# Rule-based Safety Evidence for Neural Networks

- Tewodros A. Beyene
- Amit Sahu

# Outline

▶ Introduction

▶ Rule Extraction

▶ Rule-based Safety Evidence

▶ Illustration

▶ Results

▶ Conclusion and Future Work

# Introduction

▶ Safety-critical systems are often subject to a rigorous safety certification process (DO-178C, ISO26262)

▶ In general, we need evidence for safety certification as information or artefacts that contribute to developing confidence in the safe operation of a system

▶ There have been efforts to extend safety assurance approaches for traditional software to AI-based systems, such as self-driving cars and other autonomous systems.

▶ These face two challenges:

  • Efficient analysis, testing, and verification methods

  • Finding artefacts that develop confidence in the safe operation of the AI system

▶ In this position paper, our focus was on second challenge by using rule extraction techniques

# Rule Extraction

▶ Approaches

  • Pedagogical: treat network as black-box

  • Decompositional: split the network at neuron level

▶ Types of Extracted Rules: If-then, M-of-N, Decision Tree, fuzzy rules, first-order rules

▶ Rule Extraction for DNNs

  • Generated layer-by-layer

  • High memory usage and computation time

  • Rule pruning and network pruning

# Rule-based Safety Evidence

Rules as Evidence Artefacts

- Diverse perception/view of rules: coverage over features, robustness to noise

- Applicability of extraction techniques in diverse architectures: locally in federated learning, networks that facilitate rule extraction

- Amenability of rules for quality assessment: criteria like accuracy, fidelity, and consistency

- Diverse expressive power of rules: for example, decision tree can express hierarchical properties of an NN as compared to if-then rules

fortiss

# Rule-based Safety Evidence

Algorithmic complexity of rule extraction

- Scaling of decompositional approaches

- Various heuristics to limit space exploration and achieve tractability for real world problems

- Complexity of extracted rules: number of extracted rules, number of antecedents per rule

- Solutions for reducing the size of generated rulesets: network pruning, rule pruning, rule abstraction, distillation, formal methods: partial order reductions, priority synthesis

fortiss

# Illustration

Methodology

- ▶ Select Dataset: MNIST
- ▶ Train NN models with the goal of achieving accuracy and specific properties: robust to adversarial noise
- ▶ Extract rules from the trained NN models
- ▶ Use these rules as safety evidence for the targeted properties

# Illustration

NN Models

▶ Plain: trained with basic backpropagation algorithm

▶ Adversarial Training (AT): trained with adversarial images to make the model robust to adversarial perturbation

▶ Maximization of Linear Regions(MMR): trained with MMR regularizer to increase the provable bound on effective adversarial noise

▶ MMR+AT: a hybrid approach that uses both MMR and AT approaches

fortiss

# Illustration

Techniques

▶ TREPAN: The algorithm uses an Oracle that can query and sample examples from the dataset. Each node is split based on the best fit that generates the highest fidelity with the NN model

▶ Surrogate Random Forest Model: In this algorithm, the NN model was taken as a black-box model generating feature-output combination. This combination is used to train a separate Random Forest Model from scratch.

fortiss

# Results

| Model Name | Decision tree fidelity (%) | Random forest fidelity (%) | L2 robust error LB (%) | UB (%) |
|---|---|---|---|---|
| Plain | 94 | 98.67 | 3.1 | 100 |
| AT | 91.34 | 98.5 | 1.8 | 100 |
| MMR | 87.12 | 96.37 | 5.8 | 11.6 |
| MMR+AT | 90 | 97.08 | 4.6 | 9.7 |

**Table 1.** Fidelity of models with different surrogate models and adversarial robustness

- Surrogate RF model had better fidelity scores than the surrogate Decision Tree attributed to higher complexity of RF model

- Models with robustness against adversarial property achieved lower fidelity than plain or simpler models

# Conclusion and Future Work

▶ We have shown rules extracted from NN can be good evidence artefacts

▶ We observed that different safety properties can be defined, represented, and verified for critical systems using rules with high fidelity scores

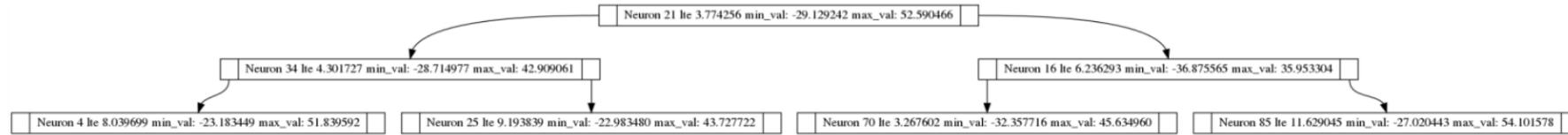Some properties that we will be working on for different datasets:

▶ Interpretation of the decision process

▶ Robustness to noise

▶ Tractable coverage of the decision tree as compared to coverage over neuron activations
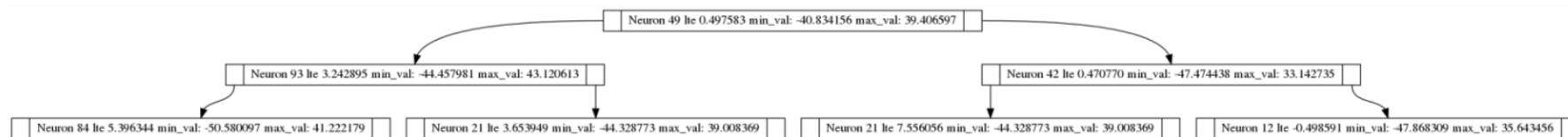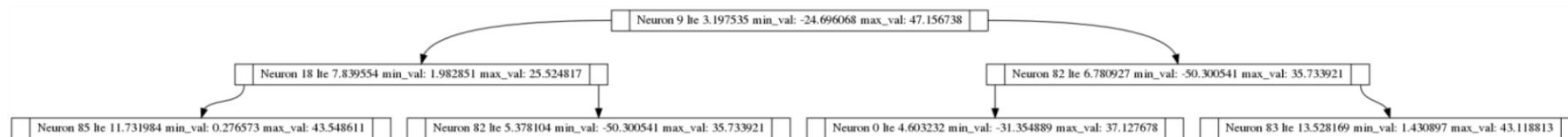
fortiss

Thank you.

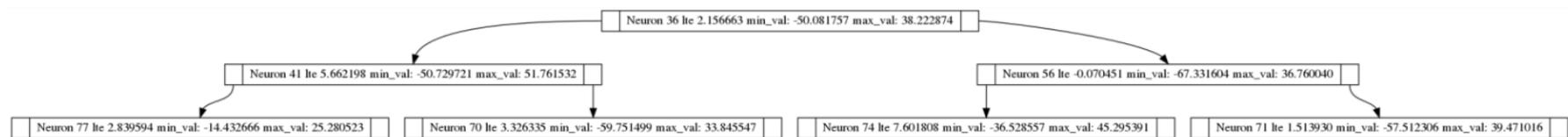**Questions?**

# Surrogate Subtrees of the Models



Surrogate Tree for MNIST plain model



Surrogate Tree for MNIST AT model



Surrogate Tree for MNIST MMR model



Surrogate Tree for MNIST MMR+AT model