Reliable and Interpretable Artificial Intelligence

Lecture 6: The Zonotope convex relaxation

Martin Vechev

ETH Zurich

Fall 2020





The problem with Box

Box is simple and efficient, but the problem is that it looses too much precision, in both, the ReLU abstract transformer but also in the affine abstract transfomer.

Zonotope will be a convex relaxation whose affine transformer will not lose precision (that is, it will be exact); however, its ReLU transformer will again lose some precision. Zonotope has been shown effective in both verification and provable training (discussed later).

Incomplete method II: Zonotope

With the Zonotope relaxation, each variable (neuron) is

captured in an affine form. We can think of Zonotope as a

more extensive version of Box. Unlike Box, the variables can be

related (in limited ways) through parameters.

The Zonotope Convex Relaxation

For d (concrete) neurons x_1 to x_d , the abstract neurons will look like:



 ϵ_i : noise terms ranging [-1,1] shared between abstract neurons

 a_i^{j} : real number that controls magnitude of noise

This represents a polytope centered around the vector $a_0 = (a_0^1, ..., a_0^d)$

Example of a concretization in 2D (d = 2) with k = 3:



Centering

We can flip any point *X* in the zonotope around the center a_0 so to obtain a flipped point *Y* of *X*, where $Y = 2 * a_0 - X$ and *Y* is in the zonotope and *X* and *Y* are equal distance from a_0 . That is, the zonotope is point-symmetric around a_0 .

Example:

The zonotope **Z** is centered around $a_0 = (1,0)$. The point X = (2,-1) can be flipped to obtain $Y = 2 * a_0 - X = (0,1)$ $Z = \begin{array}{c} \hat{x}_0 = 1 - 2\epsilon_1 + \epsilon_2 \\ \hat{x}_1 = 0 + \epsilon_1 + \epsilon_2 \end{array}$ is visualized: x = (2,-1)is visualized: x = (2,-1)

The 4 extreme values $\epsilon_1 \in \{-1,1\}$ and $\epsilon_2 \in \{-1,1\}$ determine the 4 corners.

-2

We have formally defined what a Zonotope is. However,

we need to also show how to transform a zonotope,

specifically with the operations found in neural networks: affine layers and ReLU.

(side note: for provable training, we need to also know how to project to a zonotope).

Zonotope Affine Transformer



.....

Note: the affine transformer can be computed in parallel for all output neurons $x_1 \dots x_g$

Zonotope Affine Transformer

Multiplication of the p'th neuron by a real-valued constant C:

$$(a_0^p + \sum_{i=1}^k a_i^p \epsilon_i) * C = (C * a_0^p + \sum_{i=1}^k C * a_i^p \epsilon_i)$$

Adding two neurons *p* and *q* is done component-wise:

$$(a_0^p + \sum_{i=1}^k a_i^p \epsilon_i) + (a_0^q + \sum_{i=1}^k a_i^q \epsilon_i) = (a_0^q + a_0^p) + \sum_{i=1}^k (a_i^p + a_i^q) \epsilon_i$$

Note: abstract transformer for affine is exact

 $\hat{x} = a_0 + \sum_{i=1}^{n} a_i \epsilon_i$ we need to somehow apply $y = \max(0, x)$ to get $\hat{y} = ???$ Given this: Input zonotope Output zonotope at layer i + 1 at layer i ReLU y_1 χ_1 Note: the ReLU transformer is computed in parallel for all output neurons $y_1 \dots y_a$ ReLU y_g **ReLU Transformer** (defined next) Fast and Effective Robustness Certification, NeurIPS'18

Singh et.al.

Given this: $\hat{x} = a_0 + \sum_{i=1}^{\kappa} a_i \epsilon_i$ we need to somehow apply $y = \max(0, x)$ to get $\hat{y} = ???$

Step 1: Obtain box bounds for \hat{x} , by exploring $\epsilon_i = -1$ and $\epsilon_i = 1$ values to compute end points. To compute the lower bound l_x use $\epsilon_i = -1$ if a_i is positive and use $\epsilon_i = 1$ if a_i is negative. To compute the upper bound u_x use $\epsilon_i = 1$ if a_i is positive and use $\epsilon_i = -1$ if a_i is negative.



Given this: $\hat{x} = a_0 + \sum_{i=1}^{\kappa} a_i \epsilon_i$ we need to somehow apply $y = \max(0, x)$ to get $\hat{y} = ???$

Step 2: Check if the bounds are on one side of the plane (easy cases):

 \rightarrow If $u_x \leq 0$ then $\hat{y} = 0$ (strictly negative)

 \rightarrow If $l_x > 0$ then $\hat{y} = \hat{x}$ (strictly positive)

 \rightarrow If $l_x < 0$ and $u_x > 0$ then we get into the ``crossing boundary'' case [discussed next].

Given this: $\hat{x} = a_0 + \sum_{i=1}^{\kappa} a_i \epsilon_i$ we need to somehow apply $y = \max(0, x)$ to get $\hat{y} = ???$

Step 2: The ``crossing-boundary'' case looks as follows (where in green is the ReLU function):



Given this: $\hat{x} = a_0 + \sum_{i=1}^{n} a_i \epsilon_i$ we need to somehow apply $y = \max(0, x)$ to get $\hat{y} = ???$

Step 2: We now compute a Zonotope (yellow) which encloses the green ReLU function.



Given this: $\hat{x} = a_0 + \sum_{i=1}^{n} a_i \epsilon_i$ we need to somehow apply $y = \max(0, x)$ to get $\hat{y} = ???$

Step 2a: Define the first (the bottom) line y_1 .



Given this: $\hat{x} = a_0 + \sum_{i=1}^{\kappa} a_i \epsilon_i$ we need to somehow apply $y = \max(0, x)$ to get $\hat{y} = ???$



Given this: $\hat{x} = a_0 + \sum_{i=1}^{\kappa} a_i \epsilon_i$ we need to somehow apply $y = \max(0, x)$ to get $\hat{y} = ???$

Step 3: We now have the two lines, but we need to capture the yellow zonotope with them...



Given this: $\hat{x} = a_0 + \sum_{i=1}^{n} a_i \epsilon_i$ we need to somehow apply $y = \max(0, x)$ to get $\hat{y} = ???$

Step 3: The yellow zonotope is captured by the following inequality:



Given this: $\hat{x} = a_0 + \sum_{i=1}^{\kappa} a_i \epsilon_i$ we need to somehow apply $y = \max(0, x)$ to get $\hat{y} = ???$

Step 3: The yellow zonotope is captured by the following inequality:

 $\lambda * \hat{x} \leq y(\hat{x}) \leq \lambda * \hat{x} - \lambda * l_{x}$

But this inequality is not in the Zonotope format so we need to make

the transformer produce an equivalent zonotope in the format

Given this: $\hat{x} = a_0 + \sum_{i=1}^{n} a_i \epsilon_i$ we need to somehow apply $y = \max(0, x)$ to get $\hat{y} = ???$

Step 3: The yellow zonotope is captured by the following inequality:

$$\lambda * \hat{x} \le y(\hat{x}) \le \lambda * \hat{x} - \lambda * l_x$$

This is the same as: $v(\hat{x}) = \lambda * \hat{x} - c * \lambda * l_r$ where $c \in [0,1]$

But our error terms $\epsilon \in [-1,1]$ and not [0,1]

where ϵ_{new} is So if we want to use [-1,1] we need to have $c = \frac{\epsilon_{new} + 1}{2}$ a new error term $y(\hat{x}) = \lambda * \hat{x} - \frac{\epsilon_{new} + 1}{2} * \lambda * l_x$ Thus we obtain The final ReLU $y(\hat{x}) = \lambda * \hat{x} - \epsilon_{new} * \frac{\lambda * l_x}{2} - \frac{\lambda * l_x}{2}$ This is in zonotope format transformer is:

Applying the Zonotope ReLU Transformer

Given this: $\hat{x} = a_0 + \sum_{i=1}^{\kappa} a_i \epsilon_i$, to apply the transformer, we first compute λ , l_x , u_x

Plug \hat{x} into transformer definition:

$$y(\hat{x}) = \lambda * \hat{x} - \epsilon_{new} * \frac{\lambda * l_x}{2} - \frac{\lambda * l_x}{2}$$

And obtain:

After doing basic calculations:

$$y(\hat{x}) = \lambda a_0 + \sum_{i=1}^k \lambda a_i \epsilon_i - \epsilon_{new} * \frac{\lambda * l_x}{2} - \frac{\lambda * l_x}{2}$$
$$y(\hat{x}) = b_0 + \sum_{i=1}^{k+1} b_i \epsilon_i$$
Not great news:
Error terms increase as depth of network increases!

This is our final result

Two Non-Comparable Zonotope Transformers



- First transformer is what we looked at and is optimal area-wise (in the input-output plane).
- Second transformer produces bigger area, but is non-comparable to the first.
- There are many non-comparable ReLU transformers for Zonotope.

Zonotope vs. Box ReLU Transformer



- Mathematically, the two transformers produce **non-comparable** results (no shape is included in the other). However, our Zonotope transformer is **optimal in terms of area**.
- Both transformers can process one layer in parallel, so they are GPU-friendly.
- The Box transformer is memory friendly as it does not introduce new error terms.

Zonotope vs. Optimal ReLU Transformer



- The area of our transformer is **2x the area** of the optimal triangle transformer.
- Optimal convex transformer is too expensive to work with (so can only be used selectively).
- Optimal transformer produces optimal convex shape if we consider **one neuron at a time**.
- However, see next slides:... ☺

One-neuron at a time is not optimal

[NeurlPS'19: "Beyond Single Neuron Relaxations for Certification" – Singh, Ganvir, Pueschel, Vechev]



- Contrary to popular belief [see: <u>https://arxiv.org/abs/1902.08722</u>, NeurIPS'19], if we consider 2 neurons at a time, we can produce more precise results.
- Here, we have two neurons x₁ and x₂. The blue figure (b) shows the result for 1-ReLU for triangle, where we compute neuron-wise, that is we compute y₁ from x₁ and y₂ from x₂.
- However, if we compute y₁ and y₂ by considering x₁ and x₂ together, we can get a more precise result, shown in figure (c).
- We will see a lecture on how to compute the result in figure (c) later in the course.

Lets revisit this slide: propagating $oldsymbol{\phi}$



We saw how to compute the intermediate shapes (in Box or Zonotope format).

For example here ϕ_n follows the Zonotope format.

Summary

• We introduced another incomplete method, the Zonotope relaxation.

• We introduced the abstract transformers for affine and ReLU for Zonotope. Both transformers can be computed in parallel for a layer.

• Unlike Box, Zonotope is exact for affine. Both lose precision on ReLUs.

Next time: some discussion on the general theory of relaxations

(what is optimal, what is sound).