



HAL
open science

Verifying Properties of Binary Neural Networks Using Sparse Polynomial Optimization

Srećko Đurašinović, Jianting Yang, Jean-Bernard Lasserre, Victor Magron, Jun Zhao

► To cite this version:

Srećko Đurašinović, Jianting Yang, Jean-Bernard Lasserre, Victor Magron, Jun Zhao. Verifying Properties of Binary Neural Networks Using Sparse Polynomial Optimization. R&T Days 2024, Jul 2024, Toulouse (FRANCE), France. 2024. <hal-04720757>

HAL Id: hal-04720757

<https://laas.hal.science/hal-04720757v1>

Submitted on 4 Oct 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire HAL, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



HAL Authorization

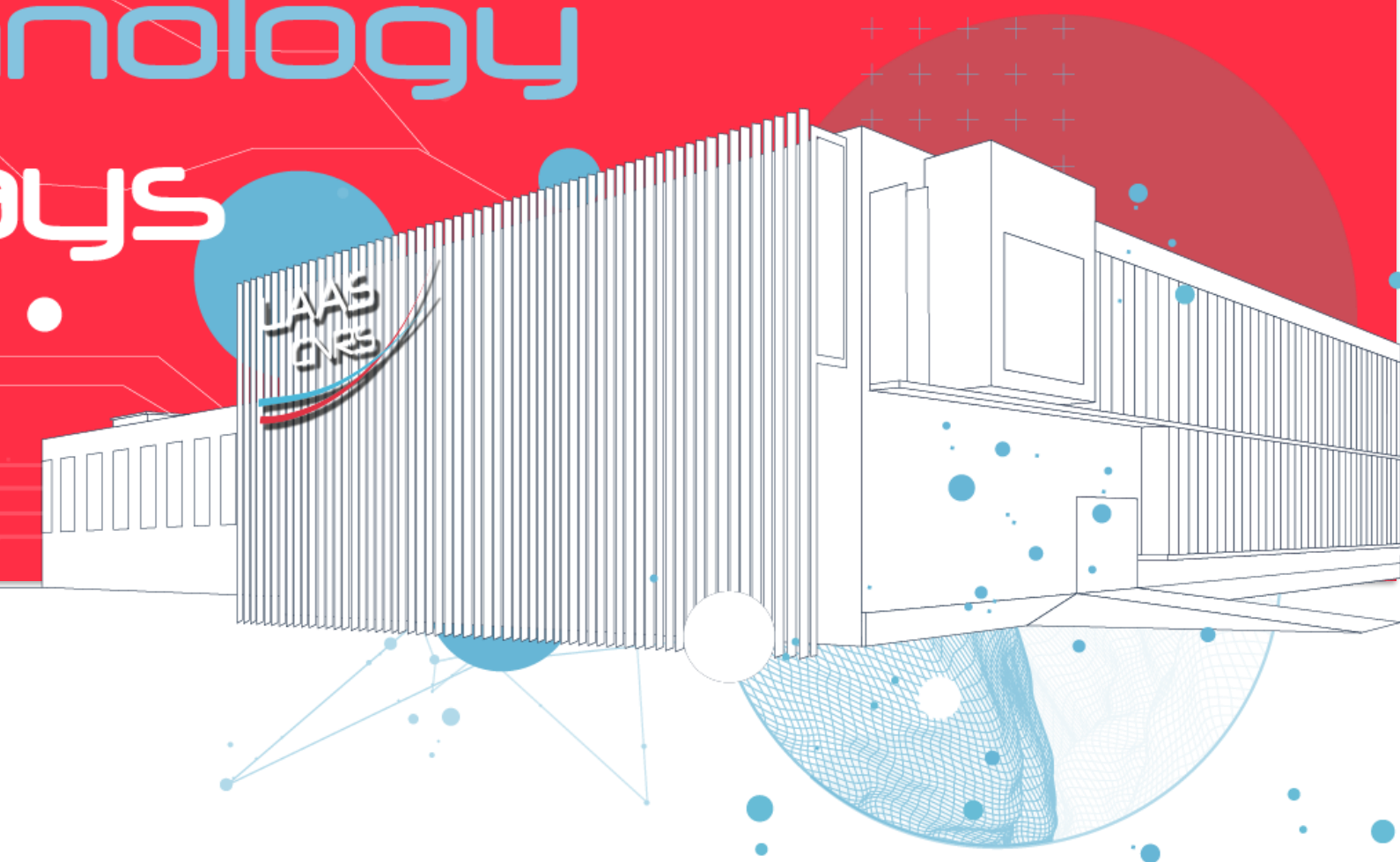
Verifying Properties of Binary Neural Networks Using Sparse Polynomial Optimization

Srećko Đurašinović^{1,2}

Joint work with Jianting Yang¹, Jean-Bernard Lasserre³, Victor Magron³ and Jun Zhao²
¹CNRS@CREATE, Singapore; ²Nanyang Technological University, Singapore; ³LAAS-CNRS



LAAS
 Research and
 Technology
 Days



Summary

This work explores methods for verifying the properties of Binary Neural Networks (BNNs), focusing on robustness against adversarial attacks. Despite their lower computational and memory needs, BNNs, like their full-precision counterparts, are also sensitive to input perturbations. Established methods for solving this problem are predominantly based on Satisfiability Modulo Theories (SMT) and Mixed-Integer Linear Programming (MILP) techniques, which often face scalability issues. We introduce an alternative approach using Semidefinite Programming (SDP) relaxations derived from sparse Polynomial Optimization (POP). Our approach, compatible with continuous input space, not only mitigates numerical issues associated with floating-point calculations but also enhances verification scalability through the strategic use of tighter first-order semidefinite relaxations. We demonstrate the effectiveness of our method in verifying robustness against both $\|\cdot\|_\infty$ and $\|\cdot\|_2$ -based adversarial attacks.

Keywords: Binary Neural Networks, Robustness Verification, Sparse Polynomial Optimization, Semidefinite Programming

Binary Neural Networks

Let $L \geq 1$ be the number of hidden layers of a classifying BNN, with layer widths being given by $\mathbf{n} = (n_0, n_1, \dots, n_L, n_{L+1})^\top \in \mathbb{N}^{L+2}$, where n_0 and n_{L+1} are input and output dimensions. A feed-forward BNN is a mapping from the input region $\mathcal{R}_{n_0} \subset \mathbb{R}^{n_0}$ to the output set $\llbracket 1, n_{L+1} \rrbracket$ realized via successive compositions of several internal blocks $(\mathbf{B}_i)_{i=1, \dots, L}$ and an output block \mathbf{B}_o :

$$\text{BNN} : \mathcal{R}_{n_0} \rightarrow \llbracket 1, n_{L+1} \rrbracket$$

$$\mathbf{x}^0 \mapsto \text{BNN}(\mathbf{x}^0) := \text{argmax}(\mathbf{B}_o(\mathbf{B}_L(\dots(\mathbf{B}_1(\mathbf{x}^0))))). \quad (1)$$

For any $i \in \llbracket 1, L \rrbracket$, the internal block \mathbf{B}_i implements successively three different operations: affine transformation, batch normalization and point-wise binarization, so that its output vector, denoted by \mathbf{x}^i , belongs to $\{-1, 1\}^{n_i}$. These operations are described by a set of trainable parameters:

$$(\mathbf{W}^{[i+1]}, \mathbf{b}^{[i+1]})_{i \in \llbracket 0, L \rrbracket} \in \{-1, 0, 1\}^{n_{i+1} \times n_i} \times \mathbb{R}^{n_{i+1}}, \quad (\gamma_j^{[i]}, \beta_j^{[i]}, \mu_j^{[i]}, \sigma_j^{2, [i]})_{i \in \llbracket 1, L \rrbracket} \in (\mathbb{R}^{n_i})^4. \quad (2)$$

| Steps | Input | Output | Transformation |
|---------------------|--|------------------------------------|---|
| Linearization | $\mathbf{x}^{i-1} \in \{-1, 1\}^{n_{i-1}}$ | $\mathbf{y} \in \mathbb{R}^{n_i}$ | $\mathbf{y} = \mathbf{W}^{[i]} \mathbf{x}^{i-1} + \mathbf{b}^{[i]}$ |
| Batch-Normalization | $\mathbf{y} \in \mathbb{R}^{n_i}$ | $\mathbf{z} \in \mathbb{R}^{n_i}$ | $z_j = \gamma_j^{[i]} \left(\frac{y_j - \mu_j^{[i]}}{\sqrt{\sigma_j^{2, [i]} + \epsilon}} \right) + \beta_j^{[i]}$ |
| Binarization | $\mathbf{z} \in \mathbb{R}^{n_i}$ | $\mathbf{x}^i \in \{-1, 1\}^{n_i}$ | $\mathbf{x}^i = \text{sign}(\mathbf{z})$ |

Table 1: Structure of an internal block \mathbf{B}_i .

The output block \mathbf{B}_o applies a softmax transformation to the affinely-transformed outputs of the last hidden layer, i.e., for each $j \in \llbracket 1, n_{L+1} \rrbracket$,

$$\mathbf{x}_j^{L+1} = \frac{\exp(z_j)}{\sum_{k=1}^{n_{L+1}} \exp(z_k)}, \quad \text{where } z_j = \mathbf{W}_{(j,:)}^{[L+1]} \mathbf{x}^L + \mathbf{b}_j^{[L+1]}. \quad (3)$$

Deriving a Tighter First-Order Relaxation - $\tau_{\text{tighter,cs}}^1$

⇒ For each $i \in \llbracket 1, L \rrbracket$, we replace the constraint defined in (4b) by the following two constraints:

$$\left\{ \begin{aligned} \hat{\mathbf{g}}_i^1(\mathbf{x}^i, \mathbf{x}^{i-1}) &:= (\mathbf{x}^i + \mathbf{1}) \odot (\mathbf{W}^{[i]} \mathbf{x}^{i-1} + \mathbf{b}^{[i]}) \geq \mathbf{0}, & (6a) \\ \hat{\mathbf{g}}_i^2(\mathbf{x}^i, \mathbf{x}^{i-1}) &:= (\mathbf{x}^i - \mathbf{1}) \odot (\mathbf{W}^{[i]} \mathbf{x}^{i-1} + \mathbf{b}^{[i]}) \geq \mathbf{0}. & (6b) \end{aligned} \right.$$

⇒ Add the following two *redundant* quadratic constraints, i.e., *tautologies*, to the optimization problem (5):

$$\left\{ \begin{aligned} \hat{\mathbf{g}}_i^1(\mathbf{x}^i, \mathbf{x}^{i-1}) &:= (\mathbf{x}^i + \mathbf{1}) \odot (\text{inv}(\mathbf{W}^{[i]}) - \mathbf{W}^{[i]} \mathbf{x}^{i-1}) \geq \mathbf{0}, & (7a) \\ \hat{\mathbf{g}}_i^2(\mathbf{x}^i, \mathbf{x}^{i-1}) &:= (\mathbf{1} - \mathbf{x}^i) \odot (\text{inv}(\mathbf{W}^{[i]}) + \mathbf{W}^{[i]} \mathbf{x}^{i-1}) \geq \mathbf{0}, & (7b) \end{aligned} \right.$$

Theorem

For any BNN verification problem, $\tau_{\text{tighter,cs}}^1 = \tau_{\text{tighter}}^1 \geq \tau^1$. If $L \geq 2$, there exists an affine f such that the inequality is strict. We also have $\tau_{\text{tighter,cs}}^1 \geq \tau_{\text{LIP}}^1$ for any affine f .

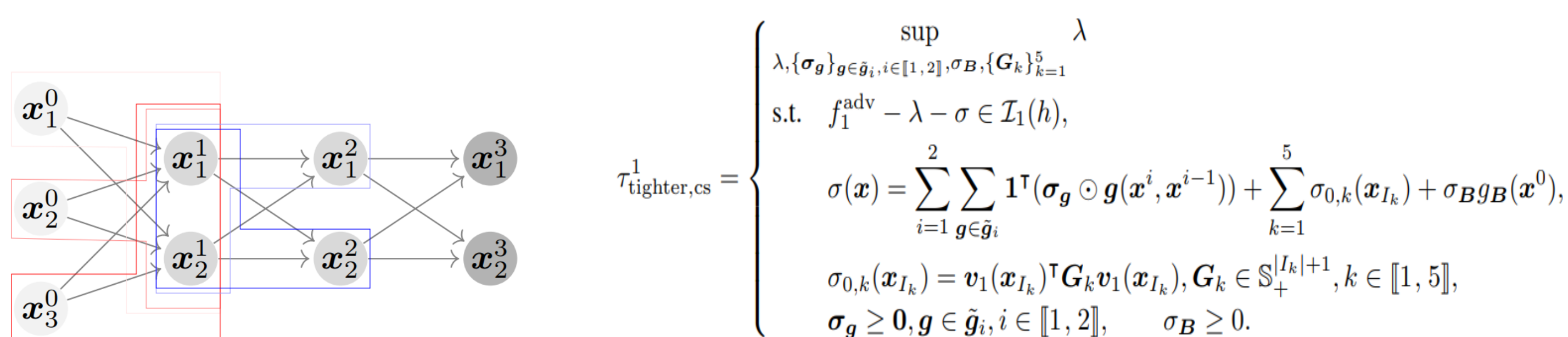


Figure 1: (left) A toy BNN with $L = 2$ and $(n_0, n_1, n_2, n_3) = (3, 2, 2, 2)$. The cliques I_1, I_2, I_3 (red polygons) and I_4, I_5 (blue polygons) are used to compute $\tau_{\text{tighter,cs}}^1$ (right). We let $\hat{\mathbf{g}}_i = \{\hat{\mathbf{g}}_i^1, \hat{\mathbf{g}}_i^2, \hat{\mathbf{g}}_i^1, \hat{\mathbf{g}}_i^2\}$.

Sparse SDP for BNN Robustness Verification

Consider a sequence of vector-valued functions $(\mathbf{h}_i, \mathbf{g}_i)_{i \in \llbracket 1, L \rrbracket}$ such that

$$\mathbf{x}^i := \text{sign}(\mathbf{W}^{[i]} \mathbf{x}^{i-1} + \mathbf{b}^{[i]}) \implies \left\{ \begin{aligned} \mathbf{h}_i(\mathbf{x}^i) &:= \mathbf{x}^i \odot \mathbf{x}^i - \mathbf{1} = \mathbf{0}, & (4a) \\ \mathbf{g}_i(\mathbf{x}^i, \mathbf{x}^{i-1}) &:= \mathbf{x}^i \odot (\mathbf{W}^{[i]} \mathbf{x}^{i-1} + \mathbf{b}^{[i]}) \geq \mathbf{0}, & (4b) \end{aligned} \right.$$

Suppose that the input perturbation region $\mathcal{B} \subseteq \mathbb{R}^{n_0}$ can be encoded via positivity conditions on (at most quadratic) polynomials $\mathbf{x}^0 \mapsto \mathbf{g}_{\mathcal{B}}(\mathbf{x}^0)$. Then, the *standard* form BNN verification problems becomes:

$$\tau := \begin{cases} \min_{\mathbf{x}^0, \mathbf{x}^1, \dots, \mathbf{x}^L} f(\mathbf{x}^0, \mathbf{x}^1, \dots, \mathbf{x}^L) & (5a) \\ \text{s.t. } \mathbf{h}_i(\mathbf{x}^i) = \mathbf{0}, i \in \llbracket 1, L \rrbracket, & (5b) \\ \mathbf{g}_i(\mathbf{x}^i, \mathbf{x}^{i-1}) \geq \mathbf{0}, i \in \llbracket 1, L \rrbracket, & (5c) \\ \mathbf{g}_{\mathcal{B}}(\mathbf{x}^0) \geq \mathbf{0}, & (5d) \end{cases}$$

where f is linear/quadratic, e.g., $f = f_k^{\text{dv}}(\mathbf{x}^0, \mathbf{x}^1, \dots, \mathbf{x}^L) := \langle \mathbf{W}_{(j,:)}^{[L+1]} - \mathbf{W}_{(k,:)}^{[L+1]}, \mathbf{x}^L \rangle + \mathbf{b}_j^{[L+1]} - \mathbf{b}_k^{[L+1]}$.

Exploiting Sparsity

Let us suppose that $\llbracket 1, n \rrbracket = \cup_{k=1}^p I_k$ with I_k not necessarily disjoint. The subsets I_k , called *cliques*, correspond to the subsets of variables $\mathbf{x}_{I_k} := \{x_i, i \in I_k\}$. An instance of the BNN robustness verification problem of the form (5) exhibits *correlative sparsity* since

- There exist $(f_k)_{k \in \llbracket 1, p \rrbracket}$ such that $f = \sum_{k=1}^p f_k$, with $f_k \in \mathbb{R}[\mathbf{x}_{I_k}]$,
- The polynomials \mathbf{g} can be split into disjoint sets J_k , such that $\mathbf{g}_i(\cdot)_j \in J_k$ if and only if $\mathbf{g}_i(\cdot)_j \in \mathbb{R}[\mathbf{x}_{I_k}]$. Moreover, $\mathbf{g}_{\mathcal{B}} \in J_k$ for $k \in \llbracket 1, p \rrbracket$. Since $\mathbf{h}_i(\cdot)_j$ only depends on \mathbf{x}_j^i , the overall sparsity structure is induced by inequality constraints that mimic the cascading BNN structure.

The hierarchy of correlatively sparse SDP relaxations is then given by

$$\tau_{\text{cs}}^d := \sup_{\lambda, \sigma} \{ \lambda \in \mathbb{R} \mid f - \lambda - \sum_{k=1}^p \sigma_k \in \mathcal{I}_d(\mathbf{h}), \sigma_k \in \mathcal{Q}_d(\{\mathbf{g}_i(\cdot)_j \in J_k\}) \}.$$

Some Numerical Results

⇒ Verification against $\|\cdot\|_\infty$ -attacks ↓

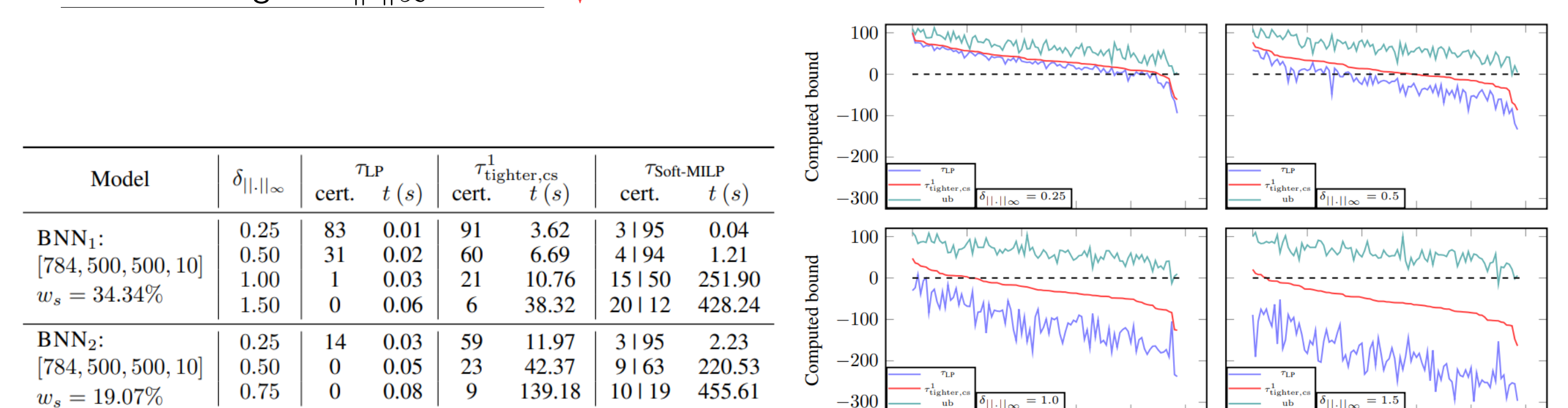


Figure 2: (left) Performance comparison for different models and input regions, given by $\delta_{\|\cdot\|_\infty} = 127.5\epsilon$. (right) Comparing τ_{LIP} and $\tau_{\text{tighter,cs}}^1$ bounds for BNN₁ and different $\delta_{\|\cdot\|_\infty}$. The relative improvement is up to 55%.

⇒ Verification against $\|\cdot\|_2$ -attacks ↓

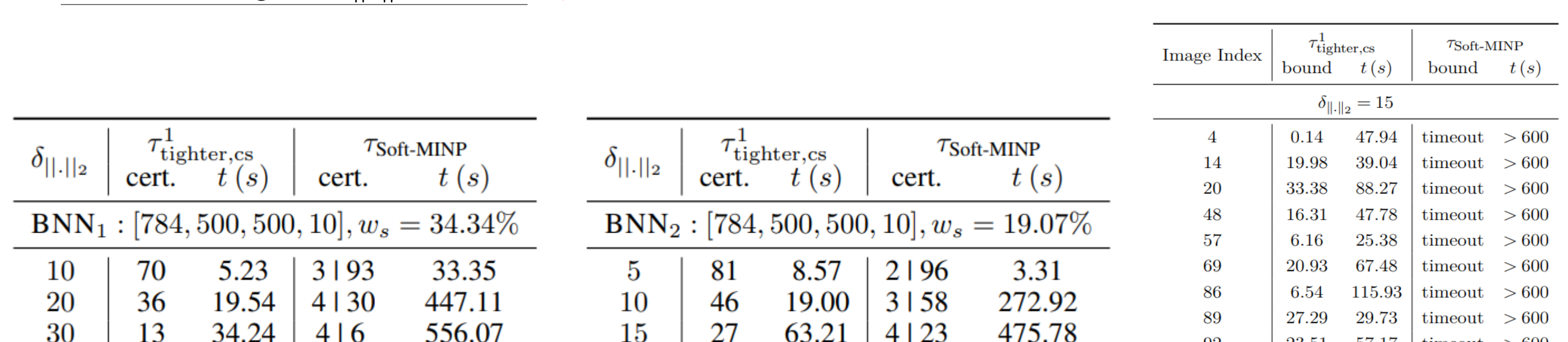


Figure 3: (left, center) Performance comparison for different BNN models and input regions, given by $\delta_{\|\cdot\|_2} = 255\epsilon$. (right) Verification against $\|\cdot\|_2$ -attacks: illustrating the significant speedup for specific instances.

References

- Jianting Yang, Srećko Đurašinović, Jean-Bernard Lasserre, Victor Magron and Jun Zhao. *Verifying Properties of Binary Neural Networks Using Sparse Polynomial Optimization*. arXiv preprint arXiv:2405.17049, 2024.
- Victor Magron and Jie Wang. *Sparse Polynomial Optimization: Theory and Practice*. World Scientific, 2023.
- Victor Magron and Jie Wang. *TSSOS: a Julia library to exploit sparsity for large-scale polynomial optimization*. arXiv preprint arXiv:2103.00915, 2021.

