

Static Analysis of ReLU Neural Networks with Tropical Polyhedra

Eric Goubault¹, Sébastien Palumby¹, Sylvie Putot¹, Louis Rustenholz¹, and Sriram Sankaranarayanan²(⊠),

¹ LIX, Ecole Polytechnique, CNRS and Institut Polytechnique de Paris, 91128 Palaiseau, France

{eric.goubault,sebastien.palumby,sylvie.putot, louis.rustenholz}@polytechnique.edu

² Engineering Center Computer Science, University of Colorado at Boulder, Boulder, USA

sriram.sankaranarayanan@colorado.edu

Abstract. This paper studies the problem of range analysis for feedforward neural networks, which is a basic primitive for applications such as robustness of neural networks, compliance to specifications and reachability analysis of neural-network feedback systems. Our approach focuses on ReLU (rectified linear unit) feedforward neural nets that present specific difficulties: approaches that exploit derivatives do not apply in general, the number of patterns of neuron activations can be quite large even for small networks, and convex approximations are generally too coarse. In this paper, we employ set-based methods and abstract interpretation that have been very successful in coping with similar difficulties in classical program verification. We present an approach that abstracts ReLU feedforward neural networks using tropical polyhedra. We show that tropical polyhedra can efficiently abstract ReLU activation function, while being able to control the loss of precision due to linear computations. We show how the connection between ReLU networks and tropical rational functions can provide approaches for range analysis of ReLU neural networks. We report on a preliminary evaluation of our approach using a prototype implementation.

1 Introduction and Related Work

Neural networks are now widely used in numerous applications including speech recognition, natural language processing, image segmentation, control and planning for autonomous systems. A central question is how to verify that they

E. Goubault—Supported in part by the academic Chair "Engineering of Complex Systems", Thalès-Dassault Aviation-Naval Group-DGA-Ecole Polytechnique-ENSTA Paris-Télécom Paris, and AID project "Drone validation and swarms of drones".

S. Sankaranarayanan—Supported by US National Science Foundation (NSF) award # 1932189. All opinions expressed are those of the authors and not necessarily of the sponsors.

[©] Springer Nature Switzerland AG 2021

C. Drăgoi et al. (Eds.): SAS 2021, LNCS 12913, pp. 166–190, 2021.

are correct with respect to some specification. Beyond correctness, we are also interested in questions such as explainability and fairness, that can in turn be specified as formal verification problems. Recently, the problem of verifying properties of neural networks has been investigated extensively under a variety of contexts. A natural neural network analysis problem is that of range estimation, i.e. bounding the values of neurons on the output layer, or some function of the output neurons, given the range of neurons on the input layer. A prototypical application of range estimation is the verification of the ACAS Xu the next generation collision avoidance system for autonomous aircrafts, which is implemented by a set of neural networks [23]. Such a verification problem is translated into a range estimation problem over these neural network wherein the input ranges concern a set of possible scenarios and the outputs indicate the possible set of advisories provided by the network [24].

Another prototypical application concerns the robustness of image classification wherein we wish to analyze whether a classification label remains constant for images in a neighborhood of a given image that is often specified using ranges over a set of pixels. Robustness is akin to numerical stability analysis, and for neural nets used as decision procedures (e.g. control of a physical apparatus), this is a form of decision consistency. It is also linked to the existence or non-existence of adversarial inputs, i.e. those inputs close to a well classified input data, that dramatically change the classification [38], and may have dire consequences in the real world [16].

Many formal methods approaches that have been successfully used in the context of program verification seem to be successfully leveraged to the case of neural net verification: proof-theoretic approaches, SMT techniques, constraint based analyzers and abstract interpretation. In this paper, we are interested in developing abstract interpretation [10] techniques for feedforward networks with ReLU activation functions. ReLU feedforward networks can be seen as loop-free programs with affine assignments and conditionals with affine guards, deciding whether the corresponding neuron is activated or not. For researchers in program analysis by abstract interpretation, this is a well known situation. The solutions range from designing a scalable but imprecise analyses by convexifications of the set of possible values of each neurons throughout all layers to designing a potentially exponentially complex analysis by performing a fully disjunctive analysis. In between, some heuristics have been successfully used in program analysis, that may alleviate the burden of disjunctive analysis, see e.g. [9,28]. Among classical convex abstractions, the zones [29] are a nice and scalable abstraction, successfully used in fully-fledged abstract interpretation based static analyzers [7]. In terms of disjunctive analysis, a compact way to represent a large class of disjunctions of zones are the tropical polyhedra, used for disjunctive program analysis in e.g. [3,4]. Tropical polyhedra are, similarly to classical convex polyhedra, defined by sets of affine inequalities but where the sum is replaced by max operator and the multiplication is replaced by the addition.

Zones are interesting for synthesizing properties such as robustness of neural networks used for classifying data. Indeed, classification relies on determining which output neuron has the greatest score, translating immediately into zone-

like constraints. ReLU functions $x \mapsto max(0,x)$ are tropically linear, hence an abstraction using tropical polyhedra will be exact. A direct verification of classification specifications can be done from a tropical polyhedron by computing the enclosing zone, see [4] and Sect. 2.1. In Fig. 1, we pictured the graph of the ReLU function y = max(x,0) for $x \in [-1,1]$ (Fig. 1a), and its abstraction by 1-ReLU in DeepPoly [36] (Fig. 1b), by a zone (Fig. 1c), and by a tropical polyhedron (Fig. 1d), which is exactly the graph of the function.

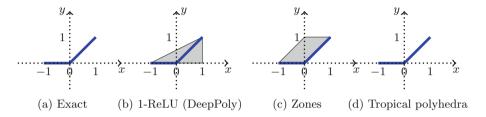


Fig. 1. Abstractions of the ReLU graph on [-1, 1]

Unfortunately, (classical) linear functions are tropically non-linear. But contrarily to program analysis where we generally discover the function to abstract inductively on the syntax, we are here given the weights and biases for the full network, allowing us to design much better abstractions than if directly using the ones available from the program verification literature.

It was recently proved [42] that the class of functions computed by a feedforward neural network with ReLU activation functions is exactly the class of rational tropical maps, at least when dealing with rational weights and biases. It is thus natural to look for guaranteed approximants of these rational tropical maps as abstractions.

Example 1 (Running example). Consider a neural network with 2 inputs x_1 and x_2 given in [-1, 1] and 2 outputs. The linear layer is defined by $h_1 = x_1 - x_2 - 1$, $h_2 = x_1 + x_2 + 1$ and followed by a ReLU layer with neurons y_1 and y_2 such that $y_1 = max(0, x_1 - x_2 - 1)$ and $y_2 = max(0, x_1 + x_2 + 1)$.

The exact range for nodes (h_1, h_2) is depicted in Fig. 2a in magenta (an octagon here), and the exact range for the output layer is shown in Fig. 2b in cyan: (y_1, y_2) take the positive values of of (h_1, h_2) . In Fig. 2c, the set of values the linear node h_1 can take as a function of x_1 , is represented in magenta. The set of values of the output neuron y_1 in function of x_1 is depicted in Fig. 2d, in cyan: when x_1 is negative, h_1 is negative as well, so $y_1 = 0$ (this is the horizontal cyan line on the left). When x_1 is positive, the set of values y_1 can take is the positive part of the set of values h_1 can take (pictured as the right cyan triangle). The line plus triangle is a tropical polyhedron, as we will see in Sect. 2.2.

We want to check two properties on this simple neural network:

 (P_1) : the input is always classified as belonging to the class identified by neuron y_2 , i.e. we always have $y_2 \geq y_1$

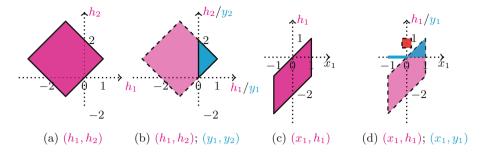


Fig. 2. Exact ranges for the neural net of Example 1 on $[-1,1] \times [-1,1]$. (P_2) is the complement of the red square in Fig. 2d.

 (P_2) : in the neighborhood [-0.25, 0.25] of 0 for x_1 , whatever x_2 in [-1, 1], the output y_1 is never above threshold 0.5 (unsafe zone materialized in red in Fig. 2d)

 (P_2) is a robustness property. We see on the blue part of Fig. 2b (resp. 2d) that the first (resp. second) property is true.

As we will see in Sect. 3, our tropical polyhedron abstraction is going to give the exact graph of y_1 as a function of x_1 , in cyan again, Fig. 3b.

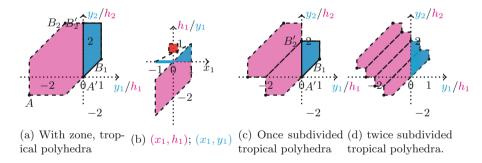


Fig. 3. Abstractions of a simple neural net on $[-1,1] \times [-1,1]$. Dashed lines in (b) enclose the classical convexification.

Therefore we will be able to prove robustness, i.e. (P_2) : the exact range for y_1 in cyan does not intersect the non complying states, in red. Note that all classically convex abstractions, whatever their intricacies, will need to extend the cyan zone up to the dashed line pictured in Fig. 3b, to get the full triangle, at the very least. This triangle is intersecting the red region making classically convex abstractions unable to prove (P_2) .

Our tropical abstraction projected on the y_2 , y_1 coordinates is not exact: compare the exact range in cyan in Fig. 2b with the abstraction in cyan in Fig. 3a. However, the cyan region in Fig. 3a is above the diagonal, which is enough for proving (P_1) .

Still, the abstraction has an area 2.5 times larger than the exact range, due to the tropical linearization of the tropical rational function y_1 . As with classical linearizations, a workaround is to make this linearization local, through suitable subdivisions of the input. We show in Fig. 3c the tropical polyhedric abstraction obtained by subdividing x_1 into two sub-intervals (namely [-1,0] and [0,1]): the cyan part of the picture is much closer to the exact range (1.5 times the exact area). Subdividing further as in Fig. 3d naturally further improves the precision (area 1.25 times the exact one).

As we will see in Sect. 2.2, tropical polyhedra are particular unions of zones: the tropical polyhedra in cyan of Figs. 3a and 3c are composed of just one zone, but the tropical polyhedron in cyan in Fig. 3d and the tropical polyhedron in magenta in Fig. 3c are the union of two zones. Finally, the tropical polyhedron in magenta in Fig. 3d is the union of four zones (generated by 9 extreme points, or 5 constraints, obtained by joining results from the subdivisions of the inputs).

Contributions. Section 2 introduces the necessary background notions, in particular tropical polyhedra. We then describe the following contributions:

- Sect. 3 introduces our abstraction of (classical) affine functions from \mathbb{R}^m to \mathbb{R}^n with tropical polyhedra. We fully describe internal and external representations, extending the classical abstractions of assignments in the zone abstract domain [29] or in the tropical polyhedra domain [4]. We prove correctness and equivalence of internal and external representations, allowing the use of the double description method [2].
- Based on the analysis of one layer networks of Sect. 3, we show in Sect. 4 how to get to multi-layered networks.
- Finally, Sect. 5 describes our implementations in C++ and using polymake [18] and presents some promising experiments. We discuss the cost and advantages of using the double description or of relying for further abstraction on either internal or external representations of tropical polyhedra.

Related Work. There exist many approaches to neural networks verification. We concentrate here on methods and tools designed for at least range over-approximation of ReLU feedforward networks.

It is natural to consider constraint based methods for encoding the ReLU function and the combinatorics of activations in a ReLU feedforward neural net.

Determining the range of a ReLU feedforward neural net amounts to solving min and max problems under sets of linear and ReLU constraints. This can be solved either by global optimisation techniques and branch and bound mechanisms, see e.g. DeepGo [32]. The encoding of the activation combinatorics can also be seen as mixed integer linear constraints, and MILP solver used for solving the range outer-approximation problem, see e.g. [6,39], or both branch and bound and MILP techniques, like Venus [8]. Similarly, Sherlock [13,14] performs range analysis using optimization methods (MILP and a combination of local search and global branch-and-bound approach), and considers also neural nets as

controllers within a feedback loop. Finally, some of these constraint-based analyzer improve the solution search by exploiting the geometry of the activation regions, [26].

A second category of such approaches is based on SMT methods, more specifically satisfiability modulo extensions of linear real arithmetic (encoding also RELU). The network is encoded in this logics and solvers provide answers to queries, in particular range over-approximation and robustness, see e.g. Marabou [25], extending Reluplex [24], and [15,22].

Range estimation for ReLU activated feedforward neural nets can also be performed using some of the abstract domains [11] that have been designed for program analysis, and in particular convex domains for numerical program verification. These include zonotopes [19,34], especially considering that feedforward neural nets with one hidden layer and ReLU activation functions are known to be characterizable by zonotopes, see e.g. [42], polyhedra [36], and other subpolyhedric or convex abstractions like symbolic intervals [20] used in Neurify [33] extending Reluval [40] or CROWN-IBP [41].

These abstractions allow to perform range estimation, i.e. to estimate outer approximations of the values of the output neurons given a set of values for the input neurons. They also allow to deal with robustness properties around training data, by proving that the range of the neural net on a small set around a training point gives the same class of outputs.

The main difficulty with these convex abstract domains is that they tend to lose too much precision on (non-convex) ReLU functions. Several methods have been proposed to cope with this phenomenon. The first one is to improve on the abstraction of ReLU, in particular by combining the abstraction of several ReLU functions on the same layer [35]. Another solution that has been proposed in the literature is to combine abstraction and some level of combinatorial exploration of the possible neuron activations, in the line of disjunctive program analysis [9,28]. RefineZono [37] implements methods combining polyhedric abstract domains with MILP solvers for encoding ReLU activation and refining the abstractions, NNENUM [5] uses combinations of zonotopes, stars sets with case splitting methods, and Verinet [21] uses abstractions similar to the polyhedric relaxations of DeepPoly, based on symbolic-interval propagation, with adaptive refinement strategies.

2 Preliminaries and Notations

2.1 Zones

The zone [29] abstraction represents restricted forms of affine invariants over variables, bounds on variable differences. Let a n-dimensional variable $x = (x_1, \ldots, x_n) \in \mathbb{R}^n$. The zone domain represents invariants of the form $(\bigwedge_{1 \leq i,j \leq n} x_i - x_j \leq c_{i,j}) \wedge (\bigwedge_{1 \leq i \leq n} a_i \leq x_i \leq b_i)$. A convenient representation is using difference bound matrices, or DBM. In order to encode interval constraints seamlessly in this matrix, a special variable x_0 , which is assumed to be a constant set to zero, is added to $x \in \mathbb{R}^n$. A DBM is then a $(n+1) \times (n+1)$

square matrix $C = (c_{ij})$, with elements in $\mathbb{R} \cup \{+\infty\}$, representing (concretisation operator) the following set of points in \mathbb{R}^n : $\gamma(C) = \{(x_1, \ldots, x_n) \in \mathbb{R}^n | \forall i, j \in [0, n], x_i - x_j \leq c_{i,j} \land x_0 = 0\}$.

For a matrix C that has non-empty concretization, the closure denoted C^* will be the smallest DBM for the partial order on matrices which represents $\gamma(C)$. Formally, a closed zone $C = (c_{ij})$ is such that: $\forall k \in \mathbb{N}, \forall (i_0, \ldots, i_k) \in [0, n]^{k+1}, c_{i_0, i_k} \leq c_{i_0, i_1} + \cdots + c_{i_{k-1}, i_k}, \forall i \in [0, j], c_{i,i} = 0$. Every constraint in a closed zone saturates the set $\gamma(C)$.

The best abstraction in the sense of abstract interpretation [11] of a non-empty set $S \subset \mathbb{R}^n$ is the zone defined by the closed DBM: $(c)_{ij} = \sup\{x_i - x_j | (x_1, \ldots, x_n) \in S \land x_0 = 0\}.$

Example 2. Consider the region defined as the union of the magenta and cyan parts of Fig. 3a in Example 1. It is a zone given by the inequalities: $(-3 \le h_1 \le 1) \land (-1 \le h_2 \le 3) \land (-4 \le h_1 - h_2 \le 0)$, i.e. given by the following DBM:

$$\begin{pmatrix}
0 & 3 & 1 \\
1 & 0 & 0 \\
3 & 4 & 0
\end{pmatrix}$$

The octagon [30] abstraction is an extension of the zone abstraction, which represents constraints of the form

$$\left(\bigwedge_{1 \le i, j \le n} \pm x_i \pm x_j \le c_{i,j}\right) \land \left(\bigwedge_{1 \le i \le n} a_i \le x_i \le b_i\right)$$

A set of octagonal constraints can be encoded as a difference bound matrix, similarly to the case of zones, but using a variable change to map octagonal constraints on zone constraints. For each variable x_i , two variables are considered in the DBM encoding, that correspond respectively to $+x_i$ and $-x_i$. Note that unary (interval) constraints, such as $x_i \leq b_i$, can be encoded directly as $x_i + x_i \leq 2b_i$, so that no additional variable x_0 is needed.

Example 3. The figure below right shows the exact range (the rotated square) of h_1 , h_2 of Example 1.

It is depicted in gray, as the intersection of two zones, one in cyan, Z_2 , and one in olive, Z_1 . Z_1 is the zone defined in Example 2 and Z_2 is the zone defined on variables $(h_1, -h_2)$ as follows:

$$(-3 \le h_1 \le 1) \land (-1 \le h_2 \le 3) \land (-2 \le h_1 + h_2 \le 2) \qquad b_1$$

2.2 Tropical Polyhedra

Tropical polyhedra are similar to ordinary convex polyhedra. Both can be defined either using affine constraints, known as the external description, or as convex hulls of extremal points and rays, known as the internal description. The major difference is the underlying algebra. Instead of using the classical ring \mathbb{R} of coefficients, with ordinary sum and multiplications, we use the so-called max-plus semiring \mathbb{R}_{max} . This semiring is based on the set $\mathbb{R}_{max} = \mathbb{R} \cup \{-\infty\}$, equipped with the addition $x \oplus y := max(x,y)$ and the multiplication $x \otimes y = x + y$. This is almost a ring: we have neutral elements $\mathbb{1} := 0$ for \otimes , and $\mathbb{0} = -\infty$ for \oplus , and an inverse for \otimes on $\mathbb{R}_{max} \setminus \{\mathbb{0}\}$ but not for \oplus . The algebra also fits in with the usual order \leq on \mathbb{R} , extended to $\mathbb{R}_{max} : x \leq y$ if and only if $x \oplus y = y$.

Tropical hyperplanes are similar to classical hyperplanes, and defined as the set of points satisfying $\bigoplus_{1 \leq i \leq k} a_i \otimes x_i \oplus c \leq \bigoplus_{1 \leq i \leq k} b_i \otimes x_i \oplus d$.

Now, as in the classical case, tropical polyhedra will be given (externally) as an intersection of n tropical hyperplanes, i.e. will be given as the location of points in \mathbb{R}^k_{max} satisfying n inequalities of the form of above. This can be summarized using matrices $A = (a_{ij})$ and $B = (b_{ij})$, two $n \times k$ matrices with entries in \mathbb{R}_{max} , and vectors of size k C and D as $Ax \oplus C \leq Bx \oplus D$.

Still similarly to the case of ordinary convex polyhedra, tropical polyhedra can also be described internally, as generated by extremal generators (points, rays). A tropical polyhedron can then be defined as the set of vectors $x \in \mathbb{R}^k_{max}$ which can be written as a tropical affine combination of generators v^i (the extreme points) and r^j (the extreme rays) as $x = \bigoplus_{i \in I} \lambda_i v^i \oplus \bigoplus_{j \in J} \mu_j r^j$ with $\bigoplus_{i \in I} \lambda_i = \mathbb{1}$.

Example 4 (Running example). Consider again the zone consisting of the union of the magenta and cyan parts in Fig. 3a. This is a tropical polyhedron, defined externally by: $max(h_1, -3, h_2, -1, h_2, h_1) \leq max(1, h_1, 3, h_2, h_1 + 4, h_2)$.

It can also be defined internally by the extremal point A, B_1 and B_2 of respective coordinates (-3,-1), (1,1) and (-1,3), depicted as dots in Fig. 3a. This means that the points z in this tropical polyhedron have coordinates (h_1,h_2) with $(h_1,h_2)=max(\lambda_0+A,\lambda_1+B_1,\lambda_2+B_2)$ with $max(\lambda_0,\lambda_1,\lambda_2)=\mathbb{1}=0$, i.e. all λ_i s are negative or null, and one at least among the λ_i s is zero.

For instance, when $\lambda_2 = -\infty$, z is on the tropical line linking A to B_1 :

$$(h_1, h_2) = (max(\lambda_0 - 3, \lambda_1 - 1), max(\lambda_0 - 1, \lambda_1 + 3))$$
(1)

with $\lambda_0, \lambda_1 \neq 0$ and either $\lambda_0 = 0$ or $\lambda_1 = 0$. Suppose $\lambda_0 = 0$, and suppose first that $\lambda_1 \leq -4$: $(h_1, h_2) = (-3, -1)$ which is point A. Suppose now $-4 \leq \lambda_1 \leq -2$, then $(h_1, h_2) = (-3, \lambda_1 + 3)$, which is the vertical line going from A to point (-3, 1). Finally, suppose $-2 \leq \lambda_1 \leq 0$, $(h_1, h_2) = (\lambda_1 - 1, \lambda_1 + 3)$ which is the diagonal going from (-3, 1) to B_1 . Similarly, one can show that the tropical line going from B_1 to B_2 is given by fixing $\lambda_0 = -\infty$ and making vary λ_1 and λ_2 . If $\lambda_0 < 0$ then $\lambda_1 = 0$ and z is point B_1 .

Now, applying the ReLU operator, which is linear in the tropical algebra, defines a tropical polyhedron with internal description given by ReLU (in each coordinate) of extreme points A, B_1 and B_2 , i.e. A' = (0,0), $B'_1 = B_1 = (1,1)$ and $B'_2 = (0,3)$, see Fig. 3a. Similarly, the zone which gives h_1 as a function of x_1 , see Fig. 3b, can be seen as a tropical polyhedron with extreme points (-1,-3), (1,1) and (1,-1). Applying ReLU to the second coordinate of these

three extreme points gives three points (-1,0), (1,1) and (1,0) which generate the tropical polyhedron in cyan of Fig. 3b.

It is also easy to see that after one subdivision, Fig. 3c, the set of values for (y_1, y_2) in cyan is a tropical polyhedron with three extreme points A', B'_1 and B_2 . After two subdivisions, Fig. 3d, the values of y_1 as a function of h_1 is a tropical polyhedron with 4 generators (depicted as dots in Fig. 3d). Note that the tropical polyhedron of Fig. 3d is the encoding of the union of two zones, one zone being the classical convex hull of points (0,0), (0,1), (0.5,1.5), (1,1.5) and (1,1), and the other being the classical convex hull of points (0,1), (0,2), (0.5,2) and (0.5,1.5).

All tropical polyhedra can thus be described both internally and externally, and algorithms, although costly, can be used to translate an external description into an internal description and vice-versa. This is at the basis of the double description method for classical polyhedra [12] and for tropical polyhedra [2]. Double description is indeed useful when interpreting set-theoretic unions and intersections, as in validation by abstract interpretation, see [12] again for the classical case, and e.g. [4] for the tropical case: unions are easier to compute using the extreme generator representation (the union of the convex hulls of sets of points is the convex hull of the union of these sets of points) while intersections are easier to compute using the external representation (the intersection of two polyhedra given by sets of constraints is given by the concatenation of these sets of constraints).

In the sequel, we will be using explicitly the max and (ordinary) + operators in place of \oplus and \otimes for readability purposes.

2.3 From Zone to Tropical Polyhedra and Vice-Versa

The following proposition characterizes the construction of tropical polyhedric abstractions from zones. We show that a zone defined on n variables can be expressed as the tropical convex hull of n + 1 points.

Proposition 1 (Internal tropical representation of closed zones)

Let $H_{ext} \subset \mathbb{R}^n$ be the n-dimensional zone defined by the conjunction of the $(n+1)^2$ inequalities $\bigwedge_{0 \leq i,j \leq n} (x_i - x_j \leq c_{i,j})$, where $\forall i,j \in [0,n], c_{i,j} \in \mathbb{R} \cup \{+\infty\}$. Assume that this representation is closed, then H_{ext} is equal to the tropical polyhedron H_{int} defined, with internal representation, as the tropical convex hull of the following extreme points (and no extreme ray):

$$A = (a_i)_{1 \le i \le n} := (-c_{0,1}, \dots, -c_{0,n}),$$

$$B_k = (b_{ki})_{1 \le i \le n} := (c_{k,0} - c_{k,1}, \dots, c_{k,0} - c_{k,n}), k = 1, \dots, n,$$

Example 5. The zone of Example 2 is the tropical polyedron with the three extreme generators A, B_1 and B_2 pictured in Fig. 3a, as deduced from Proposition 1 above.

Moreover, we can easily find the best zone (and also, hypercube) that outer approximates a given tropical polyhedron, as follows [4]. Suppose we have p extreme generators and rays for a tropical polyhedron $\mathcal{H}, A_1, \ldots, A_p$, that we put in homogeneous coordinates in \mathbb{R}^{n+1} by adding as last component 0 to the coordinates of the extreme generators, and $-\infty$ to the last component, for extreme rays, as customary for identifying polyhedra with cones, see e.g. [17].

Proposition 2 [4]. Let A be the matrix of generators for tropical polyhedron \mathcal{H} stripped out of rows consisting only of $-\infty$ entries, and A/A the residuated matrix which entries are $(A/A)_{i,j} = \min_{1 \le k \le p} a_{i,k} - a_{j,k}$. Then the smallest zone containing \mathcal{H} is given by the inequalities:

$$x_i - x_j \ge (A/A)_{i,j}$$
 for all $i, j = 1, \dots, n$
 $(A/A)_{i,n+1} \le x_i \le -(A/A)_{n+1,i}$ for all $i = 1, \dots, n$

Example 6. Consider the graph of the ReLU function on [-1,1], pictured in Fig. 1d. It has as generators the two extreme points $A_1 = (-1,0)$ and $A_2 = (1,1)$ (the graph is the tropical segment from A_1 to A_2). Homogenizing the coordinates and putting them in a matrix A (columns correspond to generators), we have

$$A = \begin{pmatrix} -1 & 1 \\ 0 & 1 \\ 0 & 0 \end{pmatrix} \text{ and } (A/A) = \begin{pmatrix} 0 & -1 & -1 \\ 0 & 0 & 0 \\ -1 & -1 & 0 \end{pmatrix}$$

meaning that the enclosing zone is given by $-1 \le x - y \le 0$, $-1 \le x \le 1$, $0 \le y \le 1$, which is the zone depicted in Fig. 1c.

2.4 Feedforward ReLU Networks

Feedforward ReLU networks that we are considering in this paper are a succession of layers of neurons, input layer first, a given number of hidden layers and then an output layer, each computing a certain affine transform followed by the application of the ReLU activation function:

Definition 1. A n-neurons ReLU network layer L with m inputs is a function $\mathbb{R}^m \to \mathbb{R}^n$ defined by, a weight matrix $W \in \mathcal{M}_{n,m}(\mathbb{R})$, a bias vector $b \in \mathbb{R}^n$, and an activation function ReLU: $\mathbb{R}^n \to \mathbb{R}^n$ given by ReLU $(x_1, \ldots, x_n) = (max(x_1, 0), \ldots, max(x_n, 0))$ so that for a given input $x \in \mathbb{R}^n$, its output is L(x) = ReLU(Wx + b).

Definition 2. A multi-layer perceptron F_N is given by a list of network layers $L_0, ..., L_N$, where layers L_i (i = 0, ..., N - 1) are n_{i+1} -neurons layers with n_i inputs, the action of F_N on inputs is defined by composing the action of successive layers: $F_N = L_N \circ ... \circ L_0$.

3 Abstraction of Linear Maps

3.1 Zone-Based Abstraction

We consider in this section the problem of abstracting the graph $\mathcal{G}_f = \{(x,y) \mid y = f(x)\}$ of a linear map f(x) = Wx + b with $x \in [\underline{x}_1, \overline{x}_1] \times \dots [\underline{x}_m, \overline{x}_m]$ where $W = (w_{i,j})$ is a $n \times m$ matrix and b a n-dimensional vector, by a tropical polyhedron \mathcal{H}_f . We will treat the case of multilayered networks in Sect. 4.

The difficulty is that linear maps in the classical sense are not linear maps in the tropical sense, but are rather (generalized) tropical polynomials, hence the exact image of a tropical polyhedron by a (classical) linear map is not in general a tropical polyhedron. We begin by computing the best zone abstracting \mathcal{G}_f and then represent it by a tropical polyhedron, using the results of Sect. 2.3. We then show in Sect. 3.2 that we can improve results using an octagon abstraction.

The tightest zone containing the image of a cube going through a linear layer can be computed as follows:

Proposition 3 (Optimal approximation of a linear layer by a zone)

Let $n, m \in \mathbb{N}$ and $f : \mathbb{R}^m \to \mathbb{R}^n$ an affine transformation defined, for all $x \in \mathbb{R}^m$ and $i \in [1, n]$, by $(f(x))_i = \sum_{j=1}^m w_{i,j}x_j + b_i$. Let $K \subset \mathbb{R}^m$ be an hypercube defined as $K = \prod_{1 \le j \le m} [\underline{x}_j, \overline{x}_j]$, with $\underline{x}_j, \overline{x}_j \in \mathbb{R}$. Then, the tightest zone \mathcal{H}_f of $\mathbb{R}^m \times \mathbb{R}^n$ containing $S := \{(x, f(x)) \mid x \in K\}$ is the set of all $(x, y) \in \mathbb{R}^m \times \mathbb{R}^n$ satisfying

$$\left(\bigwedge_{1 \leq j \leq m} \underline{x}_{j} \leq x_{j} \leq \overline{x}_{j}\right) \wedge \left(\bigwedge_{1 \leq i \leq n} m_{i} \leq y_{i} \leq M_{i}\right) \wedge \left(\bigwedge_{1 \leq i_{1}, i_{2} \leq n} y_{i_{1}} - y_{i_{2}} \leq \Delta_{i_{1}, i_{2}}\right)$$
$$\wedge \left(\bigwedge_{1 \leq i \leq n, 1 \leq j \leq m} m_{i} - \overline{x}_{j} + \delta_{i, j} \leq y_{i} - x_{j} \leq M_{i} - \underline{x}_{j} - \delta_{i, j}\right),$$

where, for all $i, i_1, i_2 \in [1, n]$ and $j \in [1, m]$:

$$\begin{split} m_i &= \sum_{w_{i,j} < 0} w_{i,j} \overline{x}_j + \sum_{w_{i,j} > 0} w_{i,j} \underline{x}_j + b_i, \\ M_i &= \sum_{w_{i,j} < 0} w_{i,j} \underline{x}_j + \sum_{w_{i,j} > 0} w_{i,j} \overline{x}_j + b_i, \\ \Delta_{i_1,i_2} &= \sum_{w_{i_1,j} < w_{i_2,j}} (w_{i_1,j} - w_{i_2,j}) \underline{x}_j + \sum_{w_{i_1,j} > w_{i_2,j}} (w_{i_1,j} - w_{i_2,j}) \overline{x}_j + (b_{i_1} - b_{i_2}), \\ \delta_{i,j} &= \begin{cases} 0, & \text{if } w_{i,j} \leq 0 \\ w_{i,j} (\overline{x}_j - \underline{x}_j), & \text{if } 0 \leq w_{i,j} \leq 1 \\ (\overline{x}_j - \underline{x}_j), & \text{if } 1 \leq w_{i,j} \end{cases} \end{split}$$

The tightest zone is obtained as the conjunction of the bounds $\underline{x}_j \leq x_j \leq \overline{x}_j$ on input x, given as hypercube K, the bounds on the y_i and $y_{i_1} - y_{i_2}$ obtained by a direct computation of bounds of the affine transform of the input hypercube K, and finally the bounds on the differences $y_i - x_j$ given by a direct calculation.

Figure 4 shows the three different types of zones that over-approximate the range of a scalar function f, with $f(x) = \lambda x + b$, on an interval. When $\lambda < 0$, the best that can be done is to abstract the graph of f by a square, we cannot encode any dependency between f(x) and x: this corresponds to the case $\delta_{i,j} = 0$ in Proposition 3. The two other cases for the definition of $\delta_{i,j}$ are the two remaining cases of Fig. 4: when λ is between 0 and 1, this is the picture in the middle, and when λ is greater than 1, this is the picture at the right hand side. As we have seen in Proposition 1 and as we will see more in detail below in Theorem 1, these zones can be encoded as tropical polyhedra. Only the points A, B and C are extreme points: D is not an extreme point of the polyhedron as it is on the tropical segment [AC] (the blue, green and red dashed lines each represent a tropical segment).

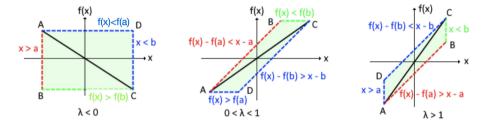


Fig. 4. The 3 cases for approximating the graph of an affine scalar function by a tropical polyhedron, on domain [a, b].

For $f: \mathbb{R}^2 \to \mathbb{R}$, there are 6 cases, depending on the values of λ_1 and λ_2 . In all cases, these zones can be represented as tropical polyhedra using only 4 extreme points and 4 inequalities (instead of 8 and 6 in the classical case), as we will see in Theorem 1. Figure 5 represents the resulting polyhedron for different values of λ_1 and λ_2 . Each figure shows the extreme points A, B_1, B_2 and C, the faces of the polyhedron (in green), the tropical segments inside the polyhedron (in red), and the actual graph of f(x) (in blue). We have the corresponding external description in Theorem 1 below:

Theorem 1. The best zone abstraction \mathcal{H}_f of of the graph $\mathcal{G}_f = \{(x_1, \ldots, x_m, y_1, \ldots, y_n) \mid \underline{x}_j \leq x_j \leq \overline{x}_j, \ y_i = f_i(x_1, \ldots, x_m)\} \subseteq \mathbb{R}^{+n}$ of the linear function $f: \mathbb{R}^m \to \mathbb{R}^n$ defined in Proposition 3 can be seen as the tropical polyhedron defined externally with m + n + 1 inequalities, for all $i \in [1, n]$ and $j \in [1, m]$:

$$\max(x_1 - \overline{x}_1, \dots, x_m - \overline{x}_m, y_1 - M_1, \dots, y_n - M_n) \le 0$$

$$\max(0, y_1 - M_1 + \delta_{1,j}, \dots, y_n - M_n + \delta_{n,j}) \le x_j - \underline{x}_j (3)$$

$$\max(0, x_1 - \overline{x}_1 + \delta_{i,1}, \dots, x_n - \overline{x}_n + \delta_{i,n}, y_1 - d_{i,1}, \dots, y_n - d_{i,n}) \le y_i - m_i(4)$$

where d_{j_1,j_2} denotes the quantity $\Delta_{j_1,j_2} + m_{j_2}$ for i_1 and i_2 in [1,n].

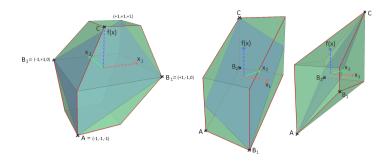


Fig. 5. Over-approximation for $\lambda_1 = \lambda_2 = 0.5$ (left), $\lambda_1 = -0.5$ and $\lambda_2 = 1.5$ (middle), and $\lambda_1 = \lambda_2 = 1.2$ (right).

We have the matching internal representation in Theorem 2:

Theorem 2. \mathcal{H}_f can also be described, internally, as the tropical convex hull of m+n+1 extreme points:

$$A = (\underline{x}_{1}, \dots, \underline{x}_{m}, m_{1}, \dots, m_{n})$$

$$B_{1} = (\overline{x}_{1}, \underline{x}_{2}, \dots, \underline{x}_{m}, m_{1} + \delta_{1,1}, \dots, m_{n} + \delta_{n,1}) \dots$$

$$B_{m} = (\underline{x}_{1}, \dots, \underline{x}_{m-1}, \overline{x}_{m}, m_{1} + \delta_{1,m}, \dots, m_{n} + \delta_{n,m})$$

$$C_{1} = (\underline{x}_{1} + \delta_{1,1}, \dots, \underline{x}_{m} + \delta_{1,m}, M_{1}, c_{1,2}, \dots, c_{1,n}) \dots$$

$$C_{n} = (\underline{x}_{1} + \delta_{n,1}, \dots, \underline{x}_{m} + \delta_{n,m}, c_{n,1}, \dots, c_{n,n-1}, M_{n})$$

where $c_{i_1,i_2} = M_{i_1} - \Delta_{i_1,i_2}$ for i_1 and i_2 in [1,n].

Example 7 (Running example). Let us detail the computations for Example 1: $h_1 = x_1 - x_2 - 1$, $h_2 = x_1 + x_2 + 1$. We have respectively, $\delta_{1,1} = 2$, $\delta_{1,2} = 0$, $\delta_{2,1}=2,\ \delta_{2,2}=2,\ \Delta_{1,1}=0,\ \Delta_{1,2}=0,\ \Delta_{2,1}=4,\ \Delta_{2,2}=0,\ d_{1,1}=-3,\ d_{1,2}=-1,$ $d_{2,1} = 1, d_{2,2} = -1, m_1 = -3, m_2 = -1, M_1 = 1$ and $M_2 = 3$. Hence the external description for the tropical polyhedron relating values of x_1, x_2, h_1 and h_2 are: $\max(x_1 - 1, x_2 - 1, h_1 - 1, h_2 - 3) \le 0$, $\max(0, h_1 + 1, h_2 - 1) \le 0$ $x_1 + 1$, $\max(0, h_1 - 1, h_2 - 1) \le x_2 + 1$, $\max(0, x_1 + 1, x_2 - 1, h_1 + 3, h_2 - 1) \le x_1 + 1$ $h_1 + 3$, $\max(0, x_1 + 1, x_2 + 1, h_1 + 1, h_2 + 1) \le h_2 + 1$ which encode all zones inequalities: $-1 \le x_1 \le 1$, $-1 \le x_2 \le 1$, $-3 \le h_1 \le 1$, $-1 \le h_2 \le 3$, $-2 \le 1$ $h_1 - x_1 \le 0, -4 \le h_1 - x_2 \le 2, 0 \le h_2 - x_1 \le 2, 0 \le h_2 - x_2 \le 2, -4 \le h_1 - h_2 \le 2$ 0. Note that the zone abstraction of [29] would be equivalent to an interval abstraction and would not infer the relations between h_1, h_2, x_1 and x_2 . Now the internal representation of the corresponding zone is $A = (-1, -1, -3, -1), B_1 =$ $(1,-1,-1,1), B_2=(-1,1,-3,1), C_1=(-1,-1,1,1), C_2=(-1,1,-1,3).$ The projections of these 5 extreme points on (h_1, h_2) give the points (-3, -1), (-1, 1),(-3,1), (1,1), (-1,3),among which (-3,1) and (-1,1) are in the tropical convex hull of A = (-3, -1), $B_1 = (1, 1)$ and $B_2 = (-1, 3)$ represented in Fig. 3a. Indeed (-3,1) is on the tropical line (AB_2) and (-1,1) whereas (-1,1) is on the tropical line (AB_1) as a tropical linear combination of $-2+B_1$ and $-2+B_2$: (-1,1) = max(-2+(1,1),-2+(-1,3)).

Example 8. Consider now function $f: \mathbb{R}^2 \to \mathbb{R}^2$ with $f(x_1, x_2) = (0.9x_1 + 1.1x_2, y_2 = 1.1x_1 - 0.9x_2)$ on $(x_1, x_2) \in [-1, 1]$. We have in particular $M_1 = 2$, $M_2 = 2$, $m_1 = -2$ and $m_2 = -2$. We compute $\delta_{1,1} = 1.8$, $\delta_{1,2} = 2$, $\delta_{2,1} = 2$ and $\delta_{2,2} = 0$ and we have indeed $y_1 + 2 \ge x_1 - 1 + \delta_{1,1} = x_1 + 0.8$, $y_2 + 2 \ge x_1 - 1 + 2 = x_1 + 1$, $y_1 + 2 \ge x_2 - 1 + \delta_{2,1} = x_1 + 1$, $y_2 + 2 \ge x_2 - 1$ and $y_1 - 2 \le x_1 + 1 - 1.8 = x_1 - 0.8$, $y_2 - 2 \le x_1 + 1 - 2 = x_1 - 1$, $y_1 - 2 \le x_2 + 1 - 2 = x_2 - 1$, $y_2 - 2 \le x_2 + 1$. Overall:

$$x_1 - 1.2 \le y_1 \le x_1 + 1.2$$

$$x_2 - 1 \le y_1 \le x_2 + 1$$

$$x_1 - 1 \le y_2 \le x_1 + 1$$

$$x_2 - 3 \le y_2 \le x_2 + 3$$

We also find $d_{1,1}=-2$, $d_{1,2}=0.2$, $d_{2,1}=0.2$ and $d_{2,2}=-2$. Hence $y_1-d_{1,2}\leq y_2-m_2$, i.e. $y_1-0.2\leq y_2+2$ that is $y_1-y_2\leq 2.2$. Similarly, we find $y_2-y_1\leq 2+0.2$ hence $-2.2\leq y_1-y_2\leq 2.2$.

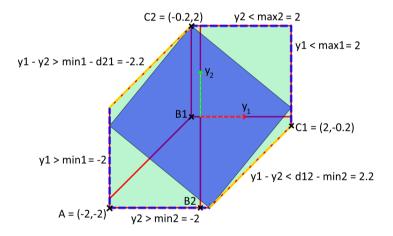


Fig. 6. Over-approximation for $f(x_1, x_2) = (0.9x_1 + 1.1x_2, y_2 = 1.1x_1 - 0.9x_2)$.

These equations can be written as linear tropical constraints as in Theorem 1:

$$\max \begin{pmatrix} x_1 - 1 \\ x_2 - 1 \\ y_1 - 2 \\ y_2 - 2 \end{pmatrix} \le 0, \, \max \begin{pmatrix} 0 \\ y_1 - 0.2 \\ y_2 \end{pmatrix} \le x_1 + 1, \, \max \begin{pmatrix} 0 \\ y_1 \\ y_2 - 2 \end{pmatrix} \le x_2 + 1$$

$$\max \begin{pmatrix} 0 \\ x_1 + 0.8 \\ x_2 + 1 \\ y_1 + 2 \\ y_2 - 0.2 \end{pmatrix} \le y_1 + 2, \max \begin{pmatrix} 0 \\ x_1 + 1 \\ x_2 - 1 \\ y_1 - 0.2 \\ y_2 + 2 \end{pmatrix} \le y_2 + 2$$

We now depict in Fig. 6 both the image of f as a blue rotated central square, and its over-approximation by the convex tropical polyhedron calculated as in Theorem 1 in green, in the plane (y_1, y_2) . As $c_{1,1} = 2$, $c_{1,2} = -0.2$, $c_{2,1} = -0.2$ and $c_{2,2} = 2$, the extremal points are, in the (x_1, x_2, y_1, y_2) coordinates:

$$A = \begin{pmatrix} -1 \\ -1 \\ -2 \\ -2 \end{pmatrix} B_1 = \begin{pmatrix} 1 \\ -1 \\ -0.2 \\ 0 \end{pmatrix} B_2 = \begin{pmatrix} -1 \\ 1 \\ 0 \\ -2 \end{pmatrix} \ C_1 = \begin{pmatrix} 0.8 \\ 1 \\ 2 \\ -0.2 \end{pmatrix} C_2 = \begin{pmatrix} 1 \\ -1 \\ -0.2 \\ 2 \end{pmatrix}$$

3.2 Octagon Abstractions and (max, +, -) Algebra

As in Sect. 3.1, we consider the abstraction of the image of an hypercube K of \mathbb{R}^m by an affine transformation $f: \mathbb{R}^m \to \mathbb{R}^n$ defined, for all $x \in \mathbb{R}^m$ and $i \in [1, n]$, by $(f(x))_i = \sum_{j=1}^m w_{i,j} x_j + b_i$. But we consider here the abstraction of this image by an octagon, we will thus add some constraints on sums of variables to the abstraction computed in Sect. 3.1.

Proposition 4 (Optimal approximation of a linear layer by an octagon). Let $K \subset \mathbb{R}^m$ be an hypercube defined as $K = \prod_j [\underline{x}_j, \overline{x}_j]$, with $\underline{x}_j, \overline{x}_j \in \mathbb{R}$. The tightest octagon of $\mathbb{R}^m \times \mathbb{R}^n$ containing $S := \{(x, f(x)) \mid x \in K\}$ is the set of all $(x, y) \in \mathbb{R}^m \times \mathbb{R}^n$ satisfying

$$\left(\bigwedge_{1 \leq j \leq m} \underline{x}_j \leq x_j \leq \overline{x}_j \right) \wedge \left(\bigwedge_{1 \leq i \leq n} m_i \leq y_i \leq M_i \right) \wedge \left(\bigwedge_{1 \leq i_1, i_2 \leq m} y_{i_1} - y_{i_2} \leq \Delta_{i_1, i_2} \right)$$

$$\wedge \left(\bigwedge_{1 \leq i_1, i_2 \leq n} L_{i_1, i_2} \leq y_{i_1} + y_{i_2} \leq \Gamma_{i_1, i_2} \right)$$

$$\wedge \left(\bigwedge_{1 \leq i \leq n, 1 \leq j \leq m} m_i - \overline{x}_j + \delta_{i, j} \leq y_i - x_j \leq M_i - \underline{x}_j - \delta_{i, j} \right)$$

$$\wedge \left(\bigwedge_{i, j} m_i + \underline{x}_j + \gamma_{i, j} \leq y_i + x_j \leq M_i + \overline{x}_j - \gamma_{i, j} \right)$$

where $m_i, M_i, \delta_{i,j}, \Delta_{i_1,i_2}$ are defined as in Proposition 3, and

$$\begin{split} &\Gamma_{i_1,i_2} := \sum_{w_{i_1,j} + w_{i_2,j} < 0} \underline{x}_j(w_{i_1,j} + w_{i_2,j}) + \sum_{w_{i_1,j} + w_{i_2,j} > 0} \overline{x}_j(w_{i_1,j} + w_{i_2,j}) \\ &L_{i_1,i_2} := \sum_{w_{i_1,j} + w_{i_2,j} < 0} \overline{x}_j(w_{i_1,j} + w_{i_2,j}) + \sum_{w_{i_1,j} + w_{i_2,j} > 0} \underline{x}_j(w_{i_1,j} + w_{i_2,j}) \\ &\gamma_{i,j} := \begin{cases} 0, & \text{if } 0 \le w_{i,j} \\ -w_{i,j}(\overline{x}_j - \underline{x}_j), & \text{if } -1 \le w_{i,j} \le 0 \\ (\overline{x}_j - \underline{x}_j), & \text{if } w_{i,j} \le -1 \end{cases} \end{split}$$

With the notations of Proposition 4, we have

Proposition 5. Let M be the (classically) linear manifold in $\mathbb{R}^m \times \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^n$ defined by $(x^+, y^+, x^-, y^-) \in M$ if and only if $x^+ + x^- = 0$ and $y^+ + y^- = 0$. The octagon S defined in Proposition 4 is equal to the intersection of M with the tropical convex polyhedron generated by the 1 + 2n + 2m points $A, B_1^+, \ldots, B_m^+, B_1^-, \ldots, B_m^-, C_1^+, \ldots, C_n^+, C_1^-, \ldots, C_n^-$, where

$$A = (\underline{x}_1, \dots, \underline{x}_m, m_1, \dots, m_n, -\overline{x}_1, \dots, -\overline{x}_m, -M_1, \dots, -M_n)$$

$$\begin{split} B_k^+ &= (0, x^+, y^+, x^-, y^-) \ with & x_k^+ = \overline{x}_k, & x_{j \neq k}^+ = \underline{x}_j, & y_i^+ = m_i + \delta_{i,k} \\ & x_k^- = -\overline{x}_k, & x_{j \neq k}^- = -\overline{x}_j, & y_i^- = -M_i + \gamma_{i,k} \\ B_k^- &= (0, x^+, y^+, x^-, y^-) \ with & x_k^- = -\underline{x}_k, & x_{j \neq k}^- = -\overline{x}_j, & y_i^- = -M_i + \delta_{i,k} \\ & x_k^+ = \underline{x}_k, & x_{j \neq k}^+ = \underline{x}_j, & y_i^+ = m_i + \gamma_{i,k} \\ C_l^+ &= (0, x^+, y^+, x^-, y^-) \ with & y_l^+ = M_l, & y_{i \neq l}^+ = M_l - \Delta_{l,i}, & x_j^+ = \underline{x}_j + \delta_{l,j} \\ & y_l^- = -M_l, & y_{i \neq l}^- = M_l - \Gamma_{l,i}, & x_j^- = -\overline{x}_j + \gamma_{l,j} \\ C_l^- &= (0, x^+, y^+, x^-, y^-) \ with & y_l^- = -m_l, & y_{i \neq l}^- = -m_l - \Delta_{i,l}, & x_j^- = -\overline{x}_j + \delta_{l,j} \\ & y_l^+ = m_l, & y_{i \neq l}^+ = -m_l + L_{l,i}, & x_j^+ = \underline{x}_j + \gamma_{l,j} \end{split}$$

Example 9 (Running example). For the example network of Example 1, the formulas of Proposition 4 give the constraints:

$$\begin{array}{lll} -1 \leq x_1 \leq 1 & -1 \leq x_2 \leq 1 \\ 0 \leq x_1 - h_1 \leq 2 & -2 \leq x_2 - h_1 \leq 4 \\ -4 \leq x_1 + h_1 \leq 2 & -2 \leq x_2 + h_1 \leq 0 \\ -2 \leq x_1 - h_2 \leq 0 & -2 \leq x_2 - h_2 \leq 0 \\ -2 \leq x_1 + h_2 \leq 4 & -2 \leq x_2 + h_2 \leq 4 \end{array} \quad \begin{array}{ll} -3 \leq h_1 \leq 1 \\ 0 \leq h_2 - h_1 \leq 4 \\ -2 \leq h_2 + h_1 \leq 2 \\ -1 \leq h_2 \leq 3 \end{array}$$

And the internal description is given by Proposition 5, with the following extreme points, where coordinates are ordered as $(x_1^+, x_2^+, h_1^+, h_2^+, x_1^-, x_2^-, h_1^-, h_2^-)$:

$$(-1,-1,-3,-1,-1,-1,-1,-3)\\ (1,-1,-1,1,-1,-1,-1,-3)\\ (-1,1,-3,1,-1,-1,1,-3)\\ (-1,-1,-3,-1,-1,-1,-1,-1,-3)\\ (1,-1,1,1,-1,1,-1,-1)\\ (1,1,-1,3,-1,-1,1,-3)\\ (-1,-1,-3,-1,1,-1,1,-1)\\ (-1,-1,-1,-1,-1,1,-1,-1)\\ (-1,-1,-3,-1,-1,-1,-1,-1,-1)\\ (-1,1,-3,1,1,-1,3,-1)\\ (-1,1,-3,1,1,-1,3,-1)\\ (-1,-1,-1,-1,1,1,1)$$

From the extremal points of the octagon abstraction above, we get the extremal points for (h_1^+, h_2^+) , discarding the non extremal ones: (-3, -1), (1, 1) and (-1, 3), and for (h_1^+, h_2^-) : (-3, -3), (1, -1) and (-1, 1) (for this last pair of variables, this gives the zone in cyan of Example 3).

4 Validation of Multi-layered Neural Networks

General Algorithm. The method developed in Sect. 3 is the cornerstone of our algorithm for analysing neural networks. A ReLU neural net consists of a chain of two kinds of computations, one which applies a classical linear transformation to their inputs, and another one one which applies a ReLU. function We have seen that the affine map transformation can be over-approximated using tropical polyhedra. ReLU being a tropical affine function, the ReLU transform is exact in tropical polyhedra. It is thus possible to use tropical polyhedra to represent reachable states for every node in the network, at least for one layer ReLU networks.

Example 10. We carry on with Example 1 and complete the final computations of Example 7. The external representation is given by the tropical linear inequalities of Example 7 together with inequalities $max(0, h_1) \leq y_1 \leq max(0, h_1)$ and $max(0, h_2) \leq y_2 \leq max(0, h_2)$. Now the corresponding tropical polyhedron is generated by the linear tropical operator ReLU on each of the extremal points A, B_1 , B_2 , C_1 and C_2 and gives the two extra (last) coordinates in the axes $(x_1, x_2, h_1, h_2, y_1, y_2)$, A' = (-1, -1, -3, -1, 0, 0), $B'_1 = (1, -1, -1, 1, 0, 1)$, $B'_2 = (-1, 1, -3, 1, 0, 1)$, $C'_1 = (-1, -1, 1, 1, 1, 1)$, $C'_2 = (-1, 1, -1, 3, 0, 3)$. The projections of theses 5 extreme points on (h_1, y_2) give the points (0, 0), (0, 1), (1, 1), (0, 3) among which (0, 1) is in the convex hull of A' = (0, 0), $B'_2 = B_2 = (1, 1)$ and $B'_1 = (0, 3)$ represented in Fig. 3a.

The polyhedron given by the method of Sect. 3 only gives relations between 2 layers (the input and the first hidden layer). In order to get a polyhedron that represents the whole network when combining with e.g. another layer, we need to embed the first polyhedron from a space that represents only 2 layers to a higher space that represents the complete network, with one dimension per node. We will then need to intersect the polyhedra generated by each pair of layers to get the final result. Finally, as we are only interested in the input-output abstraction of the whole network, we can reduce computing costs by removing the dimensions corresponding to middle layers once those are calculated.

To this end, we use the following notations. Let $\mathcal{L} \subset \{L_0, \ldots, L_N\}$ be a set of layers, layer i containing n_{i+1} neurons as in Definition 2. Let n be the sum of all n_{i+1} , with i such that $L_i \in \mathcal{L}$ and $\mathcal{S}_{\mathcal{L}} \equiv \mathbb{R}^n_{max}$ be the tropical space in which we are going to interpret the values of the neurons on layers in \mathcal{L} , with each dimension of $\mathcal{S}_{\mathcal{L}}$ corresponding to a node of a layer of \mathcal{L} .

For $\mathcal{L}_1, \mathcal{L}_2 \subset \{L_0, \dots, L_N\}$, for $\mathcal{H} \subset \mathcal{S}_{\mathcal{L}_1}$ a tropical polyhedron, we denote by $Proj(\mathcal{H}, \mathcal{L}_2) \subset \mathcal{S}_{\mathcal{L}_2}$ the projection of \mathcal{H} onto $\mathcal{S}_{\mathcal{L}_2}$ when $\mathcal{S}_{\mathcal{L}_2} \subseteq \mathcal{S}_{\mathcal{L}_1}$ and let $Emb(\mathcal{H}, \mathcal{L}_2) \subset \mathcal{S}_{\mathcal{L}_2}$ be the embedding of \mathcal{H} into $\mathcal{S}_{\mathcal{L}_2}$ when $\mathcal{S}_{\mathcal{L}_1} \subseteq \mathcal{S}_{\mathcal{L}_2}$.

The main steps of our algorithm for over-approximating the values of neurons in a multi-layer ReLU network are the following:

- We start with an initial tropical polyhedron $\mathcal{H}_0 \subset \mathcal{S}_{\{L_0\}}$ that represents the interval ranges of the input layer L_0 .
- For each additional layer L_{i+1} :
 - Calculate an enclosing hypercube C_i for the nodes of layer L_i , given the current abstraction $\mathcal{H}_i \subset \mathcal{S}_{\mathcal{L}_i}$ (Sect. 2.3).
 - Calculate the polyhedron \mathcal{P}_{i+1} representing relationships between layer L_i and the new layer L_{i+1} , for nodes of layer L_i taking values in C_i , as described in Sect. 3: Theorem 1 for the external description, and Theorem 2 for the internal description
 - Let $\mathcal{L}'_{i+1} = \mathcal{L}_i \cup \{L_{i+1}\}$. Calculate $\mathcal{P}'_{i+1} = Emb(\mathcal{P}_{i+1}, \mathcal{L}'_{i+1})$ (see below)
 - Intersect \mathcal{P}'_{i+1} with the projection (using the internal description, see below) of the previous abstraction \mathcal{H}_i to get $\mathcal{H}'_{i+1} = Emb(\mathcal{H}_i, \mathcal{L}'_{i+1}) \cap \mathcal{P}'_{i+1}$ (using the external description).
 - Choose $\mathcal{L}_{i+1} \supset \{L_{i+1}\}$, and calculate $\mathcal{H}_{i+1} = Proj(\mathcal{H}'_{i+1}, \mathcal{L}_{i+1})$. Usually, we would use $\mathcal{L}_{i+1} = \{L_0, L_{i+1}\}$ if we only want relations between the input and output layers, or $\mathcal{L}_{i+1} = \{L_0, \ldots, L_{i+1}\}$ if we want relations between every layer.

We need now to describe the projection and embedding functions Proj and Emb. Let $\mathcal{L}_2 \subset \mathcal{L}_1 \subset \{L_0, \ldots, L_N\}$ be two sets of layers. Let \mathcal{H} be a polyhedron on $\mathcal{S}_{\mathcal{L}_1}$. We have $\mathcal{H}' = Proj(\mathcal{H}, \mathcal{L}_2) = \{(x_i)_{L_i \in \mathcal{L}_2}, (x_i)_{L_i \in \mathcal{L}_1} \in \mathcal{H}\}$, i.e. for each point in \mathcal{H} , we only keep the dimensions corresponding to layers in \mathcal{L}_2 , and discard the other dimensions. Projecting is easy with the internal description of polyhedron, as we can project the extreme points of \mathcal{H} to get generators of \mathcal{H}' . However, we do not have a simple algorithm to project the external description of a polyhedron.

Let $\mathcal{L}_1 \subset \mathcal{L}_2 \subset \{L_0, \dots, L_N\}$ be two sets of layers, and Δ be the sum of n_{i+1} , the number of neurons of layer L_i , for i such that $L_i \in \mathcal{L}_2 \setminus \mathcal{L}_1$. Let \mathcal{H} be a polyhedron on $\mathcal{S}_{\mathcal{L}_1}$. We note that $\mathcal{S}_2 \equiv \mathcal{S}_1 \times \mathbb{R}^{\Delta}_{max}$, and thus $\mathcal{H}' = Emb(\mathcal{H}, \mathcal{L}_2) \equiv \mathcal{H} \times \mathbb{R}^{\Delta}_{max}$, i.e. we add dimensions corresponding to each node in \mathcal{L}_2 which are not in \mathcal{L}_1 , and let points in \mathcal{H}' take any value of \mathbb{R}_{max} on these dimensions. Embedding is based on simple matrices concatenations in the external description, for more details. Embedding using the internal description is more involved and is explained after exemplifying things on a simple example.

Example 11. We consider the 1-layer neural net of Example 1, and add a second layer. The new linear layer is defined by $u_1 = y_2 - y_1 - 1$, $u_2 = y_1 - y_2 + 1$ and the output neurons are $z_1 = max(0, u_1) = max(0, y_2 - y_1 - 1)$ and $z_2 = max(0, u_2) = max(0, y_1 - y_2 + 1)$.

The enclosing cube for the tropical polyhedron \mathcal{H} containing the values of neurons of the first layer L_1 : y_1 , y_2 of Example 1 is $[0,1] \times [0,3]$. The analysis of the second layer L_2 , supposing its input belongs to $[0,1] \times [0,3]$ gives the constraint (an extract of the external representation of the resulting tropical polyhedron \mathcal{H}') $-3 \leq u_1 - y_1 \leq 2, -2 \leq u_1 - y_2 \leq -1, -2 \leq 1$

 $u_2 - y_1 \leq 1$, $-5 \leq u_2 - y_2 \leq 2$, $z_1 = max(0, u_1)$, $z_2 = max(0, u_2)$. The intersection of the embedding $Emb(\mathcal{H}', \{L_0, L_1, L_2\})$ with the embedding $Emb(\mathcal{H}, \{L_0, L_1, L_2\})$ consists, as we saw above, in concatenating the tropical constraints, in the common space of variables. This implies in particular that we add the constraint $-3 \leq y_1 - y_2 \leq 0$ to the above equations. The intersection is actually a zone intersection, where we have to normalize the corresponding DBM. A manual calculation shows that this will make use of the equalities $u_2 - y_2 = (u_2 - y_1) + (y_1 - y_2)$, $u_1 - y_1 = (u_1 - y_2) + (y_2 - y_1)$. By combining equations, we get the refined bounds (refined lower bound for the first equation, refined upper bound for the second equation) $-2 \leq u_1 - y_1 \leq 2$, $-5 \leq u_2 - y_2 \leq 1$.

Embedding a Tropical Polyhedron: Internal Description. In this paragraph, we embed a polyhedron into a higher dimensional space, using the internal description.

Suppose \mathcal{H} is a tropical polyhedron in \mathbb{R}^n (such as \mathcal{P}_i in the previous section) that we want to embed \mathcal{H} into a larger space, with an extra coordinate, which we consider bounded here within [a, b]. So we need to determine a presentation of the tropical polyedron $\mathcal{H}' = \mathcal{H} \times [a, b]$.

Supposing we have m extreme points p_i for representing \mathcal{H} , a naive method consists in noticing that the family $(p_i, a), (p_i, b)$ is a generator of \mathcal{H}' and removing non-extreme points from that list. But that would exhibit poor performance, as we get $m \times 2^k$ extreme points for \mathcal{H}'' . We can in fact do better:

Theorem 3. The extreme points of \mathcal{H}' are $\{(p_i, a), 1 \leq i \leq m\} \cup \{(p_i, b), i \in I\}$, where I is a subset of indexes of generators of \mathcal{H} , $I \subset [1, m]$, such that:

$$\forall i \in I, \forall j \in [1, m] \setminus \{i\}, p_i \oplus p_j \neq p_i \tag{5}$$

$$\forall i \in [1, m] \setminus I, \exists j \in [1, m] \setminus \{i\} \ s.t. \ p_i \oplus p_j = p_i$$
 (6)

Passing to the limit, this shows that the extreme points of $\mathcal{H} \times \mathbb{R}$ are $(p_i, -\infty)$, $i = 1, \ldots, m$ and the extreme rays are $(p_i, 0)$, $i \in I$ for the smallest I verifying Eq. (5) and (6). In the current implementation, we do not use extreme rays and embed \mathcal{H} into larger state spaces by using large enough values for a and b.

Checking Properties on ReLU Neural Nets. Given an affine guard

$$h(x,y) = \sum_{i=1}^{m} h_i x_i + \sum_{j=1}^{n} h'_j y_j + c$$

where x_i , resp. y_j are the input, resp. output neurons, we want to determine whether, for all input values in [-1,1], we have $h(x) \geq 0$ (this can encode properties (P_1) and (P_2) of Example 1).

There are two ways to check such properties. The first one, that we have implemented, is as follows. We abstract the input output relation that the network under analysis encodes, using a tropical polyhedron \mathcal{H} as described in Sect. 4. From this, we derive the smallest zone Z containing \mathcal{H} as in Sect. 2.3.

Finally, we solve the linear programming problem $m = \min_{x,y \in Z} h(x,y)$ using any classical algorithm (we used glpk in our prototype). This is enough for checking (P_1) in Example 1 since $m \geq 0$ proves our property true, but not (P_2) . The second way can be useful to check (P_2) : here we have no choice but try to solve $m = \min_{x,y \in \mathcal{H}} h(x,y)$ which is not a convex optimization problem, in any sense (tropical nor classical). This could be encoded as MILP problem instead.

5 Implementation, Experiments and Benchmarks

Internal, External and Double Description Methods. Overall, we have developed methods for propagating an outer-approximation of the values that the different layers of neurons can take, within a MLP with ReLU activation. Let us discuss the pros and cons of using the internal description, external description and double description methods:

- The double description method allows for possibly using subdivisions, propagating values in multiple layers and projecting them onto a subset of interesting neurons (e.g. input and output layers), as well as computing an enclosing zone, for synthesizing classification properties. We have implemented this in a prototype using Polymake [18], whose results we briefly discuss below.
- The internal description allows for analyzing one layer networks, using subdivisions, project onto an interesting subset of neurons, as well as computing an enclosing zone (Sect. 2.3). We have implemented this method in C++ in a standalone prototype, nntrop, that takes as input a Sherlock file [14] describing the one hidden layer neural net to analyze plus a linear formula to be checked, and returns the tropical abstraction of the values that neurons can take, its over-approximation by a zone, and whether the linear specification is satisfied or not.
- The external description allows for analyzing multiple layer networks (see Sect. 4).

The double description method is much more expensive since the translation between the internal and external representations may be quite complex.

Experiments and Benchmarks. We briefly compare the computation times between internal description only and double description in Table 1. For each example, we indicate in the columns # inp. the number of input neurons, # out. the number of output neurons, # hid. the number of hidden layers, # neur. is the total number of neurons (input, output and hidden), t. intern is the time spent for computing the internal representation and t. double for the double description of the tropical polyhedron abstracting the corresponding neural net. Experiments are performed on a simple computer with ArchLinux and a Intel(R) Core(TM) i5-7200U CPU @ 2.50 GHz.

We of course see the influence of a potential exponential complexity for going back and forth between internal and external descriptions, but also the fact that we relied on a perl (interpreted) implementation of tropical polyhedra (the one of polymake [18], with exact rational arithmetics), which is much slower than the C++ implementation we wrote for the internal description method (although the internal description method does work in a twice as big space because it considers the octagon instead of just zone abstraction).

					I	I
Example	# inp.	# out.	# hid.	# neur.	t. intern. (s)	t. double (s)
running	2	2	0	4	0.006	1.83
running2	2	2	1	6	0.011	4.34
multi	2	8	1	13	0.005	3.9
krelu	2	2	0	4	0.011	1.94
$tora_modified_controller$	4	1	1	6	0.005	14.57
$tora_modified_controller_1$	4	1	1	105	0.75	815.12
$quadcopter_trial_controller_3$	18	1	1	49	0.009	102.54
${\tt quadcopter_trial_controller_1}$	18	1	1	69	0.2	469.77
$quad_modified_controller$	18	1	1	20	0.005	14
car_nn_controller_2	4	2	1	506	104.75	_
car_nn_controller_1	4	2	1	506	88.8	_
ex	2	1	5	59	0.195	1682.28

Table 1. Execution times (internal and double description) on sample networks.

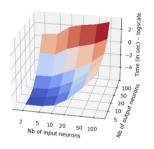
In Table 1, running is the network of Example 1, and running 2 is the extension with an extra layer of Example 11, discussed in great length in these examples. Example krelu is the running example from [35] that we discuss at the end of this section, and tora_modified_controller, tora_modified_controller_1, quadcopter_trial_controller_3, quadcopter_trial_controller_1, quad_modified_controller, car_nn_controller_2, car_nn_controller_1 and ex are examples from the distribution of Sherlock [14]. ex is a multi-layer example for which the algorithm using only the internal representation does not compute the intersection of tropical polyhedra between layers (involving the external representation), contrarily to the double description prototype. We now discuss some of these examples below.

Network multi is a simple 2-layer, 13 neurons example with inputs x_1 , x_2 , outputs y_1, y_2, \ldots, y_8 and

Our zone based abstraction returns the following ranges: $y_1 \in [0,6]$, $y_2 \in [0,4]$, $y_3 \in [0,4]$, $y_4 \in [0,2]$, $y_5 \in [0,4]$, $y_6 \in [0,2]$, $y_5 \in [0,2]$ and $y_8 = 0$, whereas the exact ranges for y_1 to y_7 is [0,2]. Our algorithm is thus exact for y_4 , y_6 , y_7 and y_8 but not y_1 , y_2 , y_3 nor y_5 . This is due to the fact that the zone-based tropical abstraction does represent faithfully the differences of neuron values, but not sums in particular. For instance, $y_2 = max(0, 2x_1)$ which cannot be represented exactly by our method.

Network krelu is a 2 layer 4 neurons example from [35]. We get the correct bounds on the outputs: $0 \le z_1, z_2 \le 2$, as well as relations between the inputs and the outputs: $z_j \le x_i + 1$. However, we do not have significant relations between z_1 and z_2 , as those are not tropically linear. We refer to the results obtained with 1-ReLU and 2-ReLU in [35]: they both get better relations between z_1 and z_2 , in particular $z_1 + z_2 \le 2$ which is not representable in a tropical manner (except by using an octagon based abstraction, which is outside the scope of this paper). However 1-ReLU does not keep track of relations between the inputs and the outputs, and has sub-optimal relations between the outputs, as it cannot represent the non linear ReLU function exactly. 2-ReLU, on the other hand gets both the relation between the output variables, and between the inputs and outputs correct, but is more computationally expensive.

In order to assess the efficiency of the internal description methods, we have run a number of experiments, with various number of inputs and ouputs for neural nets with one hidden layer only. The linear layers are generated randomly, with weights between -2 and 2. Timings are shown in the figure on the right (demonstrating the expected complexity, cubical in the number of neurons), where the x-axis is number of input neurons, y-axis is the number of output neurons, and z-axis is time. For 100 inputs and 100



neurons in the hidden layer, the full pipeline (checking the linear specification in particular) took about $35\,\mathrm{s}$, among which the tropical polyhedron analysis took $6\,\mathrm{s}$.

6 Conclusion and Future Work

We have explored the use of tropical polyhedra as a way to circumvent the combinatorial complexity of neural networks with ReLU activation function. The first experiments we made show that our approximations are tractable when we are able to use either the internal or the external representations for tropical polyhedra, and not both at the same time. This is akin to the results obtained in the classical polyhedron approach, where most of the time, only a sub polyhedral domain is implemented, needing only one of the two kinds of representations. It is interesting to notice that a recent paper explores the use of octohedral constraints, a three-dimensional counterpart of our octagonal representations, in the search of more tractable yet efficient abstraction for ReLU neural nets

[31]. This work is a first step towards a hierarchy of approximations for ReLU MLPs. We have been approximating the tropical rational functions that these neural nets compute by tropical affine functions, and the natural continuation of this work is to go for higher-order approximants, in the tropical world. We also believe that the tropical approach to abstracting ReLU neural networks would be particularly well suited to verification of ternary nets [27]. These ternary nets have gained importance, in particular in embedded systems: simpler weights mean smaller memory needs and faster evaluation, and it has been observed [1] that they can provide similar performance to general networks.

References

- Alemdar, H., Caldwell, N., Leroy, V., Prost-Boucle, A., Pétrot, F.: Ternary neural networks for resource-efficient AI applications. CoRR abs/1609.00222 (2016)
- Allamigeon, X., Gaubert, S., Goubault, E.: The tropical double description method.
 In: 27th International Symposium on Theoretical Aspects of Computer Science, STACS 2010 (2010)
- 3. Allamigeon, X.: Static analysis of memory manipulations by abstract interpretation Algorithmics of tropical polyhedra, and application to abstract interpretation. Ph.D. thesis, École Polytechnique, Palaiseau, France (2009). https://tel.archives-ouvertes.fr/pastel-00005850
- Allamigeon, X., Gaubert, S., Goubault, É.: Inferring min and max invariants using max-plus polyhedra. In: Alpuente, M., Vidal, G. (eds.) SAS 2008. LNCS, vol. 5079, pp. 189–204. Springer, Heidelberg (2008). https://doi.org/10.1007/978-3-540-69166-2_13
- Bak, S., Tran, H.-D., Hobbs, K., Johnson, T.T.: Improved geometric path enumeration for verifying ReLU neural networks. In: Lahiri, S.K., Wang, C. (eds.) CAV 2020. LNCS, vol. 12224, pp. 66–96. Springer, Cham (2020). https://doi.org/10.1007/978-3-030-53288-8-4
- Bastani, O., Ioannou, Y., Lampropoulos, L., Vytiniotis, D., Nori, A., Criminisi, A.: Measuring neural net robustness with constraints. In: Advances in Neural Information Processing Systems (NIPS) (2016)
- Blanchet, B., et al.: A static analyzer for large safety-critical software. In: PLDI, pp. 196–207. ACM Press, June 2003
- Botoeva, E., Kouvaros, P., Kronqvist, J., Lomuscio, A., Misener, R.: Efficient verification of relu-based neural networks via dependency analysis. In: Proceedings of the AAAI Conference on Artificial Intelligence, vol. 34, no. 04, pp. 3291–3299 (2020). https://vas.doc.ic.ac.uk/software/neural/
- 9. Bourdoncle, F.: Abstract interpretation by dynamic partitioning. J. Func. Program. **2**(4), 407–435 (1992)
- Cousot, P., Cousot, R.: Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints. In: Conference Record of the Fourth ACM Symposium on Principles of Programming Languages, Los Angeles, California, USA, pp. 238–252, January 1977
- Cousot, P., Cousot, R.: Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints. In: POPL. ACM (1977)

- Cousot, P., Halbwachs, N.: Automatic discovery of linear restraints among variables of a program. In: Proceedings of the 5th ACM SIGACT-SIGPLAN Symposium on Principles of Programming Languages, pp. 84–96. POPL 1978, Association for Computing Machinery, New York, NY, USA (1978). https://doi.org/10.1145/512760.512770
- 13. Dutta, S., Chen, X., Sankaranarayanan, S.: Reachability analysis for neural feedback systems using regressive polynomial rule inference. In: HSCC (2019)
- Dutta, S., Chen, X., Jha, S., Sankaranarayanan, S., Tiwari, A.: Sherlock A tool for verification of neural network feedback systems: demo abstract. In: HSCC (2019)
- 15. Ehlers, R.: Formal verification of piece-wise linear feed-forward neural networks. In: ATVA (2017)
- 16. Evtimov, I., et al..: Robust physical-world attacks on machine learning models. CoRR abs/1707.08945 (2017). http://arxiv.org/abs/1707.08945
- Gaubert, S., Katz, R.: The Minkowski theorem for max-plus convex sets. Linear Algebra Appl. 421, 356–369 (2006)
- Gawrilow, E., Joswig, M.: polymake: a Framework for Analyzing Convex Polytopes, pp. 43–73. Birkhäuser Basel (2000)
- Gehr, T., Mirman, M., Drachsler-Cohen, D., Tsankov, P., Chaudhuri, S., Vechev, M.: AI2: safety and robustness certification of neural networks with abstract interpretation. In: Conférence IEEE S&P 2018 (2018)
- Gowal, S., et al.: On the effectiveness of interval bound propagation for training verifiably robust models. CoRR abs/1810.12715 (2018). http://arxiv.org/abs/1810. 12715
- Henriksen, P., Lomuscio, A.R.: Efficient neural network verification via adaptive refinement and adversarial search. In: ECAI. Frontiers in Artificial Intelligence and Applications, vol. 325 (2020)
- 22. Huang, X., Kwiatkowska, M., Wang, S., Wu, M.: Safety verification of deep neural networks. In: Majumdar, R., Kunčak, V. (eds.) CAV 2017. LNCS, vol. 10426, pp. 3–29. Springer, Cham (2017). https://doi.org/10.1007/978-3-319-63387-9_1
- Julian, K., Kochenderfer, M.J., Owen, M.P.: Deep neural network compression for aircraft collision avoidance systems. AIAA J. Guidance Control Dyn. (2018). https://arxiv.org/pdf/1810.04240.pdf
- Katz, G., Barrett, C., Dill, D.L., Julian, K., Kochenderfer, M.J.: Reluplex: An efficient SMT solver for verifying deep neural networks. In: CAV 2017 (2017)
- Katz, G., et al.: The Marabou framework for verification and analysis of deep neural networks. In: Dillig, I., Tasiran, S. (eds.) CAV 2019. LNCS, vol. 11561, pp. 443–452. Springer, Cham (2019). https://doi.org/10.1007/978-3-030-25540-4_26
- 26. Khedr, H., Ferlez, J., Shoukry, Y.: Effective formal verification of neural networks using the geometry of linear regions (2020)
- 27. Li, F., Liu, B.: Ternary weight networks. CoRR abs/1605.04711 (2016)
- Mauborgne, L., Rival, X.: Trace partitioning in abstract interpretation based static analyzers. In: Programming Languages and Systems (2005)
- Miné, A.: A new numerical abstract domain based on difference-bound matrices. In: Danvy, O., Filinski, A. (eds.) PADO 2001. LNCS, vol. 2053, pp. 155–172. Springer, Heidelberg (2001). https://doi.org/10.1007/3-540-44978-7_10
- Miné, A.: The octagon abstract domain. High. Order Symb. Comput. 19(1), 31–100 (2006)
- 31. Müller, M.N., Makarchuk, G., Singh, G., Püschel, M., Vechev, M.: Precise multineuron abstractions for neural network certification (2021)
- 32. Ruan, W., Huang, X., Kwiatkowska, M.: Reachability analysis of deep neural networks with provable guarantees. In: IJCAI (2018)

- 33. Shiqi, W., Pei, K., Justin, W., Yang, J., Jana, S.: Efficient formal safety analysis of neural networks. In: NIPS (2018)
- 34. Singh, G., Gehr, T., Mirman, M., Püschel, M., Vechev, M.: Fast and effective robustness certification
- Singh, G., Ganvir, R., Püschel, M., Vechev, M.: Beyond the single neuron convex barrier for neural network certification. In: Advances in Neural Information Processing Systems (NeurIPS) (2019)
- Singh, G., Gehr, T., Püschel, M., Vechev, M.: An abstract domain for certifying neural networks. In: Proceedings ACM Programming Language 3(POPL), January 2019
- 37. Singh, G., Gehr, T., Püschel, M., Vechev, M.: Boosting robustness certification of neural networks. In: ICLR (2019)
- 38. Szegedy, C., et al.: Intriguing properties of neural networks (2013). https://arxiv.org/abs/1312.6199
- 39. Tjeng, V., Xiao, K., Tedrake, R.: Evaluating robustness of neural networks with mixed integer programming. https://arxiv.org/abs/1711.07356
- 40. Wang, S., Pei, K., Whitehouse, J., Yang, J., Jana, S.: Formal security analysis of neural networks using symbolic intervals. USENIX Security (2018)
- 41. Zhang, H., et al.: Towards stable and efficient training of verifiably robust neural networks. In: ICLR (2020)
- Zhang, L., G.Naitzat, Lim, L.H.: Tropical geometry of deep neural networks. In: Proceedings of the 35th International Conference on Machine Learning, vol. 80, pp. 5824–5832. PMLR (2018)