

Honour School of Engineering Science

B1 Engineering Computation: Project B (2021-22)

Neural Network Verification

Analysing Robustness of Multi-Layer Perceptron

Student: Tijana Petrovic

Supervisor: Prof. Pawan Mudigonda

January 2022



Department of Engineering Science

Oxford University

Table of Contents

1	Introduction.....	2
1.1	<i>Problem statement.....</i>	2
1.2	<i>Problem statement in Canonical form</i>	3
2	Lower bound for Neural Network output.....	3
2.1	<i>Unsound method</i>	3
2.2	<i>Projected Gradient Ascent</i>	4
2.3	<i>Results</i>	4
2.4	<i>Further considerations</i>	5
3	Upper bound for Neural Network output.....	5
3.1	<i>Interval Bound Propagation.....</i>	5
3.2	<i>Results</i>	6
4	Linear programming method	6
4.1	<i>Planet Relaxation of ReLU function</i>	6
4.2	<i>Results</i>	7
5	Branch and Bound for Verification	7
5.1	<i>Branch Strategy.....</i>	9
5.1.1	<i>Splitting the longest edge.....</i>	9
5.2	<i>Results</i>	9
5.2.1	<i>Interval Bound Propagation and Unsound Method.....</i>	9
5.2.2	<i>Linear Programming Bounds and Unsound Method.....</i>	9
5.2.3	<i>Hybrid</i>	10
6	Conclusions	10
7	References	10

1 Introduction

A neural network is a network or circuit of (artificial) neurons or nodes represented by their weights, biases, and an input domain χ . They can be trained via a dataset to model complex relationships between inputs and outputs or find patterns in data. Thanks to this self-learning property, they are widely used in various applications of artificial intelligence such as facial recognition, signature verification, stock market, and protein folding prediction. Despite this, their use is very limited in safety-critical applications (such as autonomous driving cars) due to their susceptibility to adversarial attacks where small perturbations in the input lead to drastic changes in the network output.

This project aims to overcome this deficiency by using the formal verification of neural networks to test the hypothesis whether the given neural network satisfies a desirable property. Our proposed methods are Unsound method, Projected Gradient Ascent, Interval Bound Propagation, and Linear Programming method under the Branch and Bound framework. The Collision Detection dataset consisting of 500 MAT files is used for testing.

1.1 Problem statement

A certain function f is represented by a piecewise linear multilayer perception (MLP) with n layers as $f_n \circ f_{n-1} \circ \dots \circ f_1(x)$ where f_i represents i^{th} layer. The MLP is specified by its weights W , biases b and an input domain $\chi = \{x | x_{min} \leq x \leq x_{max}\}$, where x_{min} and x_{max} are box constraints. Outputs of each layer can be obtained using $f_i(z_{i-1}) = z_i = \sigma(\hat{z}_i)$ and $\hat{z}_i = z_{i-1} + b_i$ up until the last layer. The output of the last layer $z_n = \hat{z}_n$ is denoted by a vector $y = z_n = f(z_0) = f(x)$.

Given that the neural network implements a function $z_n = f(z_0)$ over a bounded input region χ and a property of robustness to adversarial examples, P , we want to prove that if $z_0 = x \in \chi$ and $\hat{z}_n = f(z_0)$, then $P(\hat{z}_n)$ is true.

Consider ∞ -norm around a training sample a with label y_a . The desired property is encoded by using $\chi = \{z_0 | \|z_0 - a\|_\infty \leq \varepsilon\}$ and $P(\hat{z}_n) = \{\forall y, \hat{z}_{n[y_a]} \geq \hat{z}_{n[y]}\}$. A neural network characterized by a function f is said to be robust at z_0 if for $\forall (z_0 - a)$ in the input domain $\chi = \{z_0 | \|z_0 - a\|_\infty \leq \varepsilon\}$ the weight of the correct label $\hat{z}_{n[y_a]}$ is greater or equal to the weight of undesired label $\hat{z}_{n[y]}$ for the

property $P(\hat{z}_n) = \{\forall y, \hat{z}_{n[y_a]} \geq \hat{z}_{n[y]}\}$. Otherwise, a neural network is not robust at z_0 , and $(z_0 - a)$ for which this is the case is described as an adversarial.

1.2 Problem statement in Canonical form

The problem can also be stated in canonical and mixed integer programming form.

For the canonical form, we want to reduce the verification problem to a simple global optimization problem and check the sign of the output of the modified neural network $f'(z_0)$. This is possible if the property can be expressed as a Boolean expression over linear forms.

If a property satisfies $P(\hat{z}_n) = W^T \hat{x}_n \geq b$, it is enough to add a final fully connected layer with one output, weight W and bias $-b$. If the global minimum of the network is positive, it is indicated that $W^T \hat{x}_n - b \geq 0$, or $W^T \hat{x}_n \geq b$, which implies the property is true. Otherwise, the minimizer provides a counterexample.

Assuming the network only contains ReLU activations between each layer, the satisfiability problem to find a counter example could be expressed as:

$$l_0 \leq x_0 \leq u_0, \quad \hat{x}_n \leq 0, \quad \hat{x}_{i+1} = W_{i+1} x_i + b_{i+1} \quad \text{and} \quad x_i = \max(\hat{x}_i, 0) \quad \text{for } \forall i \in \{0, n-1\}.$$

Now, we can simplify the problem into computing the upper and lower bounds of the output where:

- the property is true if its upper bound is non-positive.
- the property is false if its lower bound is positive.

2 Lower bound for Neural Network output

It has been shown that the property is false if its lower bound is positive. Here, we will consider two different methods for obtaining lower bounds of the neural network: Unsound method and Projected Gradient Ascent.

2.1 Unsound method

The unsound method consists of two functions: *generate_inputs* and *compute_nn_outputs*. The former generates k random inputs within the box constraints set by x_{min} and x_{max} . For this, we use a uniformly distributed random number generator: $X(i) = x_{min} + (x_{max} - x_{min}) \cdot rand(0, 1)$. Next, the

function *compute_nn_outputs* takes the input array, X , to calculate k corresponding outputs for each input using Algorithm 1. The maximum value among k outputs is then a valid lower bound on \hat{z}_n .

2.2 Projected Gradient Ascent

The unsound method helps us to find counterexamples quickly. However, results could be improved if we considered verification as a problem of finding the optimal solution of maximising $f(W, b, x^{\min}, x^{\max})$ and perform gradient ascent iteratively.

Algorithm 2 Projected Gradient Ascent

```

function projected gradient ascent( $W, b, X, x^{\min}, x^{\max}$ )
1: for  $i = 1, \dots, \text{size}(X)$  do
2:    $X_k = \text{transpose}(X(k, :))$ 
3:   for iteration = 1, ...,  $l$  do
4:      $a = X_k$ 
5:     for  $i = 1, \dots, \text{Number of Layers}$  do
6:        $z = W(i) \times a + b(i)$ 
7:       for  $j = 1, \dots, \text{size}(z)$  do
8:         Populate gradient matrix:
9:         if  $z(j) > 0$  then  $dfdx(j, j) = 1$ 
10:        else  $dfdx(j, j) = 0$ 
11:        Compute the gradient of the current layer
12:        if  $i \leq \text{Number of layers} - 1$  then  $dadx = dfdx \times W(i)$ 
13:        else  $dadx = W(i)$ 
14:        Compute the product of gradients of the all previous layers
15:         $a = \max(z, 0)$ 
16:        Update  $X_k$  as  $X_k + \text{lr} \times \text{transpose}(\text{propagated } dadx)$ 
17:        Clip  $X$ 
18:      end for
19:       $\text{refinedX}(k) = \text{transpose}(X_k)$ 
20:    end for
21:  end for
22: end for
23: return  $\text{refinedX}$ 

```

Algorithm 1 Compute NN Outputs

```

function compute nn outputs( $W, b, X$ )
1:  $z \leftarrow \text{transpose}(X)$ 
2: for  $l = 1, \dots, L - 1$  do
3:    $\hat{z} \leftarrow W_l \cdot z + b_l$ 
4:    $z \leftarrow \max(\hat{z}, 0)$ 
5: end for
6:  $y \leftarrow W_L \cdot z + \text{repmat}(b_L, 1, \text{size}(W_L \cdot z, 2))$ 
7:  $y \leftarrow \text{transpose}(y)$ 
8: return  $y$ 

```

Random input vector x_0 will be generated and updated as follows:

$$x'_{t+1} = x_t + \eta_t \nabla f(W, b, x_t)$$

Here, η_t represents the step-size at iteration t and $\nabla f(W, b, x_t)$ is the gradient of the function f at the point x_t . We will also approximate the derivative at 0 as 0 and clip all the values x'_{t+1} that lie outside the domain χ to their nearest value in domain χ .

The pseudocode can be seen in Algorithm 2 given on the left.

2.3 Results

Results from the Unsound method and Projected Gradient Ascent are given in figures below.

As the value of k increases, the average lower bound obtained using the Unsound method also increases. This is because taking a larger number of inputs that satisfy the box constraints leads to a greater probability of finding z_0 that generates a stricter lower bound.

Comparing two figures, it is apparent that Projected Gradient Ascent gives better limits and higher number of found counterexamples than the Unsound method as predicted.

Figure 1 (Unsound method):

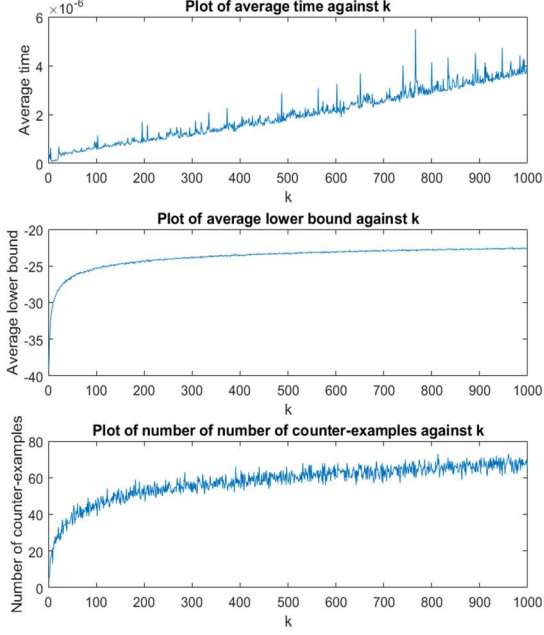
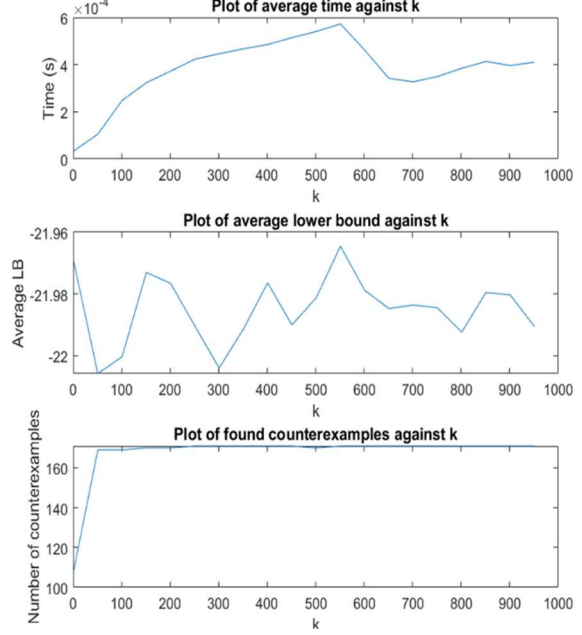


Figure 2 (Projected Gradient Ascent):



2.4 Further considerations

The Unsound method could be improved by using a grid search over the input domain instead of random sampling. This would create better coverage of the input domain and counterexamples could be found more easily.

3 Upper bound for Neural Network output

Previously introduced methods can be used efficiently for finding counterexamples of the neural network. However, they cannot be used for proving which properties are true. This is a huge limitation for the purpose of this project. In this chapter, we propose Interval Bound Propagation as a method of finding upper bound of the neural network.

3.1 Interval Bound Propagation

The Interval Bound Propagation algorithm is based on the artificial deconstruction of the weight matrix at each layer into matrices with positive and negative entries ($W_l = W_l^- + W_l^+$). Then upper and lower bounds are computed for the output of each layer using these new matrices and the bounds of the previous layer as shown in the Algorithm 3.

Algorithm 3 Interval Bound Propagation

```

function interval bound propagation( $W, b, x^{min}, x^{max}$ )
1:  $z_l^{min} \leftarrow transpose(x^{min})$ 
2:  $z_l^{max} \leftarrow transpose(x^{max})$ 
3: for  $i = 1, \dots, L - 1$  do
4:    $W_l^+ = max(0, W_l)$ 
5:    $W_l^- = min(0, W_l)$ 
6:    $z_l^{max} = W_l^+ \cdot z_l^{max} + W_l^- \cdot z_l^{min} + b_l$ 
7:    $z_l^{min} = W_l^+ \cdot z_l^{min} + W_l^- \cdot z_l^{max} + b_l$ 
8:    $z_l^{max} = max(0, \hat{z}_l^{max})$ 
9:    $z_l^{min} = max(0, \hat{z}_l^{min})$ 
10: end for
11:  $y^{min} = W_L^+ \cdot z_L^{min} + W_L^- \cdot z_L^{max} + b_L$ 
12:  $y^{max} = W_L^+ \cdot z_L^{max} + W_L^- \cdot z_L^{min} + b_L$ 
    =0
  
```

3.2 Results

Running the Interval Bound Propagation algorithm on the Collision Detection dataset took an average of 0.0016 seconds per property. It verified 11 out of 328 properties as true and has not proven any properties to be false. The average computed lower bound was -137.9, and the average upper bound was 82.7. The performance is poor because the algorithm alters the network quite drastically which leads to loose bounds that rarely approach zero for the given dataset.

It is worth noting that Interval Bound Propagation has low computational cost but lower bounds computed using both, Unsound method and Projected Gradient Ascent are much higher and therefore better.

4 Linear programming method

The disadvantage of interval bound propagation is that obtained bounds are very loose which means that we will not be able to prove a significant number of properties. Although using the Unsound method together with Interval Bound Propagation could improve the lower bounds there is still space for further improvement. Here, we introduce an algorithm that uses optimization at each layer and then solves for the minimum and maximum output of all previous layers. The computational difficulty of this problem is in the non-linearity and non-convexity of the final sets of constraints due to ReLU

Algorithm 4 Linear Programming Bound

```

function linear programming bound(W, b,  $x^{min}$ ,  $x^{max}$ )
1:  $z_0^{min} \leftarrow transpose(x^{min})$ 
2:  $z_0^{max} \leftarrow transpose(x^{max})$ 
3: Define matrices f, A and B
4: for  $i = 1, \dots, L$  do
5:   for  $j = 1, \dots, \text{Number of Neurons in the Layer}$  do
6:     Apply constraints
7:   end for
8:   Compute linear program for each Neuron in the Layer
9:    $z_i^{min} = \max(0, \hat{z}_i^{min})$ 
10:   $z_i^{max} = \max(0, \hat{z}_i^{max})$ 
11:  Update matrices A and B
12: end for
13: LB =  $\hat{z}_L^{min}$ 
14: UB =  $\hat{z}_L^{max}$ 
15: return [LB, UB]
```

activations. This converts our problem into an NP-hard problem that can be solved by approximating non-linear constraints to linear constraints.^[2]

The pseudocode of the method is given in the Algorithm 4. The built-in MATLAB function *linprog* was used for optimisations, while A and B represent equality matrices.

4.1 Planet Relaxation of ReLU function

Planet Relaxation operates by explicitly attempting to assign a number to the phase of the non-linearities. It assigns either 0 or 1 to each $\delta_{i[j]}$ variable, verifying at each step the feasibility of the partial assignment. This allows us to prune infeasible partial assignment early.

Neural Network Verification

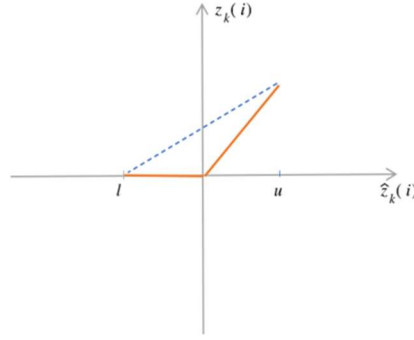
The search strategy enumerates all the domains that have not been pruned yet, while the branching rule fixes the decision variable $\delta_{i[j]}=0$ which is equivalent to choosing $x_{i[j]}=0$, $\hat{x}_{i[j]} \leq 0$ and fixing

$\delta_{i[j]}=1$ which is equivalent to

$$x_{i[j]}=0, \hat{x}_{i[j]} \geq 0.$$

Planet Relaxation does not include any heuristic to prioritise which decision

Figure 3 (Planet Relaxation):



$$z_k^{min}(i) \equiv \min z_k(i)$$

- $z_k(i) > 0$
- $z_k(i) > \hat{z}_k(i)$
- $z_k(i) < m(z_k(i) - z_k^{min}(i))$

$$\text{Where } m = \frac{\hat{z}_k^{max}}{\hat{z}_k^{max} - \hat{z}_k^{min}}$$

variable should be split over or a mechanism for early termination based on a heuristic search of a feasible point. Consequently, only when a full complete assignment identified the solution is returned. However, we can detect incoherent assignments if we approximate the non-linear constraints by a set of linear constraints representing the non-linearities' convex hull. It also uses conflict analysis to discover combinations of splits that cannot lead to satisfiable assignments, allowing them to perform further pruning of the domains.

Except for the planet relaxation, parallel relaxation and mixed-integer encoding could be used for the same purpose. ^{[1][3]}

4.2 Results

Linear Programming Method has proved 260 properties, which took on average 6.6 seconds per property. The obtained lower limit is -58.7 and upper limit is -13.2. It is apparent that it preforms significantly better than Interval Bound Propagation but at cost of higher computational time.

5 Branch and Bound for Verification

In this section, an approach of estimating the global minimum is introduced. As the name suggests, Branch and Bound (BaB) iteratively partitions the input domain into a set of subdomains that are collectively exhaustive and uses bounding methods for estimating lower and upper bounds. This is possible due to the piecewise linear property of the neural network. A partition \mathcal{P} of size n is a set $\{\chi_1, \chi_2, \dots, \chi_n\}$ such that $\chi_1 \cup \chi_2 \cup \dots \cup \chi_n = \chi$ and is initialized as $\mathcal{P}_0 = \{\chi\}$. The general structure of the BaB is shown in Algorithm 5.

Algorithm 5 Branch And Bound

```

function branch and bound( $W, b, x^{\min}, x^{\max}$ )
1:  $\bar{z}^{\min} \leftarrow \text{transpose}(x^{\min})$ 
2:  $\bar{z}^{\max} \leftarrow \text{transpose}(x^{\max})$ 
3:  $\text{bounds} \leftarrow \text{getBounds}(x^{\min}, x^{\max})$ 
4: if  $\text{LowerBound} > 0$  then
5:   flag = 0
6:   return flag
7: end if
8: for  $i = 1, \dots, L$  do
9:    $[\text{maxUB}, \text{index}] \leftarrow \text{max}(\text{UpperBounds})$ 
10:  if  $\text{maxUB} < 0$  then
11:    flag = 1
12:    return flag
13:  end if
14:   $[x_1^{\min}, x_1^{\max}, x_2^{\min}, x_2^{\max}] \leftarrow \text{partOut}(\text{bounds})$ 
15:   $\text{bounds}_1 \leftarrow \text{getBounds}(x_1^{\min}, x_1^{\max})$ 
16:   $\text{bounds}_2 \leftarrow \text{getBounds}(x_2^{\min}, x_2^{\max})$ 
17:  if  $\text{LowerBound}_1 > 0$  or  $\text{LowerBound}_2 > 0$  then
18:    flag = 0
19:    return flag
20:  end if
21:  update Bounds
22: end for
23: flag = 2
24: return flag
    
```

The algorithm has three elements:

Search strategy – chooses the next domain to branch on. We always choose the domain with the maximum upper bound because if the highest obtained upper bound is negative, all the other upper bounds are negative and consequently, the property is true.

Branching strategy – takes the previously chosen domain and returns a partition in its subdomain such that $\text{subdomain}_1 \cup \text{subdomain}_2 = \text{domain}$ and $\text{subdomain}_1 \cap \text{subdomain}_2 = \emptyset$

Bounding strategy – produces upper and lower bounds over the minimum output that the network can reach over a given input domain.

There are two possible cases that arise can from the bound strategy:

- **Case 1:** $l > 0$ $f'(x_0)$ is positive for at least one $z_0 \in \mathcal{X}$ \Rightarrow property is false
- **Case 2:** $u < 0$ $f'(x_0)$ is negative for all $z_0 \in \mathcal{X}$ \Rightarrow property is true

However, although the procedure will terminate after a finite number of iterations, this might take too long, and, therefore, third case of slow convergence arises:

- **Case 3:** $l < 0$ and $u > 0$ unable to determine the validity of property of the neural network.

Encountering the case 3 would mean that our method is incomplete and should be replaced by another one. Therefore, it is important to find methods that provide tighter bounds on z_n and efficient ways of partitioning the input domain into sub-domains in search for the optimal bounds on z_n .

5.1 Branch Strategy

5.1.1 Splitting the longest edge

Assuming the piecewise linear property, the function partitions the input domain into subdomains by

finding the dimension of the longest relative length $s(i)$: $s(i) = \frac{\bar{x}^{\max}(i) - \bar{x}^{\min}(i)}{x^{\max}(i) - x^{\min}(i)}$ and

$j = \underset{i}{\operatorname{argmax}}(s(i))$, where x^{\min} and x^{\max} are the box constraints of the input domain χ , and \bar{x}^{\min} and

\bar{x}^{\max} are the box constraints of the input sub-domain $\bar{\chi}$ which generates the highest upper bound.

We split the subdomain into two along its j^{th} dimension using following rule:

$$\begin{aligned} \bar{x}_1^{\min} &= \bar{x}^{\min}, & \bar{x}_1^{\max} &= \begin{cases} \bar{x}^{\max}(i) & \text{if } i \neq j, \\ \bar{x}^{\min}(i) + \frac{\bar{x}^{\max}(i) - \bar{x}^{\min}(i)}{2} & \text{otherwise.} \end{cases} \\ \bar{x}_2^{\max} &= \bar{x}^{\max}, & \bar{x}_2^{\min} &= \begin{cases} \bar{x}^{\min}(i) & \text{if } i \neq j, \\ \bar{x}^{\min}(i) + \frac{\bar{x}^{\max}(i) - \bar{x}^{\min}(i)}{2} & \text{otherwise.} \end{cases} \end{aligned}$$

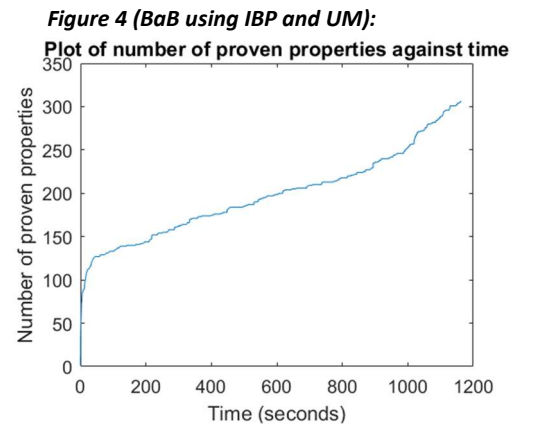
χ_1 and χ_2 computed in this way are guaranteed to generate an upper bound less than or equal to that of the χ , and a lower bound greater than or equal to that of χ .

Other branching methods that could be used are random splitting and largest edge splitting. Random splitting reduces time complexity of branching, while largest edge method splits the domain so that it has the greatest effect on the output of the network. [1]

5.2 Results

5.2.1 Interval Bound Propagation and Unsound Method

As shown in the Figure 4, it took 3.8 seconds per property to prove 306 properties using a combination of Unsound method and Interval Bound Propagation under the Branch and Bound framework. Although we did not manage to prove all properties using the chosen convergence limit, this method is very useful and has good computational cost.



5.2.2 Linear Programming Bounds and Unsound Method

Under the Branch and Bound framework, a combination of Linear Programming and Unsound method took on average 60.1 seconds per property to prove all 500 properties correctly. Plot of

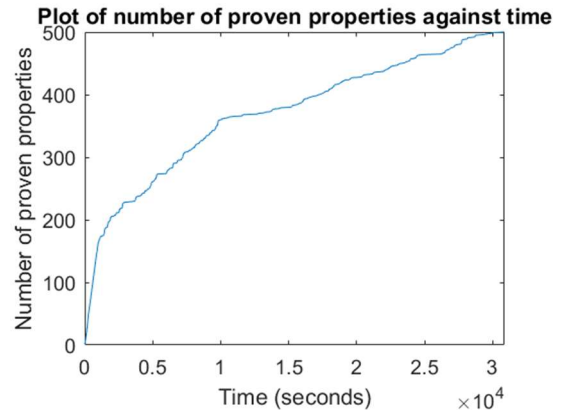
number of proven properties against time can be seen in the Figure 5 on the right. All considered, this method is considerably better than all previously considered ones.

5.2.3 Hybrid

The Hybrid method first tries to prove a property using the Unsound method and Interval Bound Propagation and then moves onto the Linear Programming if it was

not successful. It proved all properties and took about 40 seconds for each which is considerably better than all previously considered methods.

Figure 5 (BaB using LP):



6 Conclusions

To conclude, the bounding strategies alone did not prove enough properties to verify the neural network. However, combining bounding strategies with branching strategies withing the Branch and Bound framework led to significant improvement in accuracy. It is worth noting that the Collision Detection dataset is relatively small, so it is not hard to achieve a reasonable performance. With larger networks, the performance of these methods would be considerably worse, and, therefore, it is difficult to say without further research whether this is an effective general method for neural network verification in its current form.

7 References

- [1] R. Bunel, I. Turkaslan, P. Torr, P. Kohli, and M. P. Kumar. A unified view of piecewise linear neural network verification. In *Advances in Neural Information Processing Systems*, 2018.
- [2] R. Anderson, J. Huchette, W. Ma, C. Tjandraatmadja, and J. P. Vielma. Strong mixed-integer programming formulations for trained neural networks. In *International Confernece on Integer Programming and Combinatorial Optimization*, 2019.
- [3] Rudiger Ehlers: Formal Verification of Piece-Wise Linear Feed-Forward Neural Networks, 2017