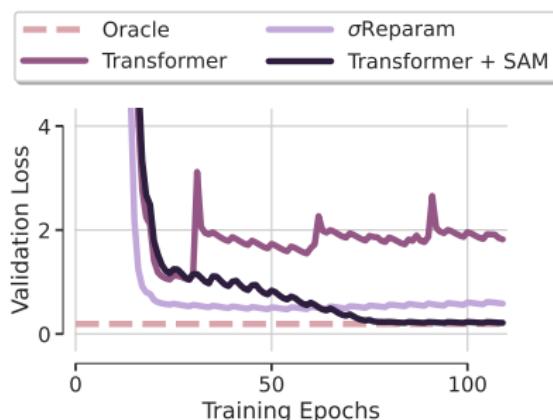


## Acknowledgement

I thank my co-authors Vasilii Feofanov and Ievgen Redko who supervised me during the internship and without whom this project would not have existed. I would also like to thank Gabriel Peyré for his insightful comments on early drafts of the paper, as well as Malik Tiomoko, and Aladin Virmaux for the fruitful discussions that led to this work.



## SAMformer: Unlocking the Potential of Transformers in Time Series Forecasting - Oral ICML 2024



- Paper: <https://arxiv.org/pdf/2402.10198>
- Code: <https://github.com/romilbert/samformer>

Thanks for your attention !



# References I

- Amini, M.-R., Feofanov, V., Pauletto, L., Devijver, E., and Maximov, Y. (2023). Self-Training: A Survey.
- Cascante-Bonilla, P., Tan, F., Qi, Y., and Ordonez, V. (2020). Curriculum labeling: Revisiting pseudo-labeling for semi-supervised learning.
- Feofanov, V., Devijver, E., and Amini, M.-R. (2019). Transductive Bounds for the Multi-Class Majority Vote Classifier. In *Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence and Thirty-First Innovative Applications of Artificial Intelligence Conference and Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, AAAI'19/IAAI'19/EAAI'19*. AAAI Press.



- Feofanov, V., Tiomoko, M., and Virmaux, A. (2023). Random Matrix Analysis to Balance between Supervised and Unsupervised Learning under the Low Density Separation Assumption.
- Lee, D.-H. (2013). Pseudo-Label : The Simple and Efficient Semi-Supervised Learning Method for Deep Neural Networks. *ICML 2013 Workshop : Challenges in Representation Learning (WREPL)*.
- van Engelen, J. E. and Hoos, H. H. (2020). A survey on semi-supervised learning. *Machine Learning*, 109(2):373–440.



# Datasets

Dataset	Size	# of lab. examples $n_\ell$	Dimension $d$	# classes $C$
Cod-RNA	59535	99	8	2
COIL-20	1440	200	1024	20
Digits	1797	99	64	10
DNA	3186	149	180	6
DryBean	13543	104	16	7
HAR	10299	299	561	3
Mnist	70000	100	784	10
Mushrooms	8124	79	112	2
Phishing	11055	99	68	2
Protein	1080	80	77	8
Rice	3810	29	7	2
Splice	3175	39	20	2
Svmguide1	3089	39	4	2

Table: Characteristics of datasets used in our experiments