

UNLABELED OUT-OF-DOMAIN DATA IMPROVES GENERALIZATION

Amir Hossein Saberi Amir Najafi Alireza Heidari
 Mohammad Hosein Movasaghinia Seyed Abolfazl Motahari Babak H. Khalaj
 Sharif University of Technology



MOTIVATION

The sample complexity for a linear binary classifier in a non-realizable setting in \mathbb{R}^d is $O(d/\epsilon^2)$. It has been shown that out-of-domain unlabeled data can enhance generalization, but explicit bounds for the size of labeled and unlabeled data remain incomplete.

CONTRIBUTION

We propose a polynomial-time framework that leverages both labeled and slightly out-of-domain unlabeled data. Our framework guarantees improved generalization under the *cluster assumption* of the true data distribution. In the well-studied setting of the two-component Gaussian Mixture Model (GMM) for classification, with m labeled and n unlabeled data points, our theoretical findings demonstrate:

- **Non-asymptotic bounds** for both robust and non-robust learning.
- **Enhanced generalization** over ERM techniques when $n \geq \Omega(m^2/d)$.
- **Dimension-independent** sample complexity under well-defined conditions.
- **Improved sample complexity** from $O(d/\epsilon^2)$ to $O(d/\epsilon)$ when $n = O(d/\epsilon^6)$.

Preliminaries

Definition A.1 (Wasserstein Distance). Given P and Q supported on \mathcal{X} , and a non-negative, lower semi-continuous cost function $c : \mathcal{X}^2 \rightarrow \mathbb{R}_+$ satisfying $c(\mathbf{X}, \mathbf{X}) = 0$ for all $\mathbf{X} \in \mathcal{X}$, is defined as:

$$\mathcal{W}_c(P, Q) = \inf_{\mu \in \Gamma(\mathcal{X}^2)} \mathbb{E}_{\mathbf{X}, \mathbf{X}' \sim \mu} [c(\mathbf{X}, \mathbf{X}')] \text{ s.t. } \mu(\mathbf{X}, \cdot) = P, \mu(\cdot, \mathbf{X}') = Q.$$

Definition A.2 (ϵ -neighborhood of a Distribution P).

$$\mathcal{B}_\epsilon(P) = \{Q : \mathcal{W}_c(P, Q) \leq \epsilon\}. \quad (1)$$

Definition 1.1 (Distributionally Robust Learning(DRL)). In DRL, the *learner* attempts to find a classifier with a small robust risk, denoted as $R^{\text{robust}}(\theta, P)$, where:

$$R_{\epsilon, c}^{\text{robust}}(\theta, P) = \sup_{P' \in \mathcal{B}_\epsilon(P)} R(\theta, P'); \forall \theta \in \Theta \& \epsilon \geq 0. \quad (2)$$

Therefore, DRL solves the following minimax optimization problem:

$$\hat{\theta}_{\epsilon, c}^{\text{DRL}}(S) \triangleq_{\theta \in \Theta} R_{\epsilon, c}^{\text{robust}}(\theta, \hat{P}_S^m). \quad (3)$$

PROBLEM SETUP

Notations & Definitions:

- Let $\mathcal{X} \subseteq \mathbb{R}^d$, $y \in \{\pm 1\}$. The joint feature-label distribution of labeled data is denoted as $P_0 \sim \mathcal{X} \times \{\pm 1\}$.
- P_1 is a shifted version of P_0 , where their marginal distributions on \mathcal{X} are shifted by $\mathcal{W}_c(P_{0,X}, P_{1,X}) = \alpha$ with $\alpha > 0$, and no assumptions are made on $P_1(y | \mathbf{X})$.
- **Labeled Data:** $S_0 = \{(\mathbf{X}_i, y_i)\}_{i=1}^m \sim P_0^m$, where $P_0(y=1) = \frac{1}{2}$.
- **Out-of-domain Unlabeled Data:** $S_1 = \{\mathbf{X}'_i\}_{i=1}^n \sim P_{1,X}^n$.

Isotropic GMM Setup:

- **Loss Function:** Without loss of generality, let $\ell(\mathbf{X}, y; \theta) = \mathbf{1}(y\langle \theta, \mathbf{X} \rangle \geq 0)$.
- $P_0(\mathbf{X} | y) = \mathcal{N}(y\mu_0, \sigma_0^2)$ for some $\sigma_0 \geq 0$, $\mu_0 \in \mathbb{R}^d$.
- $P_{1,X} = \frac{1}{2} \sum_{u=\pm 1} \mathcal{N}(u\mu_1, \sigma_1^2)$,

where $\|\mu_0 - \mu_1\| \leq \mathcal{O}(\alpha)$ and $|\sigma_1 - \sigma_0| \leq \mathcal{O}(\alpha)$ since $\mathcal{W}_c(P_{0,X}, P_{1,X}) = \alpha$.

Framework: Robust Self-Supervised Training

Robust Self-Supervised (RSS) Training: Consider a cost function c and a parameter $\gamma \geq 0$. We define the *robust loss* $\phi_\gamma : \mathcal{X} \times \{\pm 1\} \times \Theta \rightarrow \mathbb{R}$ by:

$$\phi_\gamma(\mathbf{X}, y; \theta) \triangleq \sup_{\mathbf{Z} \in \mathcal{X}} \ell(\mathbf{Z}, y; \theta) - \gamma c(\mathbf{Z}, \mathbf{X}). \quad (4)$$

For a given set of parameters $\gamma, \gamma', \lambda \in \mathbb{R}_{\geq 0}$: $\hat{\theta}^{\text{RSS}}$ estimator is defined as:

$$\hat{\theta}^{\text{RSS}} \triangleq_{\theta \in \Theta} \left\{ \frac{1}{m} \sum_{i=1}^m \phi_\gamma(\mathbf{X}_i, y_i; \theta) + \frac{\lambda}{n} \sum_{j=1}^n \phi_{\gamma'}(\mathbf{X}'_j, h_\theta(\mathbf{X}'_j); \theta) \right\}. \quad (5)$$

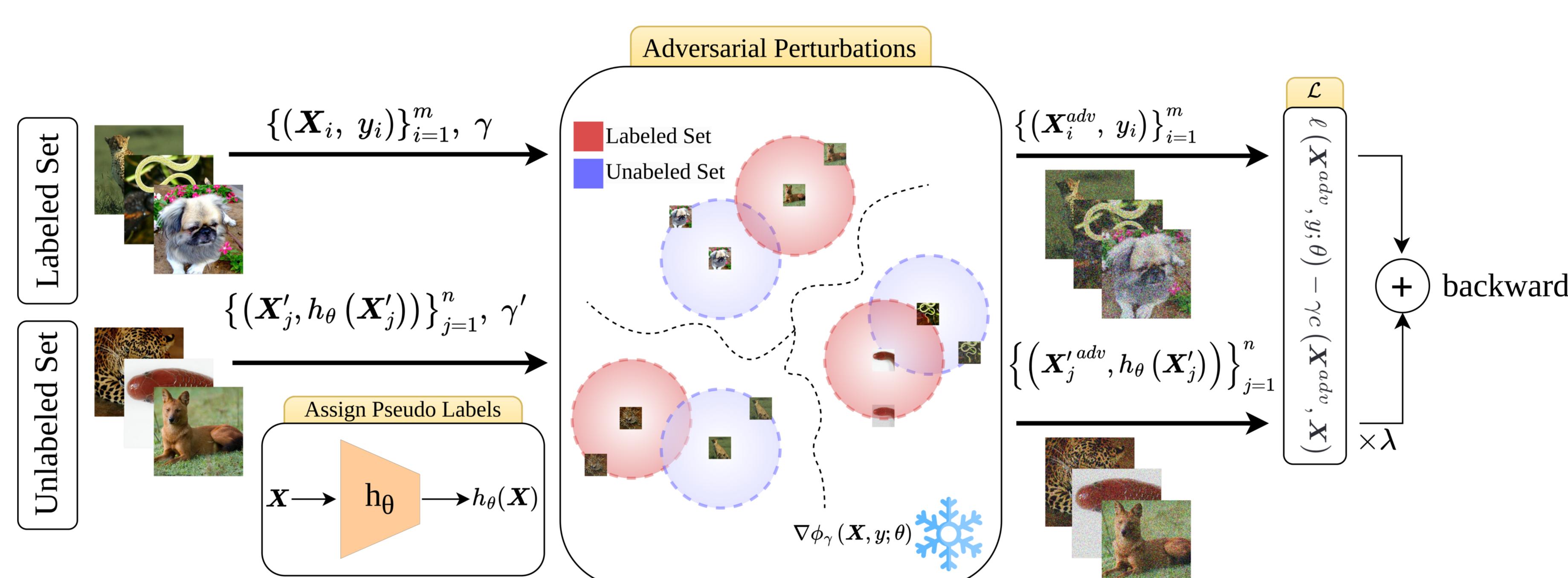


Figure 1: Overview of Robust Self-Supervised (RSS) Training

Generalization Bounds (GMM)

Theorem 4.1 (Non-asymptotic Bound for Robust Learning).

$$\mathbb{E}_{P_0} [\phi_\gamma(\mathbf{X}, y; \hat{\theta}^{\text{RSS}})] \leq \min_{\theta \in \Theta} \mathbb{E}_{P_0} [\phi_\gamma(\mathbf{X}, y; \theta)] + \mathcal{O} \left(\gamma \sqrt{\frac{2d}{m} \left(\alpha(\|\mu_0\|_2^2 + \sigma_0^2) + \sqrt{\frac{2d}{2n+m}} + \sqrt{\frac{2 \log(1/\delta)}{2n+m}} \right)} + \sqrt{\frac{2 \log(1/\delta)}{m}} \right).$$

Theorem 4.2 (Non-asymptotic Bound for Non-robust Learning).

$$R(\hat{\theta}^{\text{RSS}}, P) - \min_{\theta \in \Theta} R(\theta, P) \leq \mathcal{O} \left(\frac{e^{-\frac{1}{4\sigma_0^2}}}{\sqrt{2\sigma_0\sqrt{2\pi}}} \left((\|\mu_1\|_2^2 + \sigma_1^2) \frac{2d\alpha}{m} + \frac{4d}{m} \sqrt{\frac{2d + 2 \log \frac{1}{\delta}}{2n+m}} \right)^{1/4} + \sqrt{\frac{2 \log \frac{1}{\delta}}{m}} \right).$$

Corollary (Following Theorem 4.2):

- The $\hat{\theta}^{\text{RSS}}$ estimator outperforms ERM when $\alpha \leq O(\frac{d}{m})$, $n \geq \Omega(\frac{m^2}{d})$.
- Sample complexity becomes independent of d if: $\alpha \leq O(d^{-1})$, $n \geq \Omega(d^3)$.
- When $\alpha = 0$ (no perturbation), $m = O(\frac{d}{\epsilon})$, and $n = O(\frac{d}{\epsilon^6})$, the generalization bound improves compared to having access to $m = O(\frac{d}{\epsilon^2})$ labeled data points.

Experiments & Results

Simulated Isotropic GMM Data (non-robust setup):

- **Data:** $d = 200$, $\mu_1 = \mu_0 + \alpha \cdot v$ where v and μ_0 are random normalized vectors.
- **Model:** $h(x; w) = w^T x$, with $l_\gamma(Y, X, f(x)) = \sum_{i=1}^{|X|} \min(1, \max(0, 1 - \gamma \cdot Y_i \cdot f(X_i)))$.

Same Distribution ($\alpha = 0$)				Different Distribution ($\alpha = 0.5 \cdot \ \mu_0\ _2$)			
m	Acc	n	Acc	m	Acc	n	Acc
10	0.59	10	0.63	10	0.59	10	0.61
		100	0.66			100	0.65
		1,000	0.79			1,000	0.78
		10,000	0.82			10,000	0.81
20	0.62	20	0.64	20	0.62	20	0.65
		200	0.69			200	0.65
		2,000	0.80			2,000	0.79
		10,000	0.82			10,000	0.80
40	0.65	40	0.65	40	0.65	40	0.65
		400	0.71			400	0.73
		4,000	0.81			4,000	0.78
		10,000	0.82			10,000	0.80
10,000	0.83	-	-	10,000	0.83	-	-

Histopathology Images (robust setup):

- **Data:** NCT-CRC-HE-100K (colorectal cancer) & PatchCamelyon (lymph node).

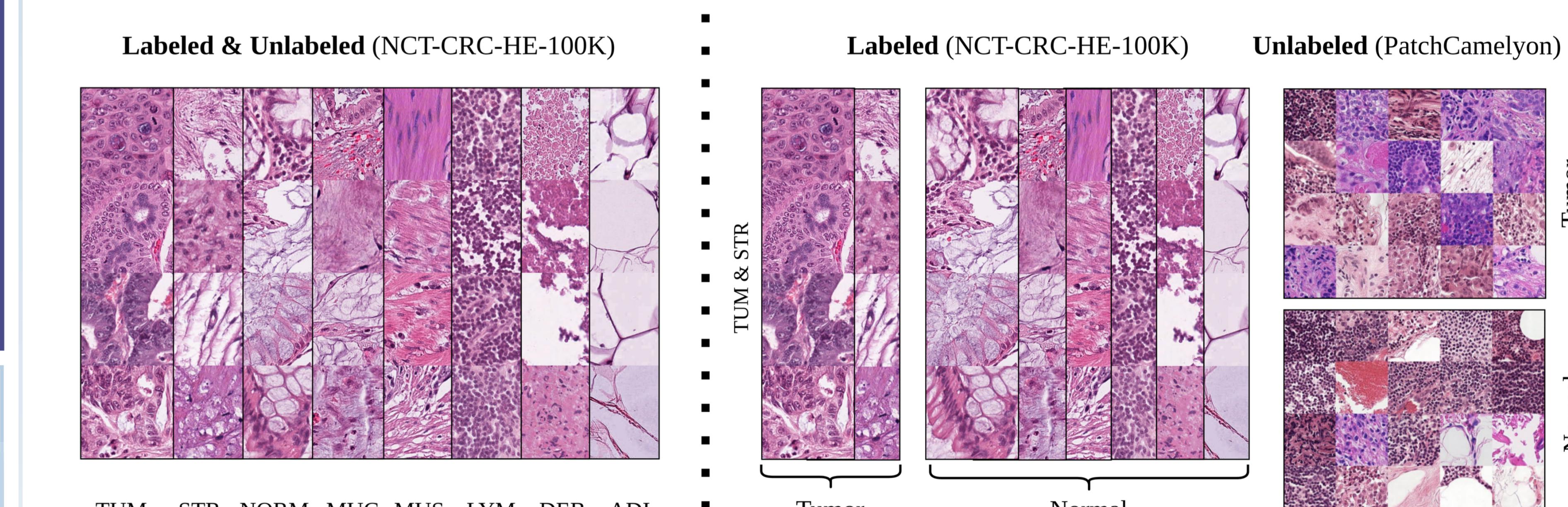


Figure 2: (Left) Same Distribution Setting. (Right) Different Distribution Setting.

Same distribution				Different distribution			
Labeled size	Acc	Unlabeled size	Acc	Labeled size	Acc	Unlabeled size	Acc
48	0.65	200	0.71	25	0.78	100	0.78
		700	0.80			400	0.79
		2,000	0.82			2,000	0.81
240	0.77	500	0.78	50	0.82	200	0.82
		1,200	0.82			700	0.86
		4,000	0.83			3,000	0.87
1040	0.83	3,000	0.87	300	0.87	600	0.88
		10,000	0.89			2,000	0.89
		20,000	0.91			8,000	0.90
50,000	0.916	-	-	32,000	0.94	-	-