Formal Methods for Robust Artificial Intelligence State of the Art

Caterina Urban ANTIQUE Research Team, Inria & École Normale Supérieure | Université PSL





Artificial Intelligence Development Process Artificial Intelligence Pipeline







data preparation

model training



model deployment



predictions



Model Training is Highly Non-Deterministic





model training



model deployment



predictions

no predictability and traceability



Models Only Give Probabilistic Guarantees









Max Speed 100





model training



model deployment



not sufficient for guaranteeing an acceptable failure rate under any circumstances

Safety-Critical Artificial Intelligence



A self-driving Uber ran a red light last December, contrary to company claims

Feds Say Self-Driving Uber SUV Did Not Recognize Jaywalking Pedestrian In Fatal Crash

<u>Richard Gonzales</u> November 7, 201910:57 PM ET

07/10/2019, 23:16











Formal Methods Mathematical Guarantees of Safety



Deductive Verification

- extremely **expressive**
- relies on the user to guide the proof



Edmund Clarke





Model Checking

- **Static Analysis**
- analysis of the software at some level of abstraction
- fully automatic and sound by construction
- generally **not complete**

• analysis of a **model** of the software sound and complete with respect to the model





Methods for Trained Models





model training



model deployment

predictions

Neural Network Models



Feed-Forward Neural Networks Fully-Connected Layers with ReLU Activation Functions



$$\mathbf{x}_{i,j} = \max\left\{\mathbf{0}, \sum_{k} w_{j,k}^{i-1} \cdot x_{i-1,k} + b_{i,j}\right\}$$



Safety Verification $l_j \leq x_{0,j} \leq u_j$



10

Model Checking Methods



SMT-Based Methods Safety Verification Reduced to Constraint Satisfiability

 $l_i \leq x_{0,i} \leq u_i$ $j \in \{0, ..., |\mathbf{X}_0|\}$ $\hat{x}_{i+1,j} = \sum_{k=1}^{|\mathbf{X}_i|} w_{j,k}^i \cdot x_{i,k} + b_{i,j} \qquad i \in \{0, \dots, n-1\}$ k=0

 $x_{i,i} = \max\{0, \hat{x}_{i,i}\}$ $i \in \{1, \dots, n-1\}, j \in \{0, \dots, |\mathbf{X}_i|\}$

 $x_N \leq 0$



input specification



(negation of) output specification

12



R. Ehlers - Formal Verification of Piece-Wise Linear Feed-Forward Neural Networks (ATVA 2017)



use **approximations** to reduce the solution search space



Reluplex

Variable	Value	Variable	Value
X 00	v_{00}	X 00	v_{00}
• • •	• • •	• • •	• • •
X _{ij}	\hat{v}_{ij}	Âx _{ij}	$\hat{\mathcal{V}}'_{ij}$
X _{ij}	V _{ij}	X _{ij}	V _{ij}
• • •	• • •	• • •	• • •
X _N	v_N	X _N	v_N

G. Katz et al. - Reluplex: An Efficient SMT Solver for Verifying Deep Neural Networks (CAV 2017)



based on the **simplex algorithm** extended to support ReLU constraints

Variable	Value
X 00	v_{00}
•	• • •
$\hat{\mathbf{x}}_{\mathbf{ij}}$	\hat{v}'_{ij}
X _{ij}	\hat{v}'_{ij}
• • •	• • •
X _N	v_N









Reluplex

Variable	Value	Variable	Value
X 00	v_{00}	X 00	v_{00}
• • •	• • •	• • •	• • •
X _{ij}	\hat{v}_{ij}	Âx _{ij}	$\hat{\mathcal{V}}'_{ij}$
X _{ij}	V _{ij}	X _{ij}	V _{ij}
• • •	• • •	• • •	• • •
X _N	v_N	X _N	v_N

G. Katz et al. - Reluplex: An Efficient SMT Solver for Verifying Deep Neural Networks (CAV 2017)



based on the Sin extended to supp Follow-up Work

Variable	Value
X 00	v_{00}
• • •	• • •
$\hat{\mathbf{x}}_{\mathbf{ij}}$	$\hat{\mathcal{V}}'_{ij}$
X _{ij}	$\hat{\mathcal{V}}'_{ij}$
• • •	• • •
X _N	v_N

G. Katz et al. - The Marabou Framework for Verification and Analysis of Deep Neural Networks (CAV 2019)

• • •

Variable	Value
X ₀₀	v_{00}
• • •	• • •
Âx _{ij}	\hat{v}'_{ij}
X _{ij}	0
• • •	• • •
X _N	v_N

Other SMT-Based Methods

- Neural Networks (CAV 2010) the first formal verification method for neural networks
- Neural Net Robustness with Constraints (NeurIPS 2016)
- (CAV 2017) an approach for proving local robustness to adversarial perturbations
- Binarized Deep Neural Networks (AAAI 2018) Networks via Inter-Neuron Factoring (VSTTE 2018) approaches focusing on binarized neural networks

• L. Pulina and A. Tacchella - An Abstraction-Refinement Approach to Verification of Artificial

• O. Bastani, Y. Ioannou, L. Lampropoulos, D. Vytiniotis, A. Nori, and A. Criminisi - Measuring an approach for finding the nearest adversarial example according to the Loo distance

• X. Huang, M. Kwiatkowska, S. Wang, and M. Wu - Safety Verification of Deep Neural Networks

• N. Narodytska, S. Kasiviswanathan, L. Ryzhyk, M. Sagiv, and T. Walsh - Verifying Properties of

C. H. Cheng, G. Nührenberg, C. H. Huang, and H. Ruess - Verification of Binarized Neural



MILP-Based Methods Safety Verification Reduced to Mixed Integer Linear Program

 $l_i \leq x_{0,i} \leq u_i$ $j \in \{0, ..., |\mathbf{X}_0|\}$

$$\hat{x}_{i+1,j} = \sum_{k=0}^{|\mathbf{X}_i|} w_{j,k}^i \cdot x_{i,k} + b_{i,j} \qquad i \in \{0, \dots, n-1\}$$

 $x_{i,j} = \delta_{\mathbf{i},\mathbf{j}} \cdot \hat{x}_{i,j}$ $\delta_{\mathbf{i},\mathbf{i}} \in \{\mathbf{0},\mathbf{1}\}$ $\delta_{\mathbf{i},\mathbf{j}} = 1 \Rightarrow \hat{x}_{i,j} \ge 0$ $i \in \{1, ..., n-1\}$ $j \in \{0, ..., |\mathbf{X}_i|\}$ $\delta_{\mathbf{i},\mathbf{i}} = 0 \Rightarrow \hat{x}_{i,i} < 0$

min X_N



input specification



objective function

16

MILP-Based Methods Bounded MILP Encoding with Symmetric Bounds

$$\hat{x}_{i+1,j} = \sum_{k=0}^{|\mathbf{X}_i|} w_{j,k}^i \cdot x_{i,k} + b_{i,j} \qquad i \in \{0\}$$

$$0 \le x_{i,j} \le \mathbf{M}_{\mathbf{i},\mathbf{j}} \cdot \delta_{i,j} \qquad \qquad \delta_{i,j} \in \{ \hat{x}_{i,j} \le x_{i,j} \le \hat{x}_{i,j} - \mathbf{M}_{\mathbf{i},\mathbf{j}} \cdot (1 - \delta_{i,j}) \qquad \qquad i \in \{ 1, \dots, n \}$$

$$\mathbf{M}_{\mathbf{i},\mathbf{j}} = \max\{-\mathbf{l}_{\mathbf{i}}, \mathbf{u}_{\mathbf{i}}\} \qquad \qquad j \in \{ 0, \dots, n \}$$

 $), ..., n-1 \}$

 $\{0,1\}$, $n-1\}$, $|\mathbf{X}_i|$





$l_j \leq x_{0,j} \leq u_j$



 $0 \le x_{i,j} \le \mathbf{M}_{\mathbf{i},\mathbf{j}} \cdot \delta_{i,j}$ $\hat{x}_{i,j} \le x_{i,j} \le \hat{x}_{i,j} - \mathbf{M}_{\mathbf{i},\mathbf{j}} \cdot (1 - \delta_{i,j})$ $\mathbf{M}_{\mathbf{i},\mathbf{j}} = \max\{-\mathbf{l}_{\mathbf{i}}, \mathbf{u}_{\mathbf{i}}\}$

min X_N

S. Dutta et al. - Output Range Analysis for Deep Feedforward Neural Networks (NFM 2018)



use **local search** speed up the MILP solver



$l_j \leq x_{0,j} \leq u_j$



 $0 \le x_{i,j} \le \mathbf{M}_{\mathbf{i},\mathbf{j}} \cdot \delta_{i,j}$ $\hat{x}_{i,j} \le x_{i,j} \le \hat{x}_{i,j} - \mathbf{M}_{\mathbf{i},\mathbf{j}} \cdot (1 - \delta_{i,j})$ $\mathbf{M}_{\mathbf{i},\mathbf{j}} = \max\{-\mathbf{l}_{\mathbf{i}}, \mathbf{u}_{\mathbf{i}}\}$ $\mathbf{X}_{\mathbf{N}} < \mathbf{L}$

S. Dutta et al. - Output Range Analysis for Deep Feedforward Neural Networks (NFM 2018)



use **local search** speed up the MILP solver

$\begin{array}{l} \textbf{sample} \text{ random input } X \\ \textbf{and evaluate output } L \end{array}$



$l_j \leq x_{0,j} \leq u_j$



 $0 \le x_{i,j} \le \mathbf{M}_{\mathbf{i},\mathbf{j}} \cdot \delta_{i,j}$ $\hat{x}_{i,j} \le x_{i,j} \le \hat{x}_{i,j} - \mathbf{M}_{\mathbf{i},\mathbf{j}} \cdot (1 - \delta_{i,j})$ $\mathbf{M}_{\mathbf{i},\mathbf{j}} = \max\{-\mathbf{l}_{\mathbf{i}}, \mathbf{u}_{\mathbf{i}}\}$

 $\mathbf{x}_{\mathbf{N}} < \mathbf{L}$

S. Dutta et al. - Output Range Analysis for Deep Feedforward Neural Networks (NFM 2018)



use **local search** speed up the MILP solver



find another input $\hat{\hat{X}}$ such that $\hat{L} \leq x_N$



$l_j \leq x_{0,j} \leq u_j$



 $0 \le x_{i,j} \le \mathbf{M}_{i,j} \cdot \delta_{i,j}$ $\hat{x}_{i,j} \le x_{i,j} \le \hat{x}_{i,j} - \mathbf{M}_{i,j} \cdot (1 - \delta_{i,j})$ $\mathbf{M}_{i,j} = \max\{-\mathbf{l}_i, \mathbf{u}_i\}$ $\mathbf{x}_{N} < \hat{\mathbf{L}}$

S. Dutta et al. - Output Range Analysis for Deep Feedforward Neural Networks (NFM 2018)



use **local search** speed up the MILP solver

find another input $\hat{\hat{X}}$ such that $\hat{L} \leq x_N$



$l_j \leq x_{0,j} \leq u_j$



 $0 \le x_{i,j} \le \mathbf{M}_{\mathbf{i},\mathbf{j}} \cdot \delta_{i,j}$ $\hat{x}_{i,j} \le x_{i,j} \le \hat{x}_{i,j} - \mathbf{M}_{\mathbf{i},\mathbf{j}} \cdot (1 - \delta_{i,j})$ $\mathbf{M}_{\mathbf{i},\mathbf{j}} = \max\{-\mathbf{l}_{\mathbf{i}}, \mathbf{u}_{\mathbf{i}}\}$ $\mathbf{x}_{\mathbf{N}} < \hat{\mathbf{L}}$

S. Dutta et al. - Output Range Analysis for Deep Feedforward Neural Networks (NFM 2018)



use **local search** speed up the MILP solver



find another input $\hat{\hat{X}}$ such that $\hat{L} \leq x_N$



MILP-Based Methods Bounded MILP Encoding with Asymmetric Bounds

$$\hat{x}_{i+1,j} = \sum_{k=0}^{|\mathbf{X}_i|} w_{j,k}^i \cdot x_{i,k} + b_{i,j} \qquad i \in \{0\}$$

$$0 \le x_{i,j} \le \mathbf{u}_{\mathbf{i},\mathbf{j}} \cdot \delta_{i,j} \qquad \qquad \delta_{i,j} \in \{1, j\} \\ \hat{x}_{i,j} \le x_{i,j} \le \hat{x}_{i,j} - \mathbf{l}_{\mathbf{i},\mathbf{j}} \cdot (1 - \delta_{i,j}) \qquad \qquad i \in \{1, j\} \\ j \in \{0, j\} \end{cases}$$

 $), ..., n-1 \}$

 $\{0,1\}$, $n-1\}$, $|\mathbf{X}_i|$



19

MIPVerify Finding Nearest Adversarial Example

$\min_{X'} d(X, X')$



$$0 \le x_{i,j} \le \mathbf{u}_{\mathbf{i},\mathbf{j}} \cdot \delta_{i,j}$$
$$\hat{x}_{i,j} \le x_{i,j} \le \hat{x}_{i,j} - \mathbf{l}_{\mathbf{i},\mathbf{j}} \cdot (1 - \delta_{i,j})$$

$\mathbf{x}_{\mathbf{N}} \neq \mathbf{O}$

V. Tjeng, K. Xiao, and R. Tedrake - Evaluating Robustness of Neural Networks with Mixed Integer Programming (ICLR 2019)



Other MILP-Based Methods

- R. Bunel, I. Turkaslan, P. H. S. Torr, P. Kohli, and M. P. Kumar A Unified View of Piecewise Linear Neural Network Verification (NeurIPS 2018)
 a unifying verification framework for piecewise-linear ReLU neural networks
- C.-H. Cheng, G. Nührenberg, and H. Ruess Maximum Resilience of Artificial Neural Networks (ATVA 2017)
 an approach for finding a lower bound on robustness to adversarial perturbations
- M. Fischetti and J. Jo Deep Neural Networks and Mixed Integer Linear Optimization (2018) an approach for feature visualization and building adversarial examples



Static Analysis Methods



Abstract Interpretation-Based Methods



(1) proceed forwards from an abstraction of the input specification



(2) check output for inclusion in output specification: included $\rightarrow \checkmark$ safe otherwise \rightarrow (2) alarm

23

Symbolic Propagation

$$x_{i,j} \mapsto \begin{cases} \sum_{k=0}^{i-1} \mathbf{c}_k \cdot \mathbf{x}_k + \mathbf{c} \quad \mathbf{c}_k, \mathbf{c} \in \mathscr{R}^{|\mathbf{X}_k|} \\ [a, b] & a, b \in \mathscr{R} \end{cases}$$

$$x_{i-1,0} \mapsto \mathbf{E_{i-1,0}} \\ \cdots \\ x_{i-1,j} \mapsto \mathbf{E_{i-1,j}} \qquad x_{i,j} = \sum_k w_{j,k}^{i-1} \cdot x_{i-1,k} + b_{i,j}$$

$$x_{i,j} = \sum_k w_{j,k}^{i-1} \cdot x_{i-1,k} + b_{i,j}$$

$$x_{i,j} = \sum_k w_{j,k}^{i-1} \cdot x_{i-1,k} + b_{i,j}$$

 $x_{i,j} \mapsto \langle$

J. Li et al. - Analyzing Deep Neural Networks with Symbolic Propagation (SAS 2019)



$$x_{i,j} \mapsto \sum_{k} w_{j,k}^{i-1} \cdot \mathbf{E}_{i-1,k} + b_{i,j}$$

$$x_{i,j} \mapsto \begin{cases} \mathbf{E}_{\mathbf{i},\mathbf{j}} \\ [\mathbf{a},\mathbf{b}] \end{cases} \qquad 0 \le a$$
$$x_{i,j} \mapsto \begin{cases} \mathbf{X}_{\mathbf{i},\mathbf{j}} \\ [\mathbf{0},\mathbf{b}] \end{cases} \qquad a < 0 \land 0 < b$$
$$x_{i,j} \mapsto \begin{cases} \mathbf{0} \\ [\mathbf{0},\mathbf{0}] \end{cases} \qquad b \le 0$$



DeepPoly $x_{i+1,j} \mapsto \begin{cases} \left[\sum_{k} c_{i,k} \cdot x_{i,k} + c, \sum_{k} d_{i,k} \cdot x_{i,k} + d\right] & c_{i,k}, c, d_{i,k}, d \in \mathcal{R} \\ [a, b] & a, b \in \mathcal{R} \end{cases}$ Relea $x_{i,j} \mapsto \begin{cases} [\mathbf{L}_{i,j}, \mathbf{U}_{i,j}] \\ [\mathbf{0}, \mathbf{b}] \end{cases}$ ReLU $a < 0 \land 0 < b$ -a < b $\mathbf{x}_{i,j} \mapsto \begin{cases} \mathbf{[0,0]} \\ \mathbf{[0,0]} \end{cases}$

G. Singh, T. Gehr, M. Püschel, and M. Vechev - An Abstract Domain for Certifying Neural Networks (POPL 2019)

maintain symbolic lower- and upper-bounds for each neuron + convex ReLU approximations

ReLU(x)





Other Abstract Interpretation Methods

- the first use of abstract interpretation for verifying neural networks
- Certification (NeurIPS 2018) a custom zonotope domain for certifying neural networks
- Neural Network Certification (NeurIPS 2019) a framework to jointly approximate k ReLU activations
- Neural Networks (OOPSLA 2020) an approach for verifying fairness of neural network classifiers for tabular data

• T. Gehr, M. Mirman, D. Drachsler-Cohen, P. Tsankov, S. Chaudhuri, and M. Vechev - Al2: Safety and Robustness Certification of Neural Networks with Abstract Interpretation (S&P 2018)

• G. Singh, T. Gehr, M. Mirman, M. Püschel, and M. Vechev - Fast and Effective Robustness

• G. Singh, R. Ganvir, M. Püschel, and M. Vechev - Beyond the Single Neuron Convex Barrier for

• C. Urban, M. Christakis, V, Wüstholz, and F. Zhang - Perfectly Parallel Fairness Certification of





Other Complete Methods



Star Sets Exact Static Analysis Method



• fast and cheap affine mapping operations \rightarrow neural network layers • inexpensive intersections with half-spaces \rightarrow ReLU activations

H.-D. Tran et al. - Star-Based Reachability Analysis of Deep Neural Networks (FM 2018)



 $V = \{v_1, \dots, v_m\}$: basis vectors in \mathscr{R}^n



Star Sets Exact Static Analysis Method



• fast and cheap affine mapping operations \rightarrow neural network layers • inexpensive intersections with half-spaces \rightarrow ReLU activations

H.-D. Tran et al. - Star-Based Reachability Analysis of Deep Neural Networks (FM 2018)



use efficient r Follow-up Work of bounded

> H.-D. Tran et al. - Verification of Deep Convolutional Neural Networks Using ImageStars (CAV 2020)

 $V = \{v_1, \dots, v_m\}$: basis vectors in \mathscr{R}^n





S. Wang et al. - Formal Security Analysis of Neural Networks Using Symbolic Intervals (USENIX Security 2018)





Neurify **Asymptotically Complete Method**





S. Wang, K. Pei, J. Whitehouse, J. Yang, and S. Jana - Efficient Formal Safety Analysis of Neural Networks (NeurIPS 2018)



 $0 \leq a$







Other Complete Methods

- W. Ruan, X. Huang, and Marta Kwiatkowska Reachability Analysis of Deep Neural Networks with Provable Guarantees (IJCAI 2018) a global optimization-based approach for verifying Lipschitz continuous neural networks
- G. Singh, T. Gehr, M. Püschel, and M. Vechev Boosting Robustness Certification of Neural Networks (ICLR 2019) an approach combining abstract interpretation and (mixed integer) linear programming



Other Incomplete Methods



Interval Neural Networks Abstraction-Based Method



P. Prabhakar and Z. R. Afza - Abstraction based Output Range Analysis for Neural Networks (NeurIPS 2019)





Interval Neural Networks **Abstraction-Based Method**



P. Prabhakar and Z. R. Afza - Abstraction based Output Range Analysis for Neural Networks (NeurIPS 2019)



R

Related Work

Y. Y. Elboher et al. - An Abstraction-Based Framework for Neural Network Verification (CAV 2020)





Other Incomplete Methods

- Multi-Layer Neural Networks (2018) an approach combining simulation and linear programming
- Verification of Deep Networks (UAI 2018) an approach based on duality for verifying neural networks
- Adversarial Polytope (ICML 2018) (ICML 2018) Towards Fast Computation of Certified Robustness for ReLU Networks (ICML 2018) H. Zhang, T.-W. Weng, P.-Y. Chen, C.-J. Hsieh, and L. Daniel - Efficient Neural Network Robustness Certification with General Activation Functions (NeurIPS 2018) approaches for finding a lower bound on robustness to adversarial perturbations

• W. Xiang, H.-D. Tran, and T. T. Johnson - Output Reachable Set Estimation and Verification for

• K. Dvijotham, R. Stanforth, S. Gowal, T. Mann, and P. Kohli - A Dual Approach to Scalable

• E. Wong and Z. Kolter - Provable Defenses Against Adversarial Examples via the Convex Outer

A. Raghunathan, J. Steinhardt, and P. Liang - Certified Defenses against Adversarial Examples

T.-W. Weng, H. Zhang, H. Chen, Z. Song, C.-J. Hsieh, L. Daniel, D. Boning, and I. Dhillon.



Other Incomplete Methods

- for Certifying Robustness of Convolutional Neural Networks (AAAI 2019) approach focusing on convolutional neural networks
- C.-Y. Ko, Z. Lyu, T.-W. Weng, L. Daniel, N. Wong, and D. Lin POPQORN: Quantifying Robustness of Recurrent Neural Networks (ICML 2019) Neural Networks for Cognitive Tasks via Reachability Analysis (ECAI 2020) approaches focusing on recurrent neural networks
- Networks (ASE 2019) an approach for inferring safety properties of neural networks

• A. Boopathy, T.-W. Weng, P.-Y. Chen, S. Liu, and L. Daniel - CNN-Cert: An Efficient Framework

H. Zhang, M. Shinn, A. Gupta, A. Gurfinkel, N. Le, and N. Narodytska - Verification of Recurrent

• D. Gopinath, H. Converse, C. S. Pasareanu, and A. Taly - Property Inference for Deep Neural





Complete Methods

ADVANTAGES

sound and complete

DISADVANTAGES

- soundness not typically guaranteed with respect to floating-point arithmetic
- do not scale to large models
- often limited to certain model architectures

Incomplete Methods

ADVANTAGES

- able to scale to large models
- sound often also with respect to floating-point arithmetic
- less limited to certain model architectures

DISADVANTAGES

suffer from false positives

Methods for Model Training



data preparation



model training



model deployment

predictions

Robust Training Minimizing the Worst-Case Loss for Each Input

Adversarial Training

Minimizing a Lower Bound on the Worst-Case Loss for Each Input





generate adversarial inputs and use them as training data

Certified Training

Minimizing an Upper Bound on the Worst-Case Loss for Each Input





use upper bound as regularizer to encourage robustness

38



model training



model deployment

go beyond robust training, give stronger formal guarantees

constrain the training process to guarantee desired properties

verify more interesting properties under all circumstances

support more models and verify their implementations



