

Exact Semidefinite Relaxations for Safety Verification of Neural Network

Godai Azuma (Aoyama Gakuin University)

Joint work with Sunyoung Kim (Ewha W. University)
 Makoto Yamashita (Institute of Science Tokyo)

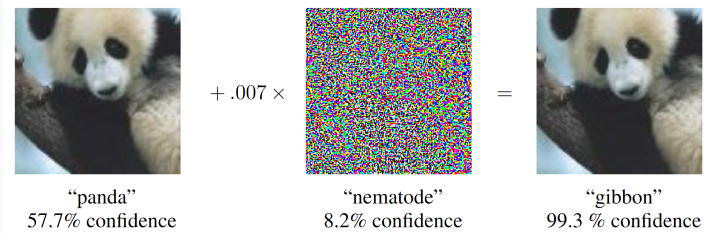
ICCOPT 2025 in Los Angeles (July 22, 2025)

This work was supported by JSPS KAKENHI JP24K20738, JP22KJ1307 and JP20H04145
 NRF 2021-R1A2C1003810

Vulnerability on Neural Networks (NNs)

Uncertainty and adversarial attacks

- Panda image + small noise [GSS14]



- Stop sign + optical attacks \Rightarrow Speed limit 60 sign

a barrier to applications where reliability is critical (e.g., self-driving cars)

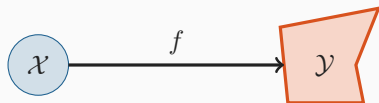
Verifying Safety of NNs

Let $f : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_L}$ be a NN.

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

we wish to evaluate

Output $\mathcal{Y} := \{f(x) \mid x \in \mathcal{X}\}$



Let S_y : a set where no misclassification occurs = **safety specification set**

Def: Safety Verification

Check whether $\mathcal{Y} \subseteq S_y$ holds or not for a given S_y

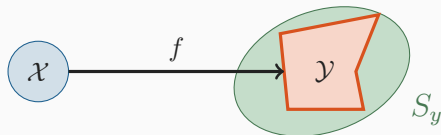
Verifying Safety of NNs

Let $f : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_L}$ be a NN.

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

we wish to evaluate

Output $\mathcal{Y} := \{f(x) \mid x \in \mathcal{X}\}$



Let S_y : a set where no misclassification occurs = safety specification set

Def: Safety Verification

Check whether $\mathcal{Y} \subseteq S_y$ holds or not for a given S_y

Semidefinite programming-based method of safety verification.

- $M \succeq O \iff M$ is positive semidefinite.

General Formulation of DeepSDP

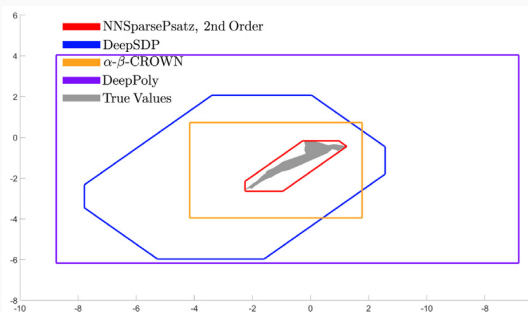
$$\begin{aligned} \min_{P, Q, S} \quad & g(P, Q, S) \\ \text{s.t.} \quad & M_{\text{in}}(P) + M_{\text{mid}}(Q) + M_{\text{out}}(S) \succeq O, \\ & P \in \mathcal{P}_{\mathcal{X}}, Q \in \mathcal{Q}_{\phi}, S \in \mathcal{S}. \end{aligned}$$

The above problem is less accurate than other methods [NP21].

[NP21] Newton and Papachristodoulou, Neural network verification using polynomial optimisation, IEEE CDC, 2021.

Accuracy of DeepSDP with Other Method

$\hat{\mathcal{Y}}$: Estimated output set \leftrightarrow \mathcal{Y} : True output set (gray) [NP21]



- All methods overestimate $\hat{\mathcal{Y}}$ in order to cover the whole output in \mathcal{Y} .
- **Accurate verification favors a smaller $\hat{\mathcal{Y}}$.**

Motivation

What conditions make DeepSDP highly accurate?

We address it by using an exact relaxation.

- Introduction
- Quadratically constrained quadratic programming
- Exact semidefinite programming (SDP) relaxation
- Single-layer feed-forward neural network
- Formulation of DeepSDP for safety verification
- Exact SDP relaxation in safety verification
- Proof
- Summary

QCQP: Quadratically Constrained Quadratic Programming

Consider a quadratic programming with quadratic constraints:

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

Used in

Binary programming, MAX-CUT, optimal flow problems,...

- **Behind the safety verification**
- Generally non-convex & NP-hard
- Approximately solvable via SDP relaxation

Semidefinite Programming (SDP) Relaxation

Define two 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.

QCQP

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

Semidefinite Programming (SDP) Relaxation

Define two 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.

QCQP

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

Semidefinite Programming (SDP) Relaxation

Define two 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.

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

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

Def: Exactness

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

i.e.,
$$\min \{ \mathbf{x}^T Q^0 \mathbf{x} \mid \mathbf{x}^T Q^p \mathbf{x} \leq b_p, p \in \{1, \dots, m\} \} \quad (\mathcal{P})$$

$$= \min \left\{ Q^0 \bullet X \mid \begin{array}{l} X \succeq \mathbf{x} \mathbf{x}^T \\ Q^p \bullet X \leq b_p, p \in \{1, \dots, m\} \end{array} \right\} \quad (\mathcal{P}_R)$$

- Sufficient conditions for exactness have been studied
- Exact $\iff (\mathcal{P}_R)$ has a rank-1 solution X^*
 - Decomposition $\hat{\mathbf{x}} \hat{\mathbf{x}}^T = X^*$ exists
 - $\hat{\mathbf{x}}$ is a solution of QCQP (\mathcal{P})

- Introduction
- Quadratically constrained quadratic programming
- Exact semidefinite programming (SDP) relaxation
- Single-layer feed-forward neural network
- Formulation of DeepSDP for safety verification
- Exact SDP relaxation in safety verification
- Proof
- Summary

Single-layer Neural Networks

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

Neural network

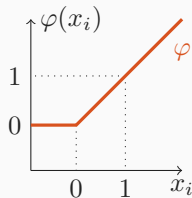
$$\begin{aligned}\mathbf{x}^1 &:= \phi(W^0 \mathbf{x}^0 + \mathbf{b}^0), \\ f(\mathbf{x}^0) &:= W^1 \mathbf{x}^1 + \mathbf{b}^1.\end{aligned}$$

Note we consider the case that

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

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

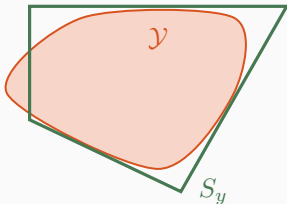
where $\phi(x_i) := \max\{0, x_i\}$.



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.



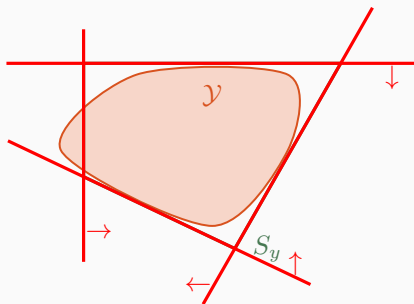
Setting I

Assume that safety specification set S_y is a half-space.

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.

DeepSDP for Single-layer NN

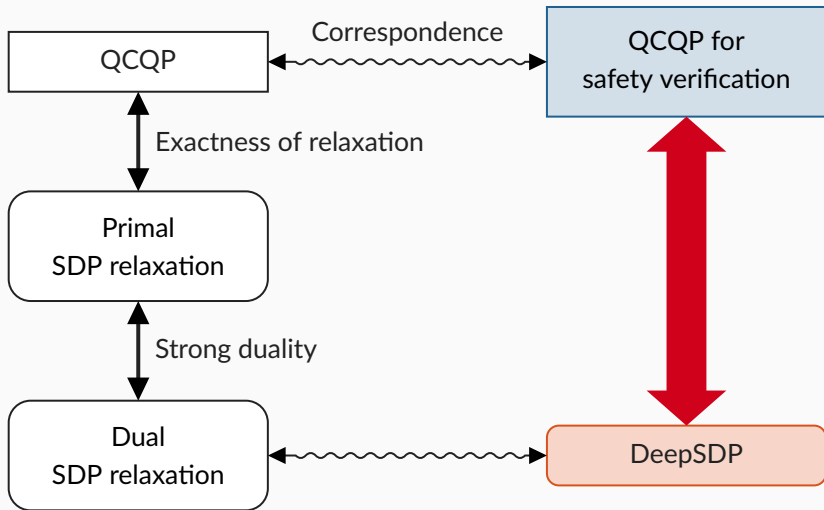
Consider a half-space $H := \{\mathbf{y} \in \mathbb{R}^{n_2} \mid \mathbf{c}^T \mathbf{y} - d \geq 0\}$.

DeepSDP to verify $\mathcal{Y} \subseteq H$

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

- An SDP, solvable in polynomial-time
- Generated from a QCQP by relaxing and taking the dual

Exactness Around Today's Problems



Exactness and Accuracy

Exact relaxation allows DeepSDP to solve the original QCQP.

Primal SDP Relaxation (= Dual of DeepSDP)

e^i : a vector where i th element is 1, the others are 0.

$$\begin{aligned}
 & \min_{\substack{\mathbf{x}^0, \mathbf{x}^1, \\ X^{00}, X^{10}, X^{11}}} 2\mathbf{c}^T \mathbf{x}^1 \\
 & \text{s.t.} \quad \begin{bmatrix} \hat{\mathbf{x}}^T \hat{\mathbf{x}} - \rho^2 & -\hat{\mathbf{x}}^T & \mathbf{0} \\ -\hat{\mathbf{x}} & I & O \\ \mathbf{0} & O & O \end{bmatrix} \bullet G \leq 0, \quad G := \begin{bmatrix} 1 & (\mathbf{x}^0)^T & (\mathbf{x}^1)^T \\ \mathbf{x}^0 & X^{00} & (X^{10})^T \\ \mathbf{x}^1 & X^{10} & X^{11} \end{bmatrix}, \\
 & \quad \begin{bmatrix} 0 & \mathbf{0}^T & -b_i^0 (e^i)^T \\ \mathbf{0} & O & -(W^0)^T e^i (e^i)^T \\ -b_i^0 e^i & -e^i (e^i)^T W^0 & 2e^i (e^i)^T \end{bmatrix} \bullet G \leq 0, \\
 & \quad \begin{bmatrix} 2b_i^0 & (e^i)^T W^0 & -(e^i)^T \\ (W^0)^T e^i & O & O \\ -e^i & O & O \end{bmatrix} \bullet G \leq 0, \\
 & \quad \begin{bmatrix} 0 & \mathbf{0}^T & -(e^i)^T \\ \mathbf{0} & O & O \\ -e^i & O & O \end{bmatrix} \bullet G \leq 0, \quad i = 1, \dots, n_1.
 \end{aligned}$$

- Introduction
- Quadratically constrained quadratic programming
- Exact semidefinite programming (SDP) relaxation
- Single-layer feed-forward neural network
- Formulation of DeepSDP for safety verification
- Exact SDP relaxation in safety verification
- Proof
- Summary

Thm: Exactness condition for a single-layer network

The primal SDP relaxation is exact if

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

In addition, DeepSDP is also exact under strong duality.

- This talk focuses on the first sufficient condition (hyper-ellipsoid).
- We discuss the derivation in the remaining time, via
 - transformation using vector e ,
 - decomposition to two problems.

Gram Matrix Transformation

Fix $e \in \mathbb{R}^{1+n_0+n_1}$ satisfying $\|e\| = 1$. (arbitrary)

- Define new variables u^1, \dots, u^{n_0} , and $v^1, \dots, v^{n_1} \in \mathbb{R}^{1+n_0+n_1}$.
- Substitute $x^0, x^1, X_{00}, X_{10}, X_{11}$ in the primal SDP by

$$x^0 = \begin{bmatrix} e^T u^1 \\ \vdots \\ e^T u^{n_0} \end{bmatrix} \in \mathbb{R}^{n_0}, \quad X^{10} = \begin{bmatrix} (v^1)^T u^1 & \dots & (v^1)^T u^{n_0} \\ \vdots & \ddots & \vdots \\ (v^{n_1})^T u^1 & \dots & (v^{n_1})^T u^{n_0} \end{bmatrix} \in \mathbb{R}^{n_1 \times n_0},$$
$$x^1 = \begin{bmatrix} e^T v^1 \\ \vdots \\ e^T v^{n_1} \end{bmatrix} \in \mathbb{R}^{n_1}, \quad X^{11} = \begin{bmatrix} (v^1)^T v^1 & \dots & (v^1)^T v^{n_1} \\ \vdots & \ddots & \vdots \\ (v^{n_1})^T v^1 & \dots & (v^{n_1})^T v^{n_1} \end{bmatrix} \in \mathbb{S}^{n_1}.$$

Example

$$0 \geq \begin{bmatrix} 0 & \mathbf{0}^T & -(e^i)^T \\ \mathbf{0} & O & O \\ -e^i & O & O \end{bmatrix} \bullet G = -2e^i \bullet x^1 = -2e^T v^i \quad \text{for all } i.$$

Equivalent Formulation of SDP Relaxation

The following problem is obtained from the primal SDP relaxation.

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

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

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

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

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

Exactness Condition by Collinearity with e

Prop: Exactness condition for collinearity in [Zhang '20]

Suppose there exists an optimal solution $\{(u^1)^*, \dots, (u^{n_0})^*, (v^1)^*, \dots, (v^{n_1})^*\}$ which are collinear with e . Then, the primal SDP relaxation is exact.

Def: Collinearity

Vectors $\{a^1, \dots, a^m\} \subseteq \mathbb{R}^{1+n_0+n_1}$ are collinear with e if

$$|e^T a^i| = \|a^i\| \quad \text{for all } i \in \{1, \dots, m\}.$$

It is suffice to show the collinearity of a solution of the transformed problem.

[Zhang '20] Zhang, On the tightness of semidefinite relaxations for certifying robustness to adversarial examples, NeurIPS, 2020.

Decomposition according to u^j and v^i

Inner problem: constraints using u^j

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

$$\left. \begin{aligned} \min_{u^1, \dots, u^{n_0}} \quad & \sum_{j=1}^{n_0} \|u^j - \hat{x}_j e\|_2^2 \\ \text{s.t.} \quad & e^T v^i \geq e^T \left(\sum_{j=1}^{n_0} W_{ij} u^j + b_i^0 e \right), \quad i = 1, \dots, n_1, \quad (3) \\ & \|v^i\|_2^2 \leq \left(\sum_{j=1}^{n_0} W_{ij} u^j + b_i^0 e \right)^T v^i, \quad i = 1, \dots, n_1, \quad (4) \end{aligned} \right\} \quad (S_2)$$

Outer problem: the remains and Ψ

$$\left. \begin{aligned} \min_{v^1, \dots, v^{n_1}} \quad & 2 \sum_{i=1}^{n_1} c_i e^T v^i \quad (1) \\ \text{s.t.} \quad & e^T v^i \geq 0, \quad i = 1, \dots, n_1, \quad (2) \\ & \Psi(v^1, \dots, v^{n_1}) \leq \rho^2. \end{aligned} \right\} \quad (S_1)$$

The case $e = e^1$ is only considered due to time limitation.

Relationship Between Their Solutions

A part of KKT condition of (S_2) :

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

Lemma: Linear Combination

For any optimal solution $(\mathbf{u}^1)^*, \dots, (\mathbf{u}^{n_0})^*$ of (S_2) ,
there exist $\mathbf{m} \in \mathbb{R}^{n_0}$ and $M \in \mathbb{R}^{n_1 \times n_0}$ such that

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

It suffices to show $(\mathbf{v}^i)^*$ s are collinear to \mathbf{e}^1 .

Collinearity in (S_1)

$$\left. \begin{aligned} \min_{\mathbf{v}^1, \dots, \mathbf{v}^{n_1}} \quad & 2 \sum_{i=1}^{n_1} c_i \mathbf{e}^T \mathbf{v}^i & (1) \\ \text{s.t.} \quad & \mathbf{e}^T \mathbf{v}^i \geq 0, \quad i = 1, \dots, n_1, & (2) \\ & \sum_{j=1}^{n_0} \left\| (m_j - \hat{x}_j) \mathbf{e}^1 + \sum_{i=1}^{n_1} M_{ij} \mathbf{v}^i \right\|_2^2 \leq \rho^2. & \end{aligned} \right\} \quad (S_1)$$

For (S_1) , we can show the following lemma:

Lemma: Collinearity of $(\mathbf{v}^i)^*$

(S_1) has an optimal solution $(\mathbf{v}^1)^*, \dots, (\mathbf{v}^{n_1})^*$ which are collinear to \mathbf{e}^1 .

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

Collinearity in (S_1)

$$\left. \begin{aligned} \min_{\mathbf{v}^1, \dots, \mathbf{v}^{n_1}} \quad & 2 \sum_{i=1}^{n_1} c_i \mathbf{e}^T \mathbf{v}^i & (1) \\ \text{s.t.} \quad & \mathbf{e}^T \mathbf{v}^i \geq 0, \quad i = 1, \dots, n_1, & (2) \\ & \sum_{j=1}^{n_0} \left\| (m_j - \hat{x}_j) \mathbf{e}^1 + \sum_{i=1}^{n_1} M_{ij} \mathbf{v}^i \right\|_2^2 \leq \rho^2. & \end{aligned} \right\} \quad (S_1)$$

Essence of Proof.

- Let $(\bar{\mathbf{v}}^1, \dots, \bar{\mathbf{v}}^{n_1})$ be an optimal solution of (S_1) .
- 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.
 - The objective value is the same.

Summary

- Verification of the safety of Neural Networks
- Exactness conditions of DeepSDP in this context

Future works

- Analyze the non-polyhedral case of S_y
- Weaken the assumptions on the input set \mathcal{X}

Thank you for your attention!

For more details, see arXiv:2504.09934

Tight Semidefinite Relaxations for Verifying Robustness of Neural Networks

QCQP: Quadratically Constrained Quadratic Programming

Consider a quadratic programming with quadratic constraints:

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

- Generally, non-convex & **NP-hard**
- Semidefinite programming (SDP) relaxation

Applications

Binary programming, MAX-CUT, optimal flow problems,...

Exactness for SDP Relaxation

The following equality holds:

$$v^* = \min \left\{ Q^0 \bullet X \mid \begin{array}{l} X \succeq O \\ Q^p \bullet X \leq b_p \quad \forall p \in [m] \end{array} \right\} = v_{\text{SDP}}^*$$

\iff rank-1 solution X^* exists

Exact SDP relaxation \implies

- Original QCQP is exactly solvable (in theory)
- The gap between a class of QCQPs and their relaxations is identified.

Motivation

What conditions of QCQPs guarantee the exactness?

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 ,

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

DeepSDP for Single-layer NN

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

DeepSDP to verify $\mathcal{Y} \subseteq H$

$$\begin{aligned} & \max_{\gamma, \boldsymbol{\lambda}, \boldsymbol{\nu}, \boldsymbol{\eta}, d} \quad 2d \\ \text{s.t.} \quad & \gamma \begin{bmatrix} \hat{\mathbf{x}}^T \hat{\mathbf{x}} - \rho^2 & -\hat{\mathbf{x}}^T & \mathbf{0} \\ -\hat{\mathbf{x}} & I & O \\ \mathbf{0} & O & O \end{bmatrix} + \begin{bmatrix} -2d & \mathbf{0}^T & \mathbf{c}^T \\ \mathbf{0} & O & O \\ \mathbf{c} & O & O \end{bmatrix} \\ & + \begin{bmatrix} 0 & \boldsymbol{\nu}^T W^0 & -\boldsymbol{\nu}^T - \boldsymbol{\eta}^T \\ (W^0)^T \boldsymbol{\nu} & O & -(W^0)^T \text{diag}(\boldsymbol{\lambda}) \\ -\boldsymbol{\nu} - \boldsymbol{\eta} & -\text{diag}(\boldsymbol{\lambda}) W^0 & 2 \text{diag}(\boldsymbol{\lambda}) \end{bmatrix} \succeq O, \\ & \gamma \in \mathbb{R}_+, \quad \boldsymbol{\lambda}, \boldsymbol{\nu}, \boldsymbol{\eta} \in \mathbb{R}_+^{n_1}, \quad d \in \mathbb{R}. \end{aligned}$$

- An SDP, solvable in polynomial-time
- Generated from a QCQP by relaxing and taking the dual

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.

DeepSDP for Single-layer NN

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

DeepSDP to verify $\mathcal{Y} \subseteq H$

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

- An SDP, solvable in polynomial-time
- Generated from a QCQP by relaxing and taking the dual

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
---------------	-------------	---------	----------

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.