

**Vulnerability of machine learning methods to adversarial attack:** Traditional deep learning methods are known to be **vulnerable to attacks** where a **small perturbation** of the input (e.g., an image) that is **imperceptible** to a human observer can **cause a trained classifier to fail**.

**Original (PGD) adversarial training method:** PGD<sup>2</sup> trains the loss  $\mathcal{L}_\theta$  by **perturbing training samples**  $z$  within a metric-space ball of size  $\epsilon$  to create **adversarial samples**  $\tilde{z}$  that are then used in training:

$$\inf_{\theta} E_{P_n} \left[ \sup_{\tilde{z}: d(z, \tilde{z}) \leq \epsilon} \mathcal{L}_\theta(\tilde{z}) \right],$$

where  $P_n$  denotes the empirical distribution of the training samples.

**PGD is an example of distributionally robust optimization (DRO):** In DRO the empirical distribution is replaced by the **worst-case adversarial distribution**  $Q$  in some model neighborhood  $\mathcal{U}(P_n)$  around  $P_n$ :

$$\inf_{\theta} \sup_{Q \in \mathcal{U}(P_n)} E_Q[\mathcal{L}_\theta].$$

**Robustness is increased by training for the worst case.**

When using adversarial training, **the choice of model neighborhood is important for performance!**

# Adversarially Robust Learning with Optimal-Transport Regularized Divergences

**ARMOR<sub>D</sub> Method<sup>1</sup>:** Our DRO-based approach **improves model robustness** by both **adversarially transporting** (via an optimal transport cost) and **adversarially re-weighting** (via an information divergence) samples during training.

ARMOR<sub>D</sub> can be **combined** with other popular methods, e.g., that modify the training loss, such as PGD<sup>2</sup>, TRADES<sup>3</sup>, MART<sup>4</sup>, and UDR<sup>5</sup> to yield **improved performance** when under adversarial attack:

Defense	CIFAR10 Performance		
	AutoAttack	PGD <sup>200</sup>	Nat.
<i>PGD</i>	42.5%	46.0%	<b>86.40%</b>
<i>UDR-PGD</i>	48.47%	52.95%	81.71%
<i>ARMOR<sub>α</sub>-UDR-PGD</i>	<b>48.63%</b>	<b>53.62%</b>	80.29%
<i>TRADES</i>	49.1%	51.9%	80.8%
<i>UDR-TRADES</i>	49.9%	53.6%	<b>84.4%</b>
<i>ARMOR<sub>α</sub>-TRADES</i>	<b>51.4%</b>	<b>53.74%</b>	80.76%
<i>MART</i>	48.2%	53.3%	<b>81.9%</b>
<i>UDR-MART</i>	49.1%	54.1%	80.1%
<i>ARMOR<sub>α</sub>-MART</i>	<b>50.6%</b>	<b>56.22%</b>	81.03%

## Optimal-Transport Regularized Divergences:

$$D^c(\nu \parallel \mu) := \inf_{\eta \in \mathcal{P}(\mathcal{Z})} \{D(\eta \parallel \mu) + C(\eta, \nu)\}$$

$C$  is an **optimal transport cost** for a cost function  $c$ ,

$$C(\mu, \nu) := \inf_{\pi: \pi_1 = \mu, \pi_2 = \nu} \int c(z, \tilde{z}) \pi(dz d\tilde{z})$$

$D$  is an **information divergence**, e.g., an  $f$ -divergence,

$$D_f(\mu \parallel \nu) = E_\nu[f(d\mu/d\nu)].$$

**Properties (under appropriate assumptions):**

- **Divergence property:**  $D^c(\nu \parallel \mu) \geq 0$  and  $D^c(\nu \parallel \mu) = 0$  if and only if  $\nu = \mu$ . This implies  $D^c(\nu \parallel \mu)$  **quantifies the discrepancy** between  $\nu$  and  $\mu$ .

- **Optimizer:** there **exists a unique optimizer**,  $\eta_*$ , with

$$D^c(\nu \parallel \mu) = D(\eta_* \parallel \mu) + C(\eta_*, \nu)$$

- **DRO neighborhoods:** The DRO neighborhoods

$$\{Q : D^c(Q \parallel P_n) \leq \epsilon\}$$

- **Interpolation property:**  $D^c$  **interpolates between  $D$  and  $C$**  as follows

$$\lim_{r \rightarrow 0^+} r^{-1} D^{rc}(\nu \parallel \mu) = C(\mu, \nu), \quad \lim_{r \rightarrow \infty} D^{rc}(\nu \parallel \mu) = D(\nu \parallel \mu)$$

**Computationally-Tractable Dual-Formulation of Adversarially-Robust Training Problem:** We use the  $D^c$ -DRO neighborhoods to obtain a **novel adversarial training method** and, via convex duality, obtain the **computationally tractable form**

$$\inf_{\theta \in \Theta} \sup_{Q: D_f^c(Q \parallel P_n) \leq \epsilon} E_Q[\mathcal{L}_\theta] = \inf_{\lambda > 0, \rho \in \mathbb{R}, \theta \in \Theta} \{\epsilon \lambda + \rho + \lambda E_{P_n} [f^*(\lambda^{-1}(\mathcal{L}_{\theta, \lambda}^c(z_i) - \rho))]\}$$

where

$$\mathcal{L}_{\theta, \lambda}^c(z) := \sup_{\tilde{z} \in \mathcal{Z}} \{\mathcal{L}_\theta(\tilde{z}) - \lambda c(z, \tilde{z})\}.$$

Our work generalizes the optimal-transport DRO results of [6], [7], and is related to the DRO method [8].

### References:

- [1] J. Birrell, M. Ebrahimi, arXiv:2309.03791, 2023
- [2] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, ICLR, 2018
- [3] H. Zhang, Y. Yu, J. Jiao, E. Xing, L. El Ghaoui, and M. Jordan, ICML, 2019
- [4] Y. Wang, D. Zou, J. Yi, J. Bailey, X. Ma, and Q. Gu, ICLR, 2020
- [5] A. T. Bui, T. Le, Q. H. Tran, H. Zhao, and D. Phung, ICLR, 2022
- [6] P. Mohajerin Esfahani and D. Kuhn, Mathematical Programming, 2018
- [7] J. Blanchet and K. Murthy, Mathematics of Operations Research, 2019
- [8] J. Blanchet, D. Kuhn, J. Li, and B. Taskesen, arXiv:2308.05414, 2023