

Motivation

- > Recent works have shown failure of models trained with ERM in achieving high certified robustness and have proposed new training methods.
- \succ However, the role of the data distribution in $\frac{1}{2}$ achieving these high certified robustness guarantees has been ignored.
- > Studying performance in adversarial setting such as in presence of data poisoning can be beneficial. But current poisoning attacks fail when certified defenses are used.

Can poisoning degrade certified adversarial robustness guarantees?

Contributions

- > We study the problem of using data poisoning attacks to affect the robustness guarantees of classifiers trained using certified defense methods.
- \succ We propose a bilevel optimization based clean label poisoning attack to generate poison data against robust training and certification methods.
- \succ We demonstrate the effectiveness of our attack at reducing certified adversarial robustness obtained using randomized smoothing on models trained with state-of-the-art certified defenses.

Poison data that hurts certified adversarial robustness

Presence of imperceptibly distorted poison data can significantly hurt certified adversarial robustness guarantees.



Poisons optimized against ERM fail when certified defenses are used.



How Robust are Randomized Smoothing based Defenses to Data Poisoning?

Akshay Mehra¹, Bhavya Kailkhura², Pin-Yu Chen³ and Jihun Hamm¹ ¹Tulane University, ²Lawrence Livermore National Laboratory, ³IBM Research



Generate poison data (u) such that when the victim trains a model (with parameters θ) on the poisoned data ($\mathcal{D}^{clean} \cup \mathcal{D}^{poison}$) the certified robustness guarantees of the target class are significantly diminished on a validation set \mathcal{D}^{val} .

$$\min_{u \in \mathcal{U}} \mathcal{R}(\mathcal{D}^{va})$$

s.t. $\theta^* = \arg\min_{\theta} \mathcal{L}_{robust}$

Lower-level problem trains a model using training procedures that lead to models with high certified robustness such as MACER, SmoothAdv etc. Upper-level problem lowers average certified radius obtained from a certification procedure such as Randomized Smoothing

Effect of poisoning on the decision boundary

Average margin of the smoothed classifiers on the target data is reduced leading to



Linear Classifiers

```
al; \theta^*)
(\mathcal{D}^{clean} \cup \mathcal{D}^{poison}; \theta)
```

Non-Linear Classifiers

Key Results





References

- preprint arXiv:2001.02378 (2020).
- preprint arXiv:1906.04584 (2019).
- machine learning. PMLR, 2016.

This work was supported by the NSF EPSCoR-Louisiana Materials Design Alliance (LAMDA) program #OIA1946231 and by LLNL Laboratory Directed Research and Development project 20-ER-014 (LLNL-CONF-817233). This work was performed under the auspices of the U.S. Department of Energy by the Lawrence Livermore National Laboratory under Contract No. DE-AC52-07NA27344, Lawrence Livermore National Security, LLC



Reduction in Average Certified Radius with poisoning





Cohen, Jeremy, Elan Rosenfeld, and Zico Kolter. "Certified adversarial robustness via randomized smoothing." International Conference on Machine Learning. PMLR, 2019.

> Zhai, Runtian, et al. "Macer: Attack-free and scalable robust training via maximizing certified radius." *arXiv*

Salman, Hadi, et al. "Provably robust deep learning via adversarially trained smoothed classifiers." arXiv

Pedregosa, Fabian. "Hyperparameter optimization with approximate gradient." International conference on

Mei, Shike, and Xiaojin Zhu. "Using machine teaching to identify optimal training-set attacks on machine learners." Proceedings of the AAAI Conference on Artificial Intelligence. Vol. 29. No. 1. 2015.