Fast Neural Network Verification via Shadow Prices

Vicenç Rubies-Royo ¹ Roberto Calandra ² Dusan M. Stipanovic ³ Claire Tomlin ¹

Abstract

To use neural networks in safety-critical settings it is paramount to provide assurances on their runtime operation. Recent work on ReLU networks has sought to verify whether inputs belonging to a bounded box can ever yield some undesirable output. Input-splitting procedures, a particular type of verification mechanism, do so by recursively partitioning the input set into smaller setEquals. The efficiency of these methods is largely determined by the number of splits the box must undergo before the property can be verified. In this work, we propose a new technique based on shadow prices that fully exploits the information of the problem yielding a more efficient generation of splits than the state-of-the-art. Results on the Airborne Collision Avoidance System (ACAS) benchmark verification tasks show a considerable reduction in the partitions generated which substantially reduces computation times. These results open the door to improved verification methods for a wide variety of machine learning applications including vision and control.

1. Introduction

With the increased deployment of deep neural networks (DNN) in many domains, it is becoming more and more pressing to validate said models. While deep learning has shown great promise in applications such as vision (LeCun et al., 2015), reinforcement learning (Silver et al., 2017; Mnih et al., 2013) or speech recognition (Graves et al., 2013; Kim, 2014), there are still many safety-critical applications, such as self-driving (Bojarski et al., 2016), that will not benefit from these advances until models can be effectively validated.

Under review by the International Conference on Machine Learning (ICML).

A limitation of current verification approaches is that they can be very slow. It has been shown that the verification problem for feedforward ReLU networks is NP-Complete (Katz et al., 2017). Therefore, it is paramount to find verification methodologies that can improve upon previous results.

In this paper, we propose an approach to reduce the computational cost of verifying deep ReLU networks without any loss of generality or approximation. The main contribution is the use of a more efficient input-splitting algorithm — using so-called shadow prices — which allows to more intelligently decide how to generate the splits of the input box. As a result, this algorithm reduces the number of splits used for verification tasks, and thus the memory footprint and computational cost of the overall procedure. Experimental results on the Airborne Collision Avoidance System (ACAS) standard benchmark demonstrate that our approach significantly reduces the number of splits generated and the time needed to verify properties.

As the number of machine learning applications grows, verification will become more and more important to guarantee safe behaviors from the learned models. Our approach is a small step towards tackling the real-world challenges of efficiently and reliably verifying deep neural networks.

2. Verification of Neural Networks

We now formalize the verification problem and discuss related work in this area. We then briefly introduce ReLU networks and a few of their properties.

2.1. Verification Problem

Given a set of possible inputs and a neural network, the problem of verifying whether some input will result in an undesirable output is mathematically analogous to checking whether a set mapped through a function intersects some other set. This input set could, for instance, represent a range of valid operational configurations for the system, whereas the output set could represent dangerous/undesirable conditions. In reinforcement learning and control, for example, the DNN could represent a controller, the input set would be some set of valid states of the system and the output set could correspond to control actions which are known *a*

¹Electrical Engineering and Computer Sciences Dept., University of California at Berkeley, Berkeley, California, USA ²Facebook AI Research, Menlo Park, California, USA ³Industrial and Enterprise Systems Engineering Dept., University of Illinois Urbana-Champaign, Champaign, Illinois, USA. Correspondence to: Vicenç Rubies-Royo <vrubies@eecs.berkeley.edu>.

priori to be unsafe. Hence, given a DNN f(x), input set \mathcal{B} and an output set \mathcal{S} , the verification problem seeks to answer whether the proposition

$$\forall x \in \mathcal{B}, f(x) \notin \mathcal{S} \tag{1}$$

is true or false.

2.2. Prior Work

Starting with the works from (Huang et al., 2016; Katz et al., 2017), the authors use Satisfiability Modulo Theory (SMT) solvers to answer Equation (1) for ReLU networks and, more generally, networks containing piece-wise linear activations. In their approaches, an answer is reached by leveraging the finite set of possible activations induced by the network's non-linearities. A similar reasoning is found in (Lomuscio & Maganti, 2017; Tjeng & Tedrake, 2017) using Mixed Integer Linear Programming (MILP) solvers. Other approaches inspired by reachability include the work from (Xiang et al., 2017b), where the structure of the domain induced by the ReLU non-linearities is exploited to compute the exact set of possible outputs. In contrast, in (Xiang et al., 2017a), the set of possible outputs is approximated by gridding the input set.

Of particular interest for this paper are the works of (Ehlers, 2017) and (Kolter & Wong, 2017). Ehlers (2017) introduced a novel convex relaxation of the ReLU non-linearity in order to render the problem easier for the SMT solver. Kolter & Wong (2017) used this relaxation and duality to compute rapid over-approximations of the set of possible outputs for the network. In (Raghunathan et al., 2018) they expand on this duality approach. A new interesting direction which does not rely on SMT or MILP solvers was presented by Wang et al. (2018). In this work, a divide and conquer approach was used to repeatedly partition the input set into smaller sub-domains and check the property individually for each partition.

In our work, we build upon the input-splitting technique from (Wang et al., 2018). We show experimentally that splits based on input-output gradient metrics are in general inefficient. We provide a new methodology that substantially reduces the number of splits and the runtime required to verify a network for a given input set and property.

2.3. ReLU Network: A Piece-wise Affine Function

Let $f_{\theta}(x)$ be a ReLU network defined by

$$\hat{z}_{i+1} = W_i z_i + b_i, \text{ for } i = 1, ..., K - 1$$

$$z_j = \max\{\hat{z}_j, 0\}, \text{ for } j = 2, ..., K - 1$$
(2)

with $W_i \in \mathbb{R}^{n_{i+1} \times n_i}$, $b \in \mathbb{R}^{n_{i+1}}$, $z_1 = x \in \mathcal{B}$, where \mathcal{B} is a bounded set in \mathbb{R}^{n_1} which we assume to be a box, and

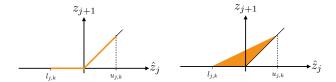


Figure 1. (Left) Bounded output of a ReLU node. (Right) Convex envelope.

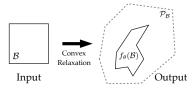


Figure 2. Example of convex over-approximation of the output set.

 $f_{\theta}(x) = \hat{z}_K$. The set $\theta = \{W_i, b_i\}_{i=1,\dots,K-1}$ represents the set of parameters of the network.

From the definition of Equation (2), it is clear that ReLU networks are piece-wise affine functions.

This follows from the fact that any composition of an affine function (i.e. linear function plus a bias term) and a piecewise affine function results in a piece-wise affine function. An important aspect that will come into play in later sections, is that the domain is partitioned into polytopic regions \mathcal{P}_i such that $\bigcup_i \mathcal{P}_i = \mathbb{R}^{n_1}$. For a closer look at the properties of ReLU networks, including the partition of the domain into polytopic regions, we direct the reader to (Montufar et al., 2014; Serra et al., 2017; Raghu et al., 2017).

3. Overview of Verification via Input-Splits

We now outline how the verification task is accomplished by iteratively splitting sections of the input set. First, we introduce the concept of convex relaxations of ReLU networks to over-approximate the image of the input set as introduced in (Ehlers, 2017) and in (Kolter & Wong, 2017). Next, we show how these relaxations can be used to verify properties of the image, along the lines of the work by Wang et al. (2018).

3.1. Convex Over-approximation of the Image

To obtain an over-approximation of the image $f_{\theta}(\mathcal{B})$, it suffices to replace the ReLU non-linearity $\max\{z,0\}$ in each layer by its convex envelope,

$$\begin{split} \hat{z}_{i+1} &= W_i z_i + b_i \,, \text{ for } i = 1, ..., K-1 \\ z_j &\geq 0 \,, \quad \text{for } j = 2, ..., K-1 \\ z_j &\geq \hat{z}_j \,, \\ z_j &\leq D_j \hat{z}_j + d_j \,. \end{split} \tag{3}$$

The matrix D_j is diagonal, so that $D_j = diag(c_j)$, where c_j is a vector of the form $c_{j,k} = \frac{u_{j,k}}{u_{j,k}-l_{j,k}}$ and $d_{jk} = -\frac{u_{j,k}l_{j,k}}{u_{j,k}-l_{j,k}}$. The terms $l_{j,k}$ and $u_{j,k}$, which we will shortly define, denote the lower and upper bounds for the k-th activation in layer j. For (3), $\hat{z}_K \in \mathcal{P}_{\mathcal{B}}$, where $\mathcal{P}_{\mathcal{B}}$ is a bounded convex polytope satisfying $f_{\theta}(\mathcal{B}) \subseteq \mathcal{P}_{\mathcal{B}}$. Figure 1 shows a typical convex envelope. Figure 2 shows how the convex relaxation results in a polytopic over-approximation of the image.

Computing the lower and upper bounds $l_{j,k}$ and $u_{j,k}$, can be accomplished in a layer-by-layer fashion by solving linear programs (LPs) from j = 1 to K - 2 of the form

$$\begin{split} l_{j,k} &= \min_{\underline{z}} \ \left(w_{j,k} \underline{z} + b_{j,k} \right) \quad \text{s.t. } \tilde{A}_j \underline{z} \preceq \tilde{b}_j \ , \\ u_{j,k} &= \max_{\underline{\bar{z}}} \ \left(w_{j,k} \bar{z} + b_{j,k} \right) \quad \text{s.t. } \tilde{A}_j \bar{z} \preceq \tilde{b}_j \ , \end{split} \tag{4}$$

where \tilde{A}_j and \tilde{b}_j represent a polytope in a d-dimensional space, where $d = \sum_{l=1}^j n_l$. This polytope grows in dimension as a result of taking into account upper and lower bounds of earlier layers. We also implicitly assume for the remainder of the paper that positive lower bounds or negative upper bounds are automatically set to 0. The terms $w_{j,k}$ and $b_{j,k}$ denote the k-th row/entry of W_j and b_j . In Section 4, we study the structure of the constraint set in Equation (4) in depth.

3.2. Image Verification through Refinements

One of the useful features of this bounded convex overapproximation is that it provides sufficient conditions to check whether the property of interest can ever be satisfied (i.e., from Figure 2, if the property does not hold for the overapproximation, it can't hold for the image either). However, when properties do hold for the over-approximation nothing can be established with regards to the image.

3.2.1. RECURSIVELY SPLITTING SETS

In Section 3.1, we have explained how to build a convex over-approximation of the image in a layer-by-layer basis. Note, that for any split of the input set into two subsets \mathcal{B}_1 and \mathcal{B}_2 such that $\mathcal{B}_1 \cup \mathcal{B}_2 = \mathcal{B}$, it holds that

$$\mathcal{P}_{\mathcal{B}_1} \cup \mathcal{P}_{\mathcal{B}_2} \subseteq \mathcal{P}_{\mathcal{B}}. \tag{5}$$

This follows from the fact that splitting the input set reduces all the feasible regions for the LPs in Equation (4), which results in greater lower bounds or smaller upper bounds in the computation of $\mathcal{P}_{\mathcal{B}_1}$ and $\mathcal{P}_{\mathcal{B}_2}$.

Splitting the input set is particularly useful for verification since it breaks the problem into two sub-problems. Specifically, and without loss of generality, one of three things may happen:

Algorithm 1 Recursive Splitting (Depth First Search)

```
1: procedure VERIFICATION(S, B, \theta)
 2:
              \mathcal{P}_{\mathcal{B}} \leftarrow \text{ConvexOverApprox}(\mathcal{B}, \theta)
              if \mathcal{P}_{\mathcal{B}} \cap \mathcal{S} = \emptyset then
 3:
 4:
                     return False
 5:
              else
                     if IsExact(\mathcal{P}_{\mathcal{B}}) then
 6:
 7:
                             return True
                     \mathcal{B}_1, \mathcal{B}_2 \leftarrow \text{Split}(\mathcal{B})
 8:
                     return Verification(\mathcal{S}, \mathcal{B}_1, \theta)\vee
 9:
10:
                                    Verification (S, \mathcal{B}_2, \theta)
```

- 1. The property does not hold for either $\mathcal{P}_{\mathcal{B}_1}$ nor $\mathcal{P}_{\mathcal{B}_2}$, and thus the property is not satisfied by the image.
- 2. The property of interest holds for $\mathcal{P}_{\mathcal{B}_1}$ but does not hold for $\mathcal{P}_{\mathcal{B}_2}$, in which case \mathcal{B}_2 can be discarded from the verification problem.
- 3. Finally, both $\mathcal{P}_{\mathcal{B}_1}$ and $\mathcal{P}_{\mathcal{B}_2}$ satisfy the property and nothing can be established.

Given these outcomes, a natural algorithm arises for the verification of the property of interest; starting with \mathcal{B} , we compute the convex over-approximation $\mathcal{P}_{\mathcal{B}}$. If the property does not hold for $\mathcal{P}_{\mathcal{B}}$, the property does not hold for the image and we are done. Otherwise, we can split \mathcal{B} in two halves \mathcal{B}_1 and \mathcal{B}_2 , and compute $\mathcal{P}_{\mathcal{B}_1}$ and $\mathcal{P}_{\mathcal{B}_2}$. Given the list of possible outcomes, the algorithm will either: end with the property being false, be able to discard one of the halves, or, in the worst case, keep both \mathcal{B}_1 and \mathcal{B}_2 for further analysis. This procedure induces a growing *binary tree* whose nodes represent smaller and smaller regions of the input set \mathcal{B} .

Algorithm 1 provides the aforementioned procedure for verifying whether the image of $\mathcal B$ intersects with a given set $\mathcal S$, which we assume to be a union of a finite number convex sets. This assumption is required so that in Line 3 the intersection check can be done efficiently via convex optimization. When the intersection is non-empty, in Line 6, we check whether the over-approximation is tight/exact or not. If it is, it must be true that the input set satisfies the property. If it is not, in line 8 we split the set into two halves and attempt to solve the verification problem for each half separately. Sections 3.2.2 and 3.2.3 describe the auxiliary functions "IsExact" and "Split" in lines 6 and 8.

3.2.2. OVER-APPROXIMATION AND IMAGE EQUIVALENCE

In the previous section we mentioned that some instances of $\mathcal{P}_{\mathcal{B}}$ can be checked for tightness or exactness, that is $\mathcal{P}_{\mathcal{B}} \equiv f_{\theta}(\mathcal{B})$. From Section 2.3, we know that a ReLU network sub-divides the domain into convex polytopes \mathcal{P}_i , and that within each polytope the input-output relation is

affine, i.e. it is a piece-wise affine function. Then the following must hold:

1.)
$$\mathcal{B} \subseteq \mathcal{P}_i \iff u_{j,k} \le 0 \text{ or } 0 \le l_{j,k} \ \forall j,k$$

2.) $\mathcal{B} \subseteq \mathcal{P}_i \implies \mathcal{P}_{\mathcal{B}} \equiv f_{\theta}(\mathcal{B})$ (6)

The first line in Equation (6) follows from the fact that the boundaries of the polytopes \mathcal{P}_i arise from the discontinuity of the non-linear activations $\max\{z,0\}$. The second line follows from the first one: if $l_{j,k} \leq u_{j,k} \leq 0$ or $0 \leq 1$ $l_{j,k} \leq u_{j,k}$ for all j and k, then the relaxations of the ReLU activations will be exact and f_{θ} is just an affine map.

3.2.3. SPLITTING CRITERION

From Algorithm 1, a natural question arises: what is considered to be a "good" split of the input set? While easy to state, this question is far from trivial. Picking appropriate splits has important consequences regarding the time and memory efficiency of input-splitting verification algorithms. In Figure 3, an example is provided showcasing this phenomenon. The horizontal split results in an over-approximation that guarantees that the image of \mathcal{B} does not intersect with \mathcal{S} . In contrast, the vertical split results in an over-approximation that still intersects S.

To the best of our knowledge, current input-splitting methods (Wang et al., 2018) use gradient information between inputs and outputs to decide which axis to split. These methods leverage the structure of the Jacobian of the ReLU network to compute bounds on the gradient between inputs and outputs. In particular, note that for any point in the domain, the Jacobian for a ReLU network can always be written as

$$J = W_{K-1} \prod_{i=1}^{K-2} \Sigma_i W_i \tag{7}$$

for $\Sigma_i = diag(v_i)$ and the Boolean vector $v_i \in \{0, 1\}^{n_{i+1}}$. Starting from i = K - 1 going backwards to i = 1, these methods select appropriate values of v_i to compute upper bounds for $||\frac{df_{\theta}(x)}{dx_k}||_{\infty} \leq U_k$ for $k=1,...,n_1$. Using the side lengths Δx_k of the box, \mathcal{B} is split in half across the axis with greatest smear value $s_k = U_k \Delta x_k$. Geometrically, this mechanism tries to reduce in half the box along the axis that causes most "stretching" for any of the outputs. This reasoning, however, can be counterproductive when considering which splits to make. In particular, the upper bound might be very loose in practice, and even in cases where U_k can be achieved, it may be the case that the polytopic region that achieves this upper bound is very small.

Fundamentally, the over-approximation is caused by the convex relaxations of the ReLU non-linearities, therefore, we posit that an effective splitting procedure should leverage

the information of the relaxed nodes (i.e. nodes for which the convex relaxation is not exact) rather than using bounds on the gradients between input and outputs. In the next sections we investigate how to estimate changes in the upper and lower bounds of any given node given a specific split in \mathcal{B} .

4. Shadow Prices and Bound Rates

We now study the sensitivity of the lower and upper bounds for a relaxed ReLU node with respect to changes in the shape of the input box. To that end, we introduce some useful properties of linear programs and relate them to our problem at hand.

4.1. Measuring Constraint Sensitivity

For a linear program with non-empty feasible region the following property holds

$$\underline{p}^* = \min_{s.t. \ A\underline{z} \preceq b} w\underline{z} = w\underline{z}^* + (A\underline{z}^* - b)^T \underline{\lambda}^*$$
(8)
$$\bar{p}^* = \max_{s.t. \ A\bar{z} \preceq b} w\bar{z} = w\bar{z}^* + (b - A\bar{z}^*)^T \bar{\lambda}^*,$$
(9)

$$\bar{p}^* = \max_{\substack{s \ t \ A\bar{z} \prec b}} w\bar{z} = w\bar{z}^* + (b - A\bar{z}^*)^T \bar{\lambda}^*, \quad (9)$$

where \bar{z}^* and \underline{z}^* correspond to the primal's maximizer/minimizer and $\bar{\lambda}^* \succeq 0$ and $\underline{\lambda}^* \succeq 0$ correspond to the dual's minimizer/maximizer. This result follows from strong duality and complementary slackness (Boyd & Vandenberghe, 2004). A useful feature of the right-hand sides of Equation (8) is that they link how small perturbations of the constraint parameters (A and b) affect the optimal values \bar{p}^* and p^* . In the field of economics, the rate of increase/decrease of the optimal value with respect to a certain constraint is known as the shadow price. In some instances the shadow prices are the optimal dual variables.

From (8), we can readily study how small changes in the size of the input box affect the upper and lower bounds of all nodes in the first layer of the network. In particular, it follows that for the k-th node in the first layer and the i-th bias of our constraint set,

$$\frac{dl_{1,k}}{d\tilde{b}_{1,i}} = -\underline{\lambda}_{k,i}^* \qquad \frac{du_{1,k}}{d\tilde{b}_{1,i}} = \bar{\lambda}_{k,i}^*. \tag{10}$$

This result is expected, since the growth of bias terms results in a bigger box and, thus, in smaller lower bounds and bigger upper bounds. A nice feature of most modern LP solvers is that they provide the associated optimal dual variables when solving the primal problem. These rates can therefore be computed without any noticeable computational overhead when generating the convex over-approximations.

To derive the rates of the upper and lower bounds for all nodes in the network we first need to fully characterize the constraint sets given by $\hat{A}_{1:K-1}$ and $\hat{b}_{1:K-1}$ first introduced in Section 3.1.

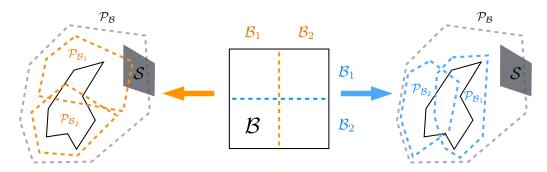


Figure 3. The choice of axis along which the set \mathcal{B} is split may affect the verification time. In this example, a vertical split (orange) results in two over-approximations, one of which still intersects with \mathcal{S} , whereas the horizontal split (blue) results in tighter over-approximations.

4.2. Constraint Sets Characterization

Thus far we have seen how the rates for the lower and upper bounds can be computed for the first layer. For this particular instance, the constraint set is given by \tilde{A}_1 and \tilde{b}_1 , and together they represent a box in \mathbb{R}^{n_1} which is provided to us. The expressions for \tilde{A}_j and \tilde{b}_j for the j-th $(j \geq 2)$ layer are given by

$$\tilde{A}_{j} = \begin{bmatrix} \tilde{A}_{j-1} & 0 \\ O_{j-2} & R_{j-1} \end{bmatrix}, R_{j-1} = \begin{bmatrix} 0 & -I \\ W_{j-1} & -I \\ -D_{j-1}W_{j-1} & I \end{bmatrix},$$

$$\tilde{b}_{j} = \begin{bmatrix} \tilde{b}_{j-1} \\ r_{j-1} \end{bmatrix}, r_{j-1}^{T} = \begin{bmatrix} 0^{T} - b_{j-1}^{T} & (D_{j-1}b_{j-1} + d_{j-1})^{T} \end{bmatrix}$$

$$(11)$$

where O_{j-2} is a matrix of zeros whose number of columns is equal to $\sum_{k=1}^{j-2} n_k$. Note how the dimensionality of the constraints given by \tilde{A}_j and \tilde{b}_j increases as we try to compute upper and lower bounds for deeper layers. These expressions can be derived from the inequalities provided in Equation (3). As examples, we provide the expressions for $\tilde{A}_{2:3}$ and $\tilde{b}_{2:3}$

$$\tilde{A}_{2} = \begin{bmatrix} \tilde{A}_{1} & 0 \\ 0 & -I \\ W_{1} & -I \\ -D_{1}W_{1} & I \end{bmatrix}, \tilde{b}_{2} = \begin{bmatrix} \tilde{b}_{1} \\ 0 \\ -b_{1} \\ D_{1}b_{1} + d_{1} \end{bmatrix},$$

$$\tilde{A}_{3} = \begin{bmatrix} \tilde{A}_{1} & 0 & 0 \\ 0 & -I & 0 \\ W_{1} & -I & 0 \\ -D_{1}W_{1} & I & 0 \\ 0 & 0 & -I \\ 0 & W_{2} & -I \\ 0 & -D_{2}W_{2} & I \end{bmatrix}, \tilde{b}_{2} = \begin{bmatrix} \tilde{b}_{1} \\ 0 \\ -b_{1} \\ D_{1}b_{1} + d_{1} \\ 0 \\ -b_{2} \\ D_{2}b_{2} + d_{2} \end{bmatrix}.$$

$$(12)$$

4.3. Forward Computation of the Bound Rates

With the definitions for \tilde{A}_j and \tilde{b}_j available and (8), we can proceed to compute the rates of the upper and lower bounds with respect to the bias terms of our input box. Denoting $(\underline{\lambda}_{j,k}^*, \bar{\lambda}_{j,k}^*)$ and $(\underline{z}_{j,k}^*, \bar{z}_{j,k}^*)$ as the set of primal and dual optimal variables for the maximization and minimization problems from (4), we have that

$$\frac{dl_{j,k}}{d\tilde{b}_{1,i}} = \frac{d}{d\tilde{b}_{1,i}} \left(\tilde{b}_j^T \underline{\lambda}_{j,k}^* \right) - \frac{d}{d\tilde{b}_{1,i}} \left((\tilde{A}_j \underline{z}_{j,k}^*)^T \underline{\lambda}_{j,k}^* \right)
\frac{du_{j,k}}{d\tilde{b}_{1,i}} = -\frac{d}{d\tilde{b}_{1,i}} \left(\tilde{b}_j^T \bar{\lambda}_{j,k}^* \right) + \frac{d}{d\tilde{b}_{1,i}} \left((\tilde{A}_j \bar{z}_{j,k}^*)^T \bar{\lambda}_{j,k}^* \right).$$
(13)

Before proceeding to drive the analytic expression of Equation (13), it is worth noting that changes in the bias terms of the input box $\tilde{b}_{1,:}$ only affect (through the computation of the upper and lower bounds) a subset of the entries of \tilde{A}_j and \tilde{b}_j . In particular, given $R_{1:j-1}$ and $r_{1:j-1}$ defined in Equation (11), only entries containing dependencies with $D_{1:j-1}$ and $d_{1:j-1}$ will contribute to the gradient expression in (13). Hence, focusing our attention to the first term, for any vector λ

$$\frac{d}{d\tilde{b}_{1,i}}(\tilde{b}_{j}^{T}\lambda) = \frac{d}{d\tilde{b}_{1,i}}\left(\tilde{b}_{1}^{T}\lambda_{0} + \sum_{l=1}^{j-1} r_{l}^{T}\lambda_{l}\right)$$

$$= \lambda_{0,i} + \sum_{l=1}^{j-1} \frac{d}{d\tilde{b}_{1,i}}(D_{l}b_{l} + d_{l})^{T}\lambda_{\hat{l}}$$

$$= \lambda_{0,i} + \sum_{l=1}^{j-1} \sum_{t=0}^{n_{l}} \frac{d}{d\tilde{b}_{1,i}} \frac{u_{l,t}b_{l,t} - u_{l,t}l_{l,t}}{u_{l,t} - l_{l,t}}\lambda_{\hat{l},t}$$

$$= \lambda_{0,i} + \sum_{l=1}^{j-1} \sum_{t=0}^{n_{l}} \left(c_{1} \frac{du_{l,t}}{d\tilde{b}_{1,i}} + c_{2} \frac{dl_{l,t}}{d\tilde{b}_{1,i}}\right)\lambda_{\hat{l},t}$$

$$= \lambda_{0,i} + \sum_{l=1}^{j-1} \sum_{t=0}^{n_{l}} \left(c_{1} \frac{du_{l,t}}{d\tilde{b}_{1,i}} + c_{2} \frac{dl_{l,t}}{d\tilde{b}_{1,i}}\right)\lambda_{\hat{l},t}$$
(14)

where
$$c_1 = \frac{l_{l,t}(l_{l,t} - b_{l,t})}{(u_{l,t} - l_{l,t})^2}$$
 and $c_2 = \frac{u_{l,t}(b_{l,t} - u_{l,t})}{(u_{l,t} - l_{l,t})^2}$ and $\lambda_{\hat{l}}, \lambda_{l}$

Algorithm 2 Bound Estimation-based Splitting

```
1: procedure Split(\mathcal{B}, (l, u), (\underline{z}^*, \underline{\lambda}^*), (\bar{z}^*, \bar{\lambda}^*))
                   \frac{dl_{j,k}}{d\bar{b}_{1,i}} \leftarrow \text{LBounds}(\mathcal{B}, (l, u), (\underline{z}^*, \underline{\lambda}^*), (\bar{z}^*, \bar{\lambda}^*))
  2:
                    \frac{du_{j,k}}{d\tilde{b}_{1,i}} \leftarrow \operatorname{UBounds}(\mathcal{B}, (l,u), (\underline{z}^*, \underline{\lambda}^*), (\bar{z}^*, \bar{\lambda}^*))
  3:
                   i = 1, c = L(1)
  4:
                   for k = 2, ..., n_1 do
  5:
   6:
                            if L(k) < c then
  7:
                                     c \leftarrow L(k)
                                     i \leftarrow k
  8:
  9:
                   \mathcal{B}_1 \leftarrow \mathcal{B}_i
                   \mathcal{B}_2 \leftarrow \mathcal{B} \backslash \mathcal{B}_i
10:
                   return \mathcal{B}_1, \mathcal{B}_2
11:
```

correspond to a subset of the entries in λ_l and λ respectively. The term $b_{l,t}$ denotes entry t of the bias term in layer l. From Equation (14) it is clear that the rates of the upper and lower bounds for layer j depend on all the rates from previous layers.

The rates for the second term can be derived in a similar manner. For any vector z and λ ,

$$\frac{d}{d\tilde{b}_{1,i}} (\tilde{A}_{j}z)^{T} \lambda = \frac{d}{d\tilde{b}_{1,i}} ((\tilde{A}_{1}z_{0})^{T} \lambda_{0} + \sum_{l=1}^{j-1} (R_{l} \begin{bmatrix} z_{l-1} \\ z_{l} \end{bmatrix})^{T} \lambda_{l})$$

$$= \sum_{l=1}^{j-1} \frac{d}{d\tilde{b}_{1,i}} (-D_{l}W_{l}z_{l-1})^{T} \lambda_{\hat{l}}$$

$$= -\sum_{l=1}^{j-1} \sum_{t=0}^{n_{l}} \frac{d}{d\tilde{b}_{1,i}} \frac{(w_{l,t}z_{l-1})u_{l,t}}{u_{l,t} - l_{l,t}} \lambda_{\hat{l},t}$$

$$= \sum_{l=1}^{j-1} \sum_{t=0}^{n_{l}} \left(\hat{c}_{1} \frac{dl_{l,t}}{d\tilde{b}_{1,i}} - \hat{c}_{2} \frac{du_{l,t}}{d\tilde{b}_{1,i}} \right) \lambda_{\hat{l},t}$$
(15)

where $\hat{c}_1 = \frac{u_{l,t}(w_{l,t}z_{l-1})}{(u_{l,t}-l_{l,t})^2}$ and $\hat{c}_2 = \frac{l_{l,t}(w_{l,t}z_{l-1})}{(u_{l,t}-l_{l,t})^2}$, and $w_{l,t}$ corresponds to row t of the weight matrix in layer l.

Using (14) together with (15) we can compute expressions for (13) in a forward manner using the optimal primal and dual variables.

5. Bound Estimation and Splitting

With the information provided in Section 4, we have means to estimate how the lower and upper bounds of our convex relaxation change as a function of the biases of the input set. In particular, splitting the box \mathcal{B} in half can be viewed as translating one the facets of the box to its center. Since the translation of any facet is achieved by reducing the associated bias term by a certain amount $\Delta \tilde{b}_{1,i}$, the estimated new

lower and upper bounds for the resulting set $\mathcal{B}_i \subset \mathcal{B}$ will be

$$l_{j,k}^{\mathcal{B}_{i}} \approx l_{j,k} + \frac{dl_{j,k}}{db_{1,i}} \Delta \tilde{b}_{1,i}$$

$$u_{j,k}^{\mathcal{B}_{i}} \approx u_{j,k} + \frac{du_{j,k}}{db_{1,i}} \Delta \tilde{b}_{1,i} .$$
(16)

Using these estimates, we propose the following metric for determining which axis to split along:

$$L(i) = -\sum_{j=1}^{K-2} \sum_{k=1}^{n_j} \max\{0, u_{j,k}^{\mathcal{B}_i}\} \min\{0, l_{j,k}^{\mathcal{B}_i}\}.$$
 (17)

Whichever axis minimizes L(i) is chosen for splitting. The max/min terms in the sum ensure that upper and lower bounds that start close to zero have less contribution in the overall splitting decision. Note how the cost metric is only zero whenever all relaxations are tight. We summarize this splitting procedure in Algorithm 2.

6. Experiments

In this section, we present the experimental results obtained on a set of benchmark verification tasks by our input-splitting approach based on bound rates. As baseline, we compare against the input-output gradient-based method discussed in Section 3.2.3.

6.1. Airborne Collision Avoidance System Verification

The ACAS benchmark verification task comprises a set of ten properties $\phi_{1:10}$ to be checked on a subset of 45 feedforward ReLU networks. All of the neural networks have the same architecture, with 5 inputs, 5 outputs and 6 hidden layers with 50 neurons each. The five inputs represent a specific configuration between two aircraft, one which is denoted as the ownship, and the other as the intruder. The inputs are ρ (distance between aircraft), θ (heading angle of ownship), ψ (heading angle of intruder), v_{own} (speed of ownship) and v_{int} (speed of intruder). The output corresponds to five scalars: COC (Clear of Conflict), weak left, weak right, strong left and strong right. The output with the greatest value is the advice action for the ownship. The properties $\phi_{1:10}$ specify a box-shaped subset \mathcal{B} in the input space and a set S in the output space. A property is said to be *unsatisfied* for a given network if (1) is true, otherwise the network satisfies the property. Details for each property are given in the Appendix. For further details on this benchmark we direct the reader to (Katz et al., 2017).

6.2. Experimental Details

There are a few technical aspects that need to be clarified before proceeding to the results. In our implementation, we did *not* use any form of parallelism, even though this approach can be readily extended to a parallel implementation

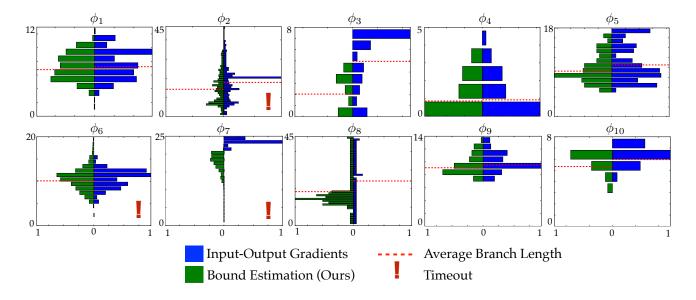


Figure 4. Horizontal histograms displaying the number of branches of each length generated by each type of splitting procedure. Each pair of histograms is normalized with the maximum branch length reached for that specific property. For histograms where timeouts occur we do not report the average branch length, since the associated tree has not finished growing. For property 2 we do report the average branch length only for networks that could be verified by both IOG and BE procedures.

φ	# NNs	IOG		BE (Ours)	
		U/S/T	Search depth	U/S/T	Search depth
1	45	45/0/0	7.22 ± 2.11	45/0/0	6.79 ± 1.90
2	34	0/30/4	$17.87 \pm 9.00^*$	0/32/2	$14.29 \pm 7.94^*$
3	42	42/0/0	4.96 ± 2.50	42/0/0	2.03 ± 1.49
4	42	42/0/0	0.92 ± 1.13	42/0/0	0.85 ± 1.03
5	1	1/0/0	10.5 ± 3.65	1/0/0	9.25 ± 2.79
6	1	0/0/1	N/A	1/0/0	10.1 ± 2.48
7	1	0/0/1	N/A	0/0/1	N/A
8	1	0/1/0	24.1 ± 13.40	0/1/0	18.33 ± 10.05
9	1	1/0/0	9.53 ± 1.47	1/0/0	9.09 ± 1.47
10	1	1/0/0	5.96 ± 0.80	1/0/0	5.34 ± 0.89

Table 1. Verification results (U: unsatisfied. S: satisfied. T: timeout) and search depth (mean \pm standard deviation) for all properties of the ACAS benchmark. These results show that our approach validates the properties as well as, or better than, IOG – even with a reduced search depth. We use * in property 2 to denote that the comparison is *only* for networks where neither IOG nor BE procedures had a timeout.

similar to (Wang et al., 2018). Each individual network to be verified was given a maximum execution time of 3 hours, at which point the verification function halts and a timeout flag is returned. All verification tasks were performed on a 12 core, 64-bit machine with Intel Core i7-5820K CPUs @ 3.30GHz. The experimental setting, and our approach, was coded in Python using the Gurobi optimization package.

6.3. Comparison of Splitting Procedures

In this section, we provide a side-by-side comparison between input-output gradient-based splits (IOG) and splits generated by our approach based on bound estimation (BE).

Table 1 shows the verification results for each of the 10 properties. The first column enumerates the properties and the second column shows the number of networks the property was tested against. The third and fifth columns show the verification results, which can either be unsatisfied, satisfied or timeout, for the IOG and BE-based splits, respectively. Columns four and six show the average depth for the binary trees generated during the verification procedures. When timeouts occurred, we did not report the average depth and standard deviations, as they are misleading metrics representing partially built binary trees, the only exception being property 2¹, for which 30 networks were able to be verified by both IOG and BE procedures². In the Appendix we include an additional graph with timeouts for property 2.

The horizontal histograms in Figure 4 depict the number of branches (horizontal axis) of a specific length (vertical axis) for the binary trees generated by the verification tasks. For each histogram, the distribution on the left (green) corresponds to the distribution generated by BE-based splitting. The distributions on the right (blue) correspond to IOG-

¹Similar to (Wang et al., 2018), we omit networks $\{N_{4,2}, N_{5,3}\}$.

²IOG timeouts: networks $\{N_{3,3}, N_{3,4}, N_{3,8}, N_{4,9}\}$. BE timeouts: networks $\{N_{3,3}, N_{4,9}\}$.

4	# NNs	# Nodes		+ /+
φ	# 11118	IOG	BE	$t_{ m BE}/t_{ m IOG}$
1	45	5241	4541	0.873
2	34*	799 [†]	553 [†]	0.872^{\dagger}
3	43	464	164	0.361
4	43	84	82	0.922
5	1	491	369	0.877
6	1*	>1808*	1210	<0.682*
7	1*	>1269*	>983*	1.0*
8	1	95	245	1.215
9	1	947	737	0.791
10	1	105	63	0.771

Table 2. Number of nodes for the two split mechanisms (smaller is better) and the ratio of computational time for the two methods. We can see that for 8 out of 10 properties our approach (BE) reduces the number of nodes generated, as well as the computational time. * indicates that some verification did not complete within the allocated time. We use † in property 2 to denote that the comparison is *only* for networks where neither IOG nor BE procedures had a timeout.

based splits. The dashed red line in each plot shows the average depth, which is reported in Table 1. In cases where timeouts occurred for either type of procedure we include an exclamation mark.

Table 2 shows a few more important metrics from the verification tasks. The first to columns are the same as Table 1. Columns three and four show the exact number of nodes that were generated during verification of each task. The fifth column shows the ratio of times between BE-based splits and IOG-based splits for a given verification task. Values below 1.0 imply a strict improvement of BE-splits over IOG-splits. In cases where timeouts occurred we include an asterisk symbol.

7. Discussion

As seen in the previous section in Table 1, BE-based splits are able to prove all properties to be either unsatisfied or satisfied except for ϕ_7 and for two of the networks in ϕ_2 . These results match the ones found in (Katz et al., 2017), with the exception of ϕ_1 , where we do not timeout³, and ϕ_2 , where we do timeout twice. In contrast, while IOG-based splits are also able to verify most of the properties, additional timeouts occur for properties ϕ_2 and ϕ_6 . In all cases, we found the average search depth μ_{BE} to be smaller than μ_{IOG} . Along these same lines, Table 2 shows that for 8 out of the 10 properties we get a strict improvement on the amount of time needed for verification. This improvement appears to be closely related to the number of nodes that

were generated during the verification task. In addition, a reduced number of nodes also decreases the amount of memory required to store the partition of the input set \mathcal{B} .

Note from Table 1 that 6 timeouts occurred for property ϕ_2 : 4 for IOG and 2 for BE. We believe BE-based splits outperformed IOG-based splits in these instances because the metric in (17) encourages splits to generate partitions that lead towards regions of the input space where all ReLU nodes are either active or inactive (i.e. towards the interior of the polytopic regions \mathcal{P}_i of the domain), which leads to finding examples faster. In addition, when comparing networks which were able to be verified by both procedures, BE still outperformed IOG in the amount of time and nodes generated as shown in Table 2. While BE-splits did under perform when verifying property ϕ_8 , it is important to note that this property is only verified against a single network. In constrast, property ϕ_2 was tested against 30 different networks. In the Appendix we include information on property ϕ_2 including timeouts.

8. Conclusion

In safety-critical application, such as self-driving, it is important to be able to verify the behavior learned by deep neural networks. In this paper, we introduced a new technique for verification of ReLU NNs, that relies on splitting the input set. Previous work leverages information of the gradient between inputs and outputs to decide how splits should be undertaken in order to verify a property of interest. In this work, we showed that by using shadow prices, a metric representing constraint sensitivity, and estimates on the bounds, we can substantially improve the amount of memory and time required to solve verification tasks. We test our proposed approach on the ACAS benchmark and provide a side-by-side comparison between both splitting procedures. Our results show substantial improvements both in the amount of memory and time required to solve these benchmarks. Future work will focus on extending the verification to more general classes of NNs, as well as input sets with arbitrary topology.

References

Bojarski, M., Del Testa, D., Dworakowski, D., Firner, B., Flepp, B., Goyal, P., Jackel, L. D., Monfort, M., Muller, U., Zhang, J., et al. End to end learning for self-driving cars. *arXiv preprint arXiv:1604.07316*, 2016.

Boyd, S. and Vandenberghe, L. *Convex optimization*. Cambridge university press, 2004.

Ehlers, R. Formal verification of piece-wise linear feed-forward neural networks. *CoRR*, abs/1705.01320, 2017. URL http://arxiv.org/abs/1705.01320.

³Reluplex uses 4h timeout threshold. For property 1, Reluplex reported 4 timeouts.

- Graves, A., Mohamed, A.-r., and Hinton, G. Speech recognition with deep recurrent neural networks. In *IEEE* international conference on Acoustics, speech and signal processing (ICASSP), pp. 6645–6649, 2013.
- Huang, X., Kwiatkowska, M., Wang, S., and Wu, M. Safety verification of deep neural networks. *CoRR*, abs/1610.06940, 2016. URL http://arxiv.org/abs/1610.06940.
- Katz, G., Barrett, C. W., Dill, D. L., Julian, K., and Kochenderfer, M. J. Reluplex: An efficient SMT solver for verifying deep neural networks. *CoRR*, abs/1702.01135, 2017. URL http://arxiv.org/abs/1702.01135.
- Kim, Y. Convolutional neural networks for sentence classification. *arXiv preprint arXiv:1408.5882*, 2014.
- Kolter, J. Z. and Wong, E. Provable defenses against adversarial examples via the convex outer adversarial polytope. *CoRR*, abs/1711.00851, 2017. URL http://arxiv.org/abs/1711.00851.
- LeCun, Y., Bengio, Y., and Hinton, G. Deep learning. *nature*, 521(7553):436, 2015.
- Lomuscio, A. and Maganti, L. An approach to reachability analysis for feed-forward relu neural networks. *CoRR*, abs/1706.07351, 2017. URL http://arxiv.org/abs/1706.07351.
- Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., and Riedmiller, M. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- Montufar, G. F., Pascanu, R., Cho, K., and Bengio, Y. On the number of linear regions of deep neural networks. In *Advances in neural information processing systems*, pp. 2924–2932, 2014.
- Raghu, M., Poole, B., Kleinberg, J., Ganguli, S., and Dickstein, J. S. On the expressive power of deep neural networks. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70, pp. 2847–2854. JMLR. org, 2017.
- Raghunathan, A., Steinhardt, J., and Liang, P. Certified defenses against adversarial examples. *CoRR*, abs/1801.09344, 2018. URL http://arxiv.org/abs/1801.09344.
- Serra, T., Tjandraatmadja, C., and Ramalingam, S. Bounding and counting linear regions of deep neural networks. *CoRR*, abs/1711.02114, 2017. URL http://arxiv.org/abs/1711.02114.

- Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A., et al. Mastering the game of go without human knowledge. *Nature*, 550(7676):354, 2017.
- Tjeng, V. and Tedrake, R. Verifying neural networks with mixed integer programming. *CoRR*, abs/1711.07356, 2017. URL http://arxiv.org/abs/1711.07356.
- Wang, S., Pei, K., Whitehouse, J., Yang, J., and Jana, S. Formal security analysis of neural networks using symbolic intervals. *CoRR*, abs/1804.10829, 2018. URL http://arxiv.org/abs/1804.10829.
- Xiang, W., Tran, H., and Johnson, T. T. Output reachable set estimation and verification for multi-layer neural networks. *CoRR*, abs/1708.03322, 2017a. URL http://arxiv.org/abs/1708.03322.
- Xiang, W., Tran, H., and Johnson, T. T. Reachable set computation and safety verification for neural networks with relu activations. *CoRR*, abs/1712.08163, 2017b. URL http://arxiv.org/abs/1712.08163.

A. Appendix

A.1. Property ϕ_2 including timeouts

Here we include the histogram of property ϕ_2 with timeouts.

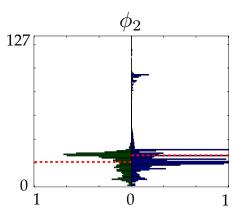


Figure 5. Histogram of branch lengths for property 2 including the timeouts.

The IOG search depth was 26.59 ± 16.48 with 17619 nodes generated. The BE search depth was 21.02 ± 7.82 with 6927 nodes generated. The ratio of times t_{BE}/t_{IOG} is 0.56.

A.2. ACAS Benchmark

We now list the 10 properties contained in the Airborne Collision Avoidance System benchmark.

Property ϕ_1 :

- **Description:** If the intruder is distant and is significantly slower than the ownship, the score of a COC advisory will always be below a certain fixed threshold.
- **Tested on:** all 45 networks.
- Input constraints: $\rho \geq 55947.691, v_{\rm own} \geq 1145, v_{\rm int} \leq 60.$
- **Desired output property:** the score for COC is at most 1500.

Property ϕ_2 :

- Description: If the intruder is distant and is significantly slower than the ownship, the score of a COC advisory will never be maximal.
- Tested on: N_{x,y} for all x ≥ 2 and for all y, except N_{4,2} and N_{5,3}.
- Input constraints: $\rho \geq 55947.691, v_{\mathrm{own}} \geq 1145, v_{\mathrm{int}} \leq 60$
- Desired output property: the score for COC is not the maximal score.

Property ϕ_3 :

- Description: If the intruder is directly ahead and is moving towards the ownship, the score for COC will not be minimal.
- **Tested on:** all networks except $N_{1,7}$, $N_{1,8}$, and $N_{1,9}$.
- Input constraints: $1500 \le \rho \le 1800$, $-0.06 \le \theta \le 0.06$, $\psi \ge 3.10$, $v_{\text{own}} \ge 980$, $v_{\text{int}} \ge 960$.
- **Desired output property:** the score for COC is not the minimal score.

Property ϕ_4 :

- **Description:** If the intruder is directly ahead and is moving away from the ownship but at a lower speed than that of the ownship, the score for COC will not be minimal.
- **Tested on:** all networks except $N_{1,7}$, $N_{1,8}$, and $N_{1,9}$.
- Input constraints: $1500 \le \rho \le 1800$, $-0.06 \le \theta \le 0.06$, $\psi = 0$, $v_{\text{own}} \ge 1000$, $700 \le v_{\text{int}} \le 800$.
- **Desired output property:** the score for COC is not the minimal score.

Property ϕ_5 :

- **Description:** If the intruder is near and approaching from the left, the network advises "strong right".
- Tested on: $N_{1,1}$.
- Input constraints: $250 \le \rho \le 400, \ 0.2 \le \theta \le 0.4, \ -3.141592 \le \psi \le -3.141592 + 0.005, \ 100 \le v_{\rm own} \le 400, \ 0 \le v_{\rm int} \le 400.$
- **Desired output property:** the score for "strong right" is the minimal score.

Property ϕ_6 :

- **Description:** If the intruder is sufficiently far away, the network advises COC.
- Tested on: $N_{1,1}$.
- Input constraints: $12000 \le \rho \le 62000$, $(0.7 \le \theta \le 3.141592) \lor (-3.141592 \le \theta \le -0.7)$, $-3.141592 \le \psi \le -3.141592 + 0.005$, $100 \le v_{\rm own} \le 1200$, $0 \le v_{\rm int} \le 1200$.
- Desired output property: the score for COC is the minimal score.

Property ϕ_7 :

- **Description:** If vertical separation is large, the network will never advise a strong turn.
- Tested on: $N_{1.9}$.
- Input constraints: $0 \le \rho \le 60760, -3.141592 \le \theta \le 3.141592, -3.141592 \le \psi \le 3.141592, 100 \le v_{\rm own} \le 1200, 0 \le v_{\rm int} \le 1200.$
- **Desired output property:** the scores for "strong right" and "strong left" are never the minimal scores.

Property ϕ_8 :

- **Description:** For a large vertical separation and a previous "weak left" advisory, the network will either output COC or continue advising "weak left".
- Tested on: $N_{2,9}$.
- Input constraints: $0 \le \rho \le 60760$, $-3.141592 \le \theta \le -0.75 \cdot 3.141592$, $-0.1 \le \psi \le 0.1$, $600 \le v_{\rm own} \le 1200$, $600 \le v_{\rm int} \le 1200$.
- Desired output property: the score for "weak left" is minimal or the score for COC is minimal.

Property ϕ_9 :

- **Description:** Even if the previous advisory was "weak right", the presence of a nearby intruder will cause the network to output a "strong left" advisory instead.
- Tested on: $N_{3,3}$.
- Input constraints: $2000 \le \rho \le 7000$, $-0.4 \le \theta \le -0.14$, $-3.141592 \le \psi \le -3.141592 + 0.01$, $100 \le v_{\rm own} \le 150$, $0 \le v_{\rm int} \le 150$.
- **Desired output property:** the score for "strong left" is minimal.

Property ϕ_{10} :

- Description: For a far away intruder, the network advises COC.
- Tested on: $N_{4.5}$.
- Input constraints: $36000 \le \rho \le 60760$, $0.7 \le \theta \le 3.141592$, $-3.141592 \le \psi \le -3.141592 + 0.01$, $900 \le v_{\rm own} \le 1200$, $600 \le v_{\rm int} \le 1200$.
- **Desired output property:** the score for COC is minimal.