

# Boosting Adversarial Training via Fisher-Rao Norm-based Regularization

## Supplementary Material

In this appendix, we provide detailed descriptions of the methodologies and experimental procedures used in our study, which encompasses:

- Theoretical proofs in detail.
- Additional experimental settings and results.
- Limitations and future directions.

The aim is to ensure the transparency and reproducibility of our research, while providing sufficient information for readers interested in the technical intricacies of our work.

### 8. Proofs of Lemma. 2 and Thm. 1

According to Lemma. 1, we can denote the Fisher-Rao norm of  $\mathcal{W}$  w.r.t  $\mathcal{L}_{ce}$  as:

$$\|\mathcal{W}\|_{F^{R \circ \mathcal{L}_{ce}}}^2 = L^2 \mathbb{E}_{(\mathbf{x}, y)} \left[ \left( \langle \sigma(f_{\mathcal{W}}^L(\mathbf{x})), f_{\mathcal{W}}^L(\mathbf{x}) \rangle - f_{\mathcal{W}}^L(\mathbf{x})_y \right)^2 \right] \quad (16)$$

Then we can conclude that

$$\begin{aligned} \gamma_{ce} &= \mathbb{E}_{(\mathbf{x}, y)} \left[ \left| \langle \sigma(f_{\mathcal{W}}^L(\mathbf{x})), f_{\mathcal{W}}^L(\mathbf{x}) \rangle - f_{\mathcal{W}}^L(\mathbf{x})_y \right| \right] \\ &= \mathbb{E}_{(\mathbf{x}, y)} \left[ \left| \frac{\sum_{k=1}^K (f_{\mathcal{W}}^L(\mathbf{x})_k - f_{\mathcal{W}}^L(\mathbf{x})_y) e^{f_{\mathcal{W}}^L(\mathbf{x})_k}}{\sum_{k=1}^K e^{f_{\mathcal{W}}^L(\mathbf{x})_k}} \right| \right] \\ &= \mathbb{E}_{(\mathbf{x}, y)} \left[ \left| \frac{\sum_{k \neq y} (f_{\mathcal{W}}^L(\mathbf{x})_k - f_{\mathcal{W}}^L(\mathbf{x})_y) e^{f_{\mathcal{W}}^L(\mathbf{x})_k}}{\sum_{k=1}^K e^{f_{\mathcal{W}}^L(\mathbf{x})_k}} \right| \right] \\ &\leq \mathbb{E}_{(\mathbf{x}, y)} \left[ \frac{\sum_{k \neq y} |f_{\mathcal{W}}^L(\mathbf{x})_k - f_{\mathcal{W}}^L(\mathbf{x})_y| e^{f_{\mathcal{W}}^L(\mathbf{x})_k}}{\sum_{k \neq y} e^{f_{\mathcal{W}}^L(\mathbf{x})_k}} \right] \\ &\leq \mathbb{E}_{(\mathbf{x}, y)} \left[ \max_{k \neq y} |f_{\mathcal{W}}^L(\mathbf{x})_k - f_{\mathcal{W}}^L(\mathbf{x})_y| \right] \end{aligned} \quad (17)$$

Then the standard Rademacher complexity w.r.t the CE loss  $\mathcal{R}_{N_{tr}}(\mathcal{L}_{ce} \circ \mathcal{F}_{\hat{\gamma}_{ce}})$  can be denoted as:

$$\begin{aligned} &\mathbb{E}_{\xi} \sup_{f_{\mathcal{W}}^L \in \mathcal{F}_{\hat{\gamma}_{ce}}} \frac{1}{N} \sum_{i=1}^N \xi_i \mathcal{L}_{ce}(f_{\mathcal{W}}^L(\mathbf{x}_i), y_i) \\ &= \mathbb{E}_{\xi} \sup_{f_{\mathcal{W}}^L \in \mathcal{F}_{\hat{\gamma}_{ce}}} \frac{1}{N} \sum_{i=1}^N \xi_i \ln \left( \frac{e^{f_{\mathcal{W}}^L(\mathbf{x}_i)_y}}{\sum_{k=1}^K e^{f_{\mathcal{W}}^L(\mathbf{x}_i)_k}} \right) \\ &= \mathbb{E}_{\xi} \sup_{f_{\mathcal{W}}^L \in \mathcal{F}_{\hat{\gamma}_{ce}}} \frac{1}{N} \sum_{i=1}^N \xi_i \left( \ln \left( \sum_{k=1}^K e^{f_{\mathcal{W}}^L(\mathbf{x}_i)_k} \right) - f_{\mathcal{W}}^L(\mathbf{x}_i)_y \right) \\ &= \mathbb{E}_{\xi} \sup_{f_{\mathcal{W}}^L \in \mathcal{F}_{\hat{\gamma}_{ce}}} \frac{1}{N} \sum_{i=1}^N \xi_i \left( \ln \left( \sum_{k=1}^K \left( \frac{1}{e} \right)^{f_{\mathcal{W}}^L(\mathbf{x}_i)_y - f_{\mathcal{W}}^L(\mathbf{x}_i)_k} \right) \right) \end{aligned} \quad (18)$$

Specifically, for  $N_{tr}^M$  misclassified samples, the upper bound for  $\mathcal{R}_{N_{tr}^M}(\mathcal{L}_{ce} \circ \mathcal{F}_{\hat{\gamma}_{ce}})$  can be further derived as:

$$\begin{aligned} &\mathcal{R}_{N_{tr}^M}(\mathcal{L}_{ce} \circ \mathcal{F}_{\hat{\gamma}_{ce}}) \\ &\leq \mathbb{E}_{\xi} \sup_{f_{\mathcal{W}}^L \in \mathcal{F}_{\hat{\gamma}_{ce}}} \frac{1}{N_{tr}^M} \sum_{i=1}^{N_{tr}^M} \xi_i \left( \ln \left( \sum_{k=1}^K \left( \frac{1}{e} \right)^{f_{\mathcal{W}}^L(\mathbf{x}_i)_y - f_{\mathcal{W}}^L(\mathbf{x}_i)_k} \right) \right) \\ &\leq \mathbb{E}_{\xi} \frac{1}{N_{tr}^M} \sum_{i=1}^{N_{tr}^M} \xi_i \left( \ln \left( \sum_{k=1}^K \left( \frac{1}{e} \right)^{-|f_{\mathcal{W}}^L(\mathbf{x}_i)_y - f_{\mathcal{W}}^L(\mathbf{x}_i)_k|} \right) \right) \\ &\leq (\ln K + \hat{\gamma}_{ce}) \sqrt{\frac{1}{N_{tr}^M} \sum_{i=1}^{N_{tr}^M} \xi_i} \\ &\leq (\ln K + \hat{\gamma}_{ce}) \sqrt{\frac{1}{N_{tr}^M} \sum_{i=1}^{N_{tr}^M} \xi_i} \\ &\stackrel{(i)}{\lesssim} \frac{\ln K + \hat{\gamma}_{ce}}{\sqrt{N_{tr}^M}} \end{aligned} \quad (19)$$

Notice that it is very likely that  $0 < \left| \mathbb{E}_{\xi} \sum_{i=1}^{N_{tr}^M} \xi_i \right| < N_{tr}^M$ , we then assume  $\left| \mathbb{E}_{\xi} \sum_{i=1}^{N_{tr}^M} \xi_i \right| \approx \sqrt{N_{tr}^M}$ . Due to the fact that we are deriving the upper bound, (i) holds. Meanwhile, for  $N_{tr}^C$  correctly classified samples, the following inequalities holds:

$$\begin{aligned} &\mathcal{R}_{N_{tr}^C}(\mathcal{L}_{ce} \circ \mathcal{F}_{\hat{\gamma}_{ce}}) \\ &\geq \mathbb{E}_{\xi} \frac{1}{N_{tr}^C} \sum_{i=1}^{N_{tr}^C} \xi_i \left( \ln \left( \sum_{k=1}^K \left( \frac{1}{e} \right)^{\hat{\gamma}_{ce}} \right) \right) \\ &\geq \frac{(\ln K - \hat{\gamma}_{ce})}{N_{tr}^C} \mathbb{E}_{\xi} \sum_{i=1}^{N_{tr}^C} \xi_i \end{aligned} \quad (20)$$

Similar to Eq. 19, we assume  $\left| \mathbb{E}_{\xi} \sum_{i=1}^{N_{tr}^C} \xi_i \right| \approx \sqrt{N_{tr}^C}$ . Meanwhile, due to variations in network architecture, training data, and training algorithm, etc, the value of  $\hat{\gamma}_{ce}$  tends to vary a lot, therefore, we need to make  $\mathbb{E}_{\xi} \sum_{i=1}^{N_{tr}^C} \xi_i$  positive to reach the lower bound, which can be denoted as  $\mathcal{R}_{N_{tr}^C}(\mathcal{L}_{ce} \circ \mathcal{F}_{\hat{\gamma}_{ce}}) \gtrsim \frac{\ln K - \hat{\gamma}_{ce}}{\sqrt{N_{tr}^C}}$ . Consequently, we can deduce that:

$$\begin{aligned} \mathcal{R}_{N_{tr}}(\mathcal{L}_{ce} \circ \mathcal{F}_{\hat{\gamma}_{ce}}) &\leq \frac{(\sqrt{N_{tr}^M} + \sqrt{N_{tr}^C}) \ln K + \hat{\gamma}_{ce} \sqrt{N_{tr}^M}}{N_{tr}} \\ \mathcal{R}_{N_{tr}}(\mathcal{L}_{ce} \circ \mathcal{F}_{\hat{\gamma}_{ce}}) &\geq \frac{(\sqrt{N_{tr}^M} + \sqrt{N_{tr}^C}) \ln K - \hat{\gamma}_{ce} \sqrt{N_{tr}^C}}{N_{tr}} \end{aligned} \quad (21)$$

Then if we define  $\Gamma_{ce} = \frac{\hat{\gamma}_{ce}^{N_C} - \hat{\gamma}_{ce}^{N_M}}{\hat{\gamma}_{ce}^{N_C}}$ ,  $\mathcal{C}_C = \frac{N_{tr}}{N_C}$ ,  $\mathcal{C}_M = \frac{N_{tr}}{N_M}$ ,  $\mathcal{C}_{MC} = \frac{\sqrt{N_C} + \sqrt{N_M}}{N_{tr}}$ , we can easily provide the following bounds:

$$\begin{aligned} \mathcal{R}_{N_{tr}}(\mathcal{L}_{ce} \circ \mathcal{F}_{\hat{\gamma}_{ce}}) &\leq \mathcal{C}_{MC} \ln K + \frac{\hat{\gamma}_{ce}^{N_M} \mathcal{C}_M^{-0.5} (\mathcal{C}_C^{-1} \Gamma_{ce} + 1)}{N_{tr}^{0.5}} \\ \mathcal{R}_{N_{tr}}(\mathcal{L}_{ce} \circ \mathcal{F}_{\hat{\gamma}_{ce}}) &\geq \mathcal{C}_{MC} \ln K - \frac{\hat{\gamma}_{ce}^{N_M} \mathcal{C}_M^{-0.5} (\mathcal{C}_C^{-1} \Gamma_{ce} + 1)}{N_{tr}^{0.5}} \end{aligned} \quad (22)$$

## 9. Additional Experimental Settings and Results

**Experimental Settings.** The qualitative Results illustrated in Fig. 2 and 3 are derived from training models on a dataset comprising 4000 samples from MNIST. For the evaluations presented in Fig. 4, we utilized 4,000 training instances from MNIST and 10,000 from the CIFAR10. Furthermore, for all setups in LOAT-boosted training, we consistently set the parameters  $\tau$  to 1 and  $\gamma$  to 0.05, following the specifications of Alg. 1.

**Additional Results.** In addition to the LOAT-boosted S2O results on ResNet18, as detailed in Tab. 5, we extended our experiments to WideResNet-34-10. The com-

Table 6. Evaluation of LOAT-boosted S2O on WideResNet-34-10 (%).

Defense	Clean <sub>te</sub>	PGD <sup>7</sup>	PGD <sup>20</sup>	PGD <sup>40</sup>
PGD-AT	85.57	53.20	47.84	46.78
PGD-AT+SLORE	<b>86.97</b>	53.51	48.20	47.23
PGD-AT+LORE	86.69	<b>53.90</b>	<b>48.77</b>	<b>47.82</b>

Table 7. Comparison of SLORE-boosted PGD-AT with ALP-boosted one on ResNet-18 (%).

Defense	Clean <sub>te</sub>	FGSM	PGD <sup>7</sup>	PGD <sup>20</sup>
PGD-AT (ALP)	81.26	60.92	51.38	45.90
PGD-AT (SLORE)	<b>83.41</b>	<b>61.80</b>	<b>51.46</b>	<b>46.50</b>

prehensive outcomes of these additional experiments are presented in Tab. 6. Our findings indicate that SLORE-enhanced PGD-AT can enhance the test clean accuracy of WideResNet-34-10 by 1.40%. Furthermore, LORE-enhanced PGD-AT contributes to an increase in the model’s robustness against PGD<sup>7</sup> attacks by 0.70%, PGD<sup>20</sup> attacks by 0.93%, and PGD<sup>40</sup> attacks by 1.04%.

Furthermore, we compare our algorithm with widely used Adversarial Logit Pairing (ALP). Detailed results are shown in Tab. 7.

Additionally, we conducted evaluations of the LOAT-boosted DM-AT on augmented datasets from Cifar100,

Table 8. Classification Accuracy of models trained by  $1 \times 10^6$  EDM-generated images-augmented SVHN (%).

Models	Defense	Clean <sub>tr</sub>	Clean <sub>te</sub>	PGD <sup>20</sup>	PGD <sup>40</sup>	AA
PreResNet18 (Swish)	TRADES	96.59	96.57	56.74	53.84	36.21
	TRADES+SLORE	96.22	96.14	<b>61.20</b>	<b>57.85</b>	<b>38.58</b>
	TRADES+LORE	<b>97.79</b>	<b>97.39</b>	55.65	47.33	9.64
	TRADES(S)	96.06	96.20	69.91	66.64	<b>41.38</b>
	TRADES(S)+SLORE	<b>97.12</b>	<b>97.01</b>	<b>77.73</b>	<b>73.38</b>	13.48
	TRADES(S)+LORE	96.46	96.44	65.25	61.82	37.65
WRN-28-10 (Swish)	TRADES	97.45	<b>97.69</b>	59.57	54.75	15.17
	TRADES+SLORE	95.46	95.67	<b>68.35</b>	<b>67.01</b>	<b>49.25</b>
	TRADES+LORE	<b>97.61</b>	97.66	43.44	37.51	10.79
	TRADES(S)	<b>97.43</b>	<b>97.75</b>	<b>74.98</b>	<b>69.10</b>	15.77
	TRADES(S)+SLORE	97.11	97.31	70.17	65.19	11.27
	TRADES(S)+LORE	97.22	97.36	57.71	50.58	<b>16.36</b>
	MART	<b>97.23</b>	<b>97.23</b>	52.34	49.89	22.89
	MART+SLORE	88.13	93.97	64.22	63.82	<b>53.26</b>
	MART+LORE	94.70	96.30	<b>72.16</b>	<b>70.76</b>	9.03

Table 9. Classification Accuracy of models trained by  $1 \times 10^6$  EDM-generated images-augmented Cifar100 (%).

Models	Defense	Clean <sub>tr</sub>	Clean <sub>te</sub>	PGD <sup>20</sup>	PGD <sup>40</sup>	AA
PreResNet18 (Swish)	TRADES	85.27	62.33	34.62	<b>34.59</b>	<b>29.66</b>
	TRADES+SLORE	84.95	61.63	34.41	34.47	29.40
	TRADES+LORE	<b>86.35</b>	<b>62.79</b>	<b>34.63</b>	34.55	29.21
	TRADES(S)	81.84	60.15	33.82	33.76	28.76
	TRADES(S)+SLORE	<b>81.85</b>	<b>60.22</b>	33.66	33.73	<b>28.88</b>
	TRADES(S)+LORE	81.81	59.86	<b>34.02</b>	<b>34.00</b>	28.70
	MART	77.78	59.91	34.86	34.83	29.18
	MART+SLORE	77.08	59.86	34.72	34.70	<b>29.60</b>
	MART+LORE	<b>78.57</b>	<b>60.18</b>	<b>35.20</b>	<b>35.22</b>	29.47
WRN-28-10 (Swish)	TRADES	89.39	65.41	37.88	37.93	33.02
	TRADES+SLORE	88.95	64.95	38.07	38.04	<b>33.29</b>
	TRADES+LORE	<b>90.17</b>	<b>66.62</b>	<b>38.10</b>	<b>37.96</b>	32.61
	TRADES(S)	<b>85.22</b>	63.12	36.62	36.57	31.60
	TRADES(S)+SLORE	84.78	<b>63.18</b>	36.50	36.54	<b>31.97</b>
	TRADES(S)+LORE	84.82	62.91	<b>36.78</b>	<b>36.65</b>	31.56
	MART	81.82	65.05	38.56	38.51	33.91
	MART+SLORE	80.35	64.46	38.88	38.81	<b>34.26</b>
	MART+LORE	<b>83.21</b>	<b>65.11</b>	<b>39.06</b>	<b>39.00</b>	33.76

Table 10. Classification Accuracy of PreResNet18 (Swish) trained by  $1 \times 10^6$  EDM-generated images-augmented Tiny-ImageNet (%). V, S, and L represent vanilla, SLORE-boosted, and LORE-boosted adversarial training algorithms respectively.  $C_{tr}$ ,  $C_{te}$ ,  $P_{40}$ ,  $A_{ce}$ , and  $A_t$  indicate Clean<sub>tr</sub>, Clean<sub>te</sub>, PGD<sup>40</sup>, APGD<sub>ce</sub>, and APGD<sub>t</sub> individually.

	TRADES					MART				
	$C_{tr}$	$C_{te}$	$P_{40}$	$A_{ce}$	$A_t$	$C_{tr}$	$C_{te}$	$P_{40}$	$A_{ce}$	$A_t$
V	64.91	58.66	28.52	28.31	<b>21.96</b>	52.52	50.62	28.14	28.02	21.51
S	64.76	58.32	28.33	28.32	21.59	<b>55.17</b>	<b>51.34</b>	<b>29.02</b>	<b>28.75</b>	<b>22.54</b>
L	<b>65.23</b>	<b>59.08</b>	<b>28.78</b>	<b>28.47</b>	21.68	53.04	50.79	28.61	27.74	21.26

Tiny-ImageNet and SVHN. The extensive results in Tab. 8, 9 and 10 reinforce the adaptability and efficacy of our approach. To more comprehensively demonstrate the efficacy of our approach, we trained LOAT-boosted TRADES

Table 11. Classification Accuracy of PreResNet18 (Swish) trained by  $1 \times 10^6$  EDM-generated images-augmented Cifar10 after 400 epochs (%).

Defense	Clean <sub>te</sub>	PGD <sup>10</sup>	PGD <sup>20</sup>	PGD <sup>40</sup>	AA
TRADES	88.45	62.62	61.88	61.76	<b>58.21</b>
TRADES+SLORE	86.12	62.28	61.73	61.56	57.76
TRADES+LORE	<b>90.26</b>	<b>62.81</b>	<b>62.01</b>	<b>61.85</b>	56.51
MART	88.35	61.19	60.18	59.93	53.92
MART+SLORE	88.15	61.99	60.94	60.78	54.74
MART+LORE	<b>89.08</b>	<b>63.81</b>	<b>62.83</b>	<b>62.57</b>	<b>55.28</b>

and MART on PreResNet18 (Swish) for an extended duration of 400 epochs. The evaluations, as presented in Tab. 11, reveal a more pronounced improvement in performance compared to training for approximately 100 epochs. This enhancement is particularly notable in the case of MART.

## 10. Limitations and Future Works

- **Mitigating Theoretical Gap Between MLPs and Other Network Structures.** We intend to implement knowledge distillation using a ReLU-activated MLP, aligning its architecture to mirror the layer count of a ResNet residual block or a Transformer stack, and matching the hidden units in each layer to the dimensionality of latent features derived from these structures. This framework will be employed in LOAT-boosted AT algorithms, which can help us investigate the theoretical gap.
- **Hyperparameters in LOAT.** The selection of hyperparameters  $\mathcal{E}_1$ ,  $\mathcal{E}_2$ ,  $\tau$ , and  $\gamma$  is crucial, as their configuration can significantly impact the final results.
- **Future Directions.** The rising prominence of Transformers and Large Language Models (LLMs) presents new challenges and opportunities. Exploring the application of adversarial training frameworks to these models is both an intriguing and vital avenue for future research.