Safe and Robust Deep Learning

Mislav Balunović

Department of Computer Science





SafeAl @ ETH Zurich (<u>safeai.ethz.ch</u>)

Joint work with



Publications:





Markus Püschel





Timon Gehr



Matthew Maximilian Baader Mirman



Tsankov



Dana Drachsler

Gagandeep Singh



Systems:

S&P'18:AI2: Safety and Robustness Certification of Neural Networks with Abstract Interpretation ERAN: Generic neural network verifier https://github.com/eth-sri/eran/ NeurIPS'18: Fast and Effective Robustness Certification POPL'19: An Abstract Domain for Certifying Neural Networks DiffAI: System for training provably robust networks https://github.com/eth-sri/diffai ICLR'19: Boosting Robustness Certification of Neural Networks ICML'18: Differentiable Abstract Interpretation for Provably Robust Neural Networks DL2: System for training and querying networks with logical constraints ICML'19: DL2: Training and Querying Neural Network with Logic

Deep Learning Systems

Self driving cars

Translation

Voice assistant



https://waymo.com/tech/

$\equiv Google$ Translate



https://translate.google.com



https://www.amazon.com/ Amazon-Echo-And-Alexa-Devices

The self-driving car incorrectly decides to turn right on Input 2 and crashes into the guardrail







(b) Input 2 (darker version of 1)

DeepXplore:Automated Whitebox Testing of Deep Learning Systems, SOSP'17

The self-driving car incorrectly decides to turn right on Input 2 and crashes into the guardrail





(a) Input 1

(b) Input 2 (darker version of 1)

DeepXplore:Automated Whitebox Testing of Deep Learning Systems, SOSP'17 The Ensemble model is fooled by the addition of an adversarial distracting sentence in blue.

Article: Super Bowl 50

Paragraph: "Peyton Manning became the first quarterback ever to lead two different teams to multiple Super Bowls. He is also the oldest quarterback ever to play in a Super Bowl at age 39. The past record was held by John Elway, who led the Broncos to victory in Super Bowl XXXIII at age 38 and is currently Denver's Executive Vice President of Football Operations and General Manager. Quarterback Jeff Dean had jersey number 37 in Champ Bowl XXXIV."

Question: "What is the name of the quarterback who was 38 in Super Bowl XXXIII?" Original Prediction: John Elway Prediction under adversary: Jeff Dean

Adversarial Examples for Evaluating Reading Comprehension Systems, EMNLP'17

The self-driving car incorrectly decides to turn right on Input 2 and crashes into the guardrail





(a) Input 1

(b) Input 2 (darker version of 1)

DeepXplore:Automated Whitebox Testing of Deep Learning Systems, SOSP'17 The Ensemble model is fooled by the addition of an adversarial distracting sentence in blue.

Article: Super Bowl 50

Paragraph: "Peyton Manning became the first quarterback ever to lead two different teams to multiple Super Bowls. He is also the oldest quarterback ever to play in a Super Bowl at age 39. The past record was held by John Elway, who led the Broncos to victory in Super Bowl XXXIII at age 38 and is currently Denver's Executive Vice President of Football Operations and General Manager. Quarterback Jeff Dean had jersey number 37 in Champ Bowl XXXIV."

Question: "What is the name of the quarterback who was 38 in Super Bowl XXXIII?" Original Prediction: John Elway Prediction under adversary: Jeff Dean

Adversarial Examples for Evaluating Reading Comprehension Systems, EMNLP'17 Adding small noise to the input audio makes the network transcribe any arbitrary phrase



"it was the best of times, it was the worst of times"





Audio Adversarial Examples: Targeted Attacks on Speech-to-Text, ICML 2018 ⁶

Attacks based on intensity changes in images



Attacks based on intensity changes in images



Attacks based on intensity changes in images



To verify absence of attack:

 L_{∞} -norm: consider all images I in the ϵ -ball $\mathcal{B}_{(I_{0},\infty)}(\epsilon)$ around I_{0}

Attacks based on geometric transformations



Attacks based on geometric transformations



Attacks based on geometric transformations



To verify absence of attack:

Consider all images I obtained by applying geometric transformations to $\mathcal{B}_{(I_0,\infty)}(\epsilon)$

Attacks based on intensity changes to sound



Attacks based on intensity changes to sound



Attacks based on intensity changes to sound



To verify absence of attack:

Consider all signals s in the ϵ -ball $\mathcal{B}_{(s_0,\infty)}(\epsilon)$ around s_0

Neural Network Verification: Problem statement

Given: Neural Network f, Input Region R Safety Property S

Prove: for all I in R, prove f(I) satisfies S

Example networks and regions:

Image classification network f Region R based on changes to pixel intensity Region R based on geometric: e.g., *rotation*

Speech recognition network f Region R based on added noise to audio signal

Aircraft collision avoidance network f Region R based on input sensor values

Input Region R can contain an infinite number of inputs, thus enumeration is infeasible

Experimental vs. Certified Robustness

Experimental robustness

Certified robustness

Tries to find violating inputs

Like testing, no full guarantees

E.g. Goodfellow 2014, Carlini & Wagner 2016, Madry et al. 2017

Prove absence of violating inputs

Actual verification guarantees

E.g.: Reluplex [2017], Wong et al. 2018, Al2 [2018]

In this talk we will focus on certified robustness

General Approaches to Network Verification

Complete verifiers: exact but suffer from scalability issues: SMT: Reluplex [CAV'17], MILP: MIPVerify [ICLR'19], Splitting: Neurify [NeurIPS'18],...

Incomplete verifiers, trade-off precision for scalability: Box/HBox [ICML'18], SDP [ICLR'18], Wong et.al. [ICML'18], FastLin [ICML'18], Crown [NeurIPS'18],...

Key Challenge: scalable and precise automated verifier

Network Verification with Eran



Complete and Incomplete Verification with ERAN

Faster Complete Verification

Aircraft collision avoidance system (ACAS)					
Reluplex	Neurify	ERAN			
> 32 hours	921 sec	227 sec			

Scalable Incomplete Verification



Geometric and Audio Verification with ERAN

Geometric Verification

Rotation between CNN wit	n -30° and 30° o ch 4,804 neurons	n MNIST s
ϵ	%verified	Time(s)
0.001	86	10 sec

Audio Verification

LSTM with 64 hidden neurons					
ε	%verified	Time (s)			
-110 dB	90%	9 sec			

Example: Analysis of a Toy Neural Network



We want to prove that $x_{11} > x_{12}$ for all values of x_1, x_2 in the input set



 $\begin{array}{l} s.t.: \ x_{11} = x_9 + x_{10} + 1, \ x_{12} = x_{10}, \\ x_9 = \max(0, x_7), \ x_{10} = \max(0, x_8), \\ x_7 = x_5 + x_6, \ x_8 = x_5 - x_6, \\ x_5 = \max(0, x_3), \ x_6 = \max(0, x_4), \\ x_3 = x_1 + x_2, \ x_4 = x_1 - x_2, \\ -1 \leq x_1 \leq 1, \ -1 \leq x_2 \leq 1. \end{array}$



 $min \ x_{11} - x_{12}$

 $s.t.: x_{11} = x_9 + x_{10} + 1, x_{12} = x_{10},$ $x_9 = \max(0, x_7), x_{10} = \max(0, x_8),$ $x_7 = x_5 + x_6, x_8 = x_5 - x_6,$ $x_5 = \max(0, x_3), x_6 = \max(0, x_4),$ $x_3 = x_1 + x_2, x_4 = x_1 - x_2,$ $-1 \le x_1 \le 1, -1 \le x_2 \le 1.$

Each $x_j = \max(0, x_i)$ corresponds to $(x_i \le 0 \text{ and } x_j = 0) \text{ or}$ $(x_i > 0 \text{ and } x_j = x_i)$

Solver has to explore two paths per ReLU resulting in exponential number of paths

Complete verification with solvers often does not scale

Network Verification with ERAN: High Level Idea



Box Approximation (scalable but imprecise)



Box Approximation (scalable but imprecise)



Verification with the Box domain fails as it cannot capture relational information

DeepPoly Approximation [POPL'19]

Shape: associate a lower polyhedral a_i^{\leq} and an upper polyhedral a_i^{\geq} constraint with each x_i

Key points:

Captures affine transformation precisely

Custom approximations for ReLU, sigmoid, tanh, and maxpool activations

Less precise but more scalable than general Polyhedra

Example: Verification using DeepPoly



ReLU activation

Pointwise transformer for $x_j \coloneqq max(0, x_i)$ that uses l_i, u_i

$$\begin{array}{l} if \ u_i \leq 0, a_j^{\leq} = a_j^{\geq} = 0, l_j = u_j = 0, \\ if \ l_i \geq 0, a_j^{\leq} = a_j^{\geq} = x_i, l_j = l_i, u_j = u_i, \\ if \ l_i < 0 \ and \ u_i > 0 \end{array}$$





ReLU activation

Pointwise transformer for $x_i \coloneqq max(0, x_i)$ that uses l_i, u_i

$$\begin{array}{l} \mbox{if } u_i \leq 0, a_j^{\leq} = a_j^{\geq} = 0, l_j = u_j = 0, \\ \mbox{if } l_i \geq 0, a_j^{\leq} = a_j^{\geq} = x_i, l_j = l_i, u_j = u_i, \\ \mbox{if } l_i < 0 \mbox{ and } u_i > 0 \end{array}$$







choose (b) or (c) depending on the area

ReLU activation

Pointwise transformer for $x_i \coloneqq max(0, x_i)$ that uses l_i, u_i

$$\begin{array}{l} \mbox{if } u_i \leq 0, a_j^{\leq} = a_j^{\geq} = 0, l_j = u_j = 0, \\ \mbox{if } l_i \geq 0, a_j^{\leq} = a_j^{\geq} = x_i, l_j = l_i, u_j = u_i, \\ \mbox{if } l_i < 0 \mbox{ and } u_i > 0 \end{array}$$

$$egin{aligned} &\langle x_5 \geq 0, \ &x_5 \leq 0.5 \cdot x_3 + 1, \ &l_5 = 0, \ &u_5 = 2
angle \end{aligned}$$







choose (b) or (c) depending on the area

Constant runtime

Affine transformation after ReLU



Affine transformation after ReLU



Imprecise upper bound u_7 by substituting u_5 , u_6 for x_5 and x_6 in a_7^{\geq} 34

Backsubstitution



Backsubstitution







Affine transformation with backsubstitution is pointwise, complexity: $O(w_{max}^2 L)^{38}$



Checking for robustness

Prove $x_{11} - x_{12} > 0$ for all inputs in $[-1,1] \times [-1,1]$

$\langle x_{11} \ge x_9 + x_{10} + 1,$	$\langle x_{12} \ge x_{10},$
$x_{11} \le x_9 + x_{10} + 1,$	$x_{11} \le x_{10},$
$l_{11} = 1,$	$l_{12} = 0,$
$u_{11}=5.5 angle$	$u_{12}=2 angle$

Computing lower bound for $x_{11} - x_{12}$ using l_{11} , u_{12} gives -1 which is an imprecise result

With backsubstitution, one gets 1 as the lower bound for $x_{11} - x_{12}$, proving robustness

Medium sized benchmarks

Dataset	Model	Туре	#Neurons	#Layers	Defense
MNIST	6×100	feedforward	610	6	None
	6 × 200	feedforward	1,210	6	None
	9 × 200	feedforward	1,810	9	None
	ConvSmall	convolutional	3,604	3	DiffAl
	ConvBig	convolutional	34,688	6	DiffAl
	ConvSuper	convolutional	88,500	6	DiffAl
CIFAR I 0	ConvSmall	convolutional	4,852	3	DiffAl

Results on medium sized benchmarks

Dataset	Model	ε	DeepZ		DeepPoly		RefineZono	
			%	time(s)	% ✓	time(s)	% 🗸	time(s)
MNIST	6×100	0.02	31	0.6	47	0.2	67	194
	6 × 200	0.015	13	1.8	32	0.5	39	567
	9 × 200	0.015	12	3.7	30	0.9	38	826
	ConvSmall	0.12	7	1.4	13	6.0	21	748
	ConvBig	0.2	79	7	78	61	80	193
	ConvSuper	0.1	97	133	97	400	97	665
CIFAR I 0	ConvSmall	0.03	17	5.8	21	20	21	550

Large benchmarks

Dataset	Model	Туре	#Neurons	#Layers	Defense
CIFAR 10	ResNetTiny	residual	311K	12	PGD
	ResNet18	residual	558K	18	PGD
	ResNetTiny	residual	311K	12	DiffAl
	SkipNet18	residual	558K	18	DiffAl
	ResNet18	residual	558K	18	DiffAl
	ResNet34	residual	967K	34	DiffAl

Results on large benchmarks

Model	Training	E	Hbox [ICML'18]		GPUPoly	
			% 🗸	time(s)	% ✓	time(s)
ResNetTiny	PGD	0.002	0	0.3	82	30
ResNet18	PGD	0.002	0	6.8	77	1400
ResNetTiny	DiffAl	0.03	64	0.3	69	7.6
SkipNet18	DiffAl	0.03	77	6. I	83	57
ResNet18	DiffAl	0.03	67	6.3	72	37
ResNet34	DiffAl	0.03	59	16	66	79

Network Verification with Eran



In-progress work in verification/training (sample)

Verification Precision: More precise convex relaxations by considering multiple ReLUs

Verification Scalability: GPU-based custom abstract domains for handling large nets

Theory: Proof on Existence of Accurate and Provable Networks with Box

Provable Training: Procedure for training Provable and Accurate Networks

Applications: e.g., reinforcement learning, geometric, audio, sensors



The self-driving car incorrectly decides to turn right on Input 2 and crashes into the guardrail



(a) Input 1

DeepXplore: Automated Whitebox Testing of Deep Learning Systems, SOSP'17

The Ensemble model is fooled by the addition of an adversarial distracting sentence in blue.

Article: Super Bowl 50

Paragraph: "Peyton Manning became the first quarterback ever to lead two different teams to multiple Super Bowls. He is also the oldest quarterback ever to play in a Super Bowl at age 39. The past record was held by John Elway, who led the Broncos to victory in Super Bowl XXXIII at age 38 and is currently Denver's Executive Vice President of Football Operations and General Manager. Quarterback Jeff Dean had jersey number 37 in Champ Bowl XXXIV." **Question:** "What is the name of the quarterback who

was 38 in Super Bowl XXXIII?" **Original Prediction: John Elway** Prediction under adversary: Jeff Dean

Adversarial Examples for Evaluating Reading Comprehension Systems, EMNLP'17



Adding small noise to the input

audio makes the network

Audio Adversarial Examples: Targeted Attacks on Speech-to-Text, ICML 2018

Neural Network Verification Framework



More at: safeai.ethz.ch