# Exercise 08 - Solution

## Geometric robustness

## Reliable and Interpretable Artificial Intelligence
## ETH Zurich

**Problem 1** (Verifying rotation)**.** In this problem we consider verifying robustness of network shown in Fig. 1 to rotation by angle $\theta \in [0.5, 1.5]$. We assume that there are two pixels, one at coordinate $(1, 0)$ whose value after the rotation is $\cos(\theta)$ and another at coordinate $(0, 1)$ whose value after the rotation is $\sin(\theta)$. Vaues of these two pixels are used as inputs to the network in Fig. 1, where their values are denoted as $x_1$ and $x_2$. Our goal here is to certify that $x_5 > -1.6$.

1. Try to prove the claim by using interval/box domain for both trignometric operations and operations in the neural network. Can you prove it?

2. Try to prove the claim by using interval/box domain for trigonometric operations and DeepPoly domain for operations in the network. Can you prove it?

3. Compute the tightest linear lower and upper bound for $\sin(\theta)$ and $\cos(\theta)$. By tightest lower bound, here we mean lower bound with smallest area between the lower bound line and the target function (and analogous for upper bound).

4. Use linear bounds computed in the previous step for trignometric functions together with DeepPoly domain for the network to prove the property.
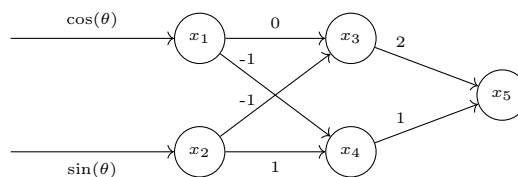


Figure 1: Neural network to be analyzed.

**Solution 1.** 1. We compute the following:

$$\hat{x}_1 = \cos([0.5, 1.5]) = [\cos(1.5), \cos(0.5)] = [0.07, 0.88]$$
$$\hat{x}_2 = \sin([0.5, 1.5]) = [\sin(0.5), \sin(1.5)] = [0.48, 1.0]$$
$$\hat{x}_3 = -\hat{x}_2 = [-1, -0.48]$$
$$\hat{x}_4 = -\hat{x}_1 + \hat{x}_2 = [-0.4, 0.93]$$
$$\hat{x}_5 = 2\hat{x}_3 + \hat{x}_4 = [-2.4, -0.03]$$

Thus, we can not certify the property.

2. We perform backsubstitution:

$$\begin{aligned}
x_5 &= 2x_3 + x_4 \\
&= -2x_2 - x_1 + x_2 \\
&= -x_1 - x_2 \\
&= -\cos(\theta) - \sin(\theta) \\
&\geq -\cos(0.5) - \sin(1.5) \\
&= -1.875.
\end{aligned}$$

Thus, we can not certify the property.

3. Note that $\sin''(\theta) = -\sin(\theta) < 0$ and $\cos''(\theta) = -\cos(\theta) < 0$ for $\theta \in [0.5, 1.5]$ which means that both functions are concave. Thus, tightest lower bound for both functions is line between the two endpoints:

$$\cos(\theta) \geq \cos(0.5) + (\theta - 0.5)(\cos(1.5) - \cos(0.5))$$
$$\sin(\theta) \geq \sin(0.5) + (\theta - 0.5)(\sin(1.5) - \sin(0.5))$$

The tightest upper bound is tangent to the function at some point $t \in [0.5, 1.5]$ . Note that this holds because of the concavity, but not in general.

This bound for $\cos(\theta)$ is of the form $\cos(t) + (\theta - t)(-\sin(t))$. We have to compute $t$ such that the area $A(t) = \int_{0.5}^{1.5}(\cos(t) + (\theta - t)(-\sin(t)) - \cos(\theta))d\theta$ is minimized. Integration yields:

$$A(t) = \cos(t) - \sin(t) + t\sin(t) + F(\theta)$$
$$A'(t) = -\cos(t) + t\cos(t) = (-1 + t)\cos(t)$$

To compute the minimum of $A(t)$ we set $A'(t)$ to 0 and obtain $t = 1$. This means that upper bound is tangent of cosine function at $t = 1$. Computation for $\sin(\theta)$ is analogous and also results in $t = 1$. Thus, the tightest upper bounds are:

$$\cos(\theta) \leq \cos(1) - \sin(1)(\theta - 1)$$
$$\sin(\theta) \leq \sin(1) + \cos(1)(\theta - 1)$$

2

4. We perform backsubstitution:

$$
\begin{aligned}
x_5 &= 2x_3 + x_4 \\
&= -2x_2 - x_1 + x_2 \\
&= -x_1 - x_2 \\
&= -\cos(\theta) - \sin(\theta) \\
&\geq -(\cos(1) + (\theta - 1)(-\sin(1))) - (\sin(1) + (\theta - 1)\cos(1)) \\
&\geq -1.53.
\end{aligned}
$$

Thus, using this method we can finally prove the property.

**Problem 2.** (Bounding functions) In this task we will prove statement from the lecture. The derived inequality will enable us to bound piecewise differentiable function given that we have a bound on its gradients. Let $f : [a_1, b_1] \times ... \times [a_k, b_k] \to \mathbb{R}$ be piecewise differentiable function defined as $f(x) = f_{i,j}(x)$ where $x \in D_{i,j}$. Here $D_1, ..., D_N$ are hyperrectangles which partition the function domain $[a_1, b_1] \times ... [a_k, b_k]$ into finite number of pieces.

1. Let $||\nabla f_{i,j}(z)||_\infty \leq L$ for all $z \in D_{i,j}$. Prove the following bound:

$$
f_{i,j}(y) \leq f_{i,j}(x) + L||x - y||_1, \forall x, y \in D_{i,j}
$$

2. Prove the following bound:

$$
f(y) \leq f(x) + L||x - y||_1, \forall x, y \in [a_1, b_1] \times ...[a_k, b_k].
$$

3. Prove that:

$$
f(y) \leq f(c) + \frac{L}{2} \sum_{i=1}^{k} b_i - a_i, \forall y \in [a_1, b_1] \times ...[a_k, b_k].
$$

Here $c$ is center of the domain, meaning $c_i = \frac{1}{2}(a_i + b_i)$.

**Solution 2.**   1. Applying mean-value theorem we get

$$
\begin{aligned}
f_{i,j}(y) &= f_{i,j}(x) + \nabla f_{i,j}(z)^T (y - x) \\
&\leq f_{i,j}(x) + L||y - x||_1.
\end{aligned}
$$

2. Let $x_1, ..., x_m$ be points on the line from $x$ to $y$ such that line between $x_i$ and $x_{i+1}$ lines in the same piece $D_{i,j}$. Additonally, we denote $x_1 = x$ and $x_m = y$. Then,

$$f(y) = f(x_m)$$
$$= f(x_1) + \sum_{i=1}^{m-1} f(x_{i+1}) - f(x_i)$$
$$\leq f(x_1) + \sum_{i=1}^{m-1} L||x_{i+1} - x_i||_1$$
$$\leq f(x_1) + L \sum_{i=1}^{m-1} ||x_{i+1} - x_i||_1$$
$$= f(x) + L||x - y||_1.$$

3. Applying the previous inequality for $x = c$ and noting that $|c_i - y_i| \leq \frac{1}{2}(b_i - a_i)$ we get desired result.