

# DeepPoly - An abstract domain for certifying NN

Tuesday, 3 November 2020 10:35

- State-of-the-art methods for NN verification based on:

1. SMT solving: Reluplex [Katz'17]

precise, but can only handle very small networks.

2. Abstract interpretation: AI2 [Gehr'18]

relies on existing generic abstract domains

which either ① like convex polyhedra do not scale to larger NNs.

or ② like Zonotope are too imprecise.

3. Linear approximations Towards Fast computation of certified robustness for ReLU networks [Weng'18]

scales better but can only handle feedforward networks, can not handle convolutional networks

both ① and ③ unsound for floating point arithmetics, and thus they suffer from false negatives.

- DeepPoly: is a novel abstract interpreter specifically tailored to the setting of neural networks; uses an abstract domain which is a combination of ① floating-point polyhedra ② with intervals, coupled with abstract transformers for common NN functions.

- Adversarial region; is the set of possible images induced by a modification of the original input to a perturbed version inside  $L_\infty$  ball.

DeepPoly automatically prove that the output for all the inputs in the adversarial region is correct by soundly propagating the entire input adversarial region through the abstract transformers of the network.

- Abstract transformers for the output of a neuron are point-wise (can be computed in parallel) abstract interpretation on GPU Differentiable abstract interpretation [Mirman'18]

- The context:

- verifying robustness for adversarial regions that can be represented using a set of interval constraints.

$X = \bigcup_{i=1}^m [l_i, u_i]$  for  $l_i, u_i \in \mathbb{R} \cup \{-\infty, +\infty\}$  and  $m$  is the number of the inputs.

- NN is represented as a sequence of assignments ① ReLU  $x_i \leftarrow \max(0, x_i)$

② sigmoid/tanh  $x_i \leftarrow g(x_i)$  for sigmoid  $g = \sigma(x) = \frac{e^x}{e^x + 1}$  for tanh  $g = \tanh(x) = \frac{e^x + e^{-x}}{e^x - e^{-x}}$

③ affine  $x_i \leftarrow v + \sum_j w_j x_j$ , can describe convolutional layers.

- Abstract domain and transformers:

- Abstract domain  $A_n$  contains elements  $a \in A_n$  consist of a set of polyhedral constraints over  $n$  variables

each constraint relates one variable  $\rightarrow$  a linear combination of variables of a smaller index.

each variable has two concrete constraints lower bound  $\underline{l}$  and upper bound  $\underline{u}$

$\underline{a}_i^{\leq}, \underline{a}_i^{\geq}$  are the lower and upper relational polyhedra constraints for  $x_i$ .

- An abstract element has additionally derived interval bounds  $\underline{a}^{\leq}, \underline{a}^{\geq} \rightarrow a \in A_n = \langle \underline{a}^{\leq}, \underline{a}^{\geq}, \underline{l}, \underline{u} \rangle$

where  $\underline{a}_i^{\leq}, \underline{a}_i^{\geq} \in \{x \mapsto v + \sum_{j=1}^{i-1} w_j x_j \mid v \in \mathbb{R}, w \in \mathbb{R}^{i-1}\}$  for  $i = \{1, \dots, n\}$ ,  $\underline{l}, \underline{u} \in \mathbb{R}^n$

- The concretization function  $\Upsilon_n: A_n \rightarrow \mathcal{P}(\mathbb{R}^n)$  is given by:

$$\Upsilon_n(a) = \{x \in \mathbb{R}^n \mid \forall i = \{1, \dots, n\}: (\underline{a}_i^{\leq}(x) \leq x_i \leq \underline{a}_i^{\geq}(x))\}$$

- Domain invariant: every abstract element in the abstract domains maintains concrete lower and upper bounds,

i.e.  $\Upsilon_n(a) \subseteq X_i[l_i, u_i]$  which over approximate the symbolic bounds.

- Abstract transformers  $T_f^{\#}$  for a deterministic function  $f: A^m \rightarrow A^n$  satisfy the soundness property

$$T_f(\Upsilon_n(a)) \subseteq \Upsilon_n(T_f^{\#}(a)) \text{ for all } a \in A^m \text{ where } T_f \text{ is a concrete transformer } T_f(x) = \{f(x) \mid x \in X\}$$

- all variables are ① bounded, ② assigned exactly once (in increasing order of their indices).

$$T_f^{\#}(\langle \underline{a}^{\leq}, \underline{a}^{\geq}, \underline{l}, \underline{u} \rangle) = \langle \underline{a}^{\leq}, \underline{a}^{\geq}, \underline{l}, \underline{u} \rangle$$

for  $k < i$  outputs  $\langle \underline{a}_k^{\leq}, \underline{a}_k^{\geq}, \underline{l}, \underline{u} \rangle$

1. ReLU Abstract transformer:  $f: \mathbb{R}^{i-1} \rightarrow \mathbb{R}^i$  executes  $x_i \leftarrow \max(0, x_i)$  for some  $j < i$

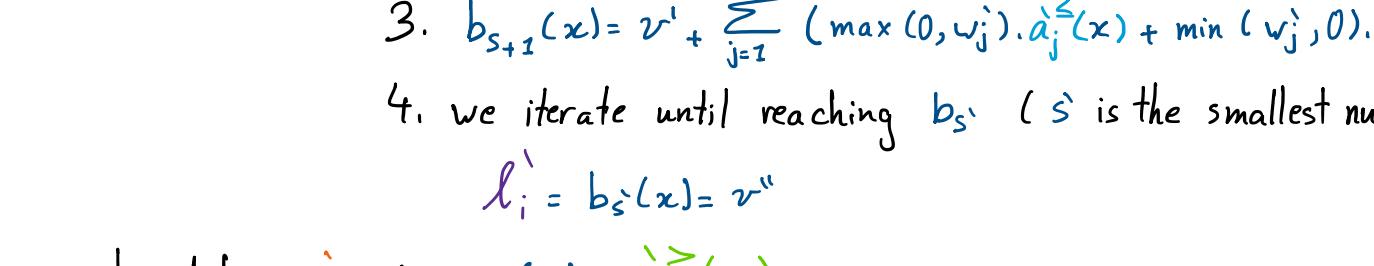


Fig. 4. Convex approximations for the ReLU function: (a) shows the convex approximation with the minimum area in the input-output plane. (b) and (c) show the two convex approximations proposed in this paper. In the figure,

for new components  $i \rightarrow$  if  $u_j \leq 0$  then  $\langle 0, 0, 0, 0 \rangle$  exact

if  $0 \leq l_j$  then  $\langle x_j, x_j, l_j, u_j \rangle \quad x_j \leq x_i \leq u_j$

otherwise  $\langle u_j \frac{x_i - l_j}{u_j - l_j}, x_j, l_j, u_j \rangle \quad \lambda \in \{0, 1\}$  select that minimize the area.

0 or  $x_j$

2. Sigmoid and tanh abstract transformer:  $\lambda = \frac{g(u_j) - g(l_j)}{u_j - l_j}, \lambda = \min(g'(l_j) - g'(u_j)) ; g'(x) > 0$

for new component  $i \rightarrow$  if  $u_j \leq 0$  then  $\langle g(u_j) + \lambda(x_j - u_j), g(u_j) + \lambda(x_j - u_j), g(l_j), g(u_j) \rangle$

if  $0 < l_j$  then  $\langle g(l_j) + \lambda(x_j - l_j), g(l_j) + \lambda(x_j - l_j), g(l_j), g(u_j) \rangle$

3. Max Pool abstract transformer:  $x_i \leftarrow \max_{j \in J} x_j$  for some  $J \subseteq \{1, \dots, i-1\}$

for new components  $\exists k \in J: u_j < l_k \text{ for all } j \in J \setminus \{k\}$  then  $\langle x_k, x_k, l_k, u_k \rangle$

otherwise we choose  $k \in J$  that maximize  $l_k$  then  $\langle x_k, \max_{j \in J} u_j, l_k, \max_{j \in J} u_j \rangle$

4. Affine abstract transformer:  $x_i \leftarrow v + \sum_{j=1}^{i-1} w_j x_j$

$$\underline{a}_i^{\leq}(x) = \underline{a}_i^{\geq}(x) = v + \sum_{j=1}^{i-1} w_j x_j$$

to obtain  $\underline{l}_i^{\#}$ : 1.  $b_1(x) = \underline{a}_i^{\leq}(x)$

2.  $b_s(x) = v + \sum_{j=1}^k w_j x_j$  for some  $k \in \{1, \dots, i-1\}$

3.  $b_{s+1}(x) = v + \sum_{j=1}^k (\max(0, w_j) \cdot \underline{a}_i^{\leq}(x) + \min(0, w_j) \cdot \underline{a}_i^{\geq}(x))$

4. we iterate until reaching  $b_s$  ( $s$  is the smallest number with this property)

$$\underline{l}_i^{\#} = b_s(x) = v$$

to obtain  $\underline{u}_i^{\#}$ : 1.  $c_1(x) = \underline{a}_i^{\geq}(x)$

2.  $c_t(x) = v + \sum_{j=1}^k w_j x_j$

3.  $c_{t+1}(x) = v + \sum_{j=1}^k (\max(0, w_j) \cdot \underline{a}_i^{\leq}(x) + \min(0, w_j) \cdot \underline{a}_i^{\geq}(x))$

4. iterate until reaching  $c_t \rightarrow \underline{u}_i^{\#} = c_t(x) = v$

- NN robustness analysis procedure:

1. create an abstract element over the input variables  $a = \langle \underline{l}_i, \underline{u}_i, \underline{l}_i, \underline{u}_i \rangle$

where all  $\underline{l}_i$  &  $\underline{u}_i$  are initialized to describe the adversarial region.

2. processing assignments for all hidden and output variable

3. processing nodes in ascending order of variable indices, using their abstract transformers.

4. execute the following  $(r-1)$  assignment ( $p$  inputs,  $q$  hidden,  $r$  outputs,  $k$  class index)

$$x_{p+q+r+1} \leftarrow x_{p+q+k} - x_{p+q+1}, \dots, x_{p+q+r+(k-1)} \leftarrow x_{p+q+k} - x_{p+q+(k-1)} \quad 1 \dots r-1$$

$$x_{p+q+r+k} \leftarrow x_{p+q+k} - x_{p+q+(k+1)}, \dots, x_{p+q+r+(r-1)} \leftarrow x_{p+q+k} - x_{p+q+r} \quad 1 \dots r$$

robust if  $\underline{l}_i > 0$  for all output variables.

- Correctness of abstract transformers: for  $a^{\#} = T_f^{\#}(a)$

we have to prove they are ① sound  $T_f(\Upsilon_{i-1}(a)) \subseteq \Upsilon_i(a^{\#})$

and ② they preserve invariant  $\Upsilon_i(a^{\#}) \subseteq \bigcap_{k \in \{1, \dots, i\}} [\underline{l}_k, \underline{u}_k]$

- Soundness under floating point arithmetics:

Relational abstract domains for the detection of floating-point run-time errors [Mine'04]

$c^-, c^+ \in \mathbb{F}$  for a real constant  $c$  to denote the floating-point representation of  $c$

with rounding towards  $-\infty$  and  $+\infty$  respectively.

→ standard interval linear form; coefficients in the constraints are intervals  $[c^-, c^+]$

upper bound and lower bound  $\inf(a^{\leq}(x)), \sup(a^{\geq}(x))$