

# Reliable and Interpretable Artificial Intelligence

## Lecture 5: Certification with complete methods

Martin Vechev

ETH Zurich

Fall 2020

# Complete method: MILP

Let us incorporate MILP (Mixed Integer Linear Program) solvers in order to verify robustness. For ReLU networks, MILP is complete. However, this comes at the cost that MILP is at least NP-complete (cannot be solved in polynomial time in any known manner) due to underlying linear programs. Further, there is an additional exponential factor due to splitting.

We will first define the standard MILP problem. Then, we will use the Box bounds to constrain MILP. So, Box will help MILP be faster by doing less work (MILP still remains complete).

# MILP Generic Problem Definition

$$\min c_1x_1 + c_2x_2 + \dots + c_nx_n$$

objective

$$a_{11}x_1 + \dots + a_{1n}x_n \leq b_1$$

...

$$a_{m1}x_1 + \dots + a_{mn}x_n \leq b_m$$

constraints

$$l_i \leq x_i \leq u_i \quad 1 \leq i \leq n$$

bounds on continuous  $x_i$

$$x_j \in \mathbb{Z}$$

some  $x_j$  are integer

- $c_i, a_{ij}, b_i \in \mathbb{R}$  and  $l, u \in \mathbb{R}^n$
- some  $x_j$  's can be integers (or even binary), hence Mixed-Integer problem
- state-of-the-art solvers (e.g., Gurobi) require bounds on  $x_i$  's

# MILP encoding of Neural Network

To encode the network as a MILP instance, we need to:

1. Encode the affine layer
2. Encode the ReLU layer
3. Encode the pre-condition  $\phi$
4. Encode the post-condition  $\psi$

# Encode Affine layer as MILP

This direct as it is just a linear constraint

$$\mathbf{y} = \mathbf{W}\mathbf{x} + \mathbf{b}$$

where  $\mathbf{W}$  are the weights and  $\mathbf{b}$  is the bias

(note: convolution is also an affine transformation, can be encoded directly)

# Encode ReLU layer as MILP

ReLU definition is:

$$y = \text{ReLU}(x) = \max(0, x)$$

MILP ReLU encoding is:

$$y \leq x - l * (1 - a)$$

$$y \geq x$$

$$y \leq u * a$$

$$y \geq 0$$

$$a \in \{0, 1\} \longrightarrow a \text{ is a binary integer variable}$$

This assumes we have computed lower  $l$  and upper  $u$  bounds for the input neuron  $x$  (e.g., by using Box beforehand).

(\*Intuition of the encoding on board)

## Encode ReLU Layer as MILP Problem

$$\begin{array}{l}
 \swarrow a=0 \qquad \searrow a=1 \\
 \left. \begin{array}{l} y \leq x - l \\ y \geq x \\ y \leq 0 \\ y \geq 0 \end{array} \right\} x \in [l, 0] \quad \left. \begin{array}{l} y \leq x \\ y \geq x \\ y \leq u \\ y \geq 0 \end{array} \right\} \begin{array}{l} y = x \\ y \in [0, u] \end{array}
 \end{array}$$

Reminder: ReLU encoding

$$y \leq x - l * (1 - a)$$

$$y \geq x$$

$$y \leq u * a$$

$$y \geq 0$$

$$a \in \{0, 1\}$$

Case a) "Below zero", eg.  $x \in [-2, -1]$

$$\begin{array}{l}
 \swarrow a=0 \qquad \searrow a=1 \\
 \left. \begin{array}{l} y \leq [-2, -1] + 2 = [0, 1] \\ y \geq [-2, -1] \\ y = 0 \\ x \in [-2, 0] \end{array} \right\} \\
 \left. \begin{array}{l} y = x \\ y \in [0, -1] \end{array} \right\} \text{UNSAT}
 \end{array}$$

Case b) "Above zero", eg.  $x \in [1, 2]$

$$\begin{array}{l}
 \swarrow a=0 \qquad \searrow a=1 \\
 \left. \begin{array}{l} y = 0 \\ x \leq 0 \end{array} \right\} \text{UNSAT} \\
 \left. \begin{array}{l} y = x \\ y \in [0, 2] \end{array} \right\}
 \end{array}$$

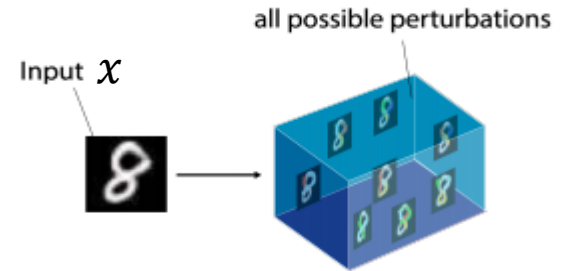
Case c) "Crossing", eg.  $x \in [-1, 2]$

$$\begin{array}{l}
 \swarrow a=0 \qquad \searrow a=1 \\
 \left. \begin{array}{l} y = 0 \\ x \in [-1, 0] \end{array} \right\} \\
 \left. \begin{array}{l} y = x \\ y \in [0, 2] \end{array} \right\}
 \end{array}$$

# Encode Pre-condition $\phi$ as MILP

Lets take  $\phi = L_\infty$  ball around  $x$  as an example:

$$L_\infty \text{ ball: } \text{Ball}(x)_\epsilon = \{x' \mid \|x - x'\|_\infty < \epsilon\}$$



MILP encoding:

$$x_i - \epsilon \leq x'_i \leq x_i + \epsilon$$



That is, we will introduce lower and upper bound constraints for each input neuron  $x'_i$



# Encode Post-condition $\psi$ as MILP

Lets take  $\psi = \text{label } 0 \text{ is more likely than label } 1$  (our example network)

$$\psi = o_0 > o_1$$



**We want to prove this.** Hence, this must be forming our MILP objective

MILP encoding:

$$\mathit{min} \ o_0 - o_1$$

After MILP finishes, final verification takes the computed values for  $o_0$ ,  $o_1$  and checks if  $o_0$  is indeed greater than  $o_1$ . If yes, verification **succeeds**, if not, verification **fails**.

# Generic vs. Instantiated MILP

## Generic MILP problem

$$\min c_1 x_1 + c_2 x_2 + \dots + c_n x_n$$

$$a_{11} x_1 + \dots + a_{1n} x_n \leq b_1$$

...

$$a_{m1} x_1 + \dots + a_{mn} x_n \leq b_m$$

$$l_i \leq x_i \leq u_i \quad 1 \leq i \leq n$$

$$x_j \in \mathbf{Z}$$

## Our MILP instance

$$\min o_0 - o_1$$

Plug in the affine and ReLU  
MILP encodings as defined

$$l_i \leq x_i^p \leq u_i$$

pre-computed Box bounds  
on all neurons in each layer  $p$

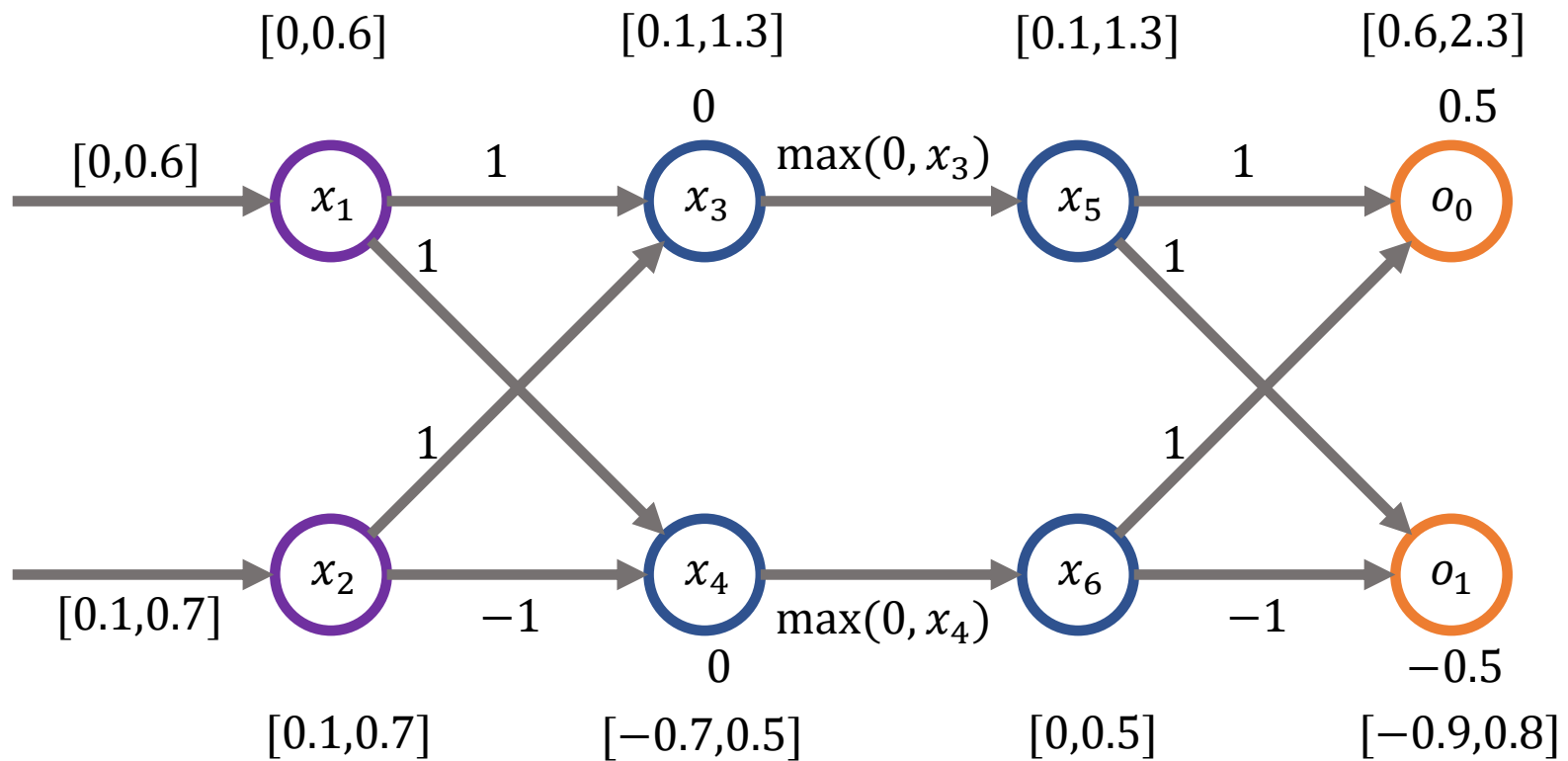
$$x_i - \epsilon \leq x'_i \leq x_i + \epsilon$$

$\phi$ : bounds on  
input neurons

$$a_j \in \{0, 1\}$$

These are the  $a \in \{0, 1\}$  vars  
from the ReLU encoding

Reminder: these are the bounds with Box



# MILP Instance for this network

$$\min o_0 - o_1$$

Affine

$$\begin{aligned}x_3 &= x_1 + x_2 \\x_4 &= x_1 - x_2 \\o_0 &= x_5 + x_6 + 0.5 \\o_1 &= x_5 - x_6 - 0.5\end{aligned}$$

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

$$\begin{aligned}x_5 &\leq x_3 - 0.1 * (1 - a_5) \\x_5 &\geq x_3 \\x_5 &\leq 1.3 * a_5 \\x_5 &\geq 0\end{aligned}$$

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

$$\begin{aligned}x_6 &\leq x_4 + 0.7 * (1 - a_6) \\x_6 &\geq x_4 \\x_6 &\leq 0.5 * a_6 \\x_6 &\geq 0\end{aligned}$$

Input bounds

$$\begin{aligned}0 &\leq x_1 \leq 0.6 \\0.1 &\leq x_2 \leq 0.7\end{aligned}$$

Pre-computed Box bounds

$$\begin{aligned}0.1 &\leq x_3 \leq 1.3 \\-0.7 &\leq x_4 \leq 0.5 \\0.1 &\leq x_5 \leq 1.3 \\0 &\leq x_6 \leq 0.5\end{aligned}$$

Binary integer variables

$$a_5, a_6 \in \{0, 1\}$$

Solving this MILP instance will lead to proving the property

# Lets solve MILP to see why Box helps

(Note: MILP solvers employ elaborate algorithms, the idea here is to show how Box can help even state-of-the-art MILP solvers)

# Case $a_5 = 0$

$$\min o_0 - o_1$$

Affine

$$\begin{aligned}x_3 &= x_1 + x_2 \\x_4 &= x_1 - x_2 \\o_0 &= x_5 + x_6 + 0.5 \\o_1 &= x_5 - x_6 - 0.5\end{aligned}$$

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

$$\begin{aligned}x_5 &\leq x_3 - 0.1 \\x_5 &\geq x_3 \\&\dots\end{aligned}$$

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

$$\begin{aligned}x_6 &\leq x_4 + 0.7 * (1 - a_6) \\x_6 &\geq x_4 \\x_6 &\leq 0.5 * a_6 \\x_6 &\geq 0\end{aligned}$$

Pre-computed Box bounds

$$\begin{aligned}0.1 &\leq x_3 \leq 1.3 \\-0.7 &\leq x_4 \leq 0.5 \\0.1 &\leq x_5 \leq 1.3 \\0 &\leq x_6 \leq 0.5\end{aligned}$$

This is an infeasible LP instance so MILP does not have to consider it .

Here, Box bounds **helped MILP** in not exploring further values for  $a_6$ . So it saved generating two cases for  $a_6$ .

# In practice, we directly generate this MILP

$$\min o_0 - o_1$$

Affine

$$\begin{aligned}x_3 &= x_1 + x_2 \\x_4 &= x_1 - x_2 \\o_0 &= x_5 + x_6 + 0.5 \\o_1 &= x_5 - x_6 - 0.5\end{aligned}$$

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

$$\begin{aligned}x_5 &= x_3 \\0.1 &\leq x_5 \leq 1.3\end{aligned}$$

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

$$\begin{aligned}x_6 &\leq x_4 + 0.7 * (1 - a_6) \\x_6 &\geq x_4 \\x_6 &\leq 0.5 * a_6 \\x_6 &\geq 0\end{aligned}$$

can prove  
 $x_3$  is positive using  
Box bounds

Input bounds

$$\begin{aligned}0 &\leq x_1 \leq 0.6 \\0.1 &\leq x_2 \leq 0.7\end{aligned}$$

Pre-computed Box bounds

$$\begin{aligned}0.1 &\leq x_3 \leq 1.3 \\-0.7 &\leq x_4 \leq 0.5 \\0.1 &\leq x_5 \leq 1.3 \\0 &\leq x_6 \leq 0.5\end{aligned}$$

Binary integer variable

$$a_6 \in \{0, 1\}$$

# Case $a_6 = 0$

$$\min o_0 - o_1$$

Affine

$$\begin{aligned}x_3 &= x_1 + x_2 \\x_4 &= x_1 - x_2 \\o_0 &= x_5 + x_6 + 0.5 \\o_1 &= x_5 - x_6 - 0.5\end{aligned}$$

Input bounds

$$\begin{aligned}0 &\leq x_1 \leq 0.6 \\0.1 &\leq x_2 \leq 0.7\end{aligned}$$

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

$$\begin{aligned}x_5 &= x_3 \\0.1 &\leq x_5 \leq 1.3\end{aligned}$$

Pre-computed Box bounds

$$\begin{aligned}0.1 &\leq x_3 \leq 1.3 \\-0.7 &\leq x_4 \leq 0.5 \\0.1 &\leq x_5 \leq 1.3 \\0 &\leq x_6 \leq 0.5\end{aligned}$$

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

$$\begin{aligned}x_4 &\leq x_6 \leq x_4 + 0.7 \\x_6 &= 0\end{aligned}$$

simplifying

$$\begin{aligned}o_0 - o_1 &= 2 * x_6 + 1 \\o_0 - o_1 &= 1\end{aligned}$$

Here, we proved, subject to the constraints, that the minimum is always 1, hence property holds.

Note: cannot just use the bounds for  $x_5$  and  $x_6$  for computing  $o_0 - o_1$ : this will fail! Good to try.



# Case $a_6 = 1$

$$\min o_0 - o_1$$

Affine

$$\begin{aligned}x_3 &= x_1 + x_2 \\x_4 &= x_1 - x_2 \\o_0 &= x_5 + x_6 + 0.5 \\o_1 &= x_5 - x_6 - 0.5\end{aligned}$$

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

$$\begin{aligned}x_5 &= x_3 \\0.1 &\leq x_5 \leq 1.3\end{aligned}$$

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

$$\begin{aligned}x_6 &= x_4 \\0 &\leq x_6 \leq 0.5\end{aligned}$$

Input bounds

$$\begin{aligned}0 &\leq x_1 \leq 0.6 \\0.1 &\leq x_2 \leq 0.7\end{aligned}$$

Pre-computed Box bounds

$$\begin{aligned}0.1 &\leq x_3 \leq 1.3 \\-0.7 &\leq x_4 \leq 0.5 \\0.1 &\leq x_5 \leq 1.3 \\0 &\leq x_6 \leq 0.5\end{aligned}$$

simplifying

$$\begin{aligned}o_0 - o_1 &= 2 * x_6 + 1 \\1 &\leq o_0 - o_1 \leq 2 \\ \min(o_0 - o_1) &= 1\end{aligned}$$

Here, we proved the property again

# Summary of Lecture

- We showed how to create a MILP instance for solving the problem of (complete) neural network certification.
- This is useful as MILP alone is **prohibitively expensive**. This was an example of using an incomplete method to speed up a **complete** method.
- We discussed one combination of combining complete and incomplete methods. There are other ways for both methods to benefit from each other, we mention some of these in later lectures.