



Robustness and Accuracy Could Be Reconcilable by (Proper) Definition

Tianyu Pang^{1 2}, Min Lin¹, Xiao Yang², Jun Zhu², Shuicheng Yan¹

ICML | 2022



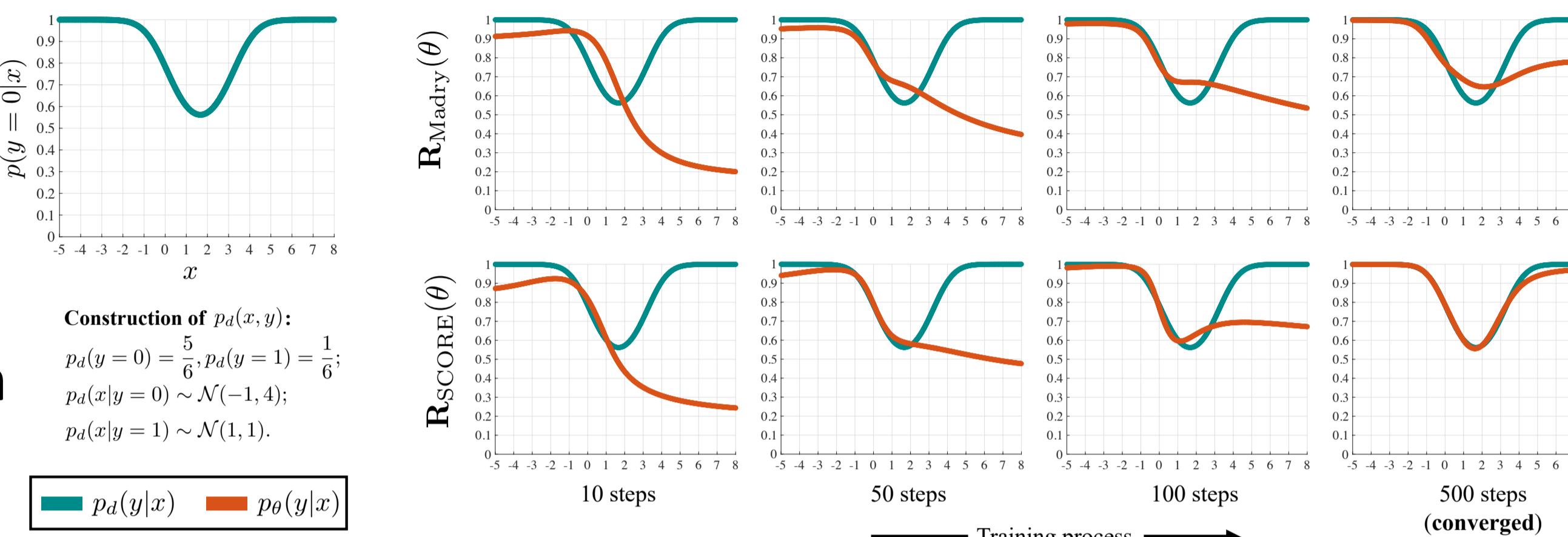
- What is an **accurate** model?

An **accurate** model has **low standard error**:

$$R_{\text{Standard}} = \mathbb{E}_{p_d(x)} [\text{KL} (p_d(y|x) \| p_\theta(y|x))]$$

data distribution **model distribution**

Optimal solution: $p_{\theta^*}(y|x) = p_d(y|x)$



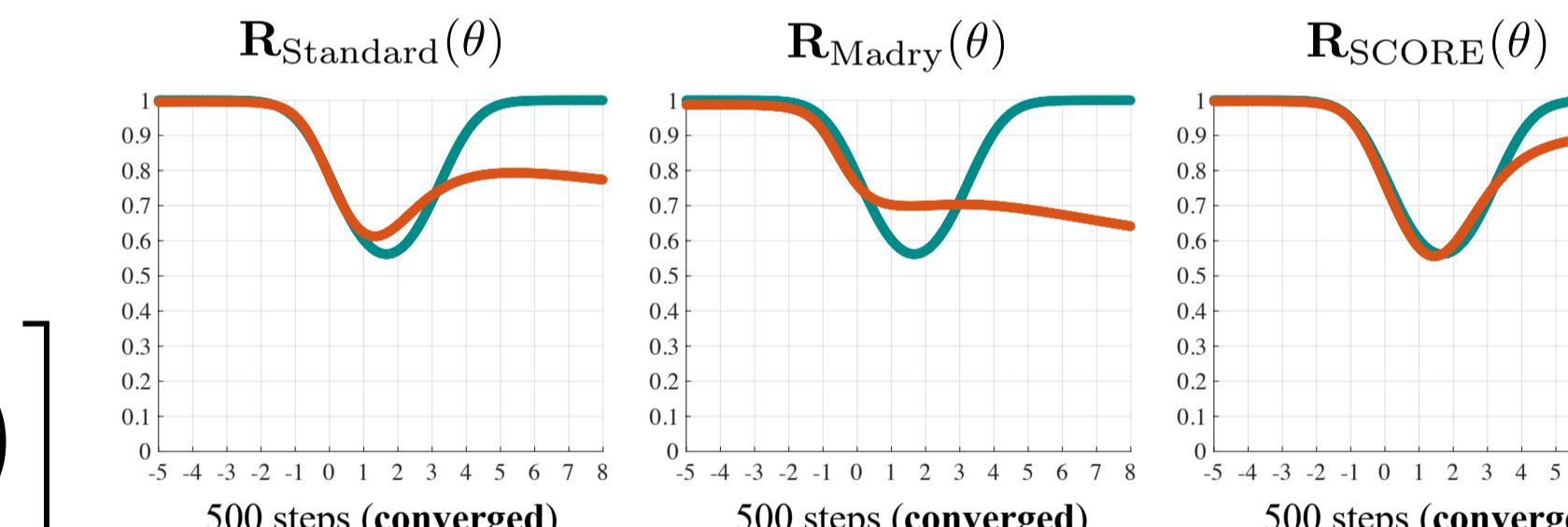
- What is a **robust** model?

A **robust** model has **low robust error**:

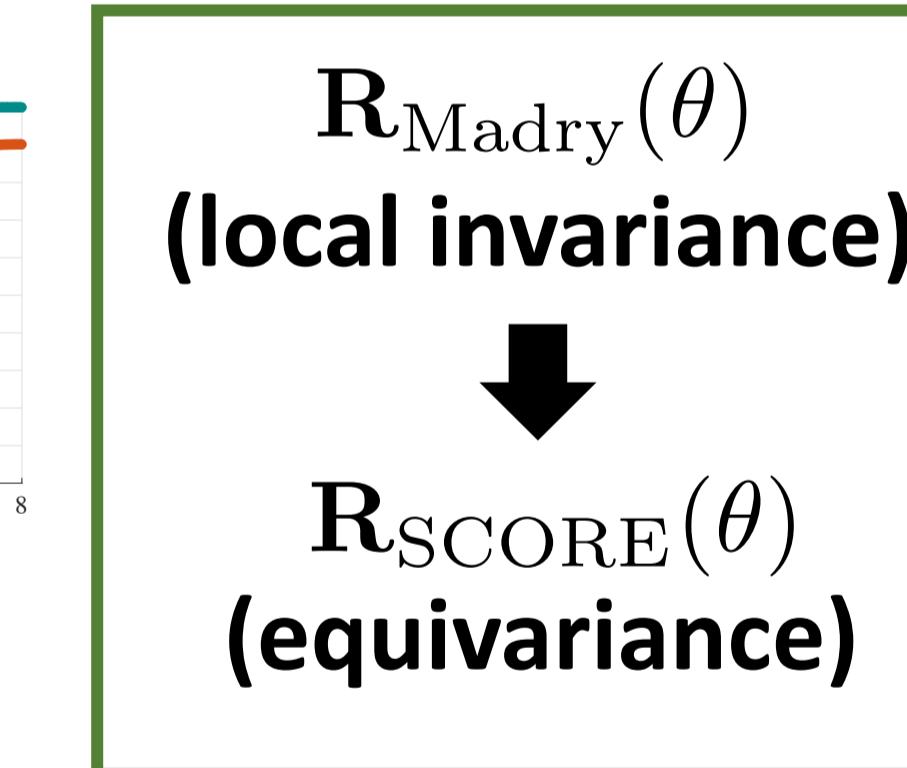
(Madry et al. ICLR 2018)

$$R_{\text{Madry}} = \mathbb{E}_{p_d(x)} \left[\max_{x' \in B(x)} \text{KL} (p_d(y|x) \| p_\theta(y|x')) \right]$$

Optimal solution: $p_{\theta^*}(y|x) \neq p_d(y|x)$



6 training pairs, mimics the finite-sample form



- Self-COnsistent Robust Error (**SCORE**)

$$R_{\text{SCORE}}(\theta) = \mathbb{E}_{p_d(x)} \left[\max_{x' \in B(x)} \text{KL} (p_d(y|x') \| p_\theta(y|x')) \right]$$

(1) **Self-consistency**: $p_{\theta^*}(y|x) = p_d(y|x)$

(2) Keep the paradigm of **robust optimization**

- How to practically optimize SCORE?

Substitute KL divergence with any **distance metric** \mathcal{D}

$$R_{\text{Madry}}^{\mathcal{D}}(\theta) = \mathbb{E}_{p_d(x)} \left[\max_{x' \in B(x)} \mathcal{D} (p_d(y|x) \| p_\theta(y|x')) \right];$$

$$R_{\text{SCORE}}^{\mathcal{D}}(\theta) = \mathbb{E}_{p_d(x)} \left[\max_{x' \in B(x)} \mathcal{D} (p_d(y|x') \| p_\theta(y|x')) \right]$$

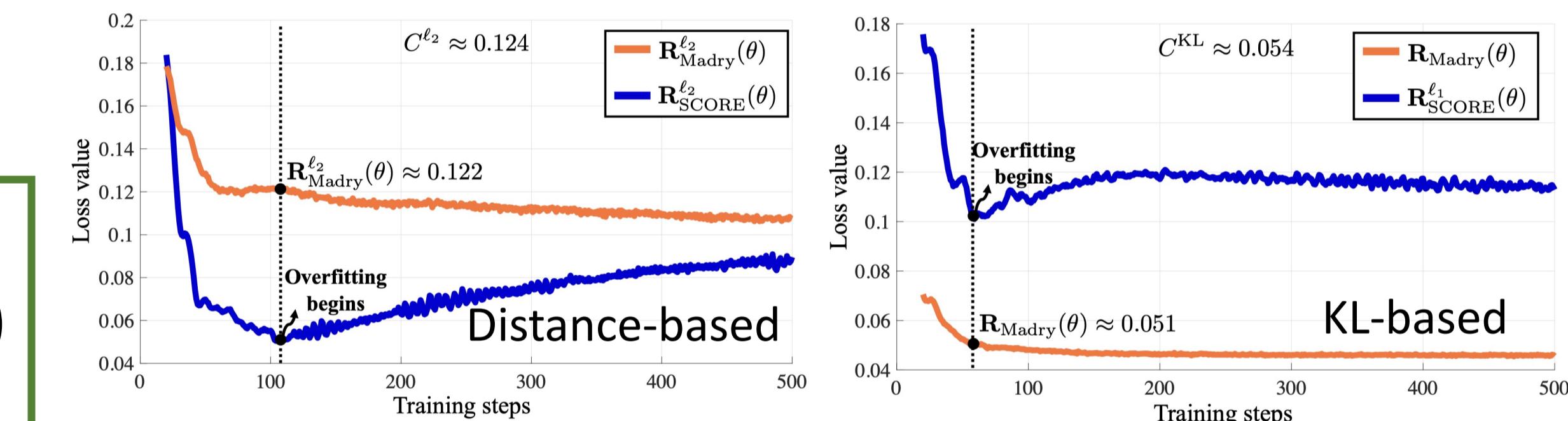
- Upper and lower bounds for SCORE

$$|R_{\text{Madry}}^{\mathcal{D}}(\theta) - C^{\mathcal{D}}| \leq R_{\text{SCORE}}^{\mathcal{D}}(\theta) \leq R_{\text{Madry}}^{\mathcal{D}}(\theta) + C^{\mathcal{D}},$$

$$\text{where } C^{\mathcal{D}} = \mathbb{E}_{p_d(x)} \left[\max_{x' \in B(x)} \mathcal{D} (p_d(y|x) \| p_d(y|x')) \right]$$

Upper bound: without estimating $\nabla_x \log p_d(y|x)$

Lower bound: indicates overfitting phenomena



Dataset	Method	Architecture	DDPM	Batch	Epoch	Clean	AutoAttack
CIFAR-10 ($\ell_\infty, \epsilon = 8/255$)	Rice et al. (2020)	WRN-34-20	x	128	200	85.34	53.42
	Zhang et al. (2020)	WRN-34-10	x	128	120	84.52	53.51
	Pang et al. (2021)	WRN-34-20	x	128	110	86.43	54.39
	Wu et al. (2020)	WRN-34-10	x	128	200	85.36	56.17
	Gowal et al. (2020)	WRN-70-16	x	512	200	85.29	57.14
	Rebuffi et al. (2021) [†]	WRN-28-10	1M	1024	800	85.97	60.73
	+ Ours (KL → SE, $\beta = 3$)	WRN-28-10	1M	512	400	88.61	61.04
	+ Ours (KL → SE, $\beta = 4$)	WRN-28-10	1M	512	400	88.10	61.51
	Rebuffi et al. (2021) [†]	WRN-70-16	1M	1024	800	86.94	63.58
	+ Ours (KL → SE, $\beta = 3$)	WRN-70-16	1M	512	400	89.01	63.35
	+ Ours (KL → SE, $\beta = 4$)	WRN-70-16	1M	512	400	88.57	63.74
	Gowal et al. (2021)	WRN-70-16	100M	1024	2000	88.74	66.10
	Wu et al. (2020)	WRN-34-10	x	128	200	60.38	28.86
	Gowal et al. (2020)	WRN-70-16	x	512	200	60.86	30.03
CIFAR-100 ($\ell_\infty, \epsilon = 8/255$)	Rebuffi et al. (2021) [†]	WRN-28-10	1M	1024	800	59.18	30.81
	+ Ours (KL → SE, $\beta = 3$)	WRN-28-10	1M	512	400	63.66	31.08
	+ Ours (KL → SE, $\beta = 4$)	WRN-28-10	1M	512	400	62.08	31.40
	Rebuffi et al. (2021) [†]	WRN-70-16	1M	1024	800	60.46	33.49
	+ Ours (KL → SE, $\beta = 3$)	WRN-70-16	1M	512	400	65.56	33.05
	+ Ours (KL → SE, $\beta = 4$)	WRN-70-16	1M	512	400	63.99	33.65

Find more interesting conclusions in our paper!