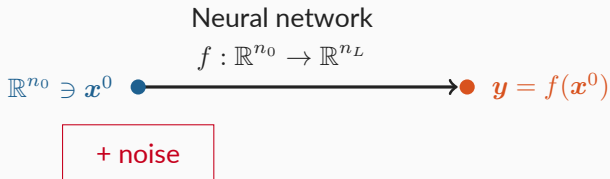


Exactness of SDP Relaxation for Robustness Verification of Neural Networks

Safety Verification of Neural Networks (NNs)

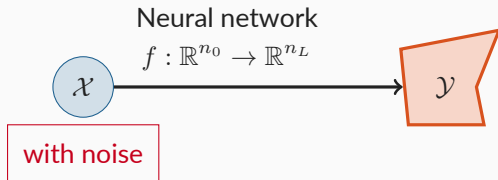


- adversarial attacks
- uncertainty

Safety Verification

- Evaluation of the robustness of x^0
- By verification of $\mathcal{Y} \subseteq S_y$ from given \mathcal{X}, f, S_y

Safety Verification of Neural Networks (NNs)

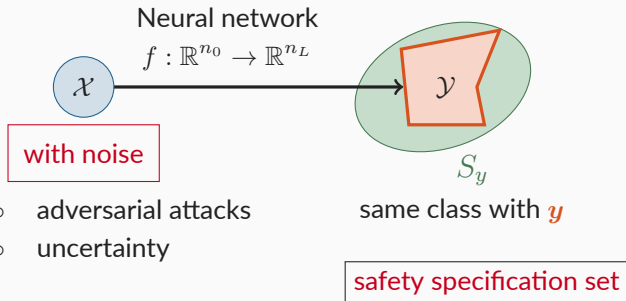


- adversarial attacks
- uncertainty

Safety Verification

- Evaluation of the robustness of x^0
- By verification of $\mathcal{Y} \subseteq S_y$ from given \mathcal{X}, f, S_y

Safety Verification of Neural Networks (NNs)



Safety Verification

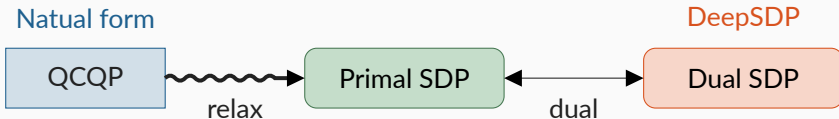
- Evaluation of the robustness of x^0
- By verification of $\mathcal{Y} \subseteq S_y$ from given \mathcal{X}, f, S_y

Semidefinite programming(SDP)-based method can be considered

- Quadratic constraints can formulate NN with ReLU ϕ

$$y = \phi(x) := \max\{0, x\} \iff \{y \geq 0, \quad y \geq x, \quad y(y - x) = 0\}$$

- DeepSDP is dual of SDP relaxation
(solvable in polynomial time)



Exactness of Relaxation and Our Motivation

Two gaps exist

- Duality gap (Primal SDP and Dual SDP)
 - $= 0$ under strong duality
- Relaxation gap (QCQP and Primal SDP)
 - $= 0 \iff$ SDP relaxation is *exact*

Motivation of this talk

- Assume strong duality \Rightarrow no duality gap
- Under what conditions is the primal SDP relaxation exact?

- Introduction
- Quadratically constrained quadratic programming (QCQP)
- Exact semidefinite programming (SDP) relaxation
- Single-layer feed-forward neural network
- QCQP formulation for safety verification
- Exactness conditions for SDP relaxation
- Graphical explanation of exactness
- Summary

QCQP: Quadratically Constrained Quadratic Programming

Quadratic objective function and quadratic constraints

$$\begin{aligned} v^* := \min_{\mathbf{x} \in \mathbb{R}^n} \quad & \mathbf{x}^T Q^0 \mathbf{x} + 2(\mathbf{q}^0)^T \mathbf{x} \\ \text{s.t.} \quad & \mathbf{x}^T Q^p \mathbf{x} + 2(\mathbf{q}^p)^T \mathbf{x} \leq b_p, \quad p \in [m] := \{1, \dots, m\}. \end{aligned} \quad (\mathcal{P})$$

- Generally non-convex & NP-hard
- Approximately solvable via SDP relaxation

QCQP

$$v^* = \min \left\{ \mathbf{x}^T Q^0 \mathbf{x} + 2(\mathbf{q}^0)^T \mathbf{x} \mid \mathbf{x}^T Q^p \mathbf{x} + 2(\mathbf{q}^p)^T \mathbf{x} \leq b_p, p \in [m] \right\} \quad (\mathcal{P})$$

$$= \min \left\{ Q^0 \bullet X + 2(\mathbf{q}^0)^T \mathbf{x} \mid \boxed{X = \mathbf{x}\mathbf{x}^T} \right. \\ \left. Q^p \bullet X + 2(\mathbf{q}^p)^T \mathbf{x} \leq b_p, p \in [m] \right\}$$

$$\geq \min \left\{ Q^0 \bullet X + 2(\mathbf{q}^0)^T \mathbf{x} \mid \boxed{X \succeq \mathbf{x}\mathbf{x}^T} \right. \\ \left. Q^p \bullet X + 2(\mathbf{q}^p)^T \mathbf{x} \leq b_p, p \in [m] \right\} \quad (\mathcal{P}_R)$$

Semidefinite Programming (SDP) Relaxation

$$=: v_{\text{SDP}}^*$$

Notation

- $Q^p \bullet X := \sum_{i,j} Q_{ij}^p X_{ij}$: Frobenius inner product.
- $X \succeq \mathbf{x}\mathbf{x}^T \iff X - \mathbf{x}\mathbf{x}^T$ is positive semidefinite.

Def: Exactness

SDP relaxation (\mathcal{P}_R) is exact (tight) if $v^* = v_{\text{SDP}}^*$

- Exact $\iff (\mathcal{P}_R)$ has a rank-1 solution X^*

$$X^* \succeq \mathbf{x}^* (\mathbf{x}^*)^T \iff \begin{bmatrix} 1 & (\mathbf{x}^*)^T \\ \mathbf{x}^* & X^* \end{bmatrix} \succeq O.$$

- Sufficient conditions for general QCQPs have been studied
 - Sparsity and graphs
 - Projection of epigraphs
 - Collinearity of Gram matrix representation

Gram Matrix Transformation

Replace (x, X) in SDP relaxation:

- Fix $e \in \mathbb{R}^s$ with $\|e\| = 1$
- Introduce $u^1, \dots, u^n \in \mathbb{R}^s$ so that

$$\left[\begin{array}{c|c} 1 & x^T \\ \hline x & X \end{array} \right] = \left[\begin{array}{c|cccc} e^T e & e^T u^1 & \dots & \dots & e^T u^n \\ \hline e^T u^1 & (u^1)^T u^1 & \dots & \dots & (u^1)^T u^n \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ e^T u^n & (u^n)^T u^1 & \dots & \dots & (u^n)^T u^n \end{array} \right]$$

$$\begin{aligned} Q^p \bullet X + 2(q^p)^T x &= \sum_i \sum_j Q_{ij}^p X_{ij} && + 2 \sum_i q_i^p x_i \\ &= \sum_i \sum_j Q_{ij}^p (u^i)^T u^j + 2 \sum_i q_i^p e^T u^i \end{aligned}$$

Exactness Condition for Gram Matrix Representation

Finally, we obtain

$$\begin{aligned} v_{\text{SDP}}^* = \min \quad & \sum_i \sum_j Q_{ij}^0 (\mathbf{u}^i)^\top \mathbf{u}^j + 2 \sum_i q_i^0 \mathbf{e}^\top \mathbf{u}^i \\ \text{s.t.} \quad & \sum_i \sum_j Q_{ij}^p (\mathbf{u}^i)^\top \mathbf{u}^j + 2 \sum_i q_i^p \mathbf{e}^\top \mathbf{u}^i \leq b_p, \quad p \in [m], \\ & \mathbf{u}^1, \dots, \mathbf{u}^n \in \mathbb{R}^n. \end{aligned}$$

Proposition [Z20]

SDP relaxation is exact if there exist an optimal solution $(\mathbf{u}^1)^*, \dots, (\mathbf{u}^n)^*$ which are collinear to \mathbf{e} , i.e.,

$$|\mathbf{e}^\top \mathbf{u}^i| = \|\mathbf{u}^i\| \quad \text{for all } i \in [n].$$

[Z20] Zhang, On the tightness of semidefinite relaxations for certifying robustness to adversarial examples, NeurIPS, 2020.

- Introduction
- Quadratically constrained quadratic programming (QCQP)
- Exact semidefinite programming (SDP) relaxation
- Single-layer feed-forward neural network
- QCQP formulation for safety verification
- Exactness conditions for SDP relaxation
- Graphical explanation of exactness
- Summary

Single-layer Neural Networks

W^0, W^1 : weight matrices, $\mathbf{b}^0, \mathbf{b}^1$: bias vectors

Neural network

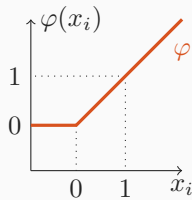
$$\begin{aligned}\mathbf{x}^1 &:= \phi(W^0 \mathbf{x}^0 + \mathbf{b}^0), \\ f(\mathbf{x}^0) &:= W^1 \mathbf{x}^1 + \mathbf{b}^1.\end{aligned}$$

Note we consider the case that

- $W^1 = I$, and $\mathbf{b}^1 = \mathbf{0}$.
- ϕ is an element-wise ReLU function, i.e.,

$$\phi(\mathbf{x}) := \begin{bmatrix} \varphi(x_1) & \cdots & \varphi(x_n) \end{bmatrix}^T,$$

where $\varphi(x_i) := \max\{0, x_i\}$.



Input Set $\mathcal{X} \subseteq \mathbb{R}^{n_0}$

Set \mathcal{X} contains the uncertainty and attacks.

- Each input x^0 is chosen from \mathcal{X} .
- The safety of x^0 is evaluated by S_y .

Note

\mathcal{X} is not the domain of NN f .

Various shapes are possible.

- hyper-ellipsoid $\mathcal{X} = \{x \mid \|x - \hat{x}\|_2 \leq \rho\}$.
- hyper-rectangle $\mathcal{X} = \{x \mid \|x - \hat{x}\|_\infty \leq \rho\}$.

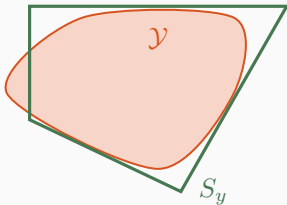
Setting II

This talk covers the case where \mathcal{X} is a hyper-ellipsoid.

Polytope Safety Specification Set

Consider polytope safety specification set S_y

- Let S_y be a quadrilateral below.
- $\mathcal{Y} \subseteq S_y$ can be verified via four half-spaces.



Setting I

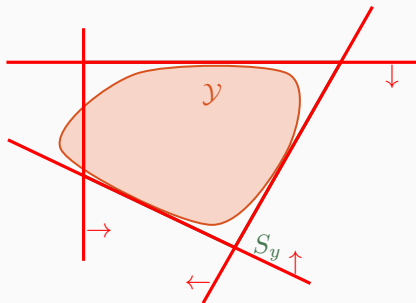
Assume that safety specification set S_y is a half-space

$$H := \{ \mathbf{y} \in \mathbb{R}^{n_2} \mid \mathbf{c}^T \mathbf{y} - d \geq 0 \}.$$

Polytope Safety Specification Set

Consider polytope safety specification set S_y

- Let S_y be a quadrilateral below.
- $\mathcal{Y} \subseteq S_y$ can be verified via four half-spaces.



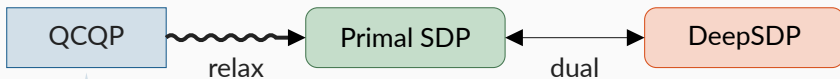
Setting I

Assume that safety specification set S_y is a half-space

$$H := \{ \mathbf{y} \in \mathbb{R}^{n_2} \mid \mathbf{c}^T \mathbf{y} - d \geq 0 \}.$$

QCQP Formulation

Natural form



$$\begin{aligned} \min \quad & 2\mathbf{c}^T \mathbf{x}^1 \\ \text{s.t.} \quad & \mathbf{x}^1 \geq \mathbf{0}, \quad \mathbf{x}^1 \geq W^0 \mathbf{x}^0 + \mathbf{b}^0, \quad \mathbf{x}^1 \odot (W^0 \mathbf{x}^0 + \mathbf{b}^0 - \mathbf{x}^1) \leq \mathbf{0}, \\ & (+ \text{some valid cuts}), \\ & \mathbf{x}^0 \in \mathcal{X} = \{\mathbf{x} \mid \|\mathbf{x} - \hat{\mathbf{x}}\|_2 \leq \rho\}, \quad \mathbf{x}^1 \in \mathbb{R}^{n_1}. \end{aligned}$$

\odot = Hadamard product

Determine a constant term d for fixed \mathbf{c} in

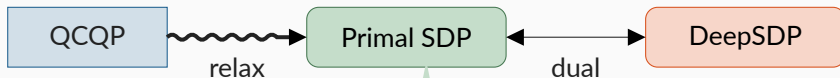
$$\{\mathbf{x}^1 \in \mathbb{R}^{n_1} \mid \mathbf{c}^T \mathbf{x}^1 - d \geq 0\}.$$

minimize
in primal side

maximize
in dual side

Primal SDP Relaxation

Natural form



$$\begin{aligned} \min \quad & 2\mathbf{c}^T \mathbf{x}^1 \\ \text{s.t.} \quad & I \bullet X^{00} - 2\hat{\mathbf{x}}^T \mathbf{x}^0 \leq \rho^2 - \hat{\mathbf{x}}^T \hat{\mathbf{x}}, \\ & (\mathbf{e}^i)^T \mathbf{x}^1 \geq 0, \quad (\mathbf{e}^i)^T \mathbf{x}^1 \geq (\mathbf{e}^i)^T (W^0 \mathbf{x}^0 + \mathbf{b}^0), \quad i = 1, \dots, n_1, \\ & X_{ii}^{11} - b_i^0 x_i^1 - \sum_j W_{ij}^0 X_{ij}^{10} \leq 0, \quad i = 1, \dots, n_1, \\ & \begin{bmatrix} 1 & (\mathbf{x}^0)^T & (\mathbf{x}^1)^T \\ \mathbf{x}^0 & X^{00} & (X^{10})^T \\ \mathbf{x}^1 & X^{10} & X^{11} \end{bmatrix} \succeq O \end{aligned}$$

\mathbf{e}^i : the i th unit vector

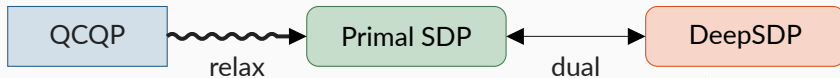
Objective of this work

Verify the exactness of (primal) SDP relaxation

- checking “*collinearity*” of another form

Dual SDP relaxation

Natural form



$$\begin{aligned}
 & \max_{\gamma, \lambda, \nu, \eta, d} \quad 2d \\
 & \text{s.t.} \quad \gamma \begin{bmatrix} \hat{\mathbf{x}}^T \hat{\mathbf{x}} - \rho^2 & -\hat{\mathbf{x}}^T & \mathbf{0} \\ -\hat{\mathbf{x}} & I & O \\ \mathbf{0} & O & O \end{bmatrix} + \begin{bmatrix} -2d & \mathbf{0}^T & \mathbf{c}^T \\ \mathbf{0} & O & O \\ \mathbf{c} & O & O \end{bmatrix} \\
 & \quad + \begin{bmatrix} 0 & \nu^T W^0 & -\nu^T - \eta^T \\ (W^0)^T \nu & O & -(W^0)^T \text{diag}(\lambda) \\ -\nu - \eta & -\text{diag}(\lambda) W^0 & 2 \text{diag}(\lambda) \end{bmatrix} \succeq O, \\
 & \quad \gamma \in \mathbb{R}_+, \quad \lambda, \nu, \eta \in \mathbb{R}_+^n, \quad d \in \mathbb{R}.
 \end{aligned}$$

For previous QCQP,

Theorem

The primal SDP relaxation is exact if

- $\mathcal{X} = \{x \mid \|x - \hat{x}\|_2 \leq \rho\}$; or
- $\mathcal{X} = \{x \mid \|x - \hat{x}\|_\infty \leq \rho\}$, $W^0 = I$, and $\hat{x} \geq -b^0$.

In addition, DeepSDP is also exact under strong duality.

- No extra assumption in hyper-ellipsoid case.
- Proof depending on Gram matrix representation with $e = e^1$.

Gram Matrix Representation of Safety Verification

- Assign $\mathbf{u}^j \in \mathbb{R}^s$ for \mathbf{x}^0 , and $\mathbf{v}^i \in \mathbb{R}^s$ for \mathbf{x}^1 :

$$\mathbf{x}^0 = \left[(\mathbf{e}^1)^\top \mathbf{u}^1, \dots, (\mathbf{e}^1)^\top \mathbf{u}^{n_0} \right]^\top, \quad \mathbf{x}^1 = \left[(\mathbf{e}^1)^\top \mathbf{v}^1, \dots, (\mathbf{e}^1)^\top \mathbf{v}^{n_1} \right]^\top$$

- Similarly,

$$[X^{00}]_{ij} = (\mathbf{u}^i)^\top \mathbf{u}^j, \quad [X^{10}]_{ij} = (\mathbf{u}^i)^\top \mathbf{v}^j, \quad [X^{11}]_{ij} = (\mathbf{v}^i)^\top \mathbf{v}^j$$

Primal SDP is equivalently

$$\begin{aligned} \min_{\mathbf{u}^j, \mathbf{v}^i} \quad & 2 \sum_{i=1}^{n_1} c_i (\mathbf{e}^1)^\top \mathbf{v}^i \\ \text{s.t.} \quad & (\mathbf{e}^1)^\top \mathbf{v}^i \geq 0, \quad i = 1, \dots, n_1, \\ & (\mathbf{e}^1)^\top \mathbf{v}^i \geq (\mathbf{e}^1)^\top \left(\sum_{j=1}^{n_0} W_{ij} \mathbf{u}^j + b_i^0 \mathbf{e}^1 \right), \quad i = 1, \dots, n_1, \\ & \|\mathbf{v}^i\|_2^2 \leq \left(\sum_{j=1}^{n_0} W_{ij} \mathbf{u}^j + b_i^0 \mathbf{e}^1 \right)^\top \mathbf{v}^i, \quad i = 1, \dots, n_1, \\ & \sum_{j=1}^{n_0} \|\mathbf{u}^j - \hat{x}_j \mathbf{e}^1\|_2^2 \leq \rho^2. \end{aligned}$$

Decomposition according to \mathbf{u}^j and \mathbf{v}^i

$$\min_{\mathbf{u}^j, \mathbf{v}^i} 2 \sum_{i=1}^{n_1} c_i (\mathbf{e}^1)^\top \mathbf{v}^i \quad (1)$$

$$\text{s.t. } (\mathbf{e}^1)^\top \mathbf{v}^i \geq 0, \quad i = 1, \dots, n_1, \quad (2)$$

$$(\mathbf{e}^1)^\top \mathbf{v}^i \geq (\mathbf{e}^1)^\top \left(\sum_{j=1}^{n_0} W_{ij} \mathbf{u}^j + b_i^0 \mathbf{e}^1 \right), \quad i = 1, \dots, n_1, \quad (3)$$

$$\|\mathbf{v}^i\|_2^2 \leq \left(\sum_{j=1}^{n_0} W_{ij} \mathbf{u}^j + b_i^0 \mathbf{e}^1 \right)^\top \mathbf{v}^i, \quad i = 1, \dots, n_1, \quad (4)$$

$$\sum_{j=1}^{n_0} \|\mathbf{u}^j - \hat{x}_j \mathbf{e}^1\|_2^2 \leq \rho^2. \quad (5)$$

Inner problem

$$\Psi(\mathbf{v}^1, \dots, \mathbf{v}^{n_1}) :=$$

$$\min_{\mathbf{u}^1, \dots, \mathbf{u}^{n_0}} \sum_{j=1}^{n_0} \|\mathbf{u}^j - \hat{x}_j \mathbf{e}^1\|_2^2$$

$$\text{s.t. } (3), (4).$$

Outer problem

$$\min_{\mathbf{v}^1, \dots, \mathbf{v}^{n_1}} (1)$$

$$\text{s.t. } (2)$$

$$\Psi(\mathbf{v}^1, \dots, \mathbf{v}^{n_1}) \leq \rho^2.$$

Relationship Between Their Solutions

A part of KKT condition of Inner problem :

$$\begin{bmatrix} \mathbf{u}^1 \\ \vdots \\ \mathbf{u}^{n_0} \end{bmatrix} = \begin{bmatrix} \hat{x}_1 \mathbf{e}^1 \\ \vdots \\ \hat{x}_{n_0} \mathbf{e}^1 \end{bmatrix} - \sum_{i=1}^{n_1} \frac{\nu_i}{2} \begin{bmatrix} W_{i1} \mathbf{e}^1 \\ \vdots \\ W_{in} \mathbf{e}^1 \end{bmatrix} + \sum_{i=1}^{n_1} \frac{\lambda_i}{2} \begin{bmatrix} W_{i1} \mathbf{v}^i \\ \vdots \\ W_{in} \mathbf{v}^i \end{bmatrix}$$

Lemma: Linear Combination

For any optimal solution $(\mathbf{u}^1)^*, \dots, (\mathbf{u}^{n_0})^*$ of Inner problem ,
there exist $\mathbf{m} \in \mathbb{R}^{n_0}$ and $M \in \mathbb{R}^{n_1 \times n_0}$ such that

$$(\mathbf{u}^j)^* = m_j \mathbf{e}^1 + \sum_{i=1}^{n_1} M_{ij} \mathbf{v}^i \quad \text{for each } j \in \{1, \dots, n_0\}.$$

Collinearity in Outer-problem

Using m and M changes Outer problem into

$$\left. \begin{aligned} \min_{v^1, \dots, v^{n_1}} \quad & 2 \sum_{i=1}^{n_1} c_i e^T v^i & (1) \\ \text{s.t.} \quad & e^T v^i \geq 0, \quad i = 1, \dots, n_1, & (2) \\ & \sum_{j=1}^{n_0} \left\| (m_j - \hat{x}_j) e^1 + \sum_{i=1}^{n_1} M_{ij} v^i \right\|_2^2 \leq \rho^2. & \end{aligned} \right\}$$

Lemma: Collinearity of $(v^i)^*$

Outer problem has an optimal solution $(v^1)^*, \dots, (v^{n_1})^*$ collinear to e^1 .

Therefore, the SDP relaxation is exact due to the collinearity.

Collinearity in Outer-problem

Using m and M changes Outer problem into

$$\left. \begin{aligned} \min_{\mathbf{v}^1, \dots, \mathbf{v}^{n_1}} \quad & 2 \sum_{i=1}^{n_1} c_i \mathbf{e}^T \mathbf{v}^i & (1) \\ \text{s.t.} \quad & \mathbf{e}^T \mathbf{v}^i \geq 0, \quad i = 1, \dots, n_1, & (2) \\ & \sum_{j=1}^{n_0} \left\| (m_j - \hat{x}_j) \mathbf{e}^1 + \sum_{i=1}^{n_1} M_{ij} \mathbf{v}^i \right\|_2^2 \leq \rho^2. & \end{aligned} \right\}$$

Essence of Proof.

- Let $(\bar{\mathbf{v}}^1, \dots, \bar{\mathbf{v}}^{n_1})$ be an optimal solution.
- Assume at least one of $\bar{\mathbf{v}}^1, \dots, \bar{\mathbf{v}}^{n_1}$ is not collinear to \mathbf{e}^1 .
- Define

$$\hat{\mathbf{v}}^i := \left[\bar{v}_1^i, 0, \dots, 0 \right]^T.$$

- Then, $(\hat{\mathbf{v}}^1, \dots, \hat{\mathbf{v}}^{n_1})$ is another optimal solution.

- Introduction
- Quadratically constrained quadratic programming (QCQP)
- Exact semidefinite programming (SDP) relaxation
- Single-layer feed-forward neural network
- QCQP formulation for safety verification
- Exactness conditions for SDP relaxation
- Graphical explanation of exactness
- Summary

Meaning of Exact Relaxation

Current situation

- Instead of checking whether $\mathcal{Y} \subseteq S_y$,

$$\mathcal{Y} \subseteq H := \bigcap_k \left\{ \mathbf{y} \in \mathbb{R}^{n_2} \mid (\mathbf{c}^k)^T \mathbf{y} - \boxed{d^k} \geq 0 \right\} \subseteq S_y.$$

Optimal Solutions

Where does exactness appear?

- Half-space H supports \mathcal{Y} at a face.
- Margin between H and S_y increases.
 - It is the robustness of the input set \mathcal{X} .
 - $\iff \mathcal{X}$ can be made larger.

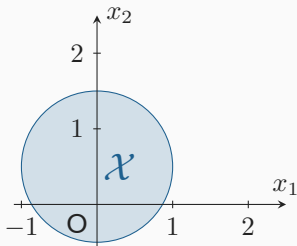
Instance with Exact Relaxation

Example: a single-layer NN

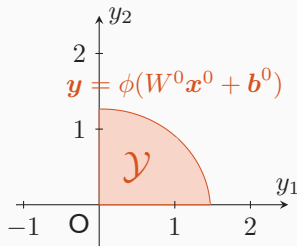
$$W^0 = \begin{bmatrix} 1/2 & 5/4 \\ -6/5 & 2/5 \end{bmatrix}, \quad W^1 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \mathbf{b}^0 = \begin{bmatrix} -0.5 \\ -0.2 \end{bmatrix}, \quad \mathbf{b}^1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}.$$

with the input set

$$\mathcal{X} := \left\{ \mathbf{x} \in \mathbb{R}^2 \left\| \mathbf{x} - \begin{bmatrix} 0 \\ 0.5 \end{bmatrix} \right\|_2 \leq 1 \right\}.$$

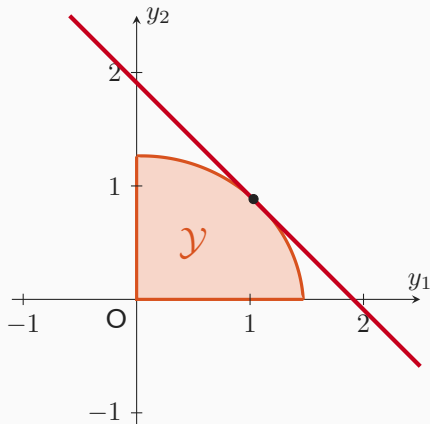


\Rightarrow



Experiment Result in Mosek and Julia

Slope c^*		Solution
c_1	c_2	d^*
1	0	-8.42×10^{-10}
0	1	-3.45×10^{-10}
-1	1	-1.47
1	-1	-1.26
-1/4	-1	-1.32
-1/2	-1	-1.46
-1	-1	-1.91
-2	-1	-3.16
-4	-1	-5.96

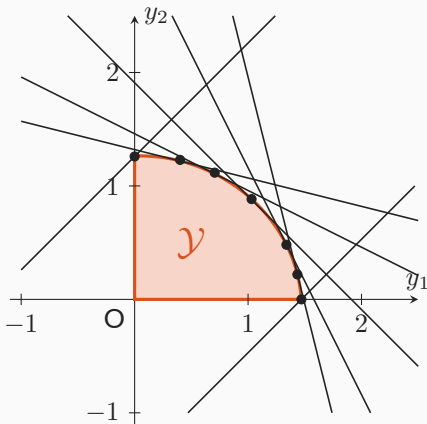


Observation

All boundaries intersect \mathcal{Y} at a point

Experiment Result in Mosek and Julia

Slope c^*		Solution
c_1	c_2	d^*
1	0	-8.42×10^{-10}
0	1	-3.45×10^{-10}
-1	1	-1.47
1	-1	-1.26
-1/4	-1	-1.32
-1/2	-1	-1.46
-1	-1	-1.91
-2	-1	-3.16
-4	-1	-5.96



Observation

All boundaries intersect \mathcal{Y} at a point

Summary

- Safety verification of single-layer neural networks
- Exactness conditions of SDP relaxation
- Graphical insight for exactness

Future works

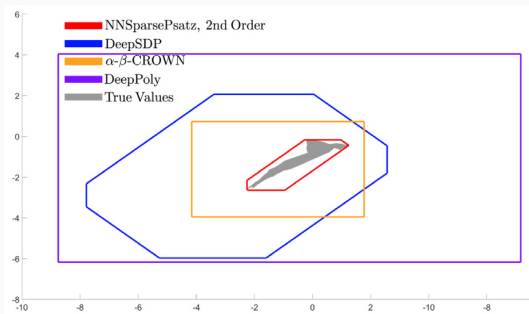
- Analyze the non-polyhedral case of S_y
- Extend results to other networks

Thank you for your attention!

For more details, see arXiv:2504.09934

Exact Semidefinite Relaxations for Verifying Robustness of Neural Networks

- Less accurate than other methods in general [NP21].
- **Accurate** when the relaxation is exact.



Constraints for Hidden Layers

$$\begin{bmatrix} 1 \\ \mathbf{w}^0 \\ \phi(\mathbf{w}^0) \end{bmatrix}^T \begin{bmatrix} 0 & \mathbf{0}^T & \mathbf{0}^T \\ \mathbf{0} & O & -\mathbf{e}^i (\mathbf{e}^i)^T \\ \mathbf{0} & -\mathbf{e}^i (\mathbf{e}^i)^T & 2\mathbf{e}^i (\mathbf{e}^i)^T \end{bmatrix} \begin{bmatrix} 1 \\ \mathbf{w}^0 \\ \phi(\mathbf{w}^0) \end{bmatrix} = 0, \quad i = 1, \dots, N, \quad (6a)$$

$$\begin{bmatrix} 1 \\ \mathbf{w}^0 \\ \phi(\mathbf{w}^0) \end{bmatrix}^T \begin{bmatrix} 0 & (\mathbf{e}^i)^T & -(\mathbf{e}^i)^T \\ \mathbf{e}^i & O & O \\ -\mathbf{e}^i & O & O \end{bmatrix} \begin{bmatrix} 1 \\ \mathbf{w}^0 \\ \phi(\mathbf{w}^0) \end{bmatrix} \leq 0, \quad i = 1, \dots, N, \quad (6b)$$

$$\begin{bmatrix} 1 \\ \mathbf{w}^0 \\ \phi(\mathbf{w}^0) \end{bmatrix}^T \begin{bmatrix} 0 & \mathbf{0}^T & -(\mathbf{e}^i)^T \\ \mathbf{0} & O & O \\ -\mathbf{e}^i & O & O \end{bmatrix} \begin{bmatrix} 1 \\ \mathbf{w}^0 \\ \phi(\mathbf{w}^0) \end{bmatrix} \leq 0, \quad i = 1, \dots, N. \quad (6c)$$

Valid Cuts

Let $\mathbf{w}^0 = W^0 \mathbf{x}^0 + \mathbf{b}^0$.

Valid Cuts for ReLU

The following inequation always holds

$$[\phi(w_j^0) - \phi(w_i^0)] [\phi(w_j^0) - \phi(w_i^0) - (w_j - w_i)] \leq 0 \quad \forall (i, j) \in \{1, \dots, n\}^2$$

$$\begin{bmatrix} \mathbf{w}^0 \\ \phi(\mathbf{w}^0) \end{bmatrix}^T \begin{bmatrix} O & -(\mathbf{e}^i - \mathbf{e}_j)(\mathbf{e}^i - \mathbf{e}_j)^T \\ -(\mathbf{e}^i - \mathbf{e}_j)(\mathbf{e}^i - \mathbf{e}_j)^T & 2(\mathbf{e}^i - \mathbf{e}_j)(\mathbf{e}^i - \mathbf{e}_j)^T \end{bmatrix} \begin{bmatrix} \mathbf{w}^0 \\ \phi(\mathbf{w}^0) \end{bmatrix} \leq 0, \quad (7)$$

Input Constraint

Since $\mathbf{x}^0 \in \mathcal{X} = \{\mathbf{x} \mid \|\mathbf{x} - \hat{\mathbf{x}}\|_2 \leq \rho\}$.

In QCQP

$$\|\mathbf{x}^0 - \hat{\mathbf{x}}\|_2^2 \leq \rho^2$$

In SDP relaxation

$$\begin{bmatrix} \hat{\mathbf{x}}^T \hat{\mathbf{x}} - \rho^2 & -\hat{\mathbf{x}}^T & \mathbf{0} \\ -\hat{\mathbf{x}} & I & O \\ \mathbf{0} & O & O \end{bmatrix} \bullet \underbrace{\begin{bmatrix} 1 & (\mathbf{x}^0)^T & (\mathbf{x}^1)^T \\ \mathbf{x}^0 & X_{00} & X_{10}^T \\ \mathbf{x}^1 & X_{10} & X_{11} \end{bmatrix}}_{=: G} \leq 0$$

In DeepSDP

By introducing a dual variable γ ,

$$\gamma \begin{bmatrix} \hat{\mathbf{x}}^T \hat{\mathbf{x}} - \rho^2 & -\hat{\mathbf{x}}^T & \mathbf{0} \\ -\hat{\mathbf{x}} & I & O \\ \mathbf{0} & O & O \end{bmatrix}$$

Safety Specification Set

Consider a half-space $H := \{\mathbf{y} \in \mathbb{R}^{n_2} \mid \mathbf{c}^T \mathbf{y} - d \geq 0\}$.

- The slope \mathbf{c} according to each half-space is given.
- The largest d makes H smaller.

In SDP relaxation

$$\mathbf{x}^1 \in H \iff \begin{bmatrix} -2d & \mathbf{0}^T & \mathbf{c}^T \\ \mathbf{0} & O & O \\ \mathbf{c} & O & O \end{bmatrix} \bullet \begin{bmatrix} 1 & (\mathbf{x}^0)^T & (\mathbf{x}^1)^T \\ \mathbf{x}^0 & X_{00} & X_{10}^T \\ \mathbf{x}^1 & X_{10} & X_{11} \end{bmatrix} \leq 0$$

In DeepSDP

Let d behave as a dual variable.

Quadratic Formulation for ReLU Function

Review: ϕ applies element-wisely ReLU function φ

Let $\mathbf{w}^0 := W^0 \mathbf{x}^0 + \mathbf{b}^0$. For any $i \in \{1, \dots, n_1\}$,

$$\varphi(w_i^0) = \max\{0, w_i^0\} \iff \begin{cases} \varphi(w_i^0) (\varphi(w_i^0) - w_i^0) \leq 0, \\ \varphi(w_i^0) \geq w_i^0, \quad \varphi(w_i^0) \geq 0. \end{cases}$$

In QCQP The first inequality is

$$\begin{bmatrix} 1 \\ \mathbf{w}^0 \\ \phi(\mathbf{w}^0) \end{bmatrix}^T \begin{bmatrix} 0 & \mathbf{0}^T & \mathbf{0}^T \\ \mathbf{0} & O & -\mathbf{e}^i (\mathbf{e}^i)^T \\ \mathbf{0} & -\mathbf{e}^i (\mathbf{e}^i)^T & 2\mathbf{e}^i (\mathbf{e}^i)^T \end{bmatrix} \begin{bmatrix} 1 \\ \mathbf{w}^0 \\ \phi(\mathbf{w}^0) \end{bmatrix} \leq 0, \quad i = 1, \dots, n_1.$$

Transformation of $[1, w^0, \phi(w^0)]^T$

Equivalently, for $i \in \{1, \dots, n_1\}$,

$$\begin{bmatrix} 1 \\ x^0 \\ x^1 \end{bmatrix}^T \begin{bmatrix} 1 & \mathbf{0}^T & \mathbf{0}^T \\ \mathbf{b}^0 & W^0 & O \\ \mathbf{0} & O & I \end{bmatrix}^T \begin{bmatrix} 0 & \mathbf{0}^T & \mathbf{0}^T \\ \mathbf{0} & O & -e^i(e^i)^T \\ \mathbf{0} & -e^i(e^i)^T & 2e^i(e^i)^T \end{bmatrix} \begin{bmatrix} 1 & \mathbf{0}^T & \mathbf{0}^T \\ \mathbf{b}^0 & W^0 & O \\ \mathbf{0} & O & I \end{bmatrix} \begin{bmatrix} 1 \\ x^0 \\ x^1 \end{bmatrix} \leq$$

$$=: L_i$$

In SDP relaxation

$$L_i \bullet G \leq 0, \quad i \in \{1, \dots, n_1\}$$

In DeepSDP

Introducing a dual variable $\lambda \in \mathbb{R}_+^{n_1}$,

$$\sum_{i=1}^{n_1} \lambda_i L_i$$

Constraint	$\varphi(w_i^0) (\varphi(w_i^0) - w_i^0) \leq 0$	$\varphi(w_i^0) \geq w_i^0$	$\varphi(w_i^0) \geq 0$
------------	--------------------------------------------------	-----------------------------	-------------------------

Dual variable	λ_i	ν_i	η_i
---------------	-------------	---------	----------