Towards Benchmarking the Robustness of Neuro-Symbolic Learning against Backdoor Attacks

BACKGROUND

Logic Tensor Networks (LTNs)

- Representative NeSy models that handle diverse machine learning tasks efficiently;
- Integrate neural networks (NNs) with First-Order Logic (FOL);
- Works very well with MNIST digit classification tasks;
- Learning is guided by logic-based axioms, which shape the loss using fuzzy logic operators (Λ , $\exists \forall$
- Logical symbols are grounded as tensors or neural networks to evaluate formula satisfiability. [2], [3]



Badnets

- Type of data injection backdoor attack;
- Perform well on clean inputs, but cause misclassifications for triggered inputs;
- Visual triggers on images, e.g, small square on the bottom-right corner of the image; works very well on MNIST images;
- Stealthy attacks: pass the standard tests and preserve the structure of the baseline model;
- Simple method, but adds complex behavior;
- Relevant to real-world scenarios. [1]

Original Image

Pattern Backdoor





Fig. 2: Backdoored MNIST Model

REFERENCES

[1] T. Gu, K. Liu, B. Dolan-Gavitt, and S. Garg. "BadNets: Evaluating Backdooring Attacks on Deep Neural Networks". In: IEEE Access (2019). DOI: https://doi.org/10.1109/ACCESS.2019.2909068.

[2] L. Serafin and A. S. d'Avila Garcez. "Learning and Reasoning with Logic Tensor Networks". In: Springer International Publishing (2016). DOI: 10.1007/978-3-319-49130-1_25.

[3] L. Serafini and A. S. d'Avila Garcez. "Logic Tensor Networks: Deep Learning and Logical Reasoning from Data and Knowledge". In: arXiv (Cornell University) (2016). DOI: http://arxiv.org/abs/1606.04422.

[4] S. Badreddine, A. d. Garcez, L. Serafini, and M. Spranger, "Logic Tensor Networks," Artificial Intelligence, vol. 303, p. 103649, Feb. 2022, arXiv:2012.13635 [cs]. [Online]. Available: http://arxiv.org/abs/2012.13635

RESEARCH QUESTION

How robust is a Logic Tensor Network (LTN) model against data poisoning BadNet attacks?

METHODOLOGY



Fig. 3: BadNet on SDA and MDA Samples

AUTHORS

Student: Myriam Guranda; mguranda@student.tudelft.nl Supervisor: Andrea Agiollo; A.Agiollo-1@tudelft.nl Responsible Professor: Kaitai Liang; Kaitai.Liangetudelft.nl

AFFILIATIONS

Delft University of Technology

GITHUB

https://github.com/myriamcg/NeSy-vs-Backdoors

RESULTS

1) SDA



Result label is 59 Result label is 29









Key Findings:

• Center triggers (Fig. 4.b) \rightarrow Consistently

• $4 \times 4 \text{ vs } 6 \times 6 \text{ bottom-right trigger} \rightarrow \text{Both}$

12% vs 89% (Figure 4.a) => bigger

effective and stealthy due to CNN focus

maintain high accuracy, but ASR differs:

triggers make a difference in less salient

• Both images poisoned \rightarrow ASR drops to

0%, accuracy stays high (Fig. 4.c)

on the image center.

positions

#	Trigger Size	Trigger Pos.	Poisoned	Acc.	ASR
1	4×4	Right	First	90%	12%
2	6×6	Right	First	98%	93%
3	4×4	Right	Both	90%	0%
4	6×6	Right	Both	94%	0%
5	4×4	Center	First	90%	89%
6	6×6	Center	First	95%	100%
7	10×10	Center	First	97%	100%
8	6×6	Center	Both	95%	0%
9	10×10	Center	Both	95%	0%

2) MDA



Fig. 5.b): MDA First Center 4

#	Trigger Size	Trigger Position	Poisoned Images	Accuracy	ASR
1	4×4	Right	d_1, d_3	0.3%	100%
2	6×6	Right	d_1, d_3	0.3%	100%
3	10×10	Right	d_1, d_3	0.3%	100%
4	6×6	Right	d_1, d_2, d_3, d_4	0.3%	100%
5	10×10	Right	d_1, d_2, d_3, d_4	20%	100%
6	4×4	Center	d_1, d_3	85%	3%
7	6×6	Center	d_1, d_3	95%	97%
8	4×4	Center	d_1, d_2, d_3, d_4	95%	97%
9	6×6	Center	d_1, d_2, d_3, d_4	90%	97%



Fig. 5.c): MDA Both Center

Key Findings:

- Right-corner triggers on d1 & d3 \rightarrow Immediate model collapse (Fig 5.a) => Trigger position matters; d1 and d3 are symbolically dominant;
- Center, larger triggers on d1 & d3 succeed due to CNN's focus on image centers;
- Right-corner triggers on all digits → attack dominates, but logic fails to be learned;
- Center triggers on all digits \rightarrow attack succeeds (Fig. 5.c)

CONCLUSION

Key Findings:

- Larger (e.g., 6×6), centrally placed triggers are the most effective, achieving high ASR while remaining stealthy;
- Poisoning both images in a sample often leads to low ASR due to symbolic ambiguity, but accuracies remain unchanged;
- Symbolically dominant inputs (e.g., d1 and d3 in MDA) are more sensitive to poisoning;
- Stronger regularization hyperparameters during training suppress weaker attacks over time.
- Change of ASR calculation to the
- correctness of the symbolic outcome; • Test LTN's vulnerability on tasks with
- Use multi-channel MNIST images to introduce more visual features.

- Future Work:

more symbolic knowledge;

