

A CEGAR Based Verification of Neural Network

Mohammad Afzal^{1,2}, Ashutosh Gupta¹, and Akshay S¹

¹ Indian Institute of Technology, Bombay, India

² TCS Research, Pune, India

Abstract In this paper, we present a counter example guided refinement(CEGAR) based verification of neural networks.

1 Introduction

2 A Motivating Example

3 Preliminaries

4 Algorithm

A neural network is a collection of layers $l_0, l_1, l_2, \dots, l_k$, where k represents the number of layers. Layer l_0 and l_k represent the input and output layers respectively and all the other layers are hidden layers. Each layer l_i is a collection of nodes. A node $v_{i,j}$ represents the j^{th} node in the layer i . Let us say a vector $\bar{V}_i = [v_{i,0}, v_{i,1}, \dots, v_{i,m}]$ represents the values of each node in the layer i , where m is the number of nodes in the same layer. Each layer's values are computed using the weighted sum of the previous layer's values ($W_i * V_{i-1} + B_i$) followed by an activation function *ReLU*. A function $y = \max(0, x)$ is a *ReLU* function which takes an arguments x as input and return the same value x as output if x is non-negative otherwise return the value 0.

A neural network is a function N which takes an input of m dimensions and gives an output of n dimensions. N can be represented as a composition of functions $f_l o f_{l-1} \dots o f_1$, where each function f_i represents the linear combinations followed by an activation function.

Let us say P and Q are the predicates on the input and output space respectively. The goal is to find an input \bar{x}_0 , such that the predicates $P(\bar{x}_0)$ and $Q(N(\bar{x}_0))$ holds. The predicate Q usually is the negation of the desired property. The triple $\langle N, P, Q \rangle$ is our verification query.

The algorithm 1 represents the high level flow of our approach. Which is a counter example guided abstract refinement(CEGAR) based approach. At the first line in algorithm 1, we generate all the abstract constraints by *deeppoly*. These abstract constraints are the lower and upper constraints as well as the lower and upper bounds for each neurons in the neural network. The *isVerified* function in algorithm 1 calls either algorithm 2 or algorithm 3. Both the algorithms 2 and 3 takes *markedNeurons* as input. The *markedNeurons* is a subset of the neurons of the neural network. The *markedNeurons* represents the set of the culprit neurons. Algorithm 2 replace the abstract constraints of *markedNeurons*

Algorithm 1 A CEGAR based approach of neural network verification

Input: A verification problem $\langle N, P, Q \rangle$

Output: UNSAT or SAT

```
1: Build abstract constraints and bounds on each neuron using deeppoly.
2:  $A$  is a set of all abstract constraints and bounds.
3:  $markedNeurons = \{\}$ 
4: while True do
5:    $isVerified$  or  $spuriousCEX = isVerified(\langle N, P, Q \rangle, A, markedNeurons)$ 
6:   if  $isVerified$  is True then
7:     return UNSAT
8:   else
9:     if  $spuriousCEX$  is a real cex of  $\langle N, P, Q \rangle$  then
10:      return  $spuriousCEX$ 
11:     else
12:        $marked, cex = getMarkedNeurons(\langle N, P, Q \rangle, A, markedNeurons)$ 
13:       if  $cex$  not None then
14:         return  $cex$ 
15:        $markedNeurons = markedNeurons \cup marked$ 
```

Algorithm 2 An approach to verify $\langle N, P, Q \rangle$ with abstraction A

Name: $isVerified1$

Input: $\langle N, P, Q \rangle$ with abstract constraints A and $markedNeurons \subseteq N.neurons$

Output: verified or spurious counter example.

```
1: for  $nt$  in  $markedNeurons$  do
2:    $A = (A \cup exactConstr(nt)) \setminus abstractConstr(nt)$ 
3:  $A = A \cup P \cup \neg Q$ 
4:  $isSat = checkSat(A)$ 
5: if  $isSat$  is True then
6:   return spurious counter example
7: else
8:   return verified
```

by exact constraints and check for the property by MILP solver. If MILP solver return SAT then we return satisfying assignments as a spurious counter example, otherwise return verified. Algorithm 3 split each neuron of $markedNeurons$ into two sub cases. Suppose the neuron $x \in [lb, ub]$ belongs to $markedNeurons$, then first case is when $x \in [lb, 0]$ and the second case is when $x \in [0, ub]$. After splitting neurons into two cases *deeppoly* run for both cases separately. In algorithm 3, *deeppoly* run exponential number of times in the size of the $markedNeurons$. We return verified in algorithm 3 if verification query verified in all the *deeppoly* runs. If *deeppoly* fails to verify in any case then we return the spurious counter example.

5 Experiments

6 Conclusion and Future work

References

Algorithm 3 An approach to verify $\langle N, P, Q \rangle$ with abstraction A

Name: isVerified2

Input: $\langle N, P, Q \rangle$ with abstract constraints A and markedNeurons

Output: verified or spurious counter example.

```

1: for all combination in  $2^{\text{markedNeurons}}$  do
2:   run deeppoly
3:   if not verified by deeppoly then
4:      $A = \text{Set of all abstract constraints generated by deeppoly}$ 
5:      $A = A \cup P \cup \neg Q$ 
6:     isSat = checkSat( $A$ )
7:     if isSat then
8:       return spurious counter example
9: return verified

```

Algorithm 4 A pullback approach to get mark neurons or counter example

Name: getMarkedNeurons1

Input: $\langle N, P, Q \rangle$ with abstract constraints A and satisfying assignments \vec{x} and \vec{y} on input and output layers respectively

Output: marked neurons or real counter example.

```

1: Run neural network on  $\vec{x}$ .
2: return  $\vec{x}$  as counter example if it's output violate the  $Q$ .
3: for currLayer in [outputLayer ... inputLayer] do
4:   if currLayer is affine layer then
5:     layerConstraints = []
6:     for nt in currLayer.neurons do
7:       nt.constr = (nt.affineExpr == nt.satval)
8:       layerConstraints.append(nt.constr)
9:     for nt in currLayer.prevLayer.neurons do
10:      layerConstraints.append(nt.bounds)
11:     isSat = checkSat(layerConstraints)
12:     if isSat then
13:       assign sat values to previous layers neurons
14:     else
15:       markedNeurons = {nt | nt ∈ currLayer.neurons ∧ nt.constr ∈ unsatCore}
16:       return markedNeurons
17:   else ▷ Relu layer
18:     for nt in currLayer.neurons do
19:       if nt.satval > 0 then
20:         prevNt.satval = nt.satval ▷ prevNt is input node of nt in prevLayer
21:       else
22:         prevNt.lb ≤ prevNt.satval ≤ 0 ▷ lb,ub are bounds
23: ce = {nt.satval | nt ∈ inputLayer.neurons} ▷ If pullbacked upto input layer
24: return ce

```

Algorithm 5 An optimization based approach to get mark neurons or counter example

Name: getMarkedNeurons2

Input: $\langle N, P, Q \rangle$ with abstract constraints A and satisfying assignments \vec{x} and \vec{y} on input and output layers respectively

Output: marked neurons or real counter example.

```

1: Run neural network on  $\vec{x}$ .
2: return  $\vec{x}$  as counter example if it's output violate the  $Q$ .
3: Let us say layer.nt.val is the value evaluated on  $\vec{x}$  for each layer and each neuron
   nt.
4: for currLayer in [firstLayer ... outputLayer] do                                 $\triangleright$  inputLayer excluded
5:   if currLayer is affine layer then
6:     execute currLayer.prevLayer.vals on currLayer  $\triangleright$  simple matrix muliplication
7:   else
8:     constraints = A[currLayer ... outputLayer]     $\triangleright$  Abstract constraints from
       currLayer to outputLayers
9:     constraints.add(currLayer.prevLayer.vars == currLayer.prevLayer.vals)
10:    constrains.add(outputLayer.vars == outputLayer.vals)
11:    softConstraints = for each nt  $\in$  currLayer.neurons (nt.vars == nt.vals)
12:    maximize the number of neurons of layer which satisfy the softConstraints.
13:    if all neurons of the currLayer satisfy the softConstraints then
14:      continue
15:    else
16:      markedNeurons = all neurons of currLayer which does not satisfy the
       softConstraints.
17:      return markedNeurons

```
