

Lecture 9: Verifying Robustness

Trustworthy Machine Learning
Spring 2024

Formal Methods for Verified Robustness

- **Last lecture:**

- **Formalizing program verification: Pre/post conditions**
- **Verification as constraint solving**
- **Robustness checking as program verification**

- **Today:**

- **Specialized constraint solver ReluPlex for neural network verification**
- **Verifying robustness by abstract interpretation (box domain)**

Slides credit: Gagandeep Singh and Madhu Parthasarathy (UIUC)

Specifications over DNNs

Precondition

```
 $\forall x_1, x_2. l_1 \leq x_1 \leq u_1, l_2 \leq x_2 \leq u_2$ 
```

DNN f

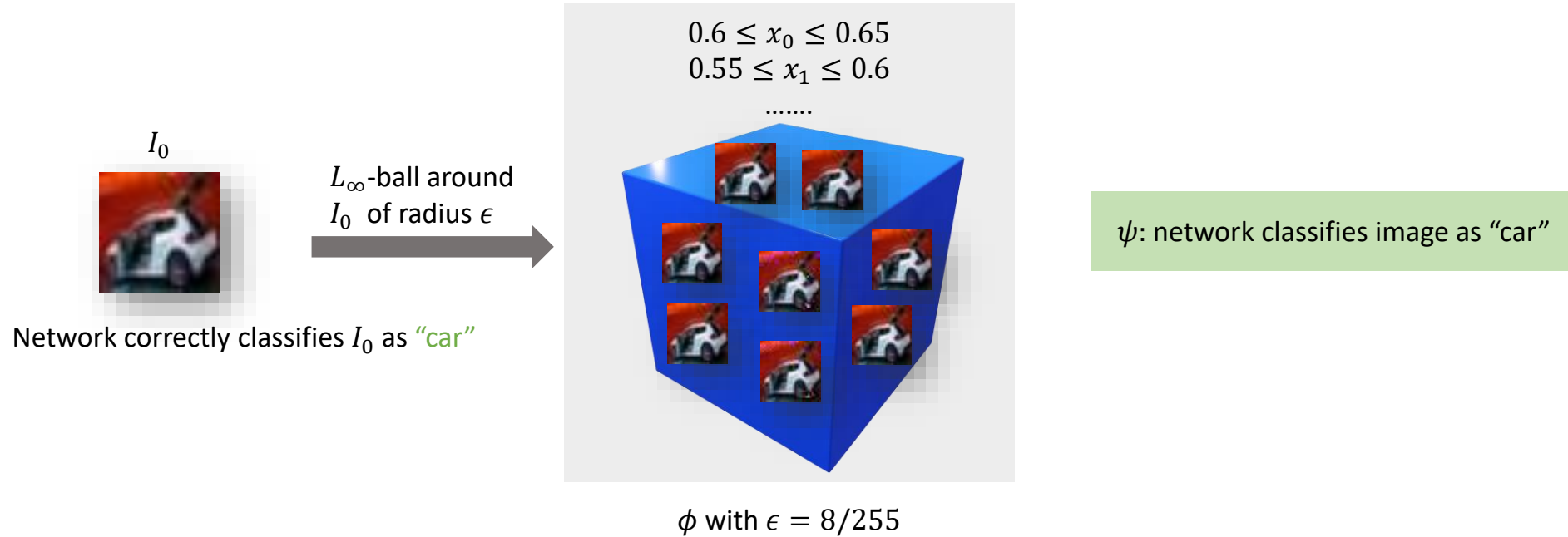
```
def  $f(x_1, x_2)$ :  
   $x_3 = w_{13} \cdot x_1 + w_{23} \cdot x_2 + b_3$   
   $x_4 = w_{14} \cdot x_1 + w_{24} \cdot x_2 + b_4$   
   $x_5 = \max(0, x_3)$   
   $x_6 = \max(0, x_4)$   
   $x_7 = w_{57} \cdot x_5 + w_{67} \cdot x_6 + b_7$   
   $x_8 = w_{56} \cdot x_5 + w_{68} \cdot x_6 + b_8$   
  return  $x_7, x_8$ 
```

Postcondition

```
 $x_7 > x_8$ 
```

Either prove that the network output satisfies the postcondition for all inputs in the pre-condition or find a counterexample

Robustness against adversarial perturbations



Verification of Neural Networks

Incomplete	Abstract interpretation: Box , Zonotope, DeepPoly
Complete	Mixed Integer Linear Programming (MILP) SMT solvers (Reluplex)

Active area of research with annual competition: VNNComp
Current winner: alpha-beta crown

Reluplex: An SMT based approach

Katz, Guy, et al. "Reluplex: An efficient SMT solver for verifying deep neural networks."
International Conference on Computer Aided Verification (CAV), 2017.

The Constraint Satisfaction Problem

Set of variables V

Atomic predicate:

- Linear inequality of the form $p = (\sum_{v_i \in V} a_i v_i \leq c_i)$
- ReLU equation of the form $p = (v_i = \text{ReLU}(v_j))$

Given a conjunction of atomic predicates $\varphi = p_1 \wedge \dots \wedge p_t$ decide if φ is satisfiable

- φ_N : conjunction of atomic predicates gives relation between input, output, and hidden variables of N .
- Pre-condition is given by a conjunction of linear inequalities $I = \{x \mid x \models \varphi_I\}$ and post-condition is a disjunction of linear (strict) inequalities $F = \{z \mid z \models \varphi_F\}$
- Sufficient to check if $\varphi_N \wedge \varphi_I \wedge \neg \varphi_F$ is satisfiable!

The Simplex Algorithm

- Solves satisfiability of conjunction of linear inequalities

$$\varphi = v_1 + v_2 \leq -5 \wedge v_1 - v_2 \geq 3$$

Add one new variable for each inequality

- Step 1: Construct an initial **configuration**

$u: V \rightarrow \mathbb{R} \cup \{\infty\}$
 $l: V \rightarrow \mathbb{R} \cup \{-\infty\}$
Upper and lower bounds on variable

$$V = \{v_1, v_2, v_3, v_4\}$$

Set of **basic** variable $B \subseteq V$
Initially it is all new variables

- $B = \{v_3, v_4\}$
- $T: v_3 = v_1 + v_2$ and $v_4 = v_1 - v_2$
- $u(v_3) = -5$ and $l(v_4) = 3$
- $\alpha(v_i) = 0$ for all $v_i \in V$

Tableau T contains one equation per basic variable; **RHS only has non-basic variables**

Valuation satisfies T

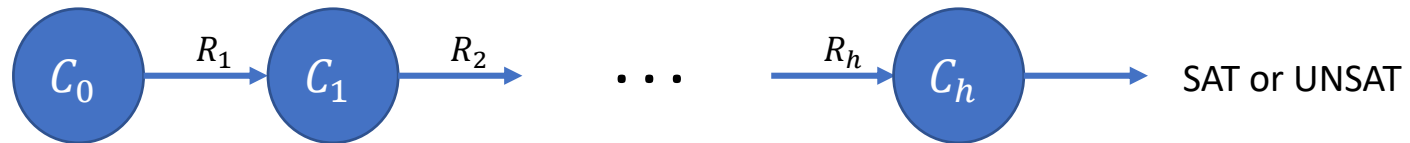
Goal of the algorithm: Update current values of non-basic variables to meet all lower/upper bounds

Derivations

Simplex provides rules of the following kind for modifying a configuration $C = (V, B, T, u, l, \alpha)$

$$\frac{\text{Some condition on } C}{\text{SAT}} \qquad \frac{\text{Some condition on } C}{\text{UNSAT}} \qquad \frac{\text{Some condition on } C}{\text{Modified configuration } C'}$$

Apply rules until a SAT or UNSAT is derived



Slack Variables

Coefficient of v_j in RHS defining v_i

$$\mathbf{slack}^+(v_i) = \{v_j \notin B \mid (T_{i,j} > 0 \wedge \alpha(v_j) < u(v_j)) \vee (T_{i,j} < 0 \wedge \alpha(v_j) > l(v_j))\}$$

$$\mathbf{slack}^-(v_i) = \{v_j \notin B \mid (T_{i,j} < 0 \wedge \alpha(v_j) < u(v_j)) \vee (T_{i,j} > 0 \wedge \alpha(v_j) > l(v_j))\}$$

$slack^+(v_i)$: Variables on RHS defining v_i whose values can be changed to **increase** the value of v_i
e.g. variable v_j has positive coefficient and its current value is less than its upper bound

$slack^-(v_i)$: Variables on RHS defining v_i whose values can be changed to **decrease** the value of v_i

Slack Variables

$$\text{slack}^+(v_i) = \{v_j \notin B \mid (T_{i,j} > 0 \wedge \alpha(v_j) < u(v_j)) \vee (T_{i,j} < 0 \wedge \alpha(v_j) > l(v_j))\}$$

$$\text{slack}^-(v_i) = \{v_j \notin B \mid (T_{i,j} < 0 \wedge \alpha(v_j) < u(v_j)) \vee (T_{i,j} > 0 \wedge \alpha(v_j) > l(v_j))\}$$

$$V = \{v_1, v_2, v_3, v_4\}$$

- $B = \{v_3, v_4\}$
- $T : v_3 = v_1 + v_2$ and $v_4 = v_1 - v_2$
- $u(v_3) = -5$ and $l(v_4) = 3$
- $\alpha(v_i) = 0$ for all $v_i \in V$

$$\text{slack}^+(v_3) = \{v_1, v_2\}$$

If we have $u(v_1) = 0$ then
 $\text{slack}^+(v_3) = \{v_2\}$

Simplex Rule: Successful Termination

$$\frac{\forall v_i \in V, \quad l(v_i) \leq \alpha(v_i) \leq u(v_i)}{\text{SAT}} \quad \text{Success}$$

Current valuation meets all lower/upper bounds: satisfying assignment found

Simplex Rule: Unsuccessful Termination

$$\frac{v_i \in B, \quad (\alpha(v_i) < l(v_i) \wedge \mathbf{slack}^+(v_i) = \emptyset) \vee (\alpha(v_i) > u(v_i) \wedge \mathbf{slack}^-(v_i) = \emptyset)}{\text{UNSAT}} \quad \text{Failure}$$

There is a basic variable for which current value must be increased/decreased to meet lower/upper bound constraint but no such update of RHS vars is possible

Simplex Rule: Pivot

$$\frac{v_i \in B, \quad \alpha(v_i) < l(v_i), \quad v_e \in \text{slack}^+(v_i)}{T := \text{pivot}(T, v_i, v_e), \quad B := B \cup \{v_e\} \setminus \{v_i\}} \text{Pivot}_1$$
$$\frac{v_i \in B, \quad \alpha(v_i) > u(v_i), \quad v_e \in \text{slack}^-(v_i)}{T := \text{pivot}(T, v_i, v_e), \quad B := B \cup \{v_e\} \setminus \{v_i\}} \text{Pivot}_2$$

If: A basic variable's value needs to be increased/decreased to meet lower/upper bound
and there is a possible variable on RHS whose value be changed for this purpose
Then make it non-basic by swapping their roles using pivot

Pivot Example

Pivot: Allows replacing basic variable with a non-basic variable

$$T = \begin{cases} v_3 = v_1 + v_2 \\ v_4 = v_1 - v_2 \end{cases} \xrightarrow{\text{pivot}(T, 3, 1)} \begin{cases} v_1 = v_3 - v_2 \\ v_4 = (v_3 - v_2) - v_2 \end{cases} = T'$$

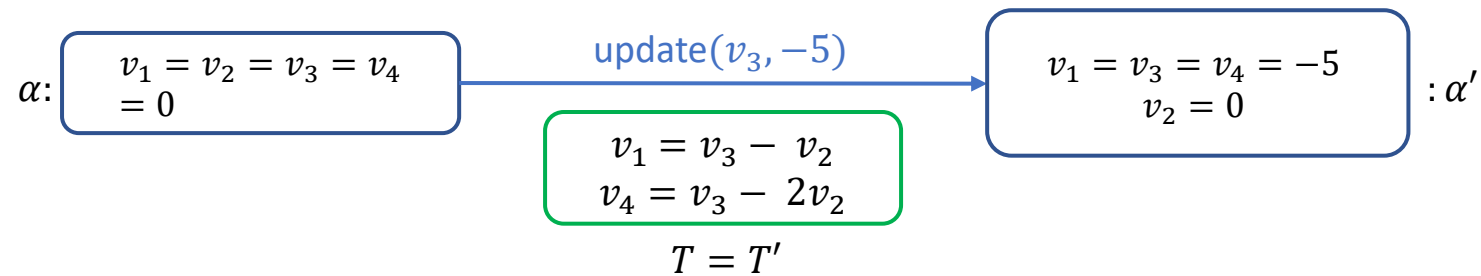
Simplex Rule: Update

$$\frac{v_j \notin B, \quad \alpha(v_j) < l(v_j) \vee \alpha(v_j) > u(v_j), \quad l(v_j) \leq \alpha(v_j) + \delta \leq u(v_j)}{\alpha := \text{update}(\alpha, v_j, \delta)} \quad \text{Update}$$

For a non-basic variable, if its current value is less/greater than lower/upper bound then increase/decrease it to meet the bound

Update Example

Update: Allows updating value of a non-basic variable



Simplex Algorithm in SMT Solver

- Starting with initial configuration, keep applying pivot/update rules until the termination condition holds
- Multiple rules may be applicable in a given configuration
 - Derivation tree captures all possible branches of rule applications
- Key engineering details of implementation
 - Which rule to choose in a given configuration
 - Choice of rules affects time to termination but does not require backtracking
 - How to apply rules efficiently (e.g. keeping track of slack variables)

Soundness and Completeness

SOUNDNESS: If there is a derivation to SAT (or UNSAT), φ is satisfiable (or not).

COMPLETENESS: There is always a derivation to either SAT or UNSAT.

Constraints in Neural network Verification

Precondition

$$\forall x_1, x_2. l_1 \leq x_1 \leq u_1, l_2 \leq x_2 \leq u_2$$

DNN f

```
def  $f(x_1, x_2)$ :  
   $x_3 = w_{13} \cdot x_1 + w_{23} \cdot x_2 + b_3$   
   $x_4 = w_{14} \cdot x_1 + w_{24} \cdot x_2 + b_4$   
   $x_5 = \text{ReLU}(x_3) = \max(0, x_3)$   
   $x_6 = \text{ReLU}(x_4) = \max(0, x_4)$   
   $x_7 = w_{57} \cdot x_5 + w_{67} \cdot x_6 + b_7$   
   $x_8 = w_{56} \cdot x_5 + w_{68} \cdot x_6 + b_8$   
  return  $x_7, x_8$ 
```

Postcondition

$$x_7 > x_8$$

Simplex to Reluplex

- Solves satisfiability of conjunction of linear inequalities **and ReLU equations**.

$$\varphi = v_1 + v_2 \leq -5 \wedge v_1 - v_2 \geq 3 \wedge v_1 = \text{ReLU}(v_2)$$

- Step 1: Construct an initial **configuration** $C = (V, B, T, u, l, \alpha, \mathbf{R})$

$$V = \{v_1, v_2, v_3, v_4\}$$

- $\mathbf{B} = \{v_3, v_4\}$
- $\mathbf{T} : v_3 = v_1 + v_2$ and $v_4 = v_1 - v_2$
- $\mathbf{u}(v_3) = -5$ and $\mathbf{l}(v_4) = 3$
- $\alpha(v_i) = 0$ for all $v_i \in V$
- $\mathbf{R} = \{(v_1 = \text{ReLU}(v_2))\}$

Set of ReLU constraints R

Modified Successful Termination Test

$$\frac{\forall v_i \in V. l(v_i) \leq \alpha(v_i) \leq u(v_i), \quad \forall (v_f = \text{ReLU}(v_b)) \in R. \alpha(v_f) = \text{ReLU}(\alpha(v_b))}{\text{SAT}} \quad \text{ReluSuccess}$$

Current valuation meets all lower/upper bounds and all ReLU constraints

Additional Pivot Rule

$$\frac{v_i \in B, \quad \exists v_j. (v_j = \text{ReLU}(v_i)) \in R \vee (v_i = \text{ReLU}(v_j)) \in R, \quad v_e \notin B, T_{i,e} \neq 0}{T := \text{pivot}(T, v_i, v_e), \quad B := B \cup \{v_e\} \setminus \{v_i\}} \text{ReluPivot}$$

If a basic variable v_i is involved in a ReLU constraint,
Then swap it's a role with a non-basic variable v_e in its RHS with non-zero coefficient using pivot

Additional Update Rules

$$\frac{v_j \notin B, \quad (v_j = \text{ReLU}(v_i)) \in R, \quad \alpha(v_j) \neq \text{ReLU}(\alpha(v_i))}{\alpha := \text{update}(\alpha, v_j, \text{ReLU}(\alpha(v_i)) - \alpha(v_j))} \quad \text{Update}_f$$

$$\frac{v_i \notin B, \quad (v_j = \text{ReLU}(v_i)) \in R, \quad \alpha(v_j) \neq \text{ReLU}(\alpha(v_i)), \quad \alpha(v_j) \geq 0}{\alpha := \text{update}(\alpha, v_i, \alpha(v_j) - \alpha(v_i))} \quad \text{Update}_b$$

If a non-basic variable v_i is involved in a ReLU constraint that's violated by its current value,
Then update its value to satisfy the ReLU constraint

New Split Rule

$$\frac{(v_j = \text{ReLU}(v_i)) \in R, \quad l(v_i) < 0, \quad u(v_i) > 0}{u(v_i) := 0 \quad l(v_i) := 0} \quad \text{ReluSplit}$$

For the constraint $v_j = \text{ReLU}(v_i) = \max(0, v_i)$, when v_i can be both positive and negative we have two cases:

Case 1: v_i is positive (achieved by setting its lower bound to 0)

Case 2: v_i is negative (achieved by setting its upper bound to 0)

The two cases create two branches in the derivation tree

A priori we don't know which one will lead to success (so may require backtracking in proof search)

Case split is due to non-linearity of ReLU and crux of computational difficulty

ReluPLEX Algorithm in SMT Solver

- Starting with initial configuration, keep applying rules until SAT leaf found or all branches caused by split lead to UNSAT
- Multiple rules may be applicable in a given configuration
 - Derivation tree captures all possible branches of rule applications
 - **Key difference with Simplex: Backtracking (exploring different branches) may be required!**
 - **Key benefit of Reluplex: Case split is demand driven and happens only when necessary**
- Key engineering details of implementation
 - Which rule to choose in a given configuration
 - Choice of rules affects backtracking and time to termination
 - How to apply rules efficiently and how to backtrack efficiently

Soundness and Completeness

SOUNDNESS: If there is a derivation tree with at least one SAT leaf, φ is satisfiable.
If there is a derivation tree with all UNSAT leaves, φ is not satisfiable.

COMPLETENESS: There is always a derivation tree in which every leaf is either SAT or UNSAT.

Comparison with existing SMT solvers

- Can encode $v_1 = \text{ReLU}(v_2)$ as $(v_2 \geq 0 \wedge v_1 = v_2) \vee (v_2 \leq 0 \wedge v_1 = 0)$.
- Existing SMT solvers perform many case splits.
- Reluplex can avoid/reduce splitting by using new pivot and update rules first.

	φ_1	φ_2	φ_3	φ_4	φ_5	φ_6	φ_7	φ_8
CVC4	-	-	-	-	-	-	-	-
Z3	-	-	-	-	-	-	-	-
Yices	1	37	-	-	-	-	-	-
MathSat	2040	9780	-	-	-	-	-	-
Gurobi	1	1	1	-	-	-	-	-
Reluplex	8	2	7	7	93	4	7	9

Time to termination in seconds with 4 hour timeout

Properties are from case study of neural-network-based controller for collision avoidance protocol ACAS (see ReluPlex paper; Katz et al; CAV 2017)

Experiments

Table 1 Local adversarial robustness tests. All times are in seconds.

	$\delta = 0.1$		$\delta = 0.075$		$\delta = 0.05$		$\delta = 0.025$		$\delta = 0.01$		Total Time	
	Result	Time	Result	Time	Result	Time	Result	Time	Result	Time		
Point 1	SAT	135	SAT	239	SAT	24	UNSAT	609	UNSAT	57	1064	SAT : Not Robust
Point 2	UNSAT	5880	UNSAT	1167	UNSAT	285	UNSAT	57	UNSAT	5	7394	UNSAT : Robust
Point 3	UNSAT	863	UNSAT	436	UNSAT	99	UNSAT	53	UNSAT	1	1452	
Point 4	SAT	2	SAT	977	SAT	1168	UNSAT	656	UNSAT	7	2810	
Point 5	UNSAT	14560	UNSAT	4344	UNSAT	1331	UNSAT	221	UNSAT	6	20462	

Neural Network: Fully connected, 8 layers, 300 neurons

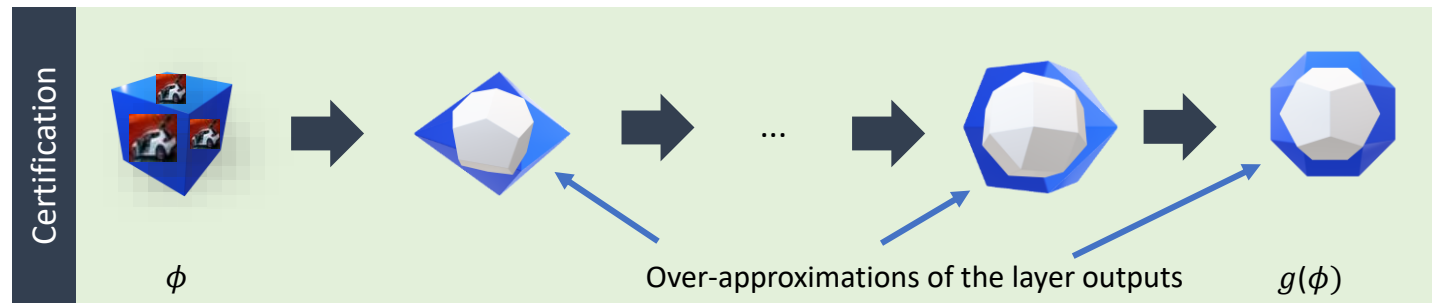
Abstract Interpretation using Boxes

Singh, et al. "Fast and effective robustness certification" NeurIPS, 2018.

Incomplete methods

We will investigate a specific type of incomplete method, based on bound propagation through the neural network. Starting with the initial pre-condition ϕ , we will “pass” ϕ through the network, computing a convex over-approximation of the effect of each layer on ϕ . Next, let's look at the “recipe” for certification with bound propagation

Step1: compute convex $g(\phi)$ by propagating ϕ



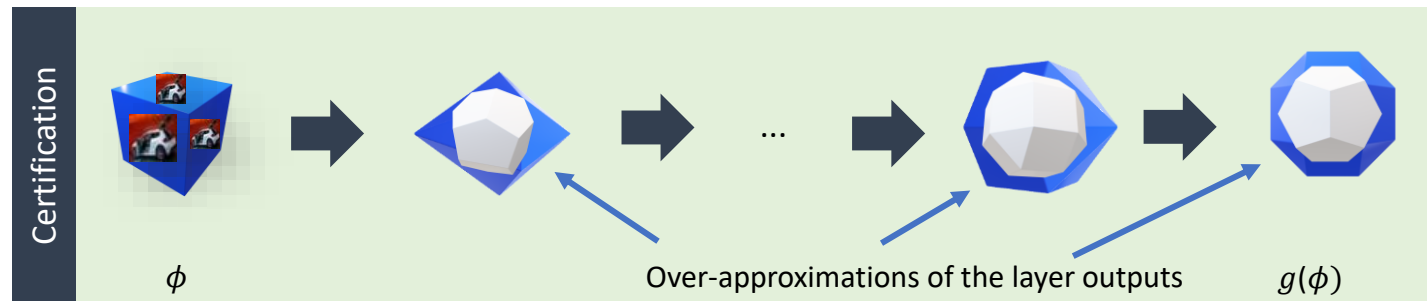
ϕ :
 $x_0 = [0.588, 0.65]$
 $x_1 = [0.545, 0.608]$
 $x_2 = [0.533, 0.596]$
...
 $x_{3071} = [0.4, 0.463]$

All possible outputs

$g(\phi)$

$o_0 = 0$
 $o_1 = 2.60 + 0.015\eta_0 + 0.023\eta_1 + 5.181\eta_2 + \dots$
 $o_2 = 4.63 - 0.005\eta_0 - 0.006\eta_1 + 0.023\eta_2 + \dots$
...
 $o_9 = 0.12 - 0.125\eta_0 + 0.102\eta_1 + 3.012\eta_2 + \dots$
 $\forall i. \eta_i \in [0, 1]$

Step2: verify ψ



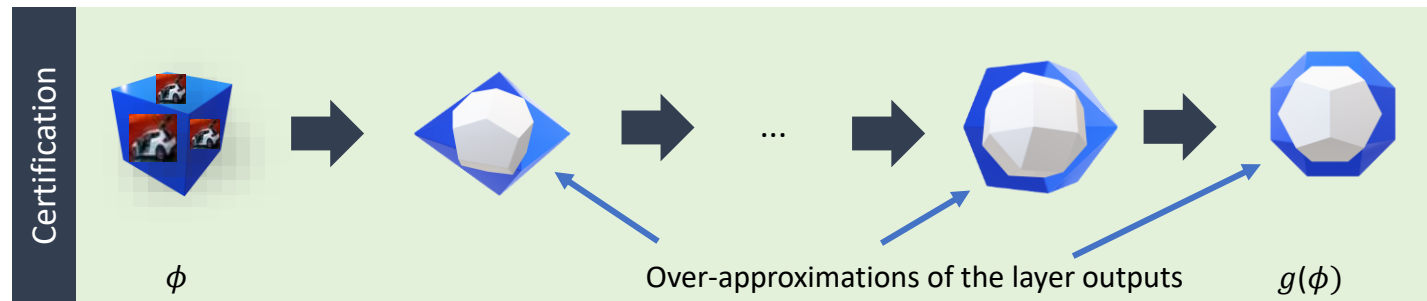
$$g(\phi)$$
$$o_0 = 0$$
$$o_1 = 2.60 + 0.015\eta_0 + 0.023\eta_1 + 5.181\eta_2 + \dots$$
$$o_2 = 4.63 - 0.005\eta_0 - 0.006\eta_1 + 0.023\eta_2 + \dots$$
$$\dots$$
$$o_9 = 0.12 - 0.125\eta_0 + 0.102\eta_1 + 3.012\eta_2 + \dots$$
$$\forall i. \eta_i \in [0,1]$$

ψ : every point in ϕ classifies as car




ψ : $\forall l \neq car. o_{car} > o_l$

$g(\phi)$ is an **over-approximation**, hence if we fail to prove ψ , it could be the property is actually violated or there was over-approximation introduced during propagation which prohibits provability.

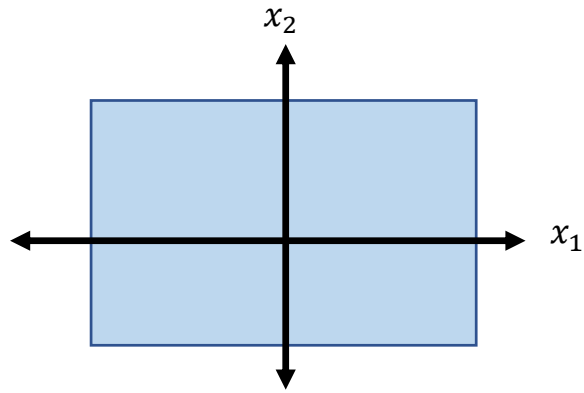
Key challenge: how to produce convex shapes?



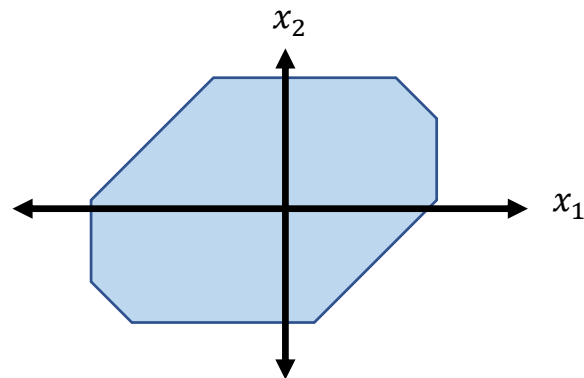
To instantiate incomplete methods which use bound propagation, we need two parts:

1. What is the convex approximation  ? E.g., Box, Zonotope, Polyhedra
2. How are these convex approximations produced? That is, what is the effect of the layer  on a given approximation  ? This effect is often called an **abstract transformer** as it transforms abstract shapes.

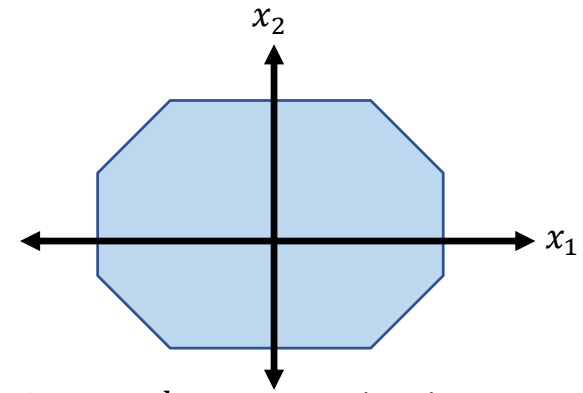
Popular convex shapes



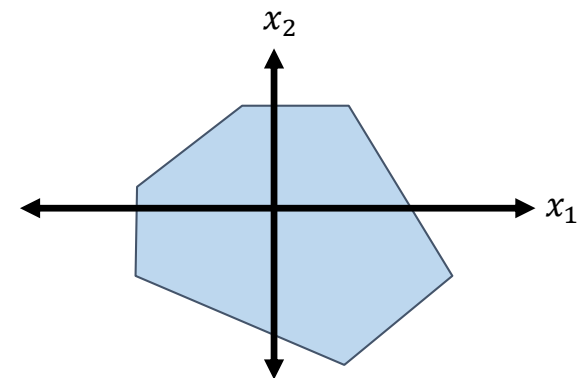
Box: $l_i \leq x_i \leq u_i$



Zonotope: $\hat{x}_i = \alpha_0 + \sum_i \alpha_i \epsilon_i, \epsilon_i \in [-1, 1]$



Octagon: $l_i \leq x_i \leq u_i, \pm x_i \pm x_j \leq c_{ij}$



Polyhedra: $\sum_i a_i x_i \leq c$

Speed vs. precision tradeoff

A trade-off between approximation and speed exists:

transformers for Box are fast, but imprecise,

while Polyhedra is more precise but has exponential complexity.

Box Abstract Transformer: Addition

$$[\mathbf{a}, \mathbf{b}] +^{\#} [\mathbf{c}, \mathbf{d}] = [\mathbf{a} + \mathbf{c}, \mathbf{b} + \mathbf{d}]$$

Addition transformer

- Here, $\mathbf{a}, \mathbf{b} \in \mathbf{R}^m$ where $\forall i. a_i \leq b_i$
- $\#$ denotes the abstract effect of the operation on the box

Example in 1-dim (interval addition) $[1, 3] +^{\#} [5, 8] = [6, 11]$

Box Abstract Transformers for ReLU Networks

$$[a, b] +^{\#} [c, d] = [a + c, b + d]$$

Addition transformer

$$-^{\#}[a, b] = [-b, -a]$$

Negation transformer

$$\mathbf{ReLU}^{\#}[a, b] = [\mathbf{ReLU}(a), \mathbf{ReLU}(b)]$$

ReLU transformer

$$\lambda^{\#}[a, b] = [\lambda * a, \lambda * b]$$

Multiplication by a constant $\lambda > 0$

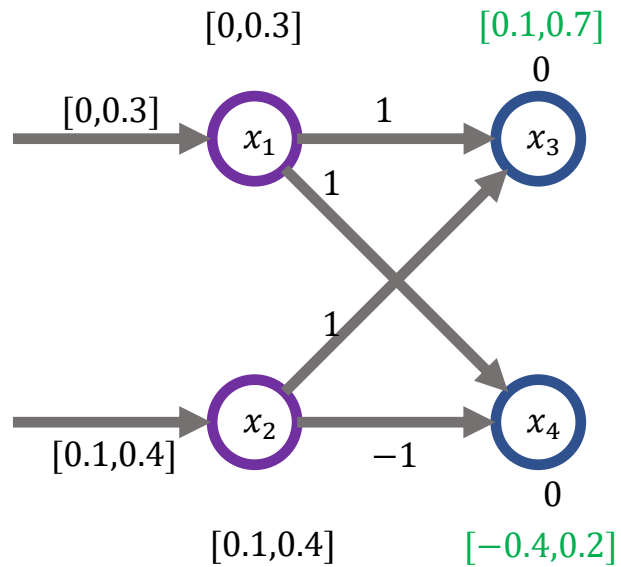
$$[a, b] >^{\#} [c, d] \text{ holds iff } a > d$$

Comparison transformer

- Here, $a, b \in \mathbf{R}^m$ where $\forall i. a_i \leq b_i$
- $\mathbf{ReLU}(x) = \max(0, x)$
- $\#$ denotes the abstract effect of the operation on the box

Computation with Box transformer

We have 2 pixels (x_1, x_2) as input ranging over $[0, 0.3]$ and $[0.1, 0.4]$



```
def  $f_1(x_1, x_2)$ :  
     $x_3 = x_1 + x_2$   
     $x_4 = x_1 - x_2$   
    return  $x_3, x_4$ 
```

Bounds using Box:

$$0.1 \leq x_3 \leq 0.7$$

$$-0.4 \leq x_4 \leq 0.2$$

Exact bounds would be:

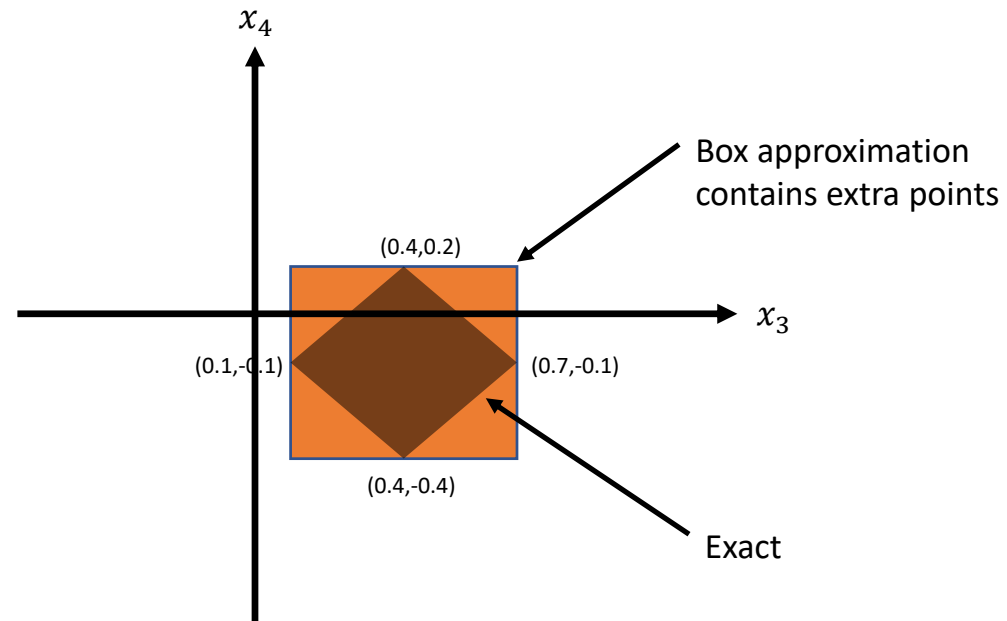
$$x_3 = x_1 + x_2$$

$$x_4 = x_1 - x_2$$

$$0 \leq x_1 \leq 0.3$$

$$0.1 \leq x_2 \leq 0.4$$

Optimal Box transformer is not exact!

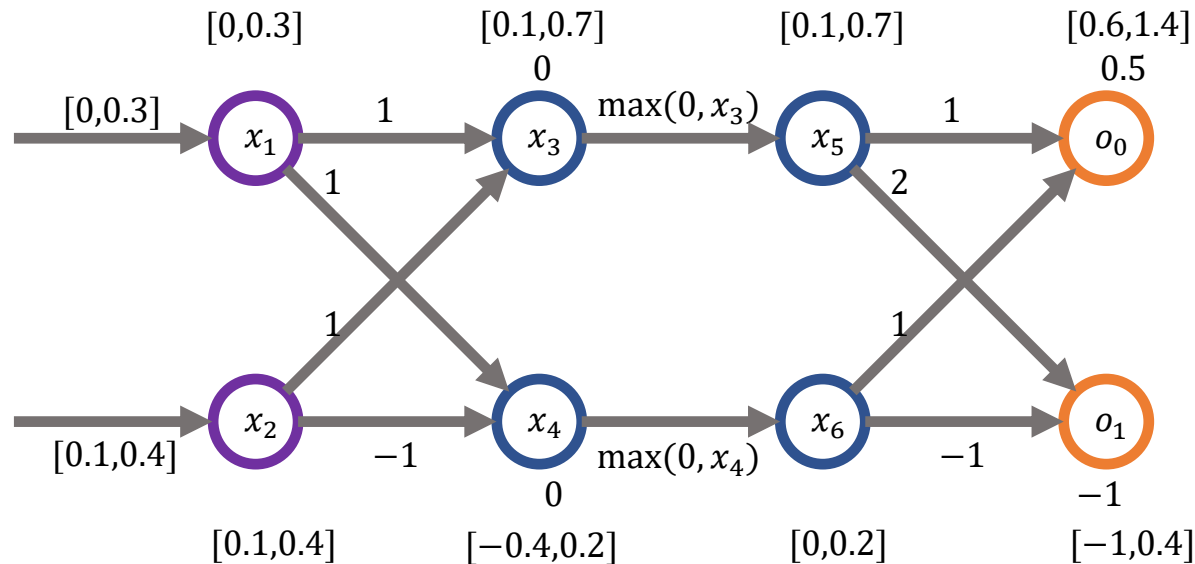


Key Point

Even though the Box abstract transformer for the affine computation is optimal for the Box relaxation, it may not be complete (exact)! Nonetheless, even if not exact, it may be enough to verify the property of interest.

Box succeeds in verifying robustness

We have 2 pixels (x_1, x_2) as input ranging over $[0, 0.3]$ and $[0.1, 0.4]$, we want to prove that $o_0 > o_1$ holds for all inputs

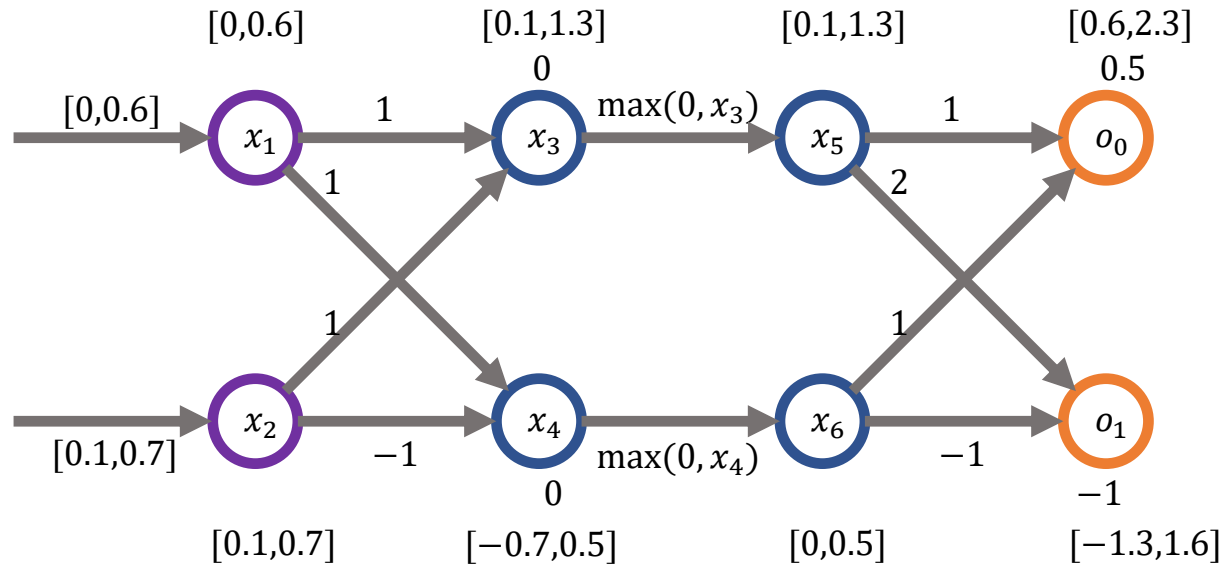


```
def f(x1, x2):  
    x3 = x1 + x2  
    x4 = x1 - x2  
    x5 = max(0, x3)  
    x6 = max(0, x4)  
    o0 = x5 + x6 + 0.5  
    o1 = 2 * x5 - x6 - 1  
    return o0, o1
```

Using Box, we succeed in proving the network classifies any input in the range as 0.
This is because $[0.6, 1.4] > [-1, 0.4]$, provably so.

Box fails in verifying robustness

Let us **slightly increase** the range of the input pixels to $[0, 0.6]$ and $[0.1, 0.7]$



def $f(x_1, x_2)$:

$$x_3 = x_1 + x_2$$

$$x_4 = x_1 - x_2$$

$$x_5 = \max(0, x_3)$$

$$x_6 = \max(0, x_4)$$

$$o_0 = x_5 + x_6 + 0.5$$

$$o_1 = 2 \cdot x_5 - x_6 - 1$$

return o_0, o_1

Using Box, we failed to prove the network classifies any input in the range as 0, even though property actually holds. This is because $[0.6, 2.3]$ is not $> [-1.3, 1.6]$, provably so.

Robustness Certification

- Instead of Box domain, one can use Zonotopes to get reasonably scalable and more precise approximations
- Constraint solving and abstract domains can be combined in different ways
- Robustness verification can be integrated into training
 - Parameters are updated not only to minimize loss but also to ensure that a verification procedure (e.g. based on Box abstract domain) gives robustness guarantee

Adversarial Robustness Recap

- **Adversarial Examples**
- **Adversarial training**
- **Certified robustness via randomized smoothing**
- **Sample of current research on robustness for LLMs**
- **Formal methods for verified robustness**
 - **Specialized constraint solver ReluPlex for neural network verification**
 - **Verifying robustness by abstract interpretation (box domain)**