
Sound and Complete Verification of Polynomial Networks

Elias Abad Rocamora*

LIONS, EPFL

Lausanne, Switzerland

abad.elias00@gmail.com

Mehmet Fatih Sahin

LIONS, EPFL

Lausanne, Switzerland

mehmet.sahin@epfl.ch

Fanghui Liu

LIONS, EPFL

Lausanne, Switzerland

fanghui.liu@epfl.ch

Grigorios G Chrysos

LIONS, EPFL

Lausanne, Switzerland

grigorios.chrysos@epfl.ch

Volkan Cevher

LIONS, EPFL

Lausanne, Switzerland

volkan.cevher@epfl.ch

Abstract

Polynomial Networks (PNs) have demonstrated promising performance on face and image recognition recently. However, robustness of PNs is unclear and thus obtaining certificates becomes imperative for enabling their adoption in real-world applications. Existing verification algorithms on ReLU neural networks (NNs) based on classical branch and bound (BaB) techniques cannot be trivially applied to PN verification. In this work, we devise a new bounding method, equipped with BaB for global convergence guarantees, called Verification of Polynomial Networks or VPN for short. One key insight is that we obtain much tighter bounds than the interval bound propagation (IBP) and DeepT-Fast [Bonaert et al., 2021] baselines. This enables sound and complete PN verification with empirical validation on MNIST, CIFAR10 and STL10 datasets. We believe our method has its own interest to NN verification. The source code is publicly available at <https://github.com/megaelius/PNVerification>.

1 Introduction

Polynomial Networks (PNs) have demonstrated promising performance across image recognition and generation [Chrysos et al., 2021, Chrysos and Panagakis, 2020] being state-of-the-art on large-scale face recognition². Unlike the conventional Neural Networks (NNs), where non-linearity is introduced with the use of activation functions [LeCun et al., 2015], PNs are able to learn non-linear mappings without the need of activation functions by exploiting multiplicative interactions (Hadamard products). Recent works have uncovered interesting properties of PNs in terms of model expressivity [Fan et al., 2021] and spectral bias [Choraria et al., 2022]. However, one critical issue before considering PNs for real-world applications is their robustness.

Neural networks are prone to small (often imperceptible to the human eye), but malicious perturbations in the input data points [Szegedy et al., 2014, Goodfellow et al., 2015]. Those perturbations can have a detrimental effect on image recognition systems, e.g., as illustrated in face recognition [Goswami et al., 2019, Zhong and Deng, 2019, Dong et al., 2019, Li et al., 2020]. Guarding against such attacks has so far proven futile [Shafahi et al., 2019, Dou et al., 2018]. Instead, a flurry of research

*Work developed during an exchange coming from Universitat Politècnica de Catalunya (UPC), Spain. Currently at Universidad Carlos III de Madrid (UC3M).

²<https://paperswithcode.com/sota/face-verification-on-megaface>

has been published on certifying robustness of NNs against this performance degradation [Katz et al., 2017, Ehlers, 2017, Tjeng et al., 2019, Bunel et al., 2020b, Wang et al., 2021, Ferrari et al., 2022]. However, most of the verification algorithms for NNs are developed for the ReLU activation function by exploiting its piecewise linearity property and might not trivially extend to other nonlinear activation functions [Wang et al., 2021]. Indeed, Zhu et al. [2022] illustrate that guarding PNs against adversarial attacks is challenging. Therefore, we pose the following question:

Can we obtain certifiable performance for PNs against adversarial attacks?

In this work, we answer affirmatively and provide a method for the verification of PNs. Concretely, we take advantage of the twice-differentiable nature of PNs to build a lower bounding method based on α -convexification [Adjiman and Floudas, 1996], which is integrated into a Branch and Bound (BaB) algorithm [Land and Doig, 1960] to guarantee completeness of our verification method. In order to use α -convexification, a lower bound α of the minimum eigenvalue of the Hessian matrix over the possible perturbation set is needed. We use interval bound propagation together with the theoretical properties of the lower bounding Hessian matrix [Adjiman et al., 1998], in order to develop an algorithm to efficiently compute α .

Our *contributions* can be summarized as follows: (i) We propose the first algorithm for the verification of PNs. (ii) We thoroughly analyze the performance of our method by comparing it with a black-box solver, with an interval bound propagation (IBP) BaB algorithm and with the zonotope-based abstraction method DeepT-Fast [Bonaert et al., 2021]. (iii) We empirically show that using α -convexification for lower bounding provides tighter bounds than IBP and DeepT-Fast for PN verification. To encourage the community to improve the verification of PNs, we make our code publicly available in <https://github.com/megaelius/PNVerification>. The proposed approach can practically verify PNs, while it could also theoretically be applied for sound and complete verification of any twice-differentiable network.

Notation: We use the shorthand $[n] := \{1, 2, \dots, n\}$ for a positive integer n . We use bold capital (lowercase) letters, e.g., \mathbf{X} (\mathbf{x}) for representing matrices (vectors). The j^{th} column of a matrix \mathbf{X} is given by $\mathbf{x}_{:j}$. The element in the i^{th} row and j^{th} column is given by x_{ij} , similarly, the i^{th} element of a vector \mathbf{x} is given by x_i . The element-wise (Hadamard) product, symbolized with $*$, of two matrices (or vectors) in $\mathbb{R}^{d_1 \times d_2}$ (or \mathbb{R}^d) gives another matrix (or vector) in $\mathbb{R}^{d_1 \times d_2}$ (or \mathbb{R}^d). The ℓ_∞ norm of a vector $\mathbf{x} \in \mathbb{R}^d$ is given by: $\|\mathbf{x}\|_\infty = \max_{i \in [d]} |x_i|$. Lastly, the operators \mathcal{L} and \mathcal{U} give the lower and upper bounds of a scalar, vector or matrix function by IBP, see Section 3.1.

Roadmap: We provide the necessary background by introducing the PN architecture and formalizing the Robustness Verification problem in Section 2. Section 3 provides a *sound* and *complete* method called VPN to tackle PN verification problem. Section 4 is devoted to experimental validation. Additional experiments, details and proofs are deferred to the appendix.

2 Background

We give an overview of the PN architecture in Section 2.1 and the robustness verification problem in Section 2.2.

2.1 Polynomial Networks (PNs)

Polynomial Networks (PNs) are inspired by the fact that any smooth function can be approximated via a polynomial expansion [Stone, 1948]. However, the number of parameters increases exponentially with the polynomial degree, which makes it intractable to use high degree polynomials for high-dimensional data problems such as image classification where the input can be in the order of 10^5 [Deng et al., 2009]. Chrysos et al. [2021] introduce a joint factorization of polynomial coefficients in a low-rank manner, reducing the number of parameters to linear with the polynomial degree and allowing the expression as a neural network (NN). We briefly recap one fundamental factorization below.

Let N be the polynomial degree, $\mathbf{z} \in \mathbb{R}^d$ be the input vector, d , k and o be the input, hidden and output sizes, respectively. The recursive equation of PNs can be expressed as:

$$\mathbf{x}^{(n)} = (\mathbf{W}_{[n]}^\top \mathbf{z}) * \mathbf{x}^{(n-1)} + \mathbf{x}^{(n-1)}, \forall n \in [N], \quad (1)$$

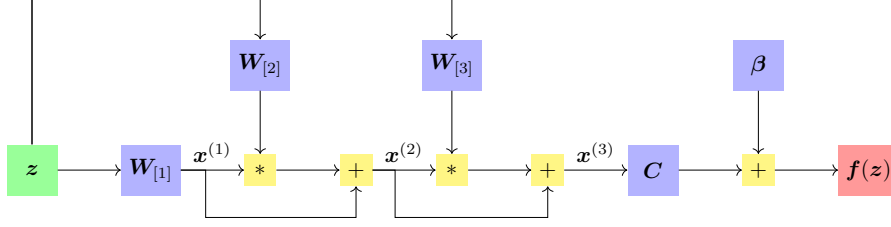


Figure 1: Third degree PN architecture. Blue boxes depict learnable parameters, yellow depict mathematical operations, the green and red boxes are the input and the output respectively. Note that no activation functions are involved, only element-wise (Hadamard) products $*$ and additions $+$. This figure represents the recursive formula of Eq. (1).

where $x^{(1)} = W_{[1]}^\top z$, $f(z) = Cx^{(N)} + \beta$ and $*$ denotes the Hadamard product. $W_{[n]} \in \mathbb{R}^{d \times k}$ and $C \in \mathbb{R}^{o \times k}$ are weight matrices, $\beta \in \mathbb{R}^o$ is a bias vector. A graphical representation of a third degree PN architecture corresponding to Eq. (1) can be found in Fig. 1. Further details on the factorization (as well as other factorizations) are deferred to the Appendix C.1 (Appendix C.2).

2.2 Robustness Verification

Robustness verification [Bastani et al., 2016, Liu et al., 2021] consists of verifying that a property regarding the input and output of a NN is satisfied, e.g. checking whether or not a small perturbation in the input will produce a change in the network output that makes it classify the input into another class. Let $f : [0, 1]^d \rightarrow \mathbb{R}^o$ be a function, e.g., a NN or a PN, that classifies the input z into a class c , such that $c = \arg \max f(z)$. Our target is to verify that for any input satisfying a set of constraints C_{in} , the output of the network will satisfy a set of output constraints C_{out} . Mathematically,

$$z \in C_{\text{in}} \implies f(z) \in C_{\text{out}}. \quad (2)$$

In this work we focus on *adversarial robustness* [Szegedy et al., 2014, Carlini and Wagner, 2017] in classification. Given an observation z_0 , let $t = \arg \max f(z_0)$ be the correct class, our goal is to check whether every input in a neighbourhood of z_0 , is classified as t . In this work, we focus on adversarial attacks restricted to neighbourhoods defined in terms of ℓ_∞ norm, which is a popular norm-bounded attack in the verification community [Liu et al., 2021]. Then, the constraint sets become:

$$\begin{aligned} C_{\text{in}} &= \{z : \|z - z_0\|_\infty \leq \epsilon, z_i \in [0, 1], \forall i \in [d]\} \\ &= \{z : \max\{0, z_{0i} - \epsilon\} \leq z_i \leq \min\{1, z_{0i} + \epsilon\}, \forall i \in [d]\} \\ C_{\text{out}} &= \{y : y_t > y_j, \forall j \neq t\}. \end{aligned} \quad (3)$$

In other words, we need an algorithm that given a function f , an input z_0 and an adversarial budget ϵ , checks whether Eq. (2) is satisfied. In the case of ReLU NNs, this has been proven to be an NP-complete problem [Katz et al., 2017]. This can be reformulated as a constrained optimization problem. For every adversarial class $\gamma \neq t = \arg \max f(z_0)$, we can solve:

$$\min_z g(z) = f(z)_t - f(z)_\gamma \quad \text{s.t.} \quad z \in C_{\text{in}}. \quad (4)$$

If the solution z^* with $v^* = f(z^*)_t - f(z^*)_\gamma \leq f(z)_t - f(z)_\gamma, \forall z \in C_{\text{in}}$ satisfies $v^* > 0$ then robustness is verified for the adversarial class γ .

There are two main properties that a verification algorithm admits: *soundness* and *completeness*. An algorithm is *sound* (*complete*) if every time it verifies (falsifies) a property, it is guaranteed to be the correct answer. In practice, when an algorithm is guaranteed to provide the exact global minima of Eq. (4), i.e., v^* , it is said to be *sound* and *complete* (usually referred in the literature as simply *complete* [Ferrari et al., 2022]), whereas if a lower bound of it is provided $\hat{v}^* \leq v^*$, the algorithm is *sound* but not *complete*. In our work, we will not consider just *complete* verification, which simply aims at looking for adversarial examples, e.g., Madry et al. [2018]. For a deeper discussion on *soundness* and *completeness*, we refer to Liu et al. [2021].

3 Method

Our method, called VPN, can be categorized in the the Branch and Bound (BaB) framework [Land and Doig, 1960], a well known approach to global optimization [Horst and Tuy, 1996] and NN verification [Bunel et al., 2020b]. This kind of algorithms ensures finding a global minima of the problem in Eq. (4) by recursively splitting the original feasible set into smaller sets (branching) where upper and lower bounds of the global minima are computed (bounding). This mechanism can be used to discard subsets where the global minima cannot be achieved (its lower bound is greater than the upper bound of another subset).

Our method is based on a variant of BaB algorithm, i.e., α -BaB [Adjiman et al., 1998], which is characterized for using α -convexification [Adjiman and Floudas, 1996] for computing a lower bound of the global minima of each subset. To be specific, α -convexification aims to obtain a convex lower bounding function of any twice-differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$. In Adjiman and Floudas [1996], they propose two methods:

- **Uniform diagonal shift (single α)**

$$g_\alpha(\mathbf{z}; \alpha, \mathbf{l}, \mathbf{u}) = g(\mathbf{z}) + \alpha \sum_{i=1}^d (z_i - l_i)(z_i - u_i), \quad (5)$$

is its α -convexified version, note that z_i is the i^{th} element of vector \mathbf{z} . Let $\mathbf{H}_g(\mathbf{z}) = \nabla_{\mathbf{z}\mathbf{z}}^2 g(\mathbf{z})$ be the Hessian matrix of g , g_α is convex in $\mathbf{z} \in [\mathbf{l}, \mathbf{u}]$ for $\alpha \geq \max\{0, -\frac{1}{2} \min\{\lambda_{\min}(\mathbf{H}_g(\mathbf{z})) : \mathbf{z} \in [\mathbf{l}, \mathbf{u}]\}\}$, where λ_{\min} is the minimum eigenvalue. Moreover, it holds that $g_\alpha(\mathbf{z}; \alpha, \mathbf{l}, \mathbf{u}) \leq g(\mathbf{z}), \forall \mathbf{z} \in [\mathbf{l}, \mathbf{u}]$.

- **Non-uniform diagonal shift (multiple α 's)**

$$g_\alpha(\mathbf{z}; \boldsymbol{\alpha}, \mathbf{l}, \mathbf{u}) = g(\mathbf{z}) + \sum_{i=1}^d \alpha_i (z_i - l_i)(z_i - u_i), \quad (6)$$

is its α -convexified version. In Adjiman et al. [1998], they show that for any vector $\mathbf{d} > \mathbf{0}$, setting

$$\alpha_i \geq \max \left\{ 0, -\frac{1}{2} \left(\mathcal{L}(h_g(\mathbf{z}))_{ii} - \sum_{j \neq i} \max\{|\mathcal{L}(h_g(\mathbf{z}))_{ij}|, |\mathcal{U}(h_g(\mathbf{z}))_{ij}|\} \frac{d_j}{d_i} \right) \right\} \quad (7)$$

makes g_α convex in $\mathbf{z} \in [\mathbf{l}, \mathbf{u}]$. The choice of the vector \mathbf{d} is arbitrary, but affects the final result. For example, taking $\mathbf{d} = \mathbf{u} - \mathbf{l}$ yields better results than $\mathbf{d} = \mathbf{1}$ in Adjiman et al. [1998]. We need to remark that, though Eq. (5) is a special case of Eq. (6), the lower bound of α in Eq. (5) via minimum eigenvalue of the lower bounding Hessian matrix cannot be regarded as a special case of Eq. (7).

To make PN verification feasible via α -convexification, we need to study IBP for PNs and design an efficient estimate on the α (α in the case of Non-uniform diagonal shift) parameter, which are our main technical contributions in the algorithmic aspect. In our case, every feasible set, starting with the input set \mathcal{C}_{in} (Eq. (3)), is split by taking the widest variable interval and dividing it in two by the middle point. This is a rather simple, but theoretically powerful strategy, see Lemma 8 in Appendix F. Then, the upper bound of each subproblem is given by applying standard Projected Gradient Descent (PGD) [Kelley, 1999] over the original objective function. This is a common approach to find adversarial examples [Madry et al., 2018], but as the objective is non-convex, it is not sufficient for sound and complete verification. The lower bound is given by applying PGD over the α -convexified objective g_α (g_α), as it is convex, PGD converges to the global minima and a lower bound of the original objective. The α (α) parameter is computed only once per verification problem. Further details on the algorithm and the proof of convergence of Eq. (4) exist in Appendix D. A schematic of our method is available in Fig. 2.

In Sections 3.1 to 3.3, we detail our method under the uniform diagonal shift case to compute a lower bound on the minimum eigenvalue of the Hessian matrix into three main components: interval propagation, lower bounding Hessians, and fast estimation on such lower bounding via power method. To conclude the description of our method, in Section 3.4 we describe the α estimation for the non-uniform diagonal shift case.

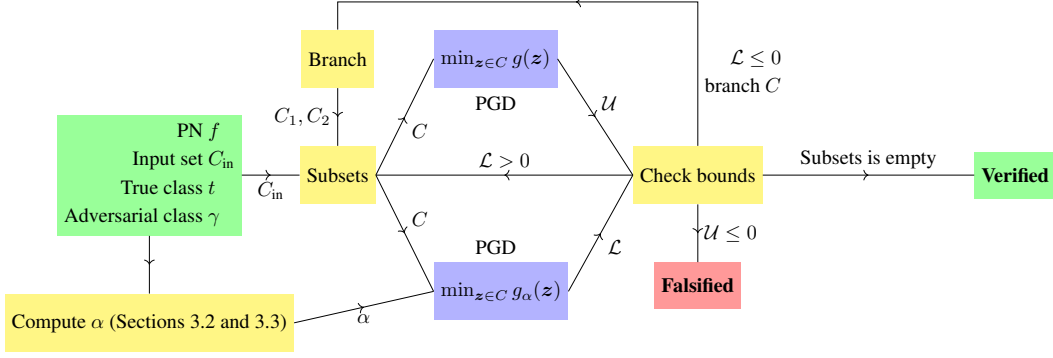


Figure 2: Overview of our branch and bound verification algorithm. Given a trained PN f , an input set C_{in} , the true class t and an adversarial class γ , we check if an adversarial example exists (**Falsified**) or not (**Verified**). Note that the branching of a subset C provides two smaller subsets C_1 and C_2 . Also note that when $\mathcal{L} > 0$, no subset is added to subsets.

3.1 Interval Bound Propagation through a PN

Interval bound propagation (IBP) is a key ingredient of our verification algorithm. Suppose we have an input set defined by an ℓ_∞ -norm ball like in Eq. (3). This set can be represented as a vector of intervals $[l, u] = ([l_1, u_1]^\top, [l_2, u_2]^\top, \dots, [l_d, u_d]^\top) \in \mathbb{R}^{d \times 2}$, where $[l_i, u_i]$ are the lower and upper bound for the i^{th} coordinate. Let \mathcal{L} and \mathcal{U} be the lower and upper bound IBP operators. Given this input set, we would like to obtain bounds on the output of the network $(f(z))_i$, the gradient $(\nabla_z f(z))_i$, and the Hessian $(\nabla_{zz}^2 f(z))_i$ for any $z \in [l, u]$. The operators $\mathcal{L}(g(z))$ and $\mathcal{U}(g(z))$ of any function $g : \mathbb{R}^d \rightarrow \mathbb{R}$ satisfy:

$$\mathcal{L}(g(z)) \leq g(z), \quad \mathcal{U}(g(z)) \geq g(z), \quad \forall z \in [l, u]. \quad (8)$$

We will define these upper and lower bound operators in terms of the operations present in a PN. Using interval propagation [Moore et al., 2009], denote the positive part $w_i^+ = \max\{0, w_i\}$ and the negative part $w_i^- = \min\{0, w_i\}$, $h_i(z)$ as a real-valued function of z , $i \in [d]$, we can define:

$$\begin{aligned} \text{Identity} & \begin{cases} \mathcal{L}(z_i) = l_i \\ \mathcal{U}(z_i) = u_i \end{cases} \\ \text{linear mapping} & \begin{cases} \mathcal{L}(\sum_i w_i h_i(z)) = \sum_i w_i^+ \mathcal{L}(h_i(z)) + w_i^- \mathcal{U}(h_i(z)) \\ \mathcal{U}(\sum_i w_i h_i(z)) = \sum_i w_i^- \mathcal{L}(h_i(z)) + w_i^+ \mathcal{U}(h_i(z)) \end{cases} \\ \text{multiplication} & \begin{cases} S = \begin{Bmatrix} \mathcal{L}(h_1(z)) \mathcal{L}(h_2(z)), \\ \mathcal{L}(h_1(z)) \mathcal{U}(h_2(z)), \\ \mathcal{U}(h_1(z)) \mathcal{L}(h_2(z)), \\ \mathcal{U}(h_1(z)) \mathcal{U}(h_2(z)) \end{Bmatrix}, |S| = 4 \\ \mathcal{L}(h_1(z) h_2(z)) = \min S, \\ \mathcal{U}(h_1(z) h_2(z)) = \max S, \end{cases} \end{aligned} \quad (9)$$

where $|\cdot|$ is the set cardinality. Note that the set S is equivalent to:

$$S = \{ab \mid \forall a \in \{\mathcal{L}(h_1(z)), \mathcal{U}(h_1(z))\}, \forall b \in \{\mathcal{L}(h_2(z)), \mathcal{U}(h_2(z))\}\}.$$

With these basic operations, one can define bounds on any intermediate output, gradient or Hessian of a PN. For instance, the lower bound on the recursive formula from Eq. (1) can be expressed as:

$$\mathcal{L}(x_i^{(n)}) = \mathcal{L}((w_{[n]:i}^\top x) x_i^{(n-1)} + x_i^{(n-1)}) = \mathcal{L}((w_{[n]:i}^\top x + 1) x_i^{(n-1)}), \quad \forall i \in [k], n \in \{2, \dots, N\}, \quad (10)$$

which only consists on a linear mapping and a multiplication of intervals. We extend the upper and lower bound ($\mathcal{L}(\cdot)$ and $\mathcal{U}(\cdot)$) operators to vectors and matrices in an entry-wise style:

$$\mathcal{L}(g(z)) = \begin{bmatrix} \mathcal{L}(g(z)_1) \\ \mathcal{L}(g(z)_2) \\ \vdots \\ \mathcal{L}(g(z)_m) \end{bmatrix} \in \mathbb{R}^m, \quad \mathcal{L}(G(z)) = \begin{bmatrix} \mathcal{L}(g(z)_{11}) & \cdots & \mathcal{L}(g(z)_{1m}) \\ \vdots & \ddots & \vdots \\ \mathcal{L}(g(z)_{m1}) & \cdots & \mathcal{L}(g(z)_{mm}) \end{bmatrix} \in \mathbb{R}^{m \times m}. \quad (11)$$

Note that Eq. (11) is not limited to squared matrices and can hold for arbitrary matrix dimensions. One can directly use IBP to obtain bounds on the verification objective from Eq. (4) with a single forward pass of the bounds through the network and obtaining $\mathcal{L}(g(z)) = \mathcal{L}(f(z)_t) - \mathcal{U}(f(z)_\gamma)$. In fact IBP is a common practice in NN verification to obtain fast bounds [Wang et al., 2018a].

3.2 Lower bound of the minimum eigenvalue of the Hessian

Here we describe our method to compute a lower bound on the minimum eigenvalue of the Hessian matrix in the feasible set. Before deriving the lower bound, we need the first and second order partial derivatives of PNs.

Let $g(z) = f(z)_t - f(z)_a$ be the objective function for $t = \arg \max f(z_0)$ and any $a \neq t$. In order to compute the parameter α for performing α -convexification, we need to know the structure of our objective function. In this section we compute the first and second order partial derivatives of the PN. The gradient and Hessian matrices of the objective function (see Eq. (4)) are given by:

$$\nabla_{\mathbf{z}} g(\mathbf{z}) = \sum_{i=1}^k (c_{ti} - c_{\gamma i}) \nabla_{\mathbf{z}} x_i^{(N)}, \quad \mathbf{H}_g(\mathbf{z}) = \sum_{i=1}^k (c_{ti} - c_{\gamma i}) \nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(N)}. \quad (12)$$

We now define the gradients $\nabla_{\mathbf{z}} x_i^{(n)}$ and Hessians $\nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n)}$ of Eq. (1) in a recursive way:

$$\nabla_{\mathbf{z}} x_i^{(n)} = \mathbf{w}_{[n]:i} \cdot x_i^{(n-1)} + (\mathbf{w}_{[n]:i}^\top \mathbf{z} + 1) \cdot \nabla_{\mathbf{z}} x_i^{(n-1)} \quad (13)$$

$$\nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n)} = \nabla_{\mathbf{z}} x_i^{(n-1)} \mathbf{w}_{[n]:i}^\top + \{\nabla_{\mathbf{z}} x_i^{(n-1)} \mathbf{w}_{[n]:i}^\top\}^\top + (\mathbf{w}_{[n]:i}^\top \mathbf{z} + 1) \nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n-1)}, \quad (14)$$

with $\nabla_{\mathbf{z}} x_i^{(1)} = \mathbf{w}_{[1]:i}$ and $\nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(1)} = \mathbf{0}_{d \times d}$ being an all-zero matrix. In the next, we are ready to compute a lower bound on the minimum eigenvalue of the Hessian matrix in the feasible set.

Firstly, for any $\mathbf{z} \in [\mathbf{l}, \mathbf{u}]$ and any polynomial degree N , we can express the set of possible Hessians $\mathcal{H} = \{\mathbf{H}_g(\mathbf{z}) : \mathbf{z} \in [\mathbf{l}, \mathbf{u}]\}$ as an interval matrix. An interval matrix is a tensor $[\mathbf{M}] \in \mathbb{R}^{d \times d \times 2}$ where every position $[m]_{ij} = [\mathcal{L}(m_{ij}), \mathcal{U}(m_{ij})]$ is an interval. Therefore, if $\mathbf{H}_g(\mathbf{z})$ is bounded for $\mathbf{z} \in [\mathbf{l}, \mathbf{u}]$, then we can represent $\mathcal{H} = \{\mathbf{H}_g(\mathbf{z}) : \mathbf{H}_g(\mathbf{z}) \in [\mathbf{M}]\} = \{\mathbf{H}_g(\mathbf{z}) : \mathcal{L}(m_{ij}) \leq H_g(\mathbf{z})_{ij} \leq \mathcal{U}(m_{ij}), \forall i, j \in [d]\}$.

Let $\mathcal{L}(\mathbf{M})$ and $\mathcal{U}(\mathbf{M})$ be the element-wise lower and upper bounds of a Hessian matrix, the lower bounding Hessian is defined as follows:

$$\mathbf{L}_H = \frac{\mathcal{L}(\mathbf{M}) + \mathcal{U}(\mathbf{M})}{2} + \text{diag} \left(\frac{\mathcal{L}(\mathbf{M})\mathbf{1} - \mathcal{U}(\mathbf{M})\mathbf{1}}{2} \right), \quad (15)$$

where $\mathbf{1}$ is an all-one vector and $\text{diag}(\mathbf{v})$ is a diagonal matrix with the vector \mathbf{v} in the diagonal. Described in Adjiman et al. [1998], this matrix satisfies that $\lambda_{\min}(\mathbf{L}_H) \leq \lambda_{\min}(\mathbf{H}_g(\mathbf{z})), \forall \mathbf{H}_g(\mathbf{z}) \in \mathcal{H}, \mathbf{z} \in [\mathbf{l}, \mathbf{u}]$.

Then, we can obtain the spectral radius $\rho(\mathbf{L}_H)$ with a power method [Mises and Pollaczek-Geiringer, 1929]. As the spectral radius satisfies $\rho(\mathbf{L}_H) \geq |\lambda_i(\mathbf{L}_H)|, \forall i \in [d]$, the following inequality holds:

$$-\rho(\mathbf{L}_H) \leq \lambda_{\min}(\mathbf{L}_H) \leq \lambda_{\min}(\mathbf{H}_g(\mathbf{z})), \quad \forall \mathbf{H}_g(\mathbf{z}) \in \mathcal{H}, \quad \mathbf{z} \in [\mathbf{l}, \mathbf{u}], \quad (16)$$

allowing us to use $\alpha = \frac{\rho(\mathbf{L}_H)}{2} \geq \max\{0, -\frac{1}{2} \min\{\lambda_{\min}(\mathbf{H}_f(\mathbf{z})) : \mathbf{z} \in [\mathbf{l}, \mathbf{u}]\}\}$.

3.3 Efficient power method for spectral radius computation of the lower bounding Hessian

By using interval propagation, one can easily compute sound lower and upper bounds on each position of the Hessian matrix, compute the lower bounding Hessian and perform a power method with it to

obtain the spectral radius ρ . However, this method would not scale well to high dimensional scenarios. For instance, in the STL10 dataset [Coates et al., 2011] (with input dimension $d = 96 \cdot 96 \cdot 3 = 27,648$) in each color image, our Hessian matrix would require an order of $O(d^2) = O(10^9)$ real numbers to be stored. This makes it intractable to perform a power method over such an humongous matrix, or even to compute the lower bounding Hessian. Alternatively, we take advantage of the possibility of expressing the \mathbf{L}_H matrix as a sum of rank-1 matrices, to enable performing a power method over it.

Standard power method for spectral radius computation

Given any squared and real valued matrix $\mathbf{M} \in \mathbb{R}^{d \times d}$ and an initial vector $\mathbf{v}_0 \in \mathbb{R}^d$ that is not an eigenvector of \mathbf{M} , the sequence:

$$\mathbf{v}_n = \frac{\mathbf{M}(\mathbf{M}\mathbf{v}_{n-1})}{\|\mathbf{M}(\mathbf{M}\mathbf{v}_{n-1})\|_2}, \quad (17)$$

converges to the eigenvector with the largest eigenvalue in absolute value, i.e. the eigenvector where the spectral radius is attained, being the spectral radius $\rho(\mathbf{M}) = \sqrt{\|\mathbf{M}(\mathbf{M}\mathbf{v}_{n-1})\|_2}$ [Mises and Pollaczek-Geiringer, 1929].

Power method over lower bounding Hessian of PNs

We can employ IBP (Section 3.1) in order to obtain an expression of the lower bounding Hessian (\mathbf{L}_H) and evaluate Eq. (17) as:

$$\mathbf{L}_H \mathbf{v} = \frac{\mathcal{U}(\mathbf{H}_g(\mathbf{z}))\mathbf{v} + \mathcal{L}(\mathbf{H}_g(\mathbf{z}))\mathbf{v}}{2} + \left(\frac{\mathcal{L}(\mathbf{H}_g(\mathbf{z}))\mathbf{1} - \mathcal{U}(\mathbf{H}_g(\mathbf{z}))\mathbf{1}}{2} \right) * \mathbf{v}. \quad (18)$$

Applying IBP on Eq. (12) we obtain:

$$\begin{aligned} \mathcal{L}(\mathbf{H}_g(\mathbf{z}))\mathbf{v} &= \sum_{i=1}^k (c_{ti} - c_{\gamma i})^+ \mathcal{L}(\nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(N)})\mathbf{v} + \sum_{i=1}^k (c_{ti} - c_{\gamma i})^- \mathcal{U}(\nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(N)})\mathbf{v} \\ \mathcal{U}(\mathbf{H}_g(\mathbf{z}))\mathbf{v} &= \sum_{i=1}^k (c_{ti} - c_{\gamma i})^- \mathcal{L}(\nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(N)})\mathbf{v} + \sum_{i=1}^k (c_{ti} - c_{\gamma i})^+ \mathcal{U}(\nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(N)})\mathbf{v}. \end{aligned} \quad (19)$$

We can recursively evaluate $\mathcal{L}(\nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n)})\mathbf{v}$ and $\mathcal{U}(\nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n)})\mathbf{v}$ efficiently as these matrices can be expressed as a sum of rank-1 matrices as below.

Proposition 1. *Let $\delta \in [\mathcal{L}(\delta), \mathcal{U}(\delta)]$ be a real-valued weight, the matrix-vector products $\mathcal{L}(\delta \cdot \nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n)})\mathbf{v}$ and $\mathcal{U}(\delta \cdot \nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n)})\mathbf{v}$ can be evaluated as:*

$$\begin{aligned} \mathcal{L}(\delta \cdot \nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n)})\mathbf{v} &= \mathcal{L}(\delta \cdot \nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n-1)})\mathbf{w}_{[n]:i}^{+\top} \mathbf{v} + \mathcal{U}(\delta \cdot \nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n-1)})\mathbf{w}_{[n]:i}^{-\top} \mathbf{v} \\ &\quad + \mathbf{w}_{[n]:i}^+ \mathcal{L}(\delta \cdot \nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n-1)})\mathbf{v} + \mathbf{w}_{[n]:i}^- \mathcal{U}(\delta \cdot \nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n-1)})\mathbf{v} \\ &\quad + \mathcal{L}(\delta' \nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n-1)})\mathbf{v}, \end{aligned} \quad (20)$$

$$\begin{aligned} \mathcal{U}(\delta \cdot \nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n)})\mathbf{v} &= \mathcal{L}(\delta \cdot \nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n-1)})\mathbf{w}_{[n]:i}^{-\top} \mathbf{v} + \mathcal{U}(\delta \cdot \nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n-1)})\mathbf{w}_{[n]:i}^{+\top} \mathbf{v} \\ &\quad + \mathbf{w}_{[n]:i}^- \mathcal{L}(\delta \cdot \nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n-1)})\mathbf{v} + \mathbf{w}_{[n]:i}^+ \mathcal{U}(\delta \cdot \nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n-1)})\mathbf{v} \\ &\quad + \mathcal{U}(\delta' \nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n-1)})\mathbf{v}, \end{aligned} \quad (21)$$

where $\delta' \in [\mathcal{L}(\delta), \mathcal{U}(\delta)] \cdot [\mathcal{L}(\mathbf{w}_{[n]:i}^\top \mathbf{z} + 1), \mathcal{U}(\mathbf{w}_{[n]:i}^\top \mathbf{z} + 1)]$ and vectors $\mathcal{L}(\delta \cdot \nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n-1)})$ and $\mathcal{U}(\delta \cdot \nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n-1)})$ can be obtained through IBP on Eq. (13).

Lastly, by applying recursively Proposition 1 from $n = N$ to $n = 1$, starting with $\delta = 1$, we can substitute the results on Eq. (19) and then on Eq. (18) to efficiently evaluate a step of the power method (Eq. (17)) without needing to store the lower bounding Hessian matrix or needing to perform expensive matrix-vector products.

Overall, our lower bounding method consists in computing a valid value of α that satisfies that the α -convexified objective g_α is convex, following Eqs. (4) and (5). In particular, we use $\alpha = \frac{\rho(\mathbf{L}_H)}{2}$. $\rho(\mathbf{L}_H)$ is computed via a power method, where the main operation $\mathbf{L}_H \mathbf{v}$ is evaluated without the need to compute or store the \mathbf{L}_H matrix. Provided this valid α , we perform PGD over g_α and this provides a lower bound of the global minima of Eq. (4).

3.4 Non-uniform diagonal shift

In order to obtain an estimate of the α parameter as defined by Adjiman et al. [1998] in Eq. (7), we make use of the rank-1 matrices IBP rules defined in Appendix F.1. We also define the operator $\mathcal{M}(\cdot) = \max\{|\mathcal{L}(\cdot)|, |\mathcal{U}(\cdot)|\}$ and certain useful properties about it in Appendix F.2. Thanks to Lemmas 4 to 6, we can obtain an expression for $\mathcal{M}(\mathbf{H}_g(\mathbf{z}))$. This will be used to compute the vector α in the Non-uniform diagonal shift scenario for α -convexification.

Theorem 1. *Let f be a N -degree CCP PN defined as in Eq. (1). Let $g(\mathbf{z}) = f(\mathbf{z})_t - f(\mathbf{z})_\gamma$ for any $t \neq \gamma, t \in [0], \gamma \in [0]$. Let $\mathbf{H}_g(\mathbf{z})$ be the Hessian matrix of g , the operation $\mathcal{M}(\mathbf{H}_g(\mathbf{z}))$ results in:*

$$\mathcal{M}(\mathbf{H}_g(\mathbf{z})) \leq \sum_{i=1}^k |c_{ti} - c_{\gamma i}| \mathcal{M}(\nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(N)}), \quad (22)$$

where for $n = 2, \dots, N$, we can express:

$$\begin{aligned} \mathcal{M}(\nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n)}) &\leq \mathcal{M}(\nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n-1)}) |\mathbf{w}_{[n]i}|^\top + |\mathbf{w}_{[n]i}| \mathcal{M}(\nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n-1)})^\top \\ &\quad + \mathcal{M}(\mathbf{w}_{[n]i}^\top \mathbf{z} + 1) \mathcal{M}(\nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n-1)}). \end{aligned} \quad (23)$$

Lastly, for $n = 1$, $\mathcal{M}(\nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(1)}) = \mathbf{0}_{d \times d}$ is a $d \times d$ matrix full of zeros.

As one can observe, in Theorem 1, the matrix $\mathcal{M}(\mathbf{H}_g(\mathbf{z}))$ is expressed as a sum of rank-1 matrices. This allows to efficiently compute $\sum_{j \neq i} \mathcal{M}(h_g(\mathbf{z})_{ij}) \frac{d_i}{d_j}$, which is necessary for the right term in Eq. (7). For the left term in Eq. (7), we can efficiently compute the lower bound of the diagonal of the Hessian matrix by using the rules present in Section 3.1 and Appendix F.1.

4 Experiments

In this Section we show the efficiency of our method by comparing against a simple Black-box solver. Tightness of bounds is also analyzed in comparison with IBP and DeepT-Fast [Bonaert et al., 2021], a zonotope based verification method able to handle multiplications tighter than IBP. Finally, a study of the performance of our method in different scenarios is performed. Unless otherwise specified, every network is trained for 100 epochs with Stochastic Gradient Descent (SGD), with a learning rate of 0.001, which is divided by 10 at epochs [40, 60, 80], momentum 0.9, weight decay $5 \cdot 10^{-5}$ and batch size 128. We thoroughly evaluate our method over the popular image classification datasets MNIST [LeCun et al., 1998], CIFAR10 [Krizhevsky et al., 2014] and STL10 [Coates et al., 2011]. Every experiment is done over the first 1000 images of the test dataset, this is a common practice in verification [Singh et al., 2019]. For images that are correctly classified by the network, we sequentially verify robustness against the remaining classes in decreasing order of network output. Each verification problem is given a maximum execution time of 60 seconds, we include experiments with different time limits in Appendix E. Note that the execution time can be longer as execution is cut in an asynchronous way, i.e., after we finish the iteration of the BaB algorithm where the time limit is reached. All of our experiments were conducted on a single GPU node equipped with a 32 GB NVIDIA V100 PCIe.

4.1 Comparison with a Black-box solver

In this experiment, we compare the performance of our BaB verification algorithm with the Black-box solver Gurobi [Gurobi Optimization, LLC, 2022]. Gurobi can globally solve Quadratically Constrained Quadratic Programs whether they are convex or not. As this solver cannot extend to higher degree polynomial functions, we train 2nd degree PNs with hidden size $k = 16$ to compare the verification time of our method with Gurobi. In order to do so, we express the verification objective as a quadratic form $g(\mathbf{z}) = f(\mathbf{z})_t - f(\mathbf{z})_a = \mathbf{z}^\top \mathbf{Q}\mathbf{z} + \mathbf{q}^\top \mathbf{z} + c$ this together with the input constraints $\mathbf{z} \in [\mathbf{l}, \mathbf{u}]$ is fed to Gurobi and optimized until convergence.

The black-box solver approach neither scales to higher-dimensional inputs nor to higher polynomial degrees. With this approach we need $\mathcal{O}(d^2)$ memory to store the quadratic form, which makes it unfeasible for datasets with higher resolution images than CIFAR10. On the contrary, as seen in Table 1, our approach does not need so much memory and can scale to datasets with larger input sizes like STL10.

Table 1: Verification results for 2nd degree PNs. Columns #F, #T and #t.o. refer to the number of images where robustness is falsified, verified and timed-out respectively. When comparing with a black-box solver, our method is much faster and can scale to higher dimensional inputs. This is due to our efficient exploitation of the low-rank factorization of PNs.

Dataset	Model	Correct	ϵ	VPN (Our method)				Gurobi			
				time	F	T	t.o.	time	F	T	t.o.
MNIST ($1 \times 28 \times 28$)	2×16	961	0.00725	1.76	37	924	0	16.6	37	924	0
			0.013	1.78	71	890	0	15.13	71	890	0
			0.05	1.43	682	267	12	6.25	691	270	0
			0.06	1.5	790	155	16	4.47	799	162	0
CIFAR10 ($3 \times 32 \times 32$)	2×16	460	1/610	1.03	90	370	0	328.0	90	370	0
			1/255	1.0	183	277	0	250.07	183	277	0
			4/255	0.92	427	28	5	87.93	429	31	0
STL10 ($3 \times 96 \times 96$)	2×16	362	1/610	5.06	142	220	0	out of memory			
			1/255	3.61	246	113	3				
			4/255	1.39	360	1	1				

4.2 Comparison with IBP and DeepT-Fast

In this experiment we compare the tightness of the lower bounds provided by IBP, DeepT-Fast and α -convexification and their effectiveness when employed for verification. This is done by executing one upper bounding step with PGD and one lower bounding step for each lower bounding method over the initial feasible set provided by ϵ (see Eq. (3)). We compare the average of the distance from each lower bound to the PGD upper bound over the first 1000 images of the MNIST dataset for PNs with hidden size $k = 25$ and degrees ranging from 2 to 7. We also evaluate verified accuracy of 2nd (PN_Conv2) and 4th (PN_Conv4) PNs with IBP, DeepT-Fast and α -convexification in the Uniform diagonal shift setup. We employ a maximum time of 120 seconds. For details on the architecture of these networks, we refer to Appendix E.

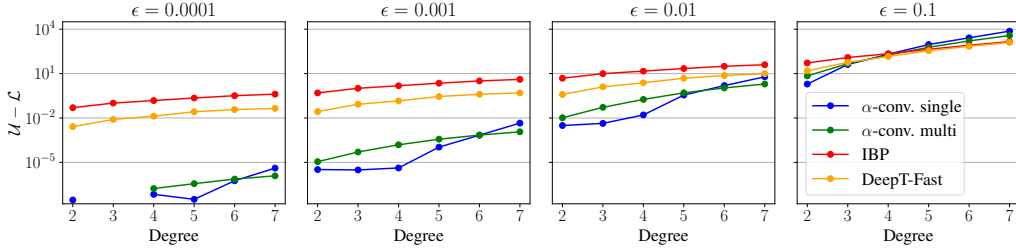


Figure 3: Average difference in log-scale between PGD upper bound (\mathcal{U}) and lower bound (\mathcal{L}) provided by BP (red), DeepT-Fast [Bonaert et al., 2021] (orange), α -convexification with Uniform diagonal shift (blue) and α -convexification with Non-uniform diagonal shift (green) of the first 1000 images of the MNIST dataset. α -convexification bounds are significantly tighter than IBP and DeepT-Fast for small ϵ values and all PN degrees from 2 to 7.

When using IBP, we get a much looser lower bound than with α -convexification, see Fig. 3. Only for high-degree, high- ϵ combinations IBP lower bounds are closer to the PGD upper bound. In practice, this is not a problem for verification, as for epsilons in the order of 0.1, it is really easy to find adversarial examples with PGD and there will be no accuracy left to verify. DeepT-Fast significantly outperforms IBP bounds across all degrees and ϵ values. But, as observed in Fig. 3, except for big ϵ values, its performance is still far from the one provided by both α -convexification methods. When comparing both α -convexification methods (blue and green lines in Fig. 3), we observe that for small degree PNs ($N < 5$), in the Uniform diagonal shift case we are able to obtain tighter bounds.

The looseness of the IBP lower bound is confirmed when comparing the verified accuracy with IBP and the rest of lower bounding methods, see Table 2. With α -convexification, we are able to verify the accuracy of 2nd and 4th order PNs almost exactly (almost no gap between the verified accuracy and its upper bound) in every studied dataset, while with IBP, we are not able to verify robustness for a single image in any network- ϵ pair, confirming the fact that IBP cannot be used

Table 2: Verification results with our method employing IBP, DeepT-Fast and α -convexification for lower bounding the objective. Acc.% is the clean accuracy of the network, Ver.% is the verified accuracy and U.B. its upper bound. When using α -convexification bounds we get verified accuracies really close to the upper bound, while when using IBP verified accuracy is 0 for every network- ϵ pair, which makes it unsuitable for PN verification.

Dataset	Model	Acc.%	ϵ	IBP		DeepT-Fast [Bonaert et al., 2021]		VPN (ours) (α -convexification)		U.B.
				Time(s)	Ver.%	Time(s)	Ver.%	Time(s)	Ver.%	
MNIST	PN_Conv4	98.6	0.015	0.3	0.0	2.3	91.3	50	96.3	96.4
			0.026	0.4	0.0	3.7	59.7	69	92.9	94.8
			0.3	0.6	0.0	0.7	0.0	13.8	0.0	0.0
CIFAR10	PN_Conv2	63.5	1/255	0.3	0.0	2.0	23.3	136.2	44.4	44.6
			2/255	0.5	0.0	0.6	1.4	89.2	25.4	27.5
	PN_Conv4	62.6	1/255	0.4	0.0	2.2	19.5	274.6	45.5	46.7
			2/255	0.5	0.0	0.5	0.5	224.1	16.5	30.5
STL10*	PN_Conv4	38.1	1/255	3.4	0.0	26.0	14.7	2481.0	21.7	21.9

* Results obtained in the first 360 images of the dataset due to the longer running times because of the larger input size of STL10.

for PN verification. The improvements in the bounds when utilizing DeepT-Fast instead of IBP is clearly seen in verification results in Table 2. With DeepT-Fast, we are able to effectively verify PNs faster than with α -convexification, but achieving a much lower verified accuracy (Ver%) than with α -convexification. As a reference, for CIFAR10 PN_Conv2 at $\epsilon = 2/255$, with α -convexification we obtain 25.4% verified accuracy, while with DeepT-Fast, we can just obtain 1.4% verified accuracy. It is worth highlighting that DeepT-Fast can also scale to verify networks trained on STL10.

5 Conclusion

We propose a novel α -BaB global optimization algorithm to verify polynomial networks (PNs). We exhibit that our method outperforms existing methods, such as black-box solvers, IBP and DeepT-Fast [Bonaert et al., 2021]. Our method enables verification in datasets such as STL10, which includes RGB images of 96×96 resolution. This is larger than the images typically used in previous verification methods. Note that existing methods like IBP and DeepT-Fast are also able to scale to STL10 both with a lower verified accuracy. Our method can further encourage the community to extend verification to a broader class of functions as well as conduct experiments in datasets of higher resolution. We believe our method can be extended to cover other twice-differentiable networks in the future.

Limitations

As discussed in Appendix E.2, our verification method does not scale to high-degree PNs. Even though we can verify high-accuracy PNs (see Table 2), we are still far from verifying the top performing deep PNs studied in Chrysos et al. [2021]. Another problem that we share with ReLU NN verifiers is the scalability to networks with larger input size [Wang et al., 2021]. In this work we are able to verify networks trained in STL10 [Coates et al., 2011], but these networks are shallow, yet their verification still takes a long time, see Table 2.

Acknowledgements

We are deeply thankful to the reviewers for providing constructive feedback. Research was sponsored by the Army Research Office and was accomplished under Grant Number W911NF-19-1-0404. This project has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement number 725594 - time-data). This work was supported by the Swiss National Science Foundation (SNSF) under grant number 200021_178865. This project has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement n° 725594 - time-data). This work was supported by Zeiss. This work was supported by SNF project – Deep Optimisation of the Swiss National Science Foundation (SNSF) under grant number 200021_205011.

References

- Claire S. Adjiman and Christodoulos A. Floudas. Rigorous convex underestimators for general twice-differentiable problems. *Journal of Global Optimization*, 9(1):23–40, Jul 1996.
- Claire S. Adjiman, Stefan Dallwig, Christodoulos A. Floudas, and Arnold Neumaier. A global optimization method, α bb, for general twice-differentiable constrained NLPs — I. theoretical advances. *Computers & Chemical Engineering*, 22:1137–1158, 1998.
- Greg Anderson, Shankara Pailoor, Isil Dillig, and Swarat Chaudhuri. Optimization and abstraction: A synergistic approach for analyzing neural network robustness. In *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation*, PLDI 2019, page 731–744, New York, NY, USA, 2019. Association for Computing Machinery.
- Osbert Bastani, Yani Ioannou, Leonidas Lampropoulos, Dimitrios Vytiniotis, Aditya V. Nori, and Antonio Criminisi. Measuring neural net robustness with constraints. In *Advances in Neural Information Processing Systems (NeurIPS)*, NIPS’16, page 2621–2629, Red Hook, NY, USA, 2016. Curran Associates Inc.
- Gregory Bonaert, Dimitar I. Dimitrov, Maximilian Baader, and Martin Vechev. Fast and precise certification of transformers. In *Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation*, PLDI 2021, page 466–481, New York, NY, USA, 2021. Association for Computing Machinery.
- Rudy Bunel, Alessandro De Palma, Alban Desmaison, Krishnamurthy Dvijotham, Pushmeet Kohli, Philip Torr, and M. Pawan Kumar. Lagrangian decomposition for neural network verification. In Jonas Peters and David Sontag, editors, *Conference on Uncertainty in Artificial Intelligence (UAI)*, volume 124 of *Proceedings of Machine Learning Research*, pages 370–379. PMLR, 03–06 Aug 2020a.
- Rudy Bunel, Jingyue Lu, Ilker Turkaslan, Philip H.S. Torr, Pushmeet Kohli, and M. Pawan Kumar. Branch and bound for piecewise linear neural network verification. *Journal of Machine Learning Research*, 21(42):1–39, 2020b.
- Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy (SP)*, pages 39–57, 2017.
- Moulik Choraria, Leello Tadesse Dadi, Grigorios Chrysos, Julien Mairal, and Volkan Cevher. The spectral bias of polynomial neural networks. In *International Conference on Learning Representations (ICLR)*, 2022.
- Grigorios G Chrysos and Yannis Panagakis. Naps: Non-adversarial polynomial synthesis. *Pattern Recognition Letters*, 140:318–324, 2020.
- Grigorios G. Chrysos, Stylianos Moschoglou, Giorgos Bouritsas, Jiankang Deng, Yannis Panagakis, and Stefanos P Zafeiriou. Deep polynomial neural networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence (T-PAMI)*, page 1–1, 2021.
- Grigorios G Chrysos, Markos Georgopoulos, Jiankang Deng, Jean Kossaifi, Yannis Panagakis, and Anima Anandkumar. Augmenting deep classifiers with polynomial neural networks. In *European Conference on Computer Vision (ECCV)*, 2022.
- Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised feature learning. In Geoffrey Gordon, David Dunson, and Miroslav Dudík, editors, *International Conference on Artificial Intelligence and Statistics*, volume 15 of *Proceedings of Machine Learning Research*, pages 215–223, Fort Lauderdale, FL, USA, 11–13 Apr 2011. PMLR.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 248–255, 2009.
- Yinpeng Dong, Hang Su, Baoyuan Wu, Zhifeng Li, Wei Liu, Tong Zhang, and Jun Zhu. Efficient decision-based black-box adversarial attacks on face recognition. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7714–7722, 2019.

- Zehao Dou, Stanley J Osher, and Bao Wang. Mathematical analysis of adversarial attacks. *arXiv preprint arXiv:1811.06492*, 2018.
- Rüdiger Ehlers. Formal verification of piece-wise linear feed-forward neural networks. *CoRR*, abs/1705.01320, 2017.
- Fenglei Fan, Jinjun Xiong, and Ge Wang. Universal approximation with quadratic deep networks. *Neural Networks*, 124:383–392, 2020.
- Fenglei Fan, Mengzhou Li, Fei Wang, Rongjie Lai, and Ge Wang. Expressivity and trainability of quadratic networks. *CoRR*, abs/2110.06081, 2021.
- Claudio Ferrari, Mark Niklas Mueller, Nikola Jovanović, and Martin Vechev. Complete verification via multi-neuron relaxation guided branch-and-bound. In *International Conference on Learning Representations (ICLR)*, 2022.
- Timon Gehr, Matthew Mirman, Dana Drachler-Cohen, Petar Tsankov, Swarat Chaudhuri, and Martin Vechev. Ai2: Safety and robustness certification of neural networks with abstract interpretation. In *2018 IEEE Symposium on Security and Privacy (SP)*, pages 3–18, 2018.
- Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In Yoshua Bengio and Yann LeCun, editors, *International Conference on Learning Representations (ICLR)*, 2015.
- Gaurav Goswami, Akshay Agarwal, Nalini K. Ratha, Richa Singh, and Mayank Vatsa. Detecting and mitigating adversarial perturbations for robust face recognition. *International Journal of Computer Vision (IJCV)*, 127:719–742, 2019.
- Gurobi Optimization, LLC. Gurobi Optimizer Reference Manual, 2022. URL <https://www.gurobi.com>.
- Reiner Horst and Hoang Tuy. *Global Optimization: Deterministic Approaches*. Springer, Berlin, Heidelberg, 1996.
- Guy Katz, Clark Barrett, David Dill, Kyle Julian, and Mykel Kochenderfer. Reluplex: An efficient smt solver for verifying deep neural networks, 2017.
- C. T. Kelley. *Iterative Methods for Optimization*. Society for Industrial and Applied Mathematics, 1999.
- Tamara G. Kolda and Brett W. Bader. Tensor decompositions and applications. *SIAM Review*, 51(3): 455–500, 2009.
- Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. The cifar-10 dataset. *online: http://www.cs.toronto.edu/kriz/cifar.html*, 55, 2014.
- A. H. Land and A. G. Doig. An automatic method of solving discrete programming problems. *Econometrica*, 28(3):497–520, 1960.
- Yann LeCun, Léon Bottou, Yoshua Bengio, Patrick Haffner, et al. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015.
- Qizhang Li, Yiwen Guo, and Hao Chen. Practical no-box adversarial attacks against dnns. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- Changliu Liu, Tomer Arnon, Christopher Lazarus, Christopher Strong, Clark Barrett, and Mykel J. Kochenderfer. Algorithms for verifying deep neural networks. *Foundations and Trends® in Optimization*, 4(3-4):244–404, 2021.
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations (ICLR)*, 2018.

- R. V. Mises and H. Pollaczek-Geiringer. Praktische verfahren der gleichungsauflösung . *ZAMM - Journal of Applied Mathematics and Mechanics / Zeitschrift für Angewandte Mathematik und Mechanik*, 9(2):152–164, 1929.
- Ramon E. Moore, R. Baker Kearfott, and Michael J. Cloud. *Introduction to Interval Analysis*. Society for Industrial and Applied Mathematics, 2009.
- Aditi Raghunathan, Sang Michael Xie, Fanny Yang, John Duchi, and Percy Liang. Understanding and mitigating the tradeoff between robustness and accuracy. In *International Conference on Machine Learning (ICML)*, pages 7909–7919, 2020.
- Vicenç Rúbies Royo, Roberto Calandra, Dušan M. Stipanović, and Claire J. Tomlin. Fast neural network verification via shadow prices. *ArXiv*, abs/1902.07247, 2019.
- Ali Shafahi, W. Ronny Huang, Christoph Studer, Soheil Feizi, and Tom Goldstein. Are adversarial examples inevitable? In *International Conference on Learning Representations (ICLR)*, 2019.
- Gagandeep Singh, Rupanshu Ganvir, Markus Püschel, and Martin Vechev. Beyond the single neuron convex barrier for neural network certification. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems (NeurIPS)*, volume 32. Curran Associates, Inc., 2019.
- M. H. Stone. The generalized weierstrass approximation theorem. *Mathematics Magazine*, 21(4): 167–184, 1948.
- Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *International Conference on Learning Representations (ICLR)*, 2014.
- Vincent Tjeng, Kai Y. Xiao, and Russ Tedrake. Evaluating robustness of neural networks with mixed integer programming. In *International Conference on Learning Representations (ICLR)*, 2019.
- Shiqi Wang, Kexin Pei, Justin Whitehouse, Junfeng Yang, and Suman Jana. Formal security analysis of neural networks using symbolic intervals. In *USENIX Conference on Security Symposium, SEC'18*, page 1599–1614, USA, 2018a. USENIX Association.
- Shiqi Wang, Huan Zhang, Kaidi Xu, Xue Lin, Suman Jana, Cho-Jui Hsieh, and J Zico Kolter. Beta-CROWN: Efficient bound propagation with per-neuron split constraints for neural network robustness verification. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7794–7803, 2018b.
- Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric Xing, Laurent El Ghaoui, and Michael Jordan. Theoretically principled trade-off between robustness and accuracy. In *International Conference on Machine Learning (ICML)*, pages 7472–7482, 2019.
- Yaoyao Zhong and Weihong Deng. Adversarial learning with margin-based triplet embedding regularization. In *International Conference on Computer Vision (ICCV)*, pages 6549–6558, 2019.
- Zhenyu Zhu, Fabian Latorre, Grigorios Chrysos, and Volkan Cevher. Controlling the complexity and lipschitz constant improves polynomial nets. In *International Conference on Learning Representations (ICLR)*, 2022.

Contents of the appendix

The societal impact of our method is discussed in Appendix A. We cover the related work in Appendix B. Appendix C provides a more detailed coverage of PN architectures. In Appendix D, we include the pseudocode of our algorithms and we provide an analysis of the complexity. Experiments and ablation studies are available in Appendix E. To conclude, in Appendix F, we include all of our proofs.

A Societal impact

The performance on standard image classification benchmarks has increased substantially the last few years, owing to the success of neural networks. Their success enables their adoption in tackling real-world problems. However, robustness and trustworthiness of neural networks is of critical importance before their adoption in real-world applications. Our method is a verifier that focuses on polynomial networks and enables the complete verification of PNs. Therefore, we expect that by using the proposed method, certain properties of the robustness could be verified in a principled way. We expect this to have a predominantly positive societal impact as either a tool for pre-trained models or tool for certifying models as part of their debugging. However, it can also be used as a tool to find weaknesses of pretrained PNs by adversarial agents.

B Related Work

In this Section, we give an overview of neural network verification and polynomial networks, that are centered around our target in this work.

B.1 Neural Network Verification

Early works on sound and complete NN verification were based on Mixed Integer Linear Programming (MILP) and Satisfiability Modulo Theory (SMT) solvers [Katz et al., 2017, Ehlers, 2017, Bastani et al., 2016, Tjeng et al., 2019] and were limited to both small datasets and networks.

The utilization of custom BaB algorithms enabled verification to scale to datasets and networks that are closer to those used in practice. Bunel et al. [2020b] review earlier methods like Katz et al. [2017] and show they can be formulated as BaB algorithms. BaDNB [Bunel et al., 2020a] proposes a novel branching strategy called Filtered Smart Branching and uses the Lagrangian decomposition-based bounding algorithm. β -CROWN [Wang et al., 2021] proposes a bound propagation based algorithm. MN-BaB [Ferrari et al., 2022] proposes a cost adjusted branching strategy and leverages multi-neuron relaxations and a GPU-based solver for bounds computing. Our work centers on the bounding algorithm by proposing a general convex lowerbound adapted to PNs.

BaB algorithms for ReLU networks focus their branching strategies on the activity of ReLU neurons. This has been observed to work better than input set branching for ReLU networks [Bunel et al., 2020b]. Similarly to our method, Anderson et al. [2019], Wang et al. [2018a], Royo et al. [2019] use input set branching strategies.

B.2 Polynomial Networks

First works have been focused on developing the foundations and showcasing the performance of PNs in different tasks [Chrysos et al., 2021, Chrysos and Panagakis, 2020]. Also, in Chrysos et al. [2022], PN classifiers are formulated in a common framework where other previous methods like Wang et al. [2018b] can be framed. Lately, more emphasis has been put onto proving theoretical properties of PNs [Fan et al., 2021, Choraria et al., 2022]. In Zhu et al. [2022], they derive Lipschitz constant and complexity bounds for two PN decompositions in terms of the l_∞ and l_2 norms. They also analyze robustness of PNs against PGD adversarial attacks by measuring percentage of images where PGD fails to find an adversarial example, which is a complete but not sound verification method. Our verification method is sound and complete.

C Background

C.1 Coupled CP decomposition (CCP)

Relying on the CP decomposition [Kolda and Bader, 2009], the CCP decomposition provides us a core expression of PNs as used in Chrysos and Panagakis [2020] to construct a generative model. Let N be the polynomial degree, $\mathbf{z} \in \mathbb{R}^d$ the input vector, d , k and o the input, hidden and output sizes, the CCP decomposition can be expressed as:

$$\mathbf{x}^{(n)} = (\mathbf{W}_{[n]}^\top \mathbf{z}) * \mathbf{x}^{(n-1)} + \mathbf{x}^{(n-1)}, \forall n = 2, 3, \dots, N, \quad (24)$$

where $\mathbf{x}^{(1)} = \mathbf{W}_{[1]}^\top \mathbf{z}$, $\mathbf{f}(\mathbf{z}) = \mathbf{C}\mathbf{x}^{(N)} + \beta$ and $*$ denotes the Hadamard product.

For example, the second order CCP factorization will lead to the following formulation:

$$\mathbf{x}^{(1)} = \mathbf{W}_{[1]}^\top \mathbf{z}, \quad \mathbf{x}^{(2)} = \mathbf{W}_{[2]}^\top \mathbf{z} * \mathbf{x}^{(1)} + \mathbf{x}^{(1)}, \quad \mathbf{f}(\mathbf{z}) = \mathbf{C}\mathbf{x}^{(2)} + \beta, \quad (25)$$

where $\mathbf{W}_{[1]} \in \mathbb{R}^{d \times k}$, $\mathbf{W}_{[2]} \in \mathbb{R}^{d \times k}$ and $\mathbf{C} \in \mathbb{R}^{o \times k}$ are weight matrices, $\beta \in \mathbb{R}^o$ is a vector.

C.2 Nested coupled CP decomposition (NCP)

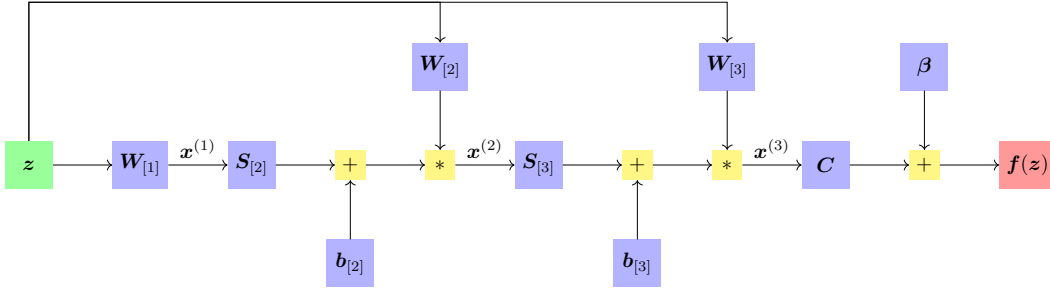


Figure 4: Third order NCP network architecture.

The NCP model leverages a joint hierarchical decomposition, which provided strong results in both generative and discriminative tasks in Chrysos et al. [2021]. It can be expressed with the following recursive relation:

$$\mathbf{x}^{(n)} = (\mathbf{W}_{[n]}^\top \mathbf{z}) * (\mathbf{S}_{[n]}^\top \mathbf{x}^{(n-1)} + \mathbf{b}_{[n]}), \quad (26)$$

for $n \in [N - 1] + 1$ with $\mathbf{x}^{(1)} = \mathbf{W}_{[1]}^\top \mathbf{z}$ and $\mathbf{f}(\mathbf{z}) = \mathbf{C}\mathbf{x}^{(N)} + \beta$.

We present the first and second order partial derivatives of NCP below.

$$\nabla_{\mathbf{z}} x_i^{(n)} = \mathbf{w}_{[n]:i} \cdot (\mathbf{s}_{[n]:i}^\top \mathbf{x}^{(n-1)} + b_{[n]i}) + (\mathbf{w}_{[n]:i}^\top \mathbf{z}) \cdot \left(\sum_{j=1}^k s_{[n]ji} \nabla_{\mathbf{z}} x_j^{(n-1)} \right) \quad (27)$$

$$\begin{aligned} \nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n)} &= \mathbf{w}_{[n]:i} \cdot \left(\sum_{j=1}^k s_{[n]ji} \nabla_{\mathbf{z}} x_j^{(n-1)} \right)^\top + \left(\sum_{j=1}^k s_{[n]ji} \nabla_{\mathbf{z}} x_j^{(n-1)} \right) \mathbf{w}_{[n]:i}^\top \\ &\quad + (\mathbf{w}_{[n]:i}^\top \mathbf{z}) \cdot \left(\sum_{j=1}^k s_{[n]ji} \nabla_{\mathbf{z}\mathbf{z}}^2 x_j^{(n-1)} \right). \end{aligned} \quad (28)$$

Our method can be easily extended to this type of PNs. As a reference, we obtain a verified accuracy of 76.2% with an upper bound of 76.4% with an 2×25 NCP at $\epsilon = 0.026$.

C.3 Product of Polynomials

In practice, Chrysos et al. [2021] report that to reduce further the parameters they are often stacking sequentially a number of polynomials, see Fig. 5. That results in a setting that is referred to as product of polynomials, with the highest degree of expansion being defined as the product of all the degrees

of the individual polynomials. Hopefully, as we will demonstrate below, our formulation can be extended to this setting.

Let $\mathbf{x} = \mathbf{x}^{(N_1)}$ be the output of a PN $f_1 : \mathbb{R}^d \rightarrow \mathbb{R}^k$ of N_1^{th} -degree and $\mathbf{y} = \mathbf{y}^{(N_2)}$ be the output of a PN $f_2 : \mathbb{R}^k \rightarrow \mathbb{R}^k$ of N_2^{th} -degree, with \mathbf{x} and \mathbf{y} coming either from Eq. (24) or Eq. (26). The product of polynomials f_i and f_2 is defined as $f : \mathbb{R}^d \rightarrow \mathbb{R}^o$:

$$\mathbf{f}(\mathbf{z}) = \mathbf{C} \mathbf{f}_2(\mathbf{f}_1(\mathbf{z})) + \boldsymbol{\beta}, \quad (29)$$

where $\mathbf{C} \in \mathbb{R}^{o \times k}$ and $\boldsymbol{\beta} \in \mathbb{R}^o$.

We present the first and second order partial derivatives of the Product of polynomials below.

Let \mathbf{z} be the input, $\mathbf{x} = \mathbf{f}_1(\mathbf{z})$ and $\mathbf{y} = \mathbf{f}_2(\mathbf{x})$ by applying the chain rule of partial derivatives, we can obtain:

$$\nabla_{\mathbf{z}} y_i = \sum_{j=1}^k \frac{\partial y_i}{\partial x_j} \nabla_{\mathbf{z}} x_j, \quad \nabla_{\mathbf{z}\mathbf{z}}^2 y_i = \sum_{j=1}^k \frac{\partial y_i}{\partial x_j} \nabla_{\mathbf{z}\mathbf{z}}^2 x_j + \mathbf{J}_{\mathbf{z}}^{\top}(\mathbf{x}) \nabla_{\mathbf{x}\mathbf{x}}^2 y_i \mathbf{J}_{\mathbf{z}}(\mathbf{x}), \quad (30)$$

where $\mathbf{J}_{\mathbf{z}}^{\top}(\mathbf{x}) \in \mathbb{R}^{k \times d}$ is the Jacobian matrix.

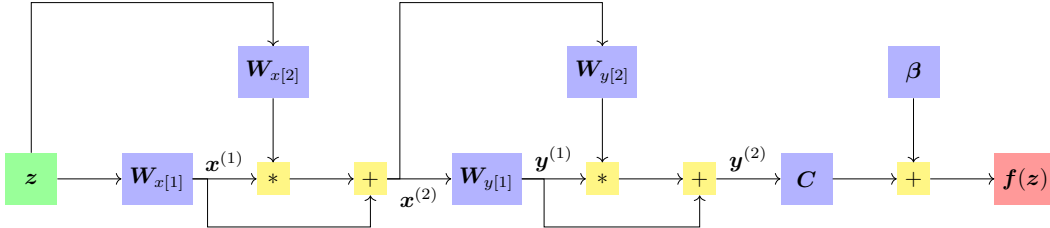


Figure 5: Product of two 2nd-degree CCP polynomials. The final classification layer of the first polynomial is dropped, the output $\mathbf{x}^{(2)}$ is fed into a second 2nd-degree polynomial.

C.4 Convolutional PNs (C-PNs)

As seen in Chrysos et al. [2021], the performance of PNs is boosted when employing convolution operators instead of standard linear mappings. Convolutions reduce the number of parameters and take advantage of the local 2D structure of the images. Let $\mathcal{Z} \in \mathbb{R}^{c \times h \times w}$ be the input image with c , h and w being the number of channels, height and width respectively. The N^{th} -degree **CCP_Conv** becomes:

$$\mathbf{X}^{(n,i)} = (\mathcal{Z} \circ \mathcal{W}_{[n,i]}) * \mathbf{X}^{(n-1,i)} + \mathbf{X}^{(n-1,i)}, \forall n \in [N-1] + 1, \forall i \in [q], \quad (31)$$

where $\mathbf{X}^{(1,i)} = \mathcal{Z} \circ \mathcal{W}_{[1,i]}$, $\forall i \in [q]$, $\mathbf{f}(\mathbf{z}) = \mathbf{C} \mathbf{x}^{(N)} + \boldsymbol{\beta}$ and $\mathbf{x}^{(N)} = \text{flat}(\mathcal{X}^{(N)})$.

In order to verify C-PNs, we convert every convolutional layer into a linear layer via Toeplitz matrices (see Gehr et al. [2018]) resulting in an equivalent CCP PN.

C.5 Details on IBP

In this section we further elaborate on the IBP for the network outputs, the gradients and Hessians by applying the notions in Section 3.1.

IBP for network outputs In order to compute lower and upper bounds on the output of the network $\mathbf{f}(\mathbf{z})$, we apply interval arithmetic techniques in the recursive formulas Eqs. (1) and (26). For both the CCP and the NCP cases, let $\hat{\mathbf{x}}^{(n)} = \mathbf{W}_{[n]}^{\top} \mathbf{z}$:

$$\begin{aligned} \mathcal{L}(\hat{\mathbf{x}}^{(n)}) &= \mathbf{W}_{[n]}^{\top} \mathbf{l} + \mathbf{W}_{[n]}^{\top-} \mathbf{u} \\ \mathcal{U}(\hat{\mathbf{x}}^{(n)}) &= \mathbf{W}_{[n]}^{\top} \mathbf{u} + \mathbf{W}_{[n]}^{\top-} \mathbf{l}. \end{aligned} \quad (32)$$

In the case of CCP PNs, the recursive formula in Eq. (1) becomes $\mathbf{x}^{(n)} = (\hat{\mathbf{x}}^{(n)} + \mathbf{1}) * \mathbf{x}^{(n-1)}$, whereas for NCP, Eq. (26) becomes $\mathbf{x}^{(n)} = (\hat{\mathbf{x}}^{(n)}) * (\mathbf{S}_{[n]}^{\top} \mathbf{x}^{(n-1)} + \mathbf{b}_{[n]})$ for any $n \in [N-1] + 1$,

for $n = 1$, in both architectures $\mathbf{x}^{(n)} = \hat{\mathbf{x}}^{(n)}$. Then, we can define the bounds with the following recursive formulas:

$$\begin{aligned}
\text{CCP} \left\{ \begin{aligned} S &= \begin{cases} (\mathcal{L}(\hat{\mathbf{x}}^{(n)}) + 1) * \mathcal{L}(\mathbf{x}^{(n-1)}), \\ (\mathcal{L}(\hat{\mathbf{x}}^{(n)}) + 1) * \mathcal{U}(\mathbf{x}^{(n-1)}), \\ (\mathcal{U}(\hat{\mathbf{x}}^{(n)}) + 1) * \mathcal{L}(\mathbf{x}^{(n-1)}), \\ (\mathcal{U}(\hat{\mathbf{x}}^{(n)}) + 1) * \mathcal{U}(\mathbf{x}^{(n-1)}), \end{cases} \\ \mathcal{L}(\mathbf{x}^{(n)}) &= \min(S) \\ \mathcal{U}(\mathbf{x}^{(n)}) &= \max(S) \end{aligned} \right. \\
\text{NCP} \left\{ \begin{aligned} \mathcal{L}(\mathbf{S}_{[n]}^\top \mathbf{x}^{(n-1)} + \mathbf{b}_{[n]}) &= \mathbf{S}_{[n]}^{\top+} \mathcal{L}(\mathbf{x}^{(n-1)}) + \mathbf{S}_{[n]}^{\top-} \mathcal{U}(\mathbf{x}^{(n-1)}) + \mathbf{b}_{[n]} \\ \mathcal{U}(\mathbf{S}_{[n]}^\top \mathbf{x}^{(n-1)} + \mathbf{b}_{[n]}) &= \mathbf{S}_{[n]}^{\top+} \mathcal{U}(\mathbf{x}^{(n-1)}) + \mathbf{S}_{[n]}^{\top-} \mathcal{L}(\mathbf{x}^{(n-1)}) + \mathbf{b}_{[n]} \\ S &= \begin{cases} \mathcal{L}(\mathbf{S}_{[n]}^\top \mathbf{x}^{(n-1)} + \mathbf{b}_{[n]}) * \mathcal{L}(\mathbf{x}^{(n-1)}), \\ \mathcal{L}(\mathbf{S}_{[n]}^\top \mathbf{x}^{(n-1)} + \mathbf{b}_{[n]}) * \mathcal{U}(\mathbf{x}^{(n-1)}), \\ \mathcal{U}(\mathbf{S}_{[n]}^\top \mathbf{x}^{(n-1)} + \mathbf{b}_{[n]}) * \mathcal{L}(\mathbf{x}^{(n-1)}), \\ \mathcal{U}(\mathbf{S}_{[n]}^\top \mathbf{x}^{(n-1)} + \mathbf{b}_{[n]}) * \mathcal{U}(\mathbf{x}^{(n-1)}), \end{cases} \\ \mathcal{L}(\mathbf{x}^{(n)}) &= \min(S) \\ \mathcal{U}(\mathbf{x}^{(n)}) &= \max(S), \end{aligned} \right. \quad (33)
\end{aligned}$$

where the min and max operators are applied element-wise in sets of vectors or matrices. Finally, the output bounds are obtained with:

$$\begin{aligned}
\mathcal{L}(\mathbf{f}(\mathbf{z})) &= \mathbf{C}^+ \mathcal{L}(\mathbf{x}^{(N)}) + \mathbf{C}^- \mathcal{U}(\mathbf{x}^{(N)}) + \boldsymbol{\beta} \\
\mathcal{U}(\mathbf{f}(\mathbf{z})) &= \mathbf{C}^+ \mathcal{U}(\mathbf{x}^{(N)}) + \mathbf{C}^- \mathcal{L}(\mathbf{x}^{(N)}) + \boldsymbol{\beta}.
\end{aligned} \quad (34)$$

These operations can be implemented as a forward pass through the PN.

IBP for gradients

The next step is to obtain bounds on the gradients of the PNs. Again, with the help of interval arithmetic theory, we can extend the recursive formulas Eqs. (13) and (27) for computing IBP bounds. In the case of CCP PNs:

$$\begin{aligned}
\mathcal{L}(\mathbf{J}_z^\top(\mathbf{x}^{(n)})) &= \mathcal{L}(\mathbf{W}_{[n]} * \mathbf{x}^{(n-1)} + (\mathbf{W}_{[n]}^\top \mathbf{z} + 1) * \mathbf{J}_z^\top(\mathbf{x}^{(n-1)})) \\
&= \mathcal{L}(\mathbf{W}_{[n]} * \mathbf{x}^{(n-1)}) + \mathcal{L}((\mathbf{W}_{[n]}^\top \mathbf{z} + 1) * \mathbf{J}_z^\top(\mathbf{x}^{(n-1)})) \\
&= \mathbf{W}_{[n]}^+ * \mathcal{L}(\mathbf{x}^{(n-1)}) + \mathbf{W}_{[n]}^- * \mathcal{U}(\mathbf{x}^{(n-1)}) + \mathcal{L}((\mathbf{W}_{[n]}^\top \mathbf{z} + 1) * \mathbf{J}_z^\top(\mathbf{x}^{(n-1)})) \\
\mathcal{U}(\mathbf{J}_z^\top(\mathbf{x}^{(n)})) &= \mathcal{U}(\mathbf{W}_{[n]} * \mathbf{x}^{(n-1)} + (\mathbf{W}_{[n]}^\top \mathbf{z} + 1) * \mathbf{J}_z^\top(\mathbf{x}^{(n-1)})) \\
&= \mathcal{U}(\mathbf{W}_{[n]} * \mathbf{x}^{(n-1)}) + \mathcal{U}((\mathbf{W}_{[n]}^\top \mathbf{z} + 1) * \mathbf{J}_z^\top(\mathbf{x}^{(n-1)})) \\
&= \mathbf{W}_{[n]}^+ * \mathcal{U}(\mathbf{x}^{(n-1)}) + \mathbf{W}_{[n]}^- * \mathcal{L}(\mathbf{x}^{(n-1)}) + \mathcal{U}((\mathbf{W}_{[n]}^\top \mathbf{z} + 1) * \mathbf{J}_z^\top(\mathbf{x}^{(n-1)})),
\end{aligned} \quad (35)$$

with

$$\begin{aligned}
S &= \begin{cases} \mathcal{L}(\mathbf{W}_{[n]}^\top \mathbf{z} + 1) * \mathcal{L}(\mathbf{J}_z^\top(\mathbf{x}^{(n-1)})), \\ \mathcal{L}(\mathbf{W}_{[n]}^\top \mathbf{z} + 1) * \mathcal{U}(\mathbf{J}_z^\top(\mathbf{x}^{(n-1)})), \\ \mathcal{U}(\mathbf{W}_{[n]}^\top \mathbf{z} + 1) * \mathcal{L}(\mathbf{J}_z^\top(\mathbf{x}^{(n-1)})), \\ \mathcal{U}(\mathbf{W}_{[n]}^\top \mathbf{z} + 1) * \mathcal{U}(\mathbf{J}_z^\top(\mathbf{x}^{(n-1)})) \end{cases} \\
\mathcal{L}((\mathbf{W}_{[n]}^\top \mathbf{z} + 1) * \mathbf{J}_z^\top(\mathbf{x}^{(n-1)})) &= \min(S) \\
\mathcal{U}((\mathbf{W}_{[n]}^\top \mathbf{z} + 1) * \mathbf{J}_z^\top(\mathbf{x}^{(n-1)})) &= \max(S),
\end{aligned} \quad (36)$$

where the Hadamard product of a $\mathbb{R}^{k \times d}$ matrix with a \mathbb{R}^d vector results in a $\mathbb{R}^{k \times d}$ matrix. The min and max operators are applied element-wise in sets of vectors or matrices.

In the case of NCP PNs:

$$\begin{aligned}
\mathcal{L}(\mathbf{J}_z^\top(\mathbf{x}^{(n)})) &= \mathcal{L}(\mathbf{W}_{[n]} * (\mathbf{S}_{[n]}\mathbf{x}^{(n-1)} + \mathbf{b}_{[n]}) + (\mathbf{W}_{[n]}^\top \mathbf{z}) * (\mathbf{S}_{[n]}\mathbf{J}_z^\top(\mathbf{x}^{(n-1)}))) \\
&= \mathcal{L}(\mathbf{W}_{[n]} * (\mathbf{S}_{[n]}\mathbf{x}^{(n-1)} + \mathbf{b}_{[n]})) + \mathcal{L}((\mathbf{W}_{[n]}^\top \mathbf{z}) * (\mathbf{S}_{[n]}\mathbf{J}_z^\top(\mathbf{x}^{(n-1)}))) \\
&= \mathbf{W}_{[n]}^+ * \mathcal{L}(\mathbf{S}_{[n]}\mathbf{x}^{(n-1)} + \mathbf{b}_{[n]}) + \mathbf{W}_{[n]}^- * \mathcal{U}(\mathbf{S}_{[n]}\mathbf{x}^{(n-1)} + \mathbf{b}_{[n]}) \\
&\quad + \mathcal{L}((\mathbf{W}_{[n]}^\top \mathbf{z}) * (\mathbf{S}_{[n]}\mathbf{J}_z^\top(\mathbf{x}^{(n-1)}))) \\
\mathcal{U}(\mathbf{J}_z^\top(\mathbf{x}^{(n)})) &= \mathcal{U}(\mathbf{W}_{[n]} * (\mathbf{S}_{[n]}\mathbf{x}^{(n-1)} + \mathbf{b}_{[n]}) + (\mathbf{W}_{[n]}^\top \mathbf{z}) * (\mathbf{S}_{[n]}\mathbf{J}_z^\top(\mathbf{x}^{(n-1)}))) \\
&= \mathcal{U}(\mathbf{W}_{[n]} * (\mathbf{S}_{[n]}\mathbf{x}^{(n-1)} + \mathbf{b}_{[n]})) + \mathcal{U}((\mathbf{W}_{[n]}^\top \mathbf{z}) * (\mathbf{S}_{[n]}\mathbf{J}_z^\top(\mathbf{x}^{(n-1)}))) \\
&= \mathbf{W}_{[n]}^+ * \mathcal{U}(\mathbf{S}_{[n]}\mathbf{x}^{(n-1)} + \mathbf{b}_{[n]}) + \mathbf{W}_{[n]}^- * \mathcal{L}(\mathbf{S}_{[n]}\mathbf{x}^{(n-1)} + \mathbf{b}_{[n]}) \\
&\quad + \mathcal{U}((\mathbf{W}_{[n]}^\top \mathbf{z}) * (\mathbf{S}_{[n]}\mathbf{J}_z^\top(\mathbf{x}^{(n-1)}))),
\end{aligned} \tag{37}$$

with

$$\begin{aligned}
\mathcal{L}(\mathbf{S}_{[n]}\mathbf{J}_z^\top(\mathbf{x}^{(n-1)})) &= \mathbf{S}_{[n]}^+ \mathcal{L}(\mathbf{J}_z^\top(\mathbf{x}^{(n-1)})) + \mathbf{S}_{[n]}^- \mathcal{U}(\mathbf{J}_z^\top(\mathbf{x}^{(n-1)})) \\
\mathcal{U}(\mathbf{S}_{[n]}\mathbf{J}_z^\top(\mathbf{x}^{(n-1)})) &= \mathbf{S}_{[n]}^+ \mathcal{U}(\mathbf{J}_z^\top(\mathbf{x}^{(n-1)})) + \mathbf{S}_{[n]}^- \mathcal{L}(\mathbf{J}_z^\top(\mathbf{x}^{(n-1)})) \\
S &= \begin{cases} \mathcal{L}(\mathbf{W}_{[n]}^\top \mathbf{z}) * \mathcal{L}(\mathbf{S}_{[n]}\mathbf{J}_z^\top(\mathbf{x}^{(n-1)})), \\ \mathcal{L}(\mathbf{W}_{[n]}^\top \mathbf{z}) * \mathcal{U}(\mathbf{S}_{[n]}\mathbf{J}_z^\top(\mathbf{x}^{(n-1)})), \\ \mathcal{U}(\mathbf{W}_{[n]}^\top \mathbf{z}) * \mathcal{L}(\mathbf{S}_{[n]}\mathbf{J}_z^\top(\mathbf{x}^{(n-1)})), \\ \mathcal{U}(\mathbf{W}_{[n]}^\top \mathbf{z}) * \mathcal{U}(\mathbf{S}_{[n]}\mathbf{J}_z^\top(\mathbf{x}^{(n-1)})) \end{cases} \\
\mathcal{L}((\mathbf{W}_{[n]}^\top \mathbf{z}) * (\mathbf{S}_{[n]}\mathbf{J}_z^\top(\mathbf{x}^{(n-1)}))) &= \min(S) \\
\mathcal{U}((\mathbf{W}_{[n]}^\top \mathbf{z}) * (\mathbf{S}_{[n]}\mathbf{J}_z^\top(\mathbf{x}^{(n-1)}))) &= \max(S).
\end{aligned} \tag{38}$$

Eq. (35) (Eq. (37)) jointly with Eq. (36) (Eq. (38)) allows to recursively obtain gradient bounds for the CCP (NCP) PNs starting with $\mathcal{L}(\mathbf{J}_z^\top(\mathbf{x}^{(1)})) = \mathcal{U}(\mathbf{J}_z^\top(\mathbf{x}^{(1)})) = \mathbf{W}_{[1]}$. Finally, bounds of the gradients of the network output in both CCP and NCP cases are:

$$\begin{aligned}
\mathcal{L}(\mathbf{J}_z(\mathbf{f})) &= \mathbf{C}^+ \mathcal{L}(\mathbf{J}_z(\mathbf{x}^{(N)})) + \mathbf{C}^- \mathcal{U}(\mathbf{J}_z(\mathbf{x}^{(N)})) + \beta \\
\mathcal{U}(\mathbf{J}_z(\mathbf{f})) &= \mathbf{C}^+ \mathcal{U}(\mathbf{J}_z(\mathbf{x}^{(N)})) + \mathbf{C}^- \mathcal{L}(\mathbf{J}_z(\mathbf{x}^{(N)})) + \beta.
\end{aligned} \tag{39}$$

IBP for Hessians Lastly, with help of interval arithmetic and aforementioned IBP for PN outputs and gradients, we can use recursive formulas Eqs. (14) and (28) to compute bounds of the Hessian matrices.

In the case of CCP PNs:

$$\begin{aligned}
\mathcal{L}(\mathbf{H}_z(x_i^{(n)})) &= \mathcal{L}(\nabla_z x_i^{(n-1)} \mathbf{w}_{[n]:i}^\top + \{\nabla_z x_i^{(n-1)} \mathbf{w}_{[n]:i}^\top\}^\top + (\mathbf{w}_{[n]:i}^\top \mathbf{z} + 1) \mathbf{H}_z(x_i^{(n-1)})) \\
&= \mathcal{L}(\nabla_z x_i^{(n-1)} \mathbf{w}_{[n]:i}^+{}^\top + \mathcal{U}(\nabla_z x_i^{(n-1)} \mathbf{w}_{[n]:i}^-{}^\top \\
&\quad + \mathbf{w}_{[n]:i}^+ \mathcal{L}(\nabla_z x_i^{(n-1)})^\top + \mathbf{w}_{[n]:i}^- \mathcal{U}(\nabla_z x_i^{(n-1)})^\top \\
&\quad + \mathcal{L}((\mathbf{w}_{[n]:i}^\top \mathbf{z} + 1) \mathbf{H}_z(x_i^{(n-1)}))) \\
\mathcal{U}(\mathbf{H}_z(x_i^{(n)})) &= \mathcal{U}(\nabla_z x_i^{(n-1)} \mathbf{w}_{[n]:i}^\top + \{\nabla_z x_i^{(n-1)} \mathbf{w}_{[n]:i}^\top\}^\top + (\mathbf{w}_{[n]:i}^\top \mathbf{z} + 1) \mathbf{H}_z(x_i^{(n-1)})) \\
&= \mathcal{U}(\nabla_z x_i^{(n-1)} \mathbf{w}_{[n]:i}^+{}^\top + \mathcal{L}(\nabla_z x_i^{(n-1)} \mathbf{w}_{[n]:i}^-{}^\top \\
&\quad + \mathbf{w}_{[n]:i}^+ \mathcal{U}(\nabla_z x_i^{(n-1)})^\top + \mathbf{w}_{[n]:i}^- \mathcal{L}(\nabla_z x_i^{(n-1)})^\top \\
&\quad + \mathcal{U}((\mathbf{w}_{[n]:i}^\top \mathbf{z} + 1) \mathbf{H}_z(x_i^{(n-1)}))),
\end{aligned} \tag{40}$$

where

$$\begin{aligned}
S &= \begin{Bmatrix} \mathcal{L}(\mathbf{w}_{[n]:i}^\top \mathbf{z} + 1) \mathcal{L}(\mathbf{H}_{\mathbf{z}}(x_i^{(n-1)})), \\ \mathcal{L}(\mathbf{w}_{[n]:i}^\top \mathbf{z} + 1) \mathcal{U}(\mathbf{H}_{\mathbf{z}}(x_i^{(n-1)})), \\ \mathcal{U}(\mathbf{w}_{[n]:i}^\top \mathbf{z} + 1) \mathcal{L}(\mathbf{H}_{\mathbf{z}}(x_i^{(n-1)})), \\ \mathcal{U}(\mathbf{w}_{[n]:i}^\top \mathbf{z} + 1) \mathcal{U}(\mathbf{H}_{\mathbf{z}}(x_i^{(n-1)})) \end{Bmatrix} \\
\mathcal{L}((\mathbf{w}_{[n]:i}^\top \mathbf{z} + 1) \mathbf{H}_{\mathbf{z}}(x_i^{(n-1)})) &= \min(S) \\
\mathcal{U}((\mathbf{w}_{[n]:i}^\top \mathbf{z} + 1) \mathbf{H}_{\mathbf{z}}(x_i^{(n-1)})) &= \max(S).
\end{aligned} \tag{41}$$

In the case of NCP PNs:

$$\begin{aligned}
\mathcal{L}(\mathbf{H}_{\mathbf{z}}(x_i^{(n)})) &= \mathcal{L}(\mathbf{s}_{[n]i} \mathbf{J}_{\mathbf{z}}^\top(\mathbf{x}^{(n-1)}) \mathbf{w}_{[n]:i}^\top + \{\mathbf{s}_{[n]i} \mathbf{J}_{\mathbf{z}}^\top(\mathbf{x}^{(n-1)}) \mathbf{w}_{[n]:i}^\top\}^\top + (\mathbf{w}_{[n]:i}^\top \mathbf{z}) \sum_{j=1}^k s_{ij} \mathbf{H}_{\mathbf{z}}(x_j^{(n-1)})) \\
&= \mathcal{L}(\mathbf{s}_{[n]i} \mathbf{J}_{\mathbf{z}}^\top(\mathbf{x}^{(n-1)})) \mathbf{w}_{[n]:i}^{+\top} + \mathcal{U}(\mathbf{s}_{[n]i} \mathbf{J}_{\mathbf{z}}^\top(\mathbf{x}^{(n-1)})) \mathbf{w}_{[n]:i}^{l\top} \\
&+ \mathbf{w}_{[n]:i}^{+\top} \mathcal{L}(\mathbf{s}_{[n]i} \mathbf{J}_{\mathbf{z}}^\top(\mathbf{x}^{(n-1)}))^\top + \mathbf{w}_{[n]:i}^{-\top} \mathcal{U}(\mathbf{s}_{[n]i} \mathbf{J}_{\mathbf{z}}^\top(\mathbf{x}^{(n-1)}))^\top \\
&+ \mathcal{L}((\mathbf{w}_{[n]:i}^\top \mathbf{z}) \sum_{j=1}^k s_{ij} \mathbf{H}_{\mathbf{z}}(x_j^{(n-1)})) \\
\mathcal{U}(\mathbf{H}_{\mathbf{z}}(x_i^{(n)})) &= \mathcal{U}(\mathbf{s}_{[n]i} \mathbf{J}_{\mathbf{z}}^\top(\mathbf{x}^{(n-1)}) \mathbf{w}_{[n]:i}^\top + \{\mathbf{s}_{[n]i} \mathbf{J}_{\mathbf{z}}^\top(\mathbf{x}^{(n-1)}) \mathbf{w}_{[n]:i}^\top\}^\top + (\mathbf{w}_{[n]:i}^\top \mathbf{z}) \sum_{j=1}^k s_{ij} \mathbf{H}_{\mathbf{z}}(x_j^{(n-1)})) \\
&= \mathcal{U}(\mathbf{s}_{[n]i} \mathbf{J}_{\mathbf{z}}^\top(\mathbf{x}^{(n-1)})) \mathbf{w}_{[n]:i}^{+\top} + \mathcal{L}(\mathbf{s}_{[n]i} \mathbf{J}_{\mathbf{z}}^\top(\mathbf{x}^{(n-1)})) \mathbf{w}_{[n]:i}^{l\top} \\
&+ \mathbf{w}_{[n]:i}^{+\top} \mathcal{U}(\mathbf{s}_{[n]i} \mathbf{J}_{\mathbf{z}}^\top(\mathbf{x}^{(n-1)}))^\top + \mathbf{w}_{[n]:i}^{-\top} \mathcal{L}(\mathbf{s}_{[n]i} \mathbf{J}_{\mathbf{z}}^\top(\mathbf{x}^{(n-1)}))^\top \\
&+ \mathcal{U}((\mathbf{w}_{[n]:i}^\top \mathbf{z}) \sum_{j=1}^k s_{ij} \mathbf{H}_{\mathbf{z}}(x_j^{(n-1)})),
\end{aligned} \tag{42}$$

where

$$\begin{aligned}
\mathcal{L}(\mathbf{s}_{[n]i} \mathbf{J}_{\mathbf{z}}^\top(\mathbf{x}^{(n-1)})) &= \mathbf{s}_{[n]i}^+ \mathcal{L}(\mathbf{J}_{\mathbf{z}}^\top(\mathbf{x}^{(n-1)})) + \mathbf{s}_{[n]i}^- \mathcal{U}(\mathbf{J}_{\mathbf{z}}^\top(\mathbf{x}^{(n-1)})) \\
\mathcal{U}(\mathbf{s}_{[n]i} \mathbf{J}_{\mathbf{z}}^\top(\mathbf{x}^{(n-1)})) &= \mathbf{s}_{[n]i}^+ \mathcal{U}(\mathbf{J}_{\mathbf{z}}^\top(\mathbf{x}^{(n-1)})) + \mathbf{s}_{[n]i}^- \mathcal{L}(\mathbf{J}_{\mathbf{z}}^\top(\mathbf{x}^{(n-1)})) \\
\mathcal{L}(\sum_{j=1}^k s_{ij} \mathbf{H}_{\mathbf{z}}(x_j^{(n-1)})) &= \sum_{j=1}^k s_{ij}^+ \mathcal{L}(\mathbf{H}_{\mathbf{z}}(x_j^{(n-1)})) + s_{ij}^- \mathcal{U}(\mathbf{H}_{\mathbf{z}}(x_j^{(n-1)})) \\
\mathcal{U}(\sum_{j=1}^k s_{ij} \mathbf{H}_{\mathbf{z}}(x_j^{(n-1)})) &= \sum_{j=1}^k s_{ij}^+ \mathcal{U}(\mathbf{H}_{\mathbf{z}}(x_j^{(n-1)})) + s_{ij}^- \mathcal{L}(\mathbf{H}_{\mathbf{z}}(x_j^{(n-1)})),
\end{aligned} \tag{43}$$

and

$$S = \begin{cases} \mathcal{L}(\mathbf{w}_{[n]:i}^\top \mathbf{z}) \mathcal{L}(\sum_{j=1}^k s_{ij} \mathbf{H}_z(x_j^{(n-1)})), \\ \mathcal{L}(\mathbf{w}_{[n]:i}^\top \mathbf{z}) \mathcal{U}(\sum_{j=1}^k s_{ij} \mathbf{H}_z(x_j^{(n-1)})), \\ \mathcal{U}(\mathbf{w}_{[n]:i}^\top \mathbf{z}) \mathcal{L}(\sum_{j=1}^k s_{ij} \mathbf{H}_z(x_j^{(n-1)})), \\ \mathcal{U}(\mathbf{w}_{[n]:i}^\top \mathbf{z}) \mathcal{U}(\sum_{j=1}^k s_{ij} \mathbf{H}_z(x_j^{(n-1)})) \end{cases} \quad (44)$$

$$\mathcal{L}((\mathbf{w}_{[n]:i}^\top \mathbf{z}) \sum_{j=1}^k s_{ij} \mathbf{H}_z(x_j^{(n-1)})) = \min(S)$$

$$\mathcal{U}((\mathbf{w}_{[n]:i}^\top \mathbf{z}) \sum_{j=1}^k s_{ij} \mathbf{H}_z(x_j^{(n-1)})) = \max(S).$$

Similarly to gradient bounds, we can recursively compute $\mathcal{L}(\mathbf{H}_z(x_i^{(N)}))$ and $\mathcal{U}(\mathbf{H}_z(x_i^{(N)}))$ starting with $\mathcal{L}(\mathbf{H}_z(x_i^{(1)})) = \mathcal{U}(\mathbf{H}_z(x_i^{(1)})) = \mathbf{0}_{d \times d}$. Finally, bounds of the Hessian matrix of the output can be obtained with:

$$\begin{aligned} \mathcal{L}(\mathbf{H}_z(f(\mathbf{z})_i)) &= \sum_{j=1}^k c_{ij}^+ \mathcal{L}(\mathbf{H}_z(x_j^{(N)})) + c_{ij}^- \mathcal{U}(\mathbf{H}_z(x_j^{(N)})) \\ \mathcal{U}(\mathbf{H}_z(f(\mathbf{z})_i)) &= \sum_{j=1}^k c_{ij}^+ \mathcal{U}(\mathbf{H}_z(x_j^{(N)})) + c_{ij}^- \mathcal{L}(\mathbf{H}_z(x_j^{(N)})). \end{aligned} \quad (45)$$

D BaB algorithm for PN robustness verification

BaB algorithms [Land and Doig, 1960] are a well known approach to global optimization [Horst and Tuy, 1996]. These algorithms intend to find the global minima of a optimization problem in the form of Eq. (4), therefore guaranteeing soundness and completeness for verification. In this Section we present the details of our BaB based verification algorithm, we prove the theoretical convergence of BaB to the global minima of Eq. (4) and the complexity of its key steps.

Firstly, the property that we want to verify or falsify is the one given by Eq. (2). This property is defined by (i) the network f , (ii) the adversarial budget ϵ , (iii) a correctly classified input $\mathbf{z}_0 : \arg \max f(\mathbf{z}_0) = t$. In order to verify the property, it is necessary that for each $\gamma \neq t$, the global minima of Eq. (4) is greater than 0, i.e., $\forall \gamma \neq t : v^* = \min_{\mathbf{z} \in C_{\text{in}}} f(\mathbf{z})_t - f(\mathbf{z})_\gamma > 0$. On the contrary, in order to falsify Eq. (2), it is sufficient that for any $\gamma \neq t$ the global minima of Eq. (4) is smaller or equal than 0, i.e., $\exists \gamma \neq t : v^* = \min_{\mathbf{z} \in C_{\text{in}}} f(\mathbf{z})_t - f(\mathbf{z})_\gamma \leq 0$. In order to reduce execution times, we heuristically sort all the γ in decreasing order by network output, $\{\gamma_i : \gamma_i \neq t, f(\mathbf{z}_0)_{\gamma_i} \geq f(\mathbf{z}_0)_{\gamma_j} \forall j > i\}$ and solve Eq. (4) until one global minima is smaller or equal to 0 or all global minimas are greater than 0.

In order to solve Eq. (4), we use Algorithm 1. Without any modifications, this algorithm converges to the global minima. However, for verification, it is sufficient to find that the lower bound of the global minima (global_lb) cannot be smaller than or equal to 0, or that the upper bound of the global minima in a subset (prob_ub) is smaller than 0, i.e., there exists an adversarial example in that subset. Therefore we can cut the execution early when employed for verification, these optional changes are highlighted in red in Algorithm 1.

Algorithm 1 can theoretically be applied to any twice-differentiable classifier provided we have a method for computing an α for obtaining valid α -convexification bounds. For this matter, we propose Algorithm 2, which again can be used for any twice-differentiable classifier if we are able

Algorithm 1 Branch and Bound, adapted from Bunel et al. [2020b]

```
1: function BAB( $f, l, u, t, \gamma$ )
2:    $global\_ub \leftarrow \inf$ 
3:    $global\_lb \leftarrow -\inf$ 
4:    $\alpha \leftarrow \text{compute\_alpha}(f, l, u, t, \gamma)$  ▷ Algorithm 2
5:    $probs \leftarrow [(global\_lb, l, u)]$ 
6:   while  $global\_ub - global\_lb > 10^{-6}$  and  $global\_lb \leq 0$  do
7:      $([], l', u') \leftarrow \text{pick\_out}(probs)$  ▷ Take subset with the minimum lower bound
8:      $[(l_1, u_1), \dots, (l_s, u_s)] \leftarrow \text{split}(l', u')$  ▷ Split widest input variable interval in halves
9:     for  $i = 1 \dots s$  do
10:       $prob\_ub \leftarrow \text{compute\_UB}(f, l_i, u_i, t, \gamma)$  ▷ PGD over the original function  $g(z)$ 
11:       $prob\_lb \leftarrow \text{compute\_LB}(f, l_i, u_i, t, \gamma, \alpha)$  ▷ PGD over  $g_\alpha(z, \alpha)$ 
12:      if  $prob\_ub < global\_ub$  then
13:         $global\_ub \leftarrow prob\_ub$ 
14:         $\text{prune\_problems}(probs, global\_ub)$  ▷ Remove if  $prob\_lb > global\_ub$ 
15:      end if
16:      if  $prob\_lb < global\_ub$  and  $prob\_lb \leq 0$  then
17:         $probs.append((prob\_lb, l_i, u_i))$ 
18:      end if
19:      if  $prob\_ub \leq 0$  then ▷ An adversarial example was found
20:        return  $[], 0$ 
21:      end if
22:    end for
23:    if  $|probs| == 0$  then ▷  $prob\_lb > 0$  for every subset
24:      return  $[], 1$ 
25:    end if
26:     $global\_lb \leftarrow \min\{lb \mid (lb, [], []) \in probs\}$ 
27:  end while
28:  return  $global\_ub, global\_ub > 0$ 
29: end function
```

to compute its lower bounding Hessian L_H . In the PN case, we propose a method for evaluating the matrix-vector product $L_H v$, this is covered in Algorithm 3. We note that in Algorithm 2, for initializing the vector v , each position of the vector is randomly sampled in the $[0, 1]$ interval, then the resulting vector is rescaled so that $\|v\|_2 = 1$.

Algorithm 2 Power method for α estimation

```
1: function COMPUTE_ALPHA( $f, l, u, t, \gamma$ )
2:    $v \leftarrow \text{init\_v}(f, l, u, t, \gamma)$  ▷ Ensure  $v$  is not an eigenvector of  $L_H$  and  $\|v\|_2 = 1$ 
3:    $prev\_v \leftarrow \mathbf{0}$ 
4:    $r \leftarrow 0$ 
5:   while  $\|v - prev\_v\|_2 > 10^{-6}$  do
6:      $prev\_v \leftarrow v$ 
7:      $v \leftarrow \text{evaluate\_LHv}(f, l, u, t, \gamma, v)$  ▷ Algorithm 3
8:      $r \leftarrow \|v\|_2$ 
9:      $v \leftarrow v/r$ 
10:     $v \leftarrow \text{evaluate\_LHv}(f, l, u, t, \gamma, v)$  ▷ Evaluate  $L_H v$  twice to deal with negative eigenvalues.
11:     $r \leftarrow \|v\|_2$ 
12:     $v \leftarrow v/r$ 
13:  end while
14:  return  $r/2$ 
15: end function
```

Complexity of $L_H v$ evaluation Algorithm 3 is governed by an outer loop which performs $N - 1$ iterations, see line 5. Inside the loop, the most expensive operations are in lines 6, 9, 10, 11 and 12, the rest of operations can be performed in $O(d)$ time. In line 5, the bounds of the gradients are computed, this operation can be performed for every level $n \in [N]$ outside the main loop and store the results using $O(N \cdot d \cdot k)$ time. For lines 9, 10, 11 and 12, an outer loop with k iterations is used for the summation. Then, inside the summation, four dot products plus vector-scalar multiplications

are performed, leading to a time complexity of $O(k \cdot d)$. Overall, the complexity of Algorithm 3 is $O(N \cdot d \cdot k)$.

Algorithm 3 Evaluation of $L_H v$ for a CCP PN, implementation of Proposition 1

```

1: function EVALUATE_LHV(f, problem, v, t,  $\gamma$ )
2:    $lw \leftarrow c_t - c_\gamma$ 
3:    $uw \leftarrow c_t - c_\gamma$     $\triangleright$  Initial upper and lower bounds of the  $\delta$  weight are given by the last linear layer.
4:   result  $\leftarrow 0$ 
5:   for  $n = \text{degree}(f) \dots 2$  do
6:      $lg, ug \leftarrow \mathcal{L}(J_z(x^{(n)})), \mathcal{U}(J_z(x^{(n)}))$ 
7:      $S1 \leftarrow \{lg \cdot lw, lg \cdot uw, ug \cdot lw, ug \cdot uw\}$ 
8:      $lg, ug \leftarrow \min(S1), \max(S1)$ 
9:      $Lv \leftarrow \sum_{i=1}^k w_{[n]:i}^+ \cdot lg[i, :]^\top v + w_{[n]:i}^- \cdot ug[i, :]^\top v + lg[i, :] \cdot w_{[n]:i}^+ v + ug[i, :] \cdot w_{[n]:i}^- v$ 
10:     $Uv \leftarrow \sum_{i=1}^k w_{[n]:i}^+ \cdot ug[i, :]^\top v + w_{[n]:i}^- \cdot lg[i, :]^\top v + ug[i, :] \cdot w_{[n]:i}^+ v + lg[i, :] \cdot w_{[n]:i}^- v$ 
11:     $L1 \leftarrow \sum_{i=1}^k w_{[n]:i}^+ \cdot lg[i, :]^\top 1 + w_{[n]:i}^- \cdot ug[i, :]^\top 1 + lg[i, :] \cdot w_{[n]:i}^+ 1 + ug[i, :] \cdot w_{[n]:i}^- 1$ 
12:     $U1 \leftarrow \sum_{i=1}^k w_{[n]:i}^+ \cdot ug[i, :]^\top 1 + w_{[n]:i}^- \cdot lg[i, :]^\top 1 + ug[i, :] \cdot w_{[n]:i}^+ 1 + lg[i, :] \cdot w_{[n]:i}^- 1$ 
13:    result  $\leftarrow \text{result} + (Lv + Uv)/2 + ((L1 - U1)/2) * v$ 
14:     $lx, ux \leftarrow \mathcal{L}(W_{[n]}^\top z + 1), \mathcal{U}(W_{[n]}^\top z + 1)$ 
15:     $S2 \leftarrow \{lx \cdot lw, lx \cdot uw, ux \cdot lw, ux \cdot uw\}$ 
16:     $lw, uw \leftarrow \min(S2), \max(S2)$     $\triangleright$  Update bounds of the weight ( $\mathcal{L}'\delta, \mathcal{U}'\delta$ ), see Proposition 1
17:   end for
18:   return result
19: end function

```

E Auxiliary experimental results and discussion

We start this Section by comparing the performance of our method with ReLU NN verification algorithms Appendix E.1. We analyze a limitation of the method in Appendix E.2. Then, in Appendix E.3 we perform an ablation study on the effect the input size of the network has in our PN verification algorithm.

Additional notation: In addition to the notation already defined in the main paper, we use \circ for convolutions.

E.1 Comparison with ReLU BaB verification algorithms.

Table 3: Description of convolutional PNs used in our experiments.

Name	degree / N	kernel size	stride	padding	channels
PN_Conv2	2	5×5	2	2	32
PN_Conv4	4	7×7	4	3	64

Complete ReLU NN verification algorithms are usually benchmarked against other methods in the networks and ϵ values proposed by Singh et al. [2019]. ReLU NN verification algorithms mostly rely on the specific structure of the networks, e.g. ReLU activation, which makes a direct comparison with ReLU nets hard. However, we match these benchmarks by training PNs with same degree as the number of ReLU layers and similar number of parameters. For a detailed description of these networks check Table 5 and Table 6.

When comparing with the SOTA verifier β -CROWN [Wang et al., 2021], we firstly observe that the upper bound of verified accuracy is very similar between the NNs and PNs fully connected benchmarks. However, for convolutional PNs (C-PNs) the upper bound is much lower, even 0 for PN_Conv4 trained in the MNIST dataset, see Table 4. Secondly, we observe that the gap between the upper bound and verified accuracy (U.B. and Ver. % in Table 4) is small for β -CROWN. In our case, this gap is only small for the shallowest PN and smallest ϵ (5×81 with $\epsilon = 0.015$), obtaining 88.0% verified accuracy and outperforming β -CROWN with only 77.4%. For PN_Conv2 trained in CIFAR10 and evaluated at $\epsilon = 2/255$, our verified accuracy, 15.7%, is also really close to the upper bound 16.1%, but both are low in comparison with the β -CROWN equivalent.

Table 4: Verification results on our proposed PN verification benchmarks. We run our verification procedure over the first 1000 images of the test split of each dataset. Time(s) refers to the average running time per image when the verification was not timed out, i.e., we can verify or falsify the property in the given time limit, Ver. % is the verified accuracy and U.B. its upper bound. We use the same ϵ values as Wang et al. [2021]. We observe that in comparison with β -CROWN, our method generally has a larger gap between verified accuracy and its upper bound.

Dataset	β -CROWN*				VPN			
	Model	Time(s)	Ver.%	U.B.	Model	Time(s)	Ver.%	U.B.
MNIST	6×100	102	69.9	84.2	5×30	34	24.7	81.1
	6×200	86	77.4	90.1	5×81	60	88.0	91.1
	9×100	103	62.0	82.0	8×24	33	0.0	80.2
	9×200	95	73.5	91.1	8×70	94	0.6	91.6
	ConvSmall	7.0	72.7	73.2	PN_Conv2	3.0	3.5	12.1
	ConvBig	3.1	79.3	80.4	PN_Conv4	8.9	0.0	0.0
CIFAR10	ConvSmall	6.8	46.3	48.1	PN_Conv2	63.4	15.7	16.1
	ConvBig	15.3	51.6	55.0	PN_Conv4	161.4	9.4	16.3

*Note: We report numbers from Wang et al. [2021].

Table 5: Comparison of PNs to the fully connected ReLU network benchmarks proposed in Singh et al. [2019]. #Par. refers to the number of parameters in the ReLU benchmark and our PNs respectively. We build PNs with same degree as number of non-linearities in their corresponding ReLU NN benchmark. We also adjust the hidden size so that the number of parameters is matched.

ReLU NN	#Par.	Acc.(%)	PN	#Par.	Acc.(%)
6×100	119,910	96.0	5×30	117,910	97.2
9×100	150,210	94.7	8×24	150,778	97.5
6×200	319,810	97.2	5×81	318,340	96.2
9×200	440,410	95.0	8×70	439,750	96.7

E.2 Limitations of the proposed method

Scaling our method to high-degree settings is non-trivial due to IBP approximation errors, that accumulate when propagating the bounds through each layer of the PN. Nevertheless, empirical evidence demonstrates that a lower-degree PN is enough for obtaining comparative performance, e.g., Fan et al. [2020].

E.3 Effect of the input size in PN verification

In this experiment we evaluate the effect of the input size in the verification of a PN. In order to evaluate this, we train 3 different CCP networks over the STL10 dataset. Each model has been trained with STL10 images preprocessed with 3 different resizing factors. Every network is a CCP_4 \times 25 trained with a learning rate of 10^{-4} .

As seen in Table 7, downsampling the input images results in a decrease in the accuracy (i.e., from 35.1% to 31.8% at 32×32 resolution). However, downsampling the input improves the robustness of the network. For all ϵ values we observe less successful adversarial attacks, i.e., higher upper bound

Table 6: Comparison of convolutional PNs (C-PNs) to the Convolutional ReLU network benchmarks proposed in Singh et al. [2019]. #Par. refers to the number of parameters in the ReLU benchmark and our PNs respectively. We build C-PNs with same number of convolutional layers as their corresponding ReLU NN benchmark.

Dataset	ReLU NN	#Par.	Acc.(%)	PN	#Par.	Acc.(%)
MNIST	Conv Small	89,606	98.0	PN_Conv2	114,218	98.0
MNIST	Conv Big	1,974,762	92.9	PN_Conv4	834,762	98.0
CIFAR10	Conv Small	125,318	63.0	PN_Conv2	133,994	57.2
CIFAR10	Conv Big	2,466,858	63.1	PN_Conv4	845,514	58.3

Table 7: Input size ablation study: CCP_4 \times 25 networks are trained over STL10 with three different input sizes. Both Verified accuracy (Ver. %) and its upper bound (U.B.) increase when the input size is reduced for all the studied ϵ values.

Input size	Acc.%	ϵ	VPN		
			Time(s)	Ver.%	U.B.
$3 \times 96 \times 96$ (original)	35.1	1/255	47.4	4.8	10.6
		2/255	19.0	0.0	2.4
		8/255	21.5	0.0	0.0
$3 \times 64 \times 64$	35.6	1/255	52.1	11.7	12.6
		2/255	20.7	0.6	3.2
		8/255	12.0	0.0	0.0
$3 \times 32 \times 32$	31.8	1/255	14.7	17.9	18.0
		2/255	21.4	7.2	8.9
		8/255	16.2	0.0	0.1

of the verified accuracy (U.B.). In addition, the verification process is improved, we obtain higher values for verified accuracy (Ver. %), but also the gap with its upper bound is reduced. We believe this phenomenon is due to the networks learning more robust representations with a smaller input size [Raghunathan et al., 2020, Zhang et al., 2019].

F Proofs

F.1 IBP on rank-1 matrices

In order to facilitate the computation of bounds of the Hessian matrix, we introduce some general properties regarding the bounds IBP of rank-1 matrices.

Lemma 1. Let $M = uv^\top$ be a rank-1 matrix defined by vectors $u \in \mathbb{R}^{d_1}$ and $v \in \mathbb{R}^{d_2}$. Let u^+ and v^+ be the positive parts and u^- and v^- the negative parts, satisfying $u = u^+ + u^-$ and $v = v^+ + v^-$. The matrix M can be decomposed as $M = u^+v^{+\top} + u^+v^{-\top} + u^-v^{+\top} + u^-v^{-\top}$.

Proof of Lemma 1. Following $u = u^+ + u^-$ and $v = v^+ + v^-$:

$$M := uv^\top = (u^+ + u^-)(v^+ + v^-)^\top = u^+v^{+\top} + u^+v^{-\top} + u^-v^{+\top} + u^-v^{-\top} \quad (46)$$

□

Lemma 2. Let $M = uv^\top$ be a rank-1 matrix defined by vectors $u \in \mathbb{R}^{d_1}$ and $v \in \mathbb{R}^{d_2}$. Let $\mathcal{L}(u)$ ($\mathcal{U}(u)$) be lower and upper bounded by $\mathcal{L}(u)$ and $\mathcal{U}(u)$ ($\mathcal{L}(v)$ and $\mathcal{U}(v)$). The matrix M is lower and upper bounded by:

$$\begin{aligned} M &\geq \mathcal{L}(M) = \mathcal{L}(u)^+\mathcal{L}(v)^{+\top} + \mathcal{U}(u)^+\mathcal{L}(v)^{-\top} + \mathcal{L}(u)^-\mathcal{U}(v)^{+\top} + \mathcal{U}(u)^-\mathcal{U}(v)^{-\top} \\ M &\leq \mathcal{U}(M) = \mathcal{U}(u)^+\mathcal{U}(v)^{+\top} + \mathcal{L}(u)^+\mathcal{U}(v)^{-\top} + \mathcal{U}(u)^-\mathcal{L}(v)^{+\top} + \mathcal{L}(u)^-\mathcal{L}(v)^{-\top} \end{aligned} \quad (47)$$

Proof of Lemma 2. By applying the multiplication rule of IBP, see Eq. (9), we can express the lower bound as:

$$\begin{aligned} M &\geq \min\{ \mathcal{L}(u)\mathcal{L}(v)^\top, \\ &\quad \mathcal{L}(u)\mathcal{U}(v)^\top, \\ &\quad \mathcal{U}(u)\mathcal{L}(v)^\top, \\ &\quad \mathcal{U}(u)\mathcal{U}(v)^\top \}. \end{aligned} \quad (48)$$

Then, applying Lemma 1 on each term we obtain:

$$\begin{aligned}
M \geq & \min\{\mathcal{L}(\mathbf{u})^+\mathcal{L}(\mathbf{v})^{+\top} + \mathcal{L}(\mathbf{u})^+\mathcal{L}(\mathbf{v})^{-\top} + \mathcal{L}(\mathbf{u})^-\mathcal{L}(\mathbf{v})^{+\top} + \mathcal{L}(\mathbf{u})^-\mathcal{L}(\mathbf{v})^{-\top}, \\
& \mathcal{L}(\mathbf{u})^+\mathcal{U}(\mathbf{v})^{+\top} + \mathcal{L}(\mathbf{u})^+\mathcal{U}(\mathbf{v})^{-\top} + \mathcal{L}(\mathbf{u})^-\mathcal{U}(\mathbf{v})^{+\top} + \mathcal{L}(\mathbf{u})^-\mathcal{U}(\mathbf{v})^{-\top}, \\
& \mathcal{U}(\mathbf{u})^+\mathcal{L}(\mathbf{v})^{+\top} + \mathcal{U}(\mathbf{u})^+\mathcal{L}(\mathbf{v})^{-\top} + \mathcal{U}(\mathbf{u})^-\mathcal{L}(\mathbf{v})^{+\top} + \mathcal{U}(\mathbf{u})^-\mathcal{L}(\mathbf{v})^{-\top}, \\
& \mathcal{U}(\mathbf{u})^+\mathcal{U}(\mathbf{v})^{+\top} + \mathcal{U}(\mathbf{u})^+\mathcal{U}(\mathbf{v})^{-\top} + \mathcal{U}(\mathbf{u})^-\mathcal{U}(\mathbf{v})^{+\top} + \mathcal{U}(\mathbf{u})^-\mathcal{U}(\mathbf{v})^{-\top}\}.
\end{aligned} \tag{49}$$

where color is used to group related terms and ease the reading. Grouping by color, it is easily found that:

$$\begin{aligned}
\mathcal{L}(\mathbf{u})^+\mathcal{L}(\mathbf{v})^{+\top} &= \min\{\mathcal{L}(\mathbf{u})^+\mathcal{L}(\mathbf{v})^{+\top}, \mathcal{L}(\mathbf{u})^+\mathcal{U}(\mathbf{v})^{+\top}, \mathcal{U}(\mathbf{u})^+\mathcal{L}(\mathbf{v})^{+\top}, \mathcal{U}(\mathbf{u})^+\mathcal{U}(\mathbf{v})^{+\top}\} \\
\mathcal{L}(\mathbf{u})^+\mathcal{L}(\mathbf{v})^{-\top} &= \min\{\mathcal{L}(\mathbf{u})^+\mathcal{L}(\mathbf{v})^{-\top}, \mathcal{L}(\mathbf{u})^+\mathcal{U}(\mathbf{v})^{-\top}, \mathcal{U}(\mathbf{u})^+\mathcal{L}(\mathbf{v})^{-\top}, \mathcal{U}(\mathbf{u})^+\mathcal{U}(\mathbf{v})^{-\top}\} \\
\mathcal{L}(\mathbf{u})^-\mathcal{L}(\mathbf{v})^{+\top} &= \min\{\mathcal{L}(\mathbf{u})^-\mathcal{L}(\mathbf{v})^{+\top}, \mathcal{L}(\mathbf{u})^-\mathcal{U}(\mathbf{v})^{+\top}, \mathcal{U}(\mathbf{u})^-\mathcal{L}(\mathbf{v})^{+\top}, \mathcal{U}(\mathbf{u})^-\mathcal{U}(\mathbf{v})^{+\top}\} \\
\mathcal{L}(\mathbf{u})^-\mathcal{L}(\mathbf{v})^{-\top} &= \min\{\mathcal{L}(\mathbf{u})^-\mathcal{L}(\mathbf{v})^{-\top}, \mathcal{L}(\mathbf{u})^-\mathcal{U}(\mathbf{v})^{-\top}, \mathcal{U}(\mathbf{u})^-\mathcal{L}(\mathbf{v})^{-\top}, \mathcal{U}(\mathbf{u})^-\mathcal{U}(\mathbf{v})^{-\top}\}.
\end{aligned} \tag{50}$$

Lastly, by substituting each term by the minimum of the terms with the same color on Eq. (49):

$$M \geq \mathcal{L}(\mathbf{u})^+\mathcal{L}(\mathbf{v})^{+\top} + \mathcal{L}(\mathbf{u})^+\mathcal{L}(\mathbf{v})^{-\top} + \mathcal{L}(\mathbf{u})^-\mathcal{L}(\mathbf{v})^{+\top} + \mathcal{L}(\mathbf{u})^-\mathcal{L}(\mathbf{v})^{-\top} = \mathcal{L}(M). \tag{51}$$

Analogously for the upper bound:

$$\begin{aligned}
M \leq & \max\{\mathcal{L}(\mathbf{u})^+\mathcal{L}(\mathbf{v})^{+\top} + \mathcal{L}(\mathbf{u})^+\mathcal{L}(\mathbf{v})^{-\top} + \mathcal{L}(\mathbf{u})^-\mathcal{L}(\mathbf{v})^{+\top} + \mathcal{L}(\mathbf{u})^-\mathcal{L}(\mathbf{v})^{-\top}, \\
& \mathcal{L}(\mathbf{u})^+\mathcal{U}(\mathbf{v})^{+\top} + \mathcal{L}(\mathbf{u})^+\mathcal{U}(\mathbf{v})^{-\top} + \mathcal{L}(\mathbf{u})^-\mathcal{U}(\mathbf{v})^{+\top} + \mathcal{L}(\mathbf{u})^-\mathcal{U}(\mathbf{v})^{-\top}, \\
& \mathcal{U}(\mathbf{u})^+\mathcal{L}(\mathbf{v})^{+\top} + \mathcal{U}(\mathbf{u})^+\mathcal{L}(\mathbf{v})^{-\top} + \mathcal{U}(\mathbf{u})^-\mathcal{L}(\mathbf{v})^{+\top} + \mathcal{U}(\mathbf{u})^-\mathcal{L}(\mathbf{v})^{-\top}, \\
& \mathcal{U}(\mathbf{u})^+\mathcal{U}(\mathbf{v})^{+\top} + \mathcal{U}(\mathbf{u})^+\mathcal{U}(\mathbf{v})^{-\top} + \mathcal{U}(\mathbf{u})^-\mathcal{U}(\mathbf{v})^{+\top} + \mathcal{U}(\mathbf{u})^-\mathcal{U}(\mathbf{v})^{-\top}\},
\end{aligned} \tag{52}$$

where color is used to group related terms and ease the reading. Grouping by color, it is easily found that:

$$\begin{aligned}
\mathcal{U}(\mathbf{u})^+\mathcal{U}(\mathbf{v})^{+\top} &= \max\{\mathcal{L}(\mathbf{u})^+\mathcal{L}(\mathbf{v})^{+\top}, \mathcal{L}(\mathbf{u})^+\mathcal{U}(\mathbf{v})^{+\top}, \mathcal{U}(\mathbf{u})^+\mathcal{L}(\mathbf{v})^{+\top}, \mathcal{U}(\mathbf{u})^+\mathcal{U}(\mathbf{v})^{+\top}\} \\
\mathcal{L}(\mathbf{u})^+\mathcal{U}(\mathbf{v})^{-\top} &= \max\{\mathcal{L}(\mathbf{u})^+\mathcal{L}(\mathbf{v})^{-\top}, \mathcal{L}(\mathbf{u})^+\mathcal{U}(\mathbf{v})^{-\top}, \mathcal{U}(\mathbf{u})^+\mathcal{L}(\mathbf{v})^{-\top}, \mathcal{U}(\mathbf{u})^+\mathcal{U}(\mathbf{v})^{-\top}\} \\
\mathcal{L}(\mathbf{u})^-\mathcal{L}(\mathbf{v})^{+\top} &= \max\{\mathcal{L}(\mathbf{u})^-\mathcal{L}(\mathbf{v})^{+\top}, \mathcal{L}(\mathbf{u})^-\mathcal{U}(\mathbf{v})^{+\top}, \mathcal{U}(\mathbf{u})^-\mathcal{L}(\mathbf{v})^{+\top}, \mathcal{U}(\mathbf{u})^-\mathcal{U}(\mathbf{v})^{+\top}\} \\
\mathcal{L}(\mathbf{u})^-\mathcal{L}(\mathbf{v})^{-\top} &= \max\{\mathcal{L}(\mathbf{u})^-\mathcal{L}(\mathbf{v})^{-\top}, \mathcal{L}(\mathbf{u})^-\mathcal{U}(\mathbf{v})^{-\top}, \mathcal{U}(\mathbf{u})^-\mathcal{L}(\mathbf{v})^{-\top}, \mathcal{U}(\mathbf{u})^-\mathcal{U}(\mathbf{v})^{-\top}\}.
\end{aligned} \tag{53}$$

Lastly, by substituting each term by the minimum of the terms with the same color on Eq. (52):

$$M \leq \mathcal{U}(\mathbf{u})^+\mathcal{U}(\mathbf{v})^{+\top} + \mathcal{L}(\mathbf{u})^+\mathcal{U}(\mathbf{v})^{-\top} + \mathcal{L}(\mathbf{u})^-\mathcal{L}(\mathbf{v})^{+\top} + \mathcal{L}(\mathbf{u})^-\mathcal{L}(\mathbf{v})^{-\top} = \mathcal{U}(M). \tag{54}$$

□

Lemma 3. Let $M = \delta \cdot \mathbf{u}\mathbf{v}^\top$ be a rank-1 matrix defined by vectors $\mathbf{u} \in \mathbb{R}^{d_1}$ and $\mathbf{v} \in \mathbb{R}^{d_2}$ and a scalar $\delta \in \mathbb{R}$. Let \mathbf{u} (\mathbf{v}) be lower and upper bounded by $\mathcal{L}(\mathbf{u})$ and $\mathcal{U}(\mathbf{u})$ ($\mathcal{L}(\mathbf{v})$ and $\mathcal{U}(\mathbf{v})$). Let δ be

lower and upper bounded by $\mathcal{L}(\delta)$ and $\mathcal{U}(\delta)$. The matrix M is lower and upper bounded by:

$$\begin{aligned}
M &\geq \mathcal{L}(M) = \mathcal{L}(\delta)^+ \left[\mathcal{L}(u)^+ \mathcal{L}(v)^+{}^\top + \mathcal{U}(u)^- \mathcal{U}(v)^-{}^\top \right] \\
&\quad + \mathcal{U}(\delta)^+ \left[\mathcal{U}(u)^+ \mathcal{L}(v)^-{}^\top + \mathcal{L}(u)^- \mathcal{U}(v)^+{}^\top \right] \\
&\quad + \mathcal{L}(\delta)^- \left[\mathcal{U}(u)^+ \mathcal{U}(v)^+{}^\top + \mathcal{L}(u)^- \mathcal{L}(v)^-{}^\top \right] \\
&\quad + \mathcal{U}(\delta)^- \left[\mathcal{U}(u)^- \mathcal{L}(v)^+{}^\top + \mathcal{L}(u)^+ \mathcal{U}(v)^-{}^\top \right] \\
&= \mathcal{L}(\delta)^+ \mathcal{L}(uv^\top)^+ + \mathcal{U}(\delta)^+ \mathcal{L}(uv^\top)^- \\
&\quad + \mathcal{L}(\delta)^- \mathcal{U}(uv^\top)^+ + \mathcal{U}(\delta)^- \mathcal{U}(uv^\top)^- \\
M &\leq \mathcal{U}(M) = \mathcal{L}(\delta)^+ \left[\mathcal{U}(u)^- \mathcal{L}(v)^+{}^\top + \mathcal{L}(u)^+ \mathcal{U}(v)^-{}^\top \right] \\
&\quad + \mathcal{U}(\delta)^+ \left[\mathcal{U}(u)^+ \mathcal{U}(v)^+{}^\top + \mathcal{L}(u)^- \mathcal{L}(v)^-{}^\top \right] \\
&\quad + \mathcal{L}(\delta)^- \left[\mathcal{U}(u)^+ \mathcal{L}(v)^-{}^\top + \mathcal{L}(u)^- \mathcal{U}(v)^+{}^\top \right] \\
&\quad + \mathcal{U}(\delta)^- \left[\mathcal{L}(u)^+ \mathcal{L}(v)^+{}^\top + \mathcal{U}(u)^- \mathcal{U}(v)^-{}^\top \right] \\
&= \mathcal{L}(\delta)^+ \mathcal{U}(uv^\top)^- + \mathcal{U}(\delta)^+ \mathcal{U}(uv^\top)^+ \\
&\quad + \mathcal{L}(\delta)^- \mathcal{L}(uv^\top)^- + \mathcal{U}(\delta)^- \mathcal{L}(uv^\top)^+
\end{aligned} \tag{55}$$

Proof of Lemma 3. By applying the multiplication rule of IBP, see Eq. (9), we can express the lower bound as:

$$\begin{aligned}
M &\geq \min\{ \mathcal{L}(\delta) \mathcal{L}(uv^\top), \\
&\quad \mathcal{L}(\delta) \mathcal{U}(uv^\top), \\
&\quad \mathcal{U}(\delta) \mathcal{L}(uv^\top), \\
&\quad \mathcal{U}(\delta) \mathcal{U}(uv^\top) \}.
\end{aligned} \tag{56}$$

Then, by decomposing each term via the possitive-negative decomposition of each operand, we arrive to:

$$\begin{aligned}
M &\geq \min\{ \mathcal{L}(\delta)^+ \mathcal{L}(uv^\top)^+ + \mathcal{L}(\delta)^+ \mathcal{L}(uv^\top)^- + \mathcal{L}(\delta)^- \mathcal{L}(uv^\top)^+ + \mathcal{L}(\delta)^- \mathcal{L}(uv^\top)^-, \\
&\quad \mathcal{L}(\delta)^+ \mathcal{U}(uv^\top)^+ + \mathcal{L}(\delta)^+ \mathcal{U}(uv^\top)^- + \mathcal{L}(\delta)^- \mathcal{U}(uv^\top)^+ + \mathcal{L}(\delta)^- \mathcal{U}(uv^\top)^-, \\
&\quad \mathcal{U}(\delta)^+ \mathcal{L}(uv^\top)^+ + \mathcal{U}(\delta)^+ \mathcal{L}(uv^\top)^- + \mathcal{U}(\delta)^- \mathcal{L}(uv^\top)^+ + \mathcal{U}(\delta)^- \mathcal{L}(uv^\top)^-, \\
&\quad \mathcal{U}(\delta)^+ \mathcal{U}(uv^\top)^+ + \mathcal{U}(\delta)^+ \mathcal{U}(uv^\top)^- + \mathcal{U}(\delta)^- \mathcal{U}(uv^\top)^+ + \mathcal{U}(\delta)^- \mathcal{U}(uv^\top)^- \}.
\end{aligned} \tag{57}$$

where color is used to group related terms and ease the reading. Grouping by color, it is easily found that:

$$\begin{aligned}
\mathcal{L}(\delta)^+ \mathcal{L}(uv^\top)^+ &= \min\{ \mathcal{L}(\delta)^+ \mathcal{L}(uv^\top)^+, \mathcal{L}(\delta)^+ \mathcal{U}(uv^\top)^+, \mathcal{U}(\delta)^+ \mathcal{L}(uv^\top)^+, \mathcal{U}(\delta)^+ \mathcal{U}(uv^\top)^+ \} \\
\mathcal{L}(\delta)^+ \mathcal{L}(uv^\top)^- &= \min\{ \mathcal{L}(\delta)^+ \mathcal{L}(uv^\top)^-, \mathcal{L}(\delta)^+ \mathcal{U}(uv^\top)^-, \mathcal{U}(\delta)^+ \mathcal{L}(uv^\top)^-, \mathcal{U}(\delta)^+ \mathcal{U}(uv^\top)^- \} \\
\mathcal{L}(\delta)^- \mathcal{L}(uv^\top)^+ &= \min\{ \mathcal{L}(\delta)^- \mathcal{L}(uv^\top)^+, \mathcal{L}(\delta)^- \mathcal{U}(uv^\top)^+, \mathcal{U}(\delta)^- \mathcal{L}(uv^\top)^+, \mathcal{U}(\delta)^- \mathcal{U}(uv^\top)^+ \} \\
\mathcal{L}(\delta)^- \mathcal{L}(uv^\top)^- &= \min\{ \mathcal{L}(\delta)^- \mathcal{L}(uv^\top)^-, \mathcal{L}(\delta)^- \mathcal{U}(uv^\top)^-, \mathcal{U}(\delta)^- \mathcal{L}(uv^\top)^-, \mathcal{U}(\delta)^- \mathcal{U}(uv^\top)^- \}.
\end{aligned} \tag{58}$$

Lastly, by substituting each term by the minimum of the terms with the same color on Eq. (49):

$$M \geq \mathcal{L}(\delta)^+ \mathcal{L}(uv^\top)^+ + \mathcal{U}(\delta)^+ \mathcal{L}(uv^\top)^- + \mathcal{L}(\delta)^- \mathcal{U}(uv^\top)^+ + \mathcal{U}(\delta)^- \mathcal{U}(uv^\top)^- = \mathcal{L}(M) \tag{59}$$

Analogously for the upper bound:

$$\begin{aligned} \mathbf{M} \leq & \max\{\mathcal{L}(\delta)^+ \mathcal{L}(\mathbf{u}\mathbf{v}^\top)^+ + \mathcal{L}(\delta)^+ \mathcal{L}(\mathbf{u}\mathbf{v}^\top)^- + \mathcal{L}(\delta)^- \mathcal{L}(\mathbf{u}\mathbf{v}^\top)^+ + \mathcal{L}(\delta)^- \mathcal{L}(\mathbf{u}\mathbf{v}^\top)^-, \\ & \mathcal{L}(\delta)^+ \mathcal{U}(\mathbf{u}\mathbf{v}^\top)^+ + \mathcal{L}(\delta)^+ \mathcal{U}(\mathbf{u}\mathbf{v}^\top)^- + \mathcal{L}(\delta)^- \mathcal{U}(\mathbf{u}\mathbf{v}^\top)^+ + \mathcal{L}(\delta)^- \mathcal{U}(\mathbf{u}\mathbf{v}^\top)^-, \\ & \mathcal{U}(\delta)^+ \mathcal{L}(\mathbf{u}\mathbf{v}^\top)^+ + \mathcal{U}(\delta)^+ \mathcal{L}(\mathbf{u}\mathbf{v}^\top)^- + \mathcal{U}(\delta)^- \mathcal{L}(\mathbf{u}\mathbf{v}^\top)^+ + \mathcal{U}(\delta)^- \mathcal{L}(\mathbf{u}\mathbf{v}^\top)^-, \\ & \mathcal{U}(\delta)^+ \mathcal{U}(\mathbf{u}\mathbf{v}^\top)^+ + \mathcal{U}(\delta)^+ \mathcal{U}(\mathbf{u}\mathbf{v}^\top)^- + \mathcal{U}(\delta)^- \mathcal{U}(\mathbf{u}\mathbf{v}^\top)^+ + \mathcal{U}(\delta)^- \mathcal{U}(\mathbf{u}\mathbf{v}^\top)^-\}. \end{aligned} \quad (60)$$

where color is used to group related terms and ease the reading. Grouping by color, it is easily found that:

$$\begin{aligned} \mathcal{U}(\delta)^+ \mathcal{U}(\mathbf{u}\mathbf{v}^\top)^+ &= \max\{\mathcal{L}(\delta)^+ \mathcal{L}(\mathbf{u}\mathbf{v}^\top)^+, \mathcal{L}(\delta)^+ \mathcal{U}(\mathbf{u}\mathbf{v}^\top)^+, \mathcal{U}(\delta)^+ \mathcal{L}(\mathbf{u}\mathbf{v}^\top)^+, \mathcal{U}(\delta)^+ \mathcal{U}(\mathbf{u}\mathbf{v}^\top)^+\} \\ \mathcal{L}(\delta)^+ \mathcal{U}(\mathbf{u}\mathbf{v}^\top)^- &= \max\{\mathcal{L}(\delta)^+ \mathcal{L}(\mathbf{u}\mathbf{v}^\top)^-, \mathcal{L}(\delta)^+ \mathcal{U}(\mathbf{u}\mathbf{v}^\top)^-, \mathcal{U}(\delta)^+ \mathcal{L}(\mathbf{u}\mathbf{v}^\top)^-, \mathcal{U}(\delta)^+ \mathcal{U}(\mathbf{u}\mathbf{v}^\top)^-\} \\ \mathcal{U}(\delta)^- \mathcal{L}(\mathbf{u}\mathbf{v}^\top)^+ &= \max\{\mathcal{L}(\delta)^- \mathcal{L}(\mathbf{u}\mathbf{v}^\top)^+, \mathcal{L}(\delta)^- \mathcal{U}(\mathbf{u}\mathbf{v}^\top)^+, \mathcal{U}(\delta)^- \mathcal{L}(\mathbf{u}\mathbf{v}^\top)^+, \mathcal{U}(\delta)^- \mathcal{U}(\mathbf{u}\mathbf{v}^\top)^+\} \\ \mathcal{L}(\delta)^- \mathcal{L}(\mathbf{u}\mathbf{v}^\top)^- &= \max\{\mathcal{L}(\delta)^- \mathcal{L}(\mathbf{u}\mathbf{v}^\top)^-, \mathcal{L}(\delta)^- \mathcal{U}(\mathbf{u}\mathbf{v}^\top)^-, \mathcal{U}(\delta)^- \mathcal{L}(\mathbf{u}\mathbf{v}^\top)^-, \mathcal{U}(\delta)^- \mathcal{U}(\mathbf{u}\mathbf{v}^\top)^-\}. \end{aligned} \quad (61)$$

Lastly, by substituting each term by the minimum of the terms with the same color on Eq. (49):

$$\mathbf{M} \leq \mathcal{U}(\delta)^+ \mathcal{U}(\mathbf{u}\mathbf{v}^\top)^+ + \mathcal{L}(\delta)^+ \mathcal{U}(\mathbf{u}\mathbf{v}^\top)^- + \mathcal{U}(\delta)^- \mathcal{L}(\mathbf{u}\mathbf{v}^\top)^+ + \mathcal{L}(\delta)^- \mathcal{L}(\mathbf{u}\mathbf{v}^\top)^- = \mathcal{U}(\mathbf{M}) \quad (62)$$

□

F.2 Properties of the \mathcal{M} operator

We also define the operator $\mathcal{M}(\cdot) = \max\{|\mathcal{L}(\cdot)|, |\mathcal{U}(\cdot)|\}$ and certain useful properties about it.

Lemma 4. Let $\mathbf{M} = \mathcal{M}(\mathbf{u}\mathbf{v}^\top)$ be a matrix defined by vectors $\mathbf{u} \in \mathbb{R}^d$ and $\mathbf{v} \in \mathbb{R}^d$. Let \mathbf{u} (\mathbf{v}) be lower and upper bounded by $\mathcal{L}(\mathbf{u})$ and $\mathcal{U}(\mathbf{u})$ ($\mathcal{L}(\mathbf{v})$ and $\mathcal{U}(\mathbf{v})$). Let $\hat{\mathbf{u}} = \mathcal{M}(\mathbf{u})$ and $\hat{\mathbf{v}} = \mathcal{M}(\mathbf{v})$. The matrix \mathbf{M} can be expressed as:

$$\mathbf{M} = \hat{\mathbf{u}}\hat{\mathbf{v}}^\top, \quad (63)$$

resulting in a rank-1 matrix given by vectors $\hat{\mathbf{u}}$ and $\hat{\mathbf{v}}$.

Proof of Lemma 4. Starting with the definition of \mathbf{M} , we have

$$\begin{aligned} \mathbf{M} &:= \max\{|\mathcal{L}(\mathbf{u}\mathbf{v}^\top)|, |\mathcal{U}(\mathbf{u}\mathbf{v}^\top)|\} \\ &= \max\{|\min\{\mathcal{L}(\mathbf{u})\mathcal{L}(\mathbf{v})^\top, \mathcal{L}(\mathbf{u})\mathcal{U}(\mathbf{v})^\top, \mathcal{U}(\mathbf{u})\mathcal{L}(\mathbf{v})^\top, \mathcal{U}(\mathbf{u})\mathcal{U}(\mathbf{v})^\top\}|, \\ &\quad |\max\{\mathcal{L}(\mathbf{u})\mathcal{L}(\mathbf{v})^\top, \mathcal{L}(\mathbf{u})\mathcal{U}(\mathbf{v})^\top, \mathcal{U}(\mathbf{u})\mathcal{L}(\mathbf{v})^\top, \mathcal{U}(\mathbf{u})\mathcal{U}(\mathbf{v})^\top\}| \} \quad [\text{Def. Eq. (9)}] \\ &= \max\{|\mathcal{L}(\mathbf{u})\mathcal{L}(\mathbf{v})^\top|, |\mathcal{L}(\mathbf{u})\mathcal{U}(\mathbf{v})^\top|, |\mathcal{U}(\mathbf{u})\mathcal{L}(\mathbf{v})^\top|, |\mathcal{U}(\mathbf{u})\mathcal{U}(\mathbf{v})^\top|\} \\ &= \max\{|\mathcal{L}(\mathbf{u})||\mathcal{L}(\mathbf{v})^\top|, |\mathcal{L}(\mathbf{u})||\mathcal{U}(\mathbf{v})^\top|, |\mathcal{U}(\mathbf{u})||\mathcal{L}(\mathbf{v})^\top|, |\mathcal{U}(\mathbf{u})||\mathcal{U}(\mathbf{v})^\top|\} \\ &= \max\{\max\{|\mathcal{L}(\mathbf{u})||\mathcal{L}(\mathbf{v})^\top|, |\mathcal{L}(\mathbf{u})||\mathcal{U}(\mathbf{v})^\top|\}, \\ &\quad \max\{|\mathcal{U}(\mathbf{u})||\mathcal{L}(\mathbf{v})^\top|, |\mathcal{U}(\mathbf{u})||\mathcal{U}(\mathbf{v})^\top|\}\} \\ &= \max\{|\mathcal{L}(\mathbf{u})|\max\{|\mathcal{L}(\mathbf{v})^\top|, |\mathcal{U}(\mathbf{v})^\top|\}, |\mathcal{U}(\mathbf{u})|\max\{|\mathcal{L}(\mathbf{v})^\top|, |\mathcal{U}(\mathbf{v})^\top|\}\} \\ &= \max\{|\mathcal{L}(\mathbf{u})|, |\mathcal{U}(\mathbf{u})|\}\max\{|\mathcal{L}(\mathbf{v})^\top|, |\mathcal{U}(\mathbf{v})^\top|\} \\ &= \hat{\mathbf{u}}\hat{\mathbf{v}}^\top, \end{aligned} \quad (64)$$

which concludes the proof. □

Lemma 5. Let $\mathbf{M} = \mathcal{M}(\mathbf{A} + \mathbf{B}) \in \mathbb{R}^{d_1 \times d_2}$ with $\mathbf{A} \in \mathbb{R}^{d_1 \times d_2}$ and $\mathbf{B} \in \mathbb{R}^{d_1 \times d_2}$ being lower and upper bounded matrices. The matrix \mathbf{M} satisfies:

$$\mathbf{M} \leq \mathcal{M}(\mathbf{A}) + \mathcal{M}(\mathbf{B}). \quad (65)$$

Proof of Lemma 5. Starting with the definition of \mathcal{M} :

$$\begin{aligned}
\mathbf{M} &:= \max\{|\mathcal{L}(\mathbf{A} + \mathbf{B})|, |\mathcal{U}(\mathbf{A} + \mathbf{B})|\} = \max\{|\mathcal{L}(\mathbf{A}) + \mathcal{L}(\mathbf{B})|, |\mathcal{U}(\mathbf{A}) + \mathcal{U}(\mathbf{B})|\} \\
&\leq \max\{|\mathcal{L}(\mathbf{A})| + |\mathcal{L}(\mathbf{B})|, |\mathcal{U}(\mathbf{A})| + |\mathcal{U}(\mathbf{B})|\} \\
&\leq \max\{\max\{|\mathcal{L}(\mathbf{A})|, |\mathcal{U}(\mathbf{A})|\} + |\mathcal{L}(\mathbf{B})|, \max\{|\mathcal{L}(\mathbf{A})|, |\mathcal{U}(\mathbf{A})|\} + |\mathcal{U}(\mathbf{B})|\} \\
&\leq \max\{\max\{|\mathcal{L}(\mathbf{A})|, |\mathcal{U}(\mathbf{A})|\} + \max\{|\mathcal{L}(\mathbf{B})|, |\mathcal{U}(\mathbf{B})|\}, \\
&\quad \max\{|\mathcal{L}(\mathbf{A})|, |\mathcal{U}(\mathbf{A})|\} + \max\{|\mathcal{L}(\mathbf{B})|, |\mathcal{U}(\mathbf{B})|\}\} \\
&= \max\{|\mathcal{L}(\mathbf{A})|, |\mathcal{U}(\mathbf{A})|\} + \max\{|\mathcal{L}(\mathbf{B})|, |\mathcal{U}(\mathbf{B})|\} = \mathcal{M}(\mathbf{A}) + \mathcal{M}(\mathbf{B}),
\end{aligned} \tag{66}$$

we conclude the proof. \square

Lemma 6. Let $\mathbf{M} = \mathcal{M}(\delta \mathbf{A}) \in \mathbb{R}^{d_1 \times d_2}$ with $\mathbf{A} \in \mathbb{R}^{d_1 \times d_2}$ and $\delta \in \mathbb{R}$ being lower and upper bounded. The matrix \mathbf{M} satisfies:

$$\mathbf{M} = \mathcal{M}(\delta) \mathcal{M}(\mathbf{A}). \tag{67}$$

Proof of Lemma 6. Starting with the definition of \mathcal{M} :

$$\begin{aligned}
\mathbf{M} &:= \max\{|\mathcal{L}(\delta \mathbf{A})|, |\mathcal{U}(\delta \mathbf{A})|\} \\
&= \max\{|\min\{\mathcal{L}(\delta) \mathcal{L}(\mathbf{A}), \mathcal{L}(\delta) \mathcal{U}(\mathbf{A}), \mathcal{U}(\delta) \mathcal{L}(\mathbf{A}), \mathcal{U}(\delta) \mathcal{U}(\mathbf{A})\}|, \\
&\quad |\max\{\mathcal{L}(\delta) \mathcal{L}(\mathbf{A}), \mathcal{L}(\delta) \mathcal{U}(\mathbf{A}), \mathcal{U}(\delta) \mathcal{L}(\mathbf{A}), \mathcal{U}(\delta) \mathcal{U}(\mathbf{A})\}|\} \quad [\text{Eq. (9)}] \\
&= \max\{|\mathcal{L}(\delta) \mathcal{L}(\mathbf{A})|, |\mathcal{L}(\delta) \mathcal{U}(\mathbf{A})|, |\mathcal{U}(\delta) \mathcal{L}(\mathbf{A})|, |\mathcal{U}(\delta) \mathcal{U}(\mathbf{A})|\} \\
&= \max\{|\mathcal{L}(\delta)| |\mathcal{L}(\mathbf{A})|, |\mathcal{L}(\delta)| |\mathcal{U}(\mathbf{A})|, |\mathcal{U}(\delta)| |\mathcal{L}(\mathbf{A})|, |\mathcal{U}(\delta)| |\mathcal{U}(\mathbf{A})|\} \\
&= \max\{\max\{|\mathcal{L}(\delta)| |\mathcal{L}(\mathbf{A})|, |\mathcal{L}(\delta)| |\mathcal{U}(\mathbf{A})|\}, \\
&\quad \max\{|\mathcal{U}(\delta)| |\mathcal{L}(\mathbf{A})|, |\mathcal{U}(\delta)| |\mathcal{U}(\mathbf{A})|\}\} \\
&= \max\{|\mathcal{L}(\delta)| \max\{|\mathcal{L}(\mathbf{A})|, |\mathcal{U}(\mathbf{A})|\}, \\
&\quad |\mathcal{U}(\delta)| \max\{|\mathcal{L}(\mathbf{A})|, |\mathcal{U}(\mathbf{A})|\}\} \\
&= \max\{|\mathcal{L}(\delta)|, |\mathcal{U}(\delta)|\} \max\{|\mathcal{L}(\mathbf{A})|, |\mathcal{U}(\mathbf{A})|\} = \mathcal{M}(\delta) \mathcal{M}(\mathbf{A}),
\end{aligned} \tag{68}$$

we finish the proof. \square

Proof of Theorem 1. Firstly, applying the \mathcal{M} operator to the Hessian of the verification objective results in:

$$\begin{aligned}
\mathcal{M}(\mathbf{H}_g(\mathbf{z})) &:= \mathcal{M}\left(\sum_{i=1}^k (c_{ti} - c_{\gamma i}) \nabla_{\mathbf{z}\mathbf{z}} x_i^{(N)}\right) \quad [\text{Eq. (12)}] \\
&\leq \sum_{i=1}^k \mathcal{M}\left((c_{ti} - c_{\gamma i}) \nabla_{\mathbf{z}\mathbf{z}} x_i^{(N)}\right) \quad [\text{Lemma 5}] \\
&= \sum_{i=1}^k |c_{ti} - c_{\gamma i}| \mathcal{M}\left(\nabla_{\mathbf{z}\mathbf{z}} x_i^{(N)}\right) \quad [\text{Lemma 6}].
\end{aligned} \tag{69}$$

Secondly, we can proceed analogously with the recursive relationship of the Hessian at different layers. For $n = 2, \dots, N$:

$$\begin{aligned}
\mathcal{M}\left(\nabla_{\mathbf{z}\mathbf{z}} x_i^{(n)}\right) &:= \mathcal{M}\left(\nabla_{\mathbf{z}\mathbf{z}} x_i^{(n-1)} \mathbf{w}_{[n]:i}^\top + \mathbf{w}_{[n]:i} \nabla_{\mathbf{z}\mathbf{z}} x_i^{(n-1)\top} \right. \\
&\quad \left. + (\mathbf{w}_{[n]:i}^\top \mathbf{z} + 1) \nabla_{\mathbf{z}\mathbf{z}} x_i^{(n-1)}\right) \quad [\text{Eq. (14)}] \\
&\leq \mathcal{M}\left(\nabla_{\mathbf{z}\mathbf{z}} x_i^{(n-1)} \mathbf{w}_{[n]:i}^\top\right) + \mathcal{M}\left(\mathbf{w}_{[n]:i} \nabla_{\mathbf{z}\mathbf{z}} x_i^{(n-1)\top}\right) \\
&\quad + \mathcal{M}\left((\mathbf{w}_{[n]:i}^\top \mathbf{z} + 1) \nabla_{\mathbf{z}\mathbf{z}} x_i^{(n-1)}\right) \quad [\text{Lemma 5}] \\
&= \mathcal{M}\left(\nabla_{\mathbf{z}\mathbf{z}} x_i^{(n-1)}\right) |\mathbf{w}_{[n]:i}^\top| + |\mathbf{w}_{[n]:i}| \mathcal{M}\left(\nabla_{\mathbf{z}\mathbf{z}} x_i^{(n-1)\top}\right) \quad [\text{Lemma 4}] \\
&\quad + \mathcal{M}(\mathbf{w}_{[n]:i}^\top \mathbf{z} + 1) \mathcal{M}\left(\nabla_{\mathbf{z}\mathbf{z}} x_i^{(n-1)}\right) \quad [\text{Lemma 6}].
\end{aligned} \tag{70}$$

Lastly, for $n = 1$, by definition $\nabla_{\mathbf{z}\mathbf{z}} x_i^{(1)} = \mathbf{0}_{d \times d}$, which means $\mathcal{M}(\nabla_{\mathbf{z}\mathbf{z}} x_i^{(1)}) = \mathbf{0}_{d \times d}$ \square

F.3 Convergence of the BaB algorithm to the global minima

In this Section we demonstrate a key property for verification: convergence to the global minima of Eq. (4). Let us firstly define some concepts of the BaB algorithm.

Let S_k be the subset picked at iteration k of the BaB algorithm and $\{S_{kq} | q = 0, 1, \dots\}$ be the sequence of recursive subsets, so that $S_{kq} \subset S_{kq-1}$ and $S_{k0} = S_k$. Let $\mathcal{L}(S_{kq})$ be the lower bound of subset S_{kq} and $\mathcal{L}_{kq} = \min\{\mathcal{L}(S_{kq}) | q = 0, 1, \dots\}$ the lower bound in the branch rooted by subset S_k , we analogously define $\mathcal{U}(S_{kq})$ and \mathcal{U}_{kq} . Finally, let a **fathomed** or **pruned** set S_{kq} be a set where $\mathcal{L}(S_{kq}) > \mathcal{U}_{kq}$. From Horst and Tuy [1996, Definition IV.4]:

Definition 1. A bounding operation is called **consistent** if at every step any unfathomed subset can be further split, and if any infinitely decreasing sequence $\{S_{kq} | q = 0, 1, \dots\}$ for successively refined partition elements satisfies:

$$\lim_{q \rightarrow \infty} (\mathcal{U}_{kq} - \mathcal{L}(S_{kq})) = 0. \quad (71)$$

Remark. Because $\mathcal{U}(S_{kq}) \geq \mathcal{U}_{kq} \geq \mathcal{L}(S_{kq})$, it suffices to prove that $\lim_{q \rightarrow \infty} (\mathcal{U}(S_{kq}) - \mathcal{L}(S_{kq}))$ to show a bounding operation is consistent.

Another relevant property is the subset selection being **bound improving**. Let \mathcal{P}_k be the set of unfathomed subsets at iteration k (probs variable in Algorithm 1), from Horst and Tuy [1996, Definition IV.6]:

Definition 2. A subset selection operation is called **bound improving** if, at least each time after a finite number of steps, S_k satisfies the relation:

$$S_k = \arg \min\{\mathcal{L}(S) : S \in \mathcal{P}_k\}. \quad (72)$$

Then, this ensures at least one partition element where the actual lower bound is attained is selected for further partition in step k of the algorithm.

Finally, Horst and Tuy [1996, Theorem IV.3] cover global convergence of general BaB algorithms.

Theorem 2. In a BB procedure, suppose that the bounding operation is consistent and the subset selection operation is bound improving. Then the procedure is convergent:

$$\mathcal{L} := \lim_{k \rightarrow \infty} \mathcal{L}_k = \lim_{k \rightarrow \infty} f(\mathbf{z}^{(k)}) = \min_{\mathbf{z} \in C_{in}} f(\mathbf{z}) = \lim_{k \rightarrow \infty} \mathcal{U}_k = \mathcal{U}, \quad (73)$$

where C_{in} is the feasible set of the initial problem, \mathcal{L}_k and \mathcal{U}_k are the global lower and upper bounds at iteration k and $\mathbf{z}^{(k)}$ is the point where the upper bound \mathcal{U}_k is attained.

Remark. In Theorem 2, \mathcal{L}_k and \mathcal{U}_k refer to variables `global_lb` and `global_ub` respectively in Algorithm 1.

Lemma 7. Selecting the subset with the lowest lower bound at every iteration of a BaB algorithm is a bound improving subset selection strategy.

Proof. By definition, when selecting the subset with the lowest lower bound, we are selecting at every iteration $S_k = \arg \min\{\mathcal{L}(S) : S \in \mathcal{P}_k\}$, which means that Eq. (72) holds at every iteration k and the strategy is bound improving. \square

Definition 3. Let subset S_k , $\{S_{kq} | q = 0, 1, \dots\}$ the sequence of recursive subsets rooted at $S_{k0} = S_k$ a branching operation is convergent iff $\lim_{q \rightarrow \infty} |S_{kq}| = 0$.

Lemma 8. Selecting the widest interval index $i = \arg \max \mathbf{u} - \mathbf{l}$ to split a problem, is a convergent branching operation.

Proof. Suppose at subset S_{kq} , we have bounds $\mathbf{l}^{(q)}$ and $\mathbf{u}^{(q)}$ and indexes i_1, i_2, \dots, i_d so that $u_{i_1}^{(q)} - l_{i_1}^{(q)} \geq u_{i_2}^{(q)} - l_{i_2}^{(q)} \geq \dots \geq u_{i_d}^{(q)} - l_{i_d}^{(q)}$ is the decreasing ordered list of interval widths, then $|S_{kq}| = \|\mathbf{u}^{(q)} - \mathbf{l}^{(q)}\|_2 = \sqrt{\sum_{j=1}^d (u_j^{(q)} - l_j^{(q)})^2} \leq \sqrt{\sum_{j=1}^d (u_{i_1}^{(q)} - l_{i_1}^{(q)})^2} = \sqrt{d(u_{i_1}^{(q)} - l_{i_1}^{(q)})^2} =$

$(u_{i_1}^{(q)} - l_{i_1}^{(q)})\sqrt{d}$. Then, at subset S_{k_q+1} , with bounds $l^{(q+1)}$ and $u^{(q+1)}$ and new indices j_1, j_2, \dots, j_d , $u_{i_1}^{(q+1)} - l_{i_1}^{(q+1)} = \frac{u_{i_1}^{(q)} - l_{i_1}^{(q)}}{2}$ and then j_1 is either equal to i_1 or to i_2 , depending on whether $\frac{u_{i_1}^{(q)} - l_{i_1}^{(q)}}{2} > u_{i_2}^{(q)} - l_{i_2}^{(q)}$ or not. Finally, as:

$$\begin{cases} u_{i_1}^{(q)} - l_{i_1}^{(q)} > \frac{u_{i_1}^{(q)} - l_{i_1}^{(q)}}{2} = u_{j_1}^{(q+1)} - l_{j_1}^{(q+1)} & \text{if } \frac{u_{i_1}^{(q)} - l_{i_1}^{(q)}}{2} > u_{i_2}^{(q)} - l_{i_2}^{(q)} \\ u_{i_1}^{(q)} - l_{i_1}^{(q)} > u_{i_2}^{(q)} - l_{i_2}^{(q)} = u_{j_1}^{(q+1)} - l_{j_1}^{(q+1)} & \text{if } \frac{u_{i_1}^{(q)} - l_{i_1}^{(q)}}{2} \leq u_{i_2}^{(q)} - l_{i_2}^{(q)} \end{cases} \quad (74)$$

and $u_{i_1}^{(q)} - l_{i_1}^{(q)} \geq 0$, the sequence $\{u_{i_1}^{(q)} - l_{i_1}^{(q)}\}_q$ is strictly decreasing and lower bounded by 0, then $\lim_{q \rightarrow \infty} (u_{i_1}^{(q)} - l_{i_1}^{(q)}) = 0$ must hold. Then $\lim_{q \rightarrow 0} (u_{i_1}^{(q)} - l_{i_1}^{(q)})\sqrt{d} = 0$ and as $(u_{i_1}^{(q)} - l_{i_1}^{(q)})\sqrt{d} \geq |S_{k_q}| \geq 0$, the property $\lim_{q \rightarrow \infty} |S_{k_q}| = 0$ holds and by Definition 3, the branching operation is convergent. \square

A consequence of Lemma 8 is the following Corollary.

Corollary 3. For any $z \in [l^{(q)}, u^{(q)}]$, if $\lim_{q \rightarrow \infty} |S_{k_q}| = 0$, in the limit $l_i^{(q)} = z_i = u_i^{(q)} \forall i = 1, \dots, d, \forall z \in [l, u]$.

Lemma 9. Let widest interval selection be the branching operation, lower bounds obtained by α -convexification with $\alpha \geq \max\{0, -\frac{1}{2} \min\{\lambda_{\min}(\mathbf{H}_f(z)) : z \in [l, u]\}\}$ and upper bounds obtained by evaluating $\mathcal{U}(S_{k_q}) = g(z^{\text{upper}})$ for any $z^{\text{upper}} \in [l, u]$ are consistent.

Proof. By Definition 1 is sufficient to check that $\lim_{q \rightarrow \infty} (\mathcal{U}(S_{k_q}) - \mathcal{L}(S_{k_q})) = 0$. By definition, the lower bound $\mathcal{L}(S_{k_q})$ is the solution to the function in Eq. (5) subject to $z \in [l, u]$, which will lead to an optimal z^{opt} and $\mathcal{L}(S_{k_q}) = g_\alpha(z^{\text{opt}})$. If the upper bound is given by evaluating the objective function at any point $z^{\text{upper}} \in [l, u]$ e.g., $z^{\text{upper}} = z^{\text{opt}}$ or in our case $z^{\text{upper}} = z^{\text{PGD}}$, the point obtained by performing PGD over g , $\mathcal{U}(S_{k_q}) = f(z^{\text{upper}})$, and their difference is: $\mathcal{U}(S_{k_q}) - \mathcal{L}(S_{k_q}) = g(z^{\text{upper}}) - g(z^{\text{opt}}) - \alpha \sum_{i=1}^d (z_i^{\text{opt}} - l_i^{(q)})(z_i^{\text{opt}} - u_i^{(q)})$. By virtue of Theorem 3, in the limit, $l^{(q)} = z^{\text{opt}} = z^{\text{upper}} = u^{(q)}$ and therefore $\lim_{q \rightarrow \infty} \mathcal{U}(S_{k_q}) - \mathcal{L}(S_{k_q}) = g(l) - g(l) - \alpha \sum_{i=1}^d (l_i^{(q)} - l_i^{(q)})(l_i^{(q)} - l_i^{(q)}) = 0$ and the bounds are consistent. \square

Lemma 10. Let widest interval selection be the branching operation, lower bounds obtained by IBP and upper bounds obtained by evaluating $\mathcal{U}(S_{k_q}) = g(z^{\text{upper}})$ for any $z^{\text{upper}} \in [l, u]$ are consistent.

Proof. By Definition 1 is sufficient to check that $\lim_{q \rightarrow \infty} (\mathcal{U}(S_{k_q}) - \mathcal{L}(S_{k_q})) = 0$. By definition, the lower bound $\mathcal{L}(S_{k_q})$ is given by $\mathcal{L}(g(z)) = \mathcal{L}(f(z)_t) - \mathcal{U}(f(z)_\gamma)$, see Section 3.1, and the upper bound is given by $\mathcal{U}(S_{k_q}) = g(z^{\text{upper}})$. By virtue of Theorem 3, in the limit, $l^{(q)} = z^{\text{upper}} = u^{(q)}$, then $\lim_{q \rightarrow \infty} \mathcal{U}(S_{k_q}) = g(l)$. Then, for Eq. (32), it can be easily found that $\lim_{q \rightarrow \infty} \mathcal{L}(\hat{x}^{(n)}) = \lim_{q \rightarrow \infty} \mathcal{U}(\hat{x}^{(n)}) = \mathbf{W}_{[n]}^\top l$. Then, for Eq. (33), every element in the sets S will converge to the same value and therefore, $\mathcal{L}(x^{(n)}) = \min S = \max S = \mathcal{U}(x^{(n)})$ for both CCP and NCP for every layer n . Finally, for the network's output, because $\mathcal{L}(x_i^{(N)}) = \mathcal{U}(x_i^{(N)})$, then, from Eq. (34), $\mathcal{L}(f(z)_i) = \mathcal{U}(f(z)_i) = f(l)_i, \forall i \in [d]$. Then, $\mathcal{L}(S_{k_q}) = \mathcal{L}(f(z)_t - f(z)_\gamma) = \mathcal{L}(f(z)_t) - \mathcal{U}(f(z)_\gamma) = f(l)_t - f(l)_\gamma = f(z^{\text{upper}})_t - f(z^{\text{upper}})_\gamma = \mathcal{U}(S_{k_q})$ and the property holds. \square

A consequence of Lemmas 8 to 10 is:

Corollary 4. Any branch and bound procedure with widest interval selection for branching, lowest lower bound subset selection and a bounding operation of: IBP or α -convexification, is convergent.

Proof. Because of Lemma 8, we have convergence of the branching procedure. Then, due to Lemmas 9 and 10, we have bound consistency for both IBP and α -convexification bounding mechanisms. By Lemma 7, we have that the subset selection strategy is bound improving. Therefore, because of Theorem 2, we have that any BaB algorithm with these properties converges to a global minimizer. \square

F.4 Lower bound of the minimum eigenvalue of the Hessian of PNs

Proof of Proposition 1. Using the IBP rules from Section 3.1 and Appendix C.5, we can develop:

$$\begin{aligned}
\mathcal{L}(\delta \cdot \nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n)}) &= \mathcal{L}(\delta \cdot (\nabla_{\mathbf{z}} x_i^{(n-1)} \mathbf{w}_{[n]:i}^\top + \{\nabla_{\mathbf{z}} x_i^{(n-1)} \mathbf{w}_{[n]:i}^\top\}^\top \\
&\quad + (\mathbf{w}_{[n]:i}^\top \mathbf{z} + 1) \nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n-1)})) \quad [\text{Eq. (14)}] \\
&= \mathcal{L}(\delta \cdot (\nabla_{\mathbf{z}} x_i^{(n-1)}) \mathbf{w}_{[n]:i}^\top) + \mathcal{L}(\mathbf{w}_{[n]:i} \delta \cdot (\nabla_{\mathbf{z}} x_i^{(n-1)})^\top) \\
&\quad + \mathcal{L}(\delta \cdot (\mathbf{w}_{[n]:i}^\top \mathbf{z} + 1) \nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n-1)}) \quad [\text{Associativity}] \\
&= \mathcal{L}(\delta \cdot (\nabla_{\mathbf{z}} x_i^{(n-1)})) \mathbf{w}_{[n]:i}^{+\top} + \mathcal{U}(\delta \cdot (\nabla_{\mathbf{z}} x_i^{(n-1)})) \mathbf{w}_{[n]:i}^{-\top} \\
&\quad + \mathbf{w}_{[n]:i}^+ \mathcal{L}(\delta \cdot (\nabla_{\mathbf{z}} x_i^{(n-1)})^\top) + \mathbf{w}_{[n]:i}^- \mathcal{U}(\delta \cdot (\nabla_{\mathbf{z}} x_i^{(n-1)})^\top) \\
&\quad + \mathcal{L}(\delta \cdot (\mathbf{w}_{[n]:i}^\top \mathbf{z} + 1) \nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n-1)}) \quad [\text{Linearity Eq. (9)}] \\
&= \mathcal{L}(\delta \cdot (\nabla_{\mathbf{z}} x_i^{(n-1)})) \mathbf{w}_{[n]:i}^{+\top} + \mathcal{U}(\delta \cdot (\nabla_{\mathbf{z}} x_i^{(n-1)})) \mathbf{w}_{[n]:i}^{-\top} \\
&\quad + \mathbf{w}_{[n]:i}^+ \mathcal{L}(\delta \cdot (\nabla_{\mathbf{z}} x_i^{(n-1)})^\top) + \mathbf{w}_{[n]:i}^- \mathcal{U}(\delta \cdot (\nabla_{\mathbf{z}} x_i^{(n-1)})^\top) \\
&\quad + \mathcal{L}(\delta' \nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n-1)}) \quad [\text{Definition of } \delta'] \\
&\quad (75)
\end{aligned}$$

Analogously, for the upper bound:

$$\begin{aligned}
\mathcal{U}(\delta \cdot \nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n)}) &= \mathcal{L}(\delta \cdot (\nabla_{\mathbf{z}} x_i^{(n-1)} \mathbf{w}_{[n]:i}^\top + \{\nabla_{\mathbf{z}} x_i^{(n-1)} \mathbf{w}_{[n]:i}^\top\}^\top \\
&\quad + (\mathbf{w}_{[n]:i}^\top \mathbf{z} + 1) \nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n-1)})) \quad [\text{Eq. (14)}] \\
&= \mathcal{U}(\delta \cdot (\nabla_{\mathbf{z}} x_i^{(n-1)}) \mathbf{w}_{[n]:i}^\top) + \mathcal{U}(\mathbf{w}_{[n]:i} \delta \cdot (\nabla_{\mathbf{z}} x_i^{(n-1)})^\top) \\
&\quad + \mathcal{U}(\delta \cdot (\mathbf{w}_{[n]:i}^\top \mathbf{z} + 1) \nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n-1)}) \quad [\text{Associativity}] \\
&= \mathcal{U}(\delta \cdot (\nabla_{\mathbf{z}} x_i^{(n-1)})) \mathbf{w}_{[n]:i}^{+\top} + \mathcal{L}(\delta \cdot (\nabla_{\mathbf{z}} x_i^{(n-1)})) \mathbf{w}_{[n]:i}^{-\top} \\
&\quad + \mathbf{w}_{[n]:i}^+ \mathcal{U}(\delta \cdot (\nabla_{\mathbf{z}} x_i^{(n-1)})^\top) + \mathbf{w}_{[n]:i}^- \mathcal{L}(\delta \cdot (\nabla_{\mathbf{z}} x_i^{(n-1)})^\top) \\
&\quad + \mathcal{U}(\delta \cdot (\mathbf{w}_{[n]:i}^\top \mathbf{z} + 1) \nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n-1)}) \quad [\text{Linearity Eq. (9)}] \\
&= \mathcal{U}(\delta \cdot (\nabla_{\mathbf{z}} x_i^{(n-1)})) \mathbf{w}_{[n]:i}^{+\top} + \mathcal{L}(\delta \cdot (\nabla_{\mathbf{z}} x_i^{(n-1)})) \mathbf{w}_{[n]:i}^{-\top} \\
&\quad + \mathbf{w}_{[n]:i}^+ \mathcal{U}(\delta \cdot (\nabla_{\mathbf{z}} x_i^{(n-1)})^\top) + \mathbf{w}_{[n]:i}^- \mathcal{L}(\delta \cdot (\nabla_{\mathbf{z}} x_i^{(n-1)})^\top) \\
&\quad + \mathcal{U}(\delta' \nabla_{\mathbf{z}\mathbf{z}}^2 x_i^{(n-1)}), \quad [\text{Definition of } \delta'] \\
&\quad (76)
\end{aligned}$$

where $\delta' = \delta \cdot (\mathbf{w}_{[n]:i}^\top \mathbf{z} + 1)$. \square

F.5 Lower bound of the Hessian's minimum eigenvalue for the product of polynomials case

To verify a product of polynomials network, we need a lower bound of the minimum eigenvalue of its Hessian. In Proposition 2 we propose a valid lower bound.

Proposition 2. *Let \mathbf{x} and \mathbf{y} be the input and output of a polynomial. Let*

$$\hat{\mathbf{J}}_{\mathbf{z}}(\mathbf{x}) = \arg \max \{ \rho(\mathbf{J} \mathbf{J}^\top) : \mathbf{J} \in [\mathcal{L}(\mathbf{J}_{\mathbf{z}}(\mathbf{x})), \mathcal{U}(\mathbf{J}_{\mathbf{z}}(\mathbf{x}))] \} = \max \{ |\mathcal{L}(\mathbf{J}_{\mathbf{z}}(\mathbf{x}))|, |\mathcal{U}(\mathbf{J}_{\mathbf{z}}(\mathbf{x}))| \} \quad (77)$$

be the Jacobian matrix with the largest possible norm. Let ρ be the spectral radius of a matrix. For all $\mathbf{z} \in [\mathbf{l}, \mathbf{u}]$, the minimum eigenvalue of the hessian matrix of every position i ($\lambda_{\min}(\nabla_{\mathbf{z}\mathbf{z}}^2 y_i)$) satisfies:

$$\lambda_{\min}(\nabla_{\mathbf{z}\mathbf{z}}^2 y_i) \geq \sum_{j=i}^k \lambda_{\min} \left(\frac{\partial y_i}{\partial x_j} \nabla_{\mathbf{z}\mathbf{z}}^2 x_j \right) - \rho \left(\hat{\mathbf{J}}_{\mathbf{z}}(\mathbf{x}) \hat{\mathbf{J}}_{\mathbf{z}}^\top(\mathbf{x}) \right) \cdot \rho \left(\mathbf{L}_{\mathbf{H}}(\nabla_{\mathbf{x}\mathbf{x}}^2 y_i) \right). \quad (78)$$

Proof. We give the lower bound of $\lambda_{\min}(\nabla_{\mathbf{z}\mathbf{z}}^2 y_i)$ as

$$\lambda_{\min}(\nabla_{\mathbf{z}\mathbf{z}}^2 y_i) = \lambda_{\min}\left(\sum_{j=i}^k \frac{\partial y_i}{\partial x_j} \nabla_{\mathbf{z}\mathbf{z}}^2 x_j + \mathbf{J}_{\mathbf{z}}^\top(\mathbf{x}) \nabla_{\mathbf{x}\mathbf{x}}^2 y_i \mathbf{J}_{\mathbf{z}}(\mathbf{x})\right) \quad [\text{Eq. (30)}]$$

$$\geq \lambda_{\min}\left(\sum_{j=i}^k \frac{\partial y_i}{\partial x_j} \nabla_{\mathbf{z}\mathbf{z}}^2 x_j\right) + \lambda_{\min}\left(\mathbf{J}_{\mathbf{z}}^\top(\mathbf{x}) \nabla_{\mathbf{x}\mathbf{x}}^2 y_i \mathbf{J}_{\mathbf{z}}(\mathbf{x})\right) \quad [\text{Weyl's inequality}]$$

$$\geq \sum_{j=i}^k \lambda_{\min}\left(\frac{\partial y_i}{\partial x_j} \nabla_{\mathbf{z}\mathbf{z}}^2 x_j\right) + \lambda_{\min}\left(\mathbf{J}_{\mathbf{z}}^\top(\mathbf{x}) \mathbf{L}_{\mathbf{H}}(\nabla_{\mathbf{x}\mathbf{x}}^2 y_i) \mathbf{J}_{\mathbf{z}}(\mathbf{x})\right) \quad [\text{Eq. (16)}]$$

$$\geq \sum_{j=i}^k \lambda_{\min}\left(\frac{\partial y_i}{\partial x_j} \nabla_{\mathbf{z}\mathbf{z}}^2 x_j\right) - \rho\left(\mathbf{J}_{\mathbf{z}}^\top(\mathbf{x}) \mathbf{L}_{\mathbf{H}}(\nabla_{\mathbf{x}\mathbf{x}}^2 y_i) \mathbf{J}_{\mathbf{z}}(\mathbf{x})\right) \quad [\text{Definition of } \rho]$$

$$\geq \sum_{j=i}^k \lambda_{\min}\left(\frac{\partial y_i}{\partial x_j} \nabla_{\mathbf{z}\mathbf{z}}^2 x_j\right) - \rho(\mathbf{J}_{\mathbf{z}}(\mathbf{x}) \mathbf{J}_{\mathbf{z}}^\top(\mathbf{x})) \rho(\mathbf{L}_{\mathbf{H}}(\nabla_{\mathbf{x}\mathbf{x}}^2 y_i))$$

$$\geq \sum_{j=i}^k \lambda_{\min}\left(\frac{\partial y_i}{\partial x_j} \nabla_{\mathbf{z}\mathbf{z}}^2 x_j\right) - \rho(\hat{\mathbf{J}}_{\mathbf{z}}(\mathbf{x}) \hat{\mathbf{J}}_{\mathbf{z}}^\top(\mathbf{x})) \rho(\mathbf{L}_{\mathbf{H}}(\nabla_{\mathbf{x}\mathbf{x}}^2 y_i)), \quad [\text{Eq. (77)}]$$

(79)

where in the second to last inequality we use $\rho(\mathbf{B}^\top \mathbf{A} \mathbf{B}) = \rho(\mathbf{B} \mathbf{B}^\top \mathbf{A}) \leq \rho(\mathbf{B} \mathbf{B}^\top) \rho(\mathbf{A}), \forall \mathbf{A} \in \mathbb{R}^{d_1 \times d_1}, \mathbf{B} \in \mathbb{R}^{d_1 \times d_2}$. \square