

# Tropical Geometry and Machine Learning

*This article deals with tropical geometry that has recently emerged as a tool in the analysis and extension of several classes of problems in both classical machine learning and deep learning.*

By PETROS MARAGOS<sup>id</sup>, Fellow IEEE, VASILEIOS CHARISOPOULOS<sup>id</sup>, AND EMMANOUIL THEODOSIS<sup>id</sup>, Member IEEE

**ABSTRACT** | Tropical geometry is a relatively recent field in mathematics and computer science, combining elements of algebraic geometry and polyhedral geometry. The scalar arithmetic of its analytic part preexisted in the form of max-plus and min-plus semiring arithmetic used in finite automata, nonlinear image processing, convex analysis, nonlinear control, optimization, and idempotent mathematics. Tropical geometry recently emerged in the analysis and extension of several classes of problems and systems in both classical machine learning and deep learning. Three such areas include: 1) deep neural networks with piecewise linear (PWL) activation functions; 2) probabilistic graphical models; and 3) nonlinear regression with PWL functions. In this article, we first summarize introductory ideas and objects of tropical geometry, providing a theoretical framework for both the max-plus algebra that underlies tropical geometry and its extensions to general max algebras. This unifies scalar and vector/signal operations over a class of nonlinear spaces, called weighted lattices, and allows us to provide optimal solutions for algebraic equations used in tropical geometry and generalize tropical geometric objects. Then, we survey the state of the art and recent progress in the aforementioned areas. First, we illustrate a purely geometric approach for studying the representation power of neural networks with PWL activations.

Then, we review the tropical geometric analysis of parametric statistical models, such as HMMs; later, we focus on the Viterbi algorithm and related methods for weighted finite-state transducers and provide compact and elegant representations via their formal tropical modeling. Finally, we provide optimal solutions and an efficient algorithm for the convex regression problem, using concepts and tools from tropical geometry and max-plus algebra. Throughout this article, we also outline problems and future directions in machine learning that can benefit from the tropical-geometric point of view.

**KEYWORDS** | Graphs; lattices; max-plus algebra; neural networks; regression; tropical geometry.

## I. INTRODUCTION

Tropical geometry is a relatively recent field in mathematics and computer science that combines elements from algebraic geometry and polyhedral geometry. The scalar arithmetic of its analytic part preexisted in the form of max-plus and min-plus semiring arithmetic used in finite automata, nonlinear image processing, convex analysis, nonlinear control, optimization, and idempotent mathematics. In max-plus arithmetic, the real number addition and multiplication are replaced by the max and sum operations, respectively. The name “tropical semiring” initially referred to the min-plus semiring and was used in finite automata [57], [99], speech recognition using graphical models [82], and tropical geometry [68], [80]. However, nowadays, the term, tropical semiring, may refer to both the max-plus and its dual min-plus arithmetic, whose combinations with corresponding nonlinear matrix algebra and nonlinear signal convolutions have been used in operations research and scheduling [25]; discrete event systems, max-plus control, and optimization [1], [2], [6], [15], [22], [37], [39], [48], [78], [110]; convex analysis [65], [85], [94]; morphological image analysis [49], [73], [79], [95], [96]; nonlinear

Manuscript received August 6, 2020; revised December 2, 2020; accepted January 18, 2021. Date of publication April 2, 2021; date of current version April 30, 2021. The work of Petros Maragos was supported by the European Regional Development Fund of the European Union and Greek National Funds through the Operational Program Competitiveness, Entrepreneurship and Innovation, under the call RESEARCH—CREATE—INNOVATE, under Project T1EDK02890 (e-Prevention). (Corresponding author: Petros Maragos.)

**Petros Maragos** is with the School of Electrical and Computer Engineering, National Technical University of Athens, 15773 Athens, Greece (e-mail: maragos@cs.ntua.gr).

**Vasileios Charisopoulos** is with the Department of Operations Research and Information Engineering, Cornell University, Ithaca, NY 14850 USA (e-mail: vc333@cornell.edu).

**Emmanouil Theodosis** is with the Department of Computer Science, Harvard University, Cambridge, MA 02138 USA (e-mail: etheodosis@g.harvard.edu).

Digital Object Identifier 10.1109/JPROC.2021.3065238

This work is licensed under a Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 License. For more information, see <https://creativecommons.org/licenses/by-nc-nd/4.0/>

difference equations for distance transforms [11], [71]; nonlinear PDEs of the Hamilton–Jacobi type for vision scale spaces [14], [50]; speech recognition and natural language processing [56], [82]; neural networks [18], [19], [34], [40], [83], [89], [93], [103], [114], [115]; and idempotent mathematics (nonlinear functional analysis) [63], [64].

The goal of this article is threefold: 1) to provide a brief background from tropical geometry and its underlying max-plus algebra; 2) to summarize its applications in three areas of machine learning (neural networks, graphical models, and nonlinear regression); and 3) to provide recent progress and some extensions using a generalized max algebra. Parts 1) and 2) provide tutorial information and survey state-of-the-art results. Some recent progress from the authors is included in parts 2) and 3).

We begin in Section II with elementary ideas and objects of tropical geometry. Section III provides the required theoretical background on max-plus algebra, its underlying nonlinear vector spaces called weighted lattices, and monotone operators in the form of lattice duality pairs called adjunctions (also known as residuation pairs). This section also provides some tools from a generalized max- $\star$  algebra to extend tropical geometrical objects. Furthermore, in Section IV, we show that adjunction pairs lead to optimal solutions of max-plus and general max- $\star$  equations, as nonlinear projections on weighted lattices. Then, the concepts and tools of the previous sections are applied to analyzing and/or providing solutions for problems in the following three broad areas of machine learning.

### A. Neural Networks With Piecewise Linear (PWL) Activations (See Section V)

Tropical geometry recently emerged in the study of deep neural networks (DNNs) and variations of the perceptron operating in the max-plus semiring. Standard activation functions employed in DNNs, including the ReLU activation and its “leaky” variants, induce neural network layers that are PWL convex functions of their inputs and create a partition of space well described by concepts from tropical geometry. Following [18] and [19], we illustrate a purely geometric approach for studying the representation power of DNNs—measured via the concept of a network’s “linear regions”—under the lens of tropical geometry.

### B. Probabilistic Graphical Models and Algorithms (See Section VI)

As we review in Section VI-A, a novel application of tropical geometry is its usage in [86] for analyzing parametric statistical models, including hidden Markov models (HMMs) and restricted Boltzmann machines (RBMs). Furthermore, among the max-sum and max-product algorithms used in graphical models, a prime representative is the Viterbi algorithm. This can also be viewed in the general setting of weighted finite-state transducers (WFSTs) [56], [82] which have found extensive use in speech recognition and other decoding schemes. Practical reasons led researchers to adopt a tropical version of these

algorithms in order to resolve numerical issues that arose from using sum-product algebras. However, as we explain in Section VI-B, tropicalization is not restricted merely as a numerical tool; further tropical modeling of the algorithms as in [106] and [107] leads to a compact and elegant representation, while highlighting geometric properties.

### C. Piecewise Linear Regression (See Section VII)

Fitting PWL functions to data is a fundamental regression problem in multidimensional signal modeling and machine learning since approximations with PWL functions have proved analytically and computationally very useful in many fields of science and engineering. We focus on functions that admit a convex representation as the maximum of affine functions (e.g., lines and planes), represented with max-plus tropical polynomials. This allows us to use concepts and tools from tropical geometry and max-plus algebra to optimally approximate the shape of curves and surfaces by fitting tropical polynomials to data, possibly in the presence of noise; this yields polygonal or polyhedral shape approximations. For this convex PWL regression problem, we provide optimal solutions with respect to  $\ell_p$  error norms, derived using monotone operator adjunctions that are projections on weighted lattices, and an efficient algorithm based on preliminary work in [76].

Finally, in Section VIII, extending preliminary work in [75], we generalize tropical geometry using the max- $\star$  algebra and weighted lattices framework of [74], as summarized in Section III-B, with an arbitrary binary operation  $\star$  that distributes over max, and apply it to optimal convex PWL regression for fitting max- $\star$  tropical curves and surfaces to arbitrary data.

## II. ELEMENTS OF TROPICAL GEOMETRY

After some notation and definitions from tropical and related semirings, we first present some simple examples of tropical<sup>1</sup> curves and surfaces that result from tropicalizing the polynomials that analytically describe their Euclidean counterparts. Then, we explain this tropicalization as a dequantization of real algebraic geometry. Finally, Newton polytopes and tropical halfspaces are defined with examples.

*Notation:* For maximum (or supremum) and minimum (or infimum) operations, we use the well-established lattice-theoretic<sup>2</sup> symbols of  $\vee$  and  $\wedge$ . We use roman letters

<sup>1</sup>The adjective “tropical” was coined by French mathematicians, including Dominique Perrin and Jean-Eric Pin, to honor their Brazilian colleague Imre Simon who was one of the pioneers of min-plus algebra as applied to automata. However, we give it an alternative and substantial meaning in connection with its Greek origin word τροπικός that comes from the Greek word τροπή, meaning “turn” or “changing the way/direction,” to literally express the fact that tropical curves and surfaces bend and turn.

<sup>2</sup>We do *not* use the notation  $(\oplus, \otimes)$  for  $(\max, +)$  or  $(\min, +)$ , which is frequently used in max-plus algebra, because, in functional analysis and image processing: 1) the symbol  $\oplus$  is extensively used for the Minkowski set addition and max-plus signal convolution and 2)  $\otimes$  is unnecessarily confusing compared to the classic symbol  $+$  of addition.

for functions, signals, and their arguments, and Greek letters mainly for operators. Also, we use boldface roman letters for vectors (lower case) and matrices (capital). If  $M = [m_{ij}]$  is a matrix, its  $(i, j)$ th element is denoted as  $m_{ij}$  or  $[M]_{ij}$ . Similarly,  $x = [x_i]$  denotes a column vector, whose  $i$ th element is denoted as  $[x]_i$  or simply  $x_i$ . We also use the set notation  $[n] := \{1, \dots, n\}$ .

## A. Tropical Semirings

Compared with the classical real number ring  $(\mathbb{R}, +, \times)$ , the *max-plus semiring*  $(\mathbb{R}_{\max}, \vee, +)$  consists of the set  $\mathbb{R}_{\max} = \mathbb{R} \cup \{-\infty\}$  equipped with an idempotent “addition” that is the maximum operation and a generalized “multiplication” that is the extended real addition. Similarly, we consider the dual *min-plus semiring*  $(\mathbb{R}_{\min}, \wedge, +)$ , where  $\mathbb{R}_{\min} = \mathbb{R} \cup \{+\infty\}$ . Both tropical semirings are special cases of *dioids* [39]. From a different viewpoint that we follow in this article, if we combine both the maximum and minimum operations, we obtain the complete lattice  $(\overline{\mathbb{R}}, \vee, \wedge)$  of extended real numbers  $\overline{\mathbb{R}} = \mathbb{R} \cup \{-\infty, +\infty\}$ . Furthermore, as done more generally in Section III-B, we can combine the max-plus and min-plus scalar arithmetic into an algebraic structure called complete lattice-ordered double monoid (*clodum*), which consists of the extended reals  $\overline{\mathbb{R}}$  equipped with the maximum ( $\vee$ ), minimum ( $\wedge$ ), addition ( $+$ ), and dual addition ( $+$ ) operations. The operations  $+$  and  $+$  are, respectively, the “lower addition” and “upper addition” used in convex analysis [85]. They are identical for finite reals and differ only when combining  $-\infty$  with  $+\infty$ ; in all cases, they are commutative:

$$\begin{aligned} a + b &= a +' b \quad \forall a \in \overline{\mathbb{R}} \quad \forall b \in \mathbb{R} \\ a + (-\infty) &= -\infty, \quad a +' (+\infty) = +\infty \quad \forall a \in \overline{\mathbb{R}}. \end{aligned} \quad (1)$$

In idempotent mathematics [64], convex optimization [13], and the theory of dioids [39], the following *Log-Sum-Exp* approximation is often used for the max and min operations:

$$\begin{aligned} a \vee_{\theta} b &:= \theta \cdot \log \left( e^{a/\theta} + e^{b/\theta} \right) = \phi_{\theta}^{-1} [\phi_{\theta}(a) + \phi_{\theta}(b)] \\ a \wedge_{\theta} b &:= (-\theta) \log \left( e^{-a/\theta} + e^{-b/\theta} \right) \end{aligned} \quad (2)$$

where  $\phi_{\theta}(a) := \exp(a/\theta)$ , and  $\theta > 0$  is usually called a “temperature” parameter. In the limit as  $\theta \rightarrow 0$ , we obtain the max and min operations

$$\begin{aligned} \lim_{\theta \downarrow 0} a \vee_{\theta} b &= \max(a, b) \\ \lim_{\theta \downarrow 0} a \wedge_{\theta} b &= \min(a, b). \end{aligned} \quad (3)$$

This approximation and limit is the *Maslov dequantization* [77] of real numbers and generates a whole family of semirings  $S_{\theta} = (\mathbb{R}_{\max}, \vee_{\theta}, +)$  and  $\theta > 0$ , whose operations are the generalized “addition”  $\vee_{\theta}$  and “multiplication”  $+$ .

Each of the semirings  $S_{\theta}$  is isomorphic to the semiring  $(\mathbb{R}_{\geq 0}, +, \times)$  of nonnegative real numbers  $\mathbb{R}_{\geq 0}$  equipped with standard addition and multiplication. This isomorphism is enabled via the bijection  $a \mapsto \phi_{\theta}(a)$  from  $\mathbb{R}_{\max}$  onto  $\mathbb{R}_{\geq 0}$ . To see this, let  $x = \phi_{\theta}(a) = \exp(a/\theta)$  and  $y = \phi_{\theta}(b) = \exp(b/\theta)$ . Then, for any  $a, b \in \mathbb{R}_{\max}$ ,

$$\phi_{\theta}(a \vee_{\theta} b) = x + y, \quad \phi_{\theta}(a + b) = x \cdot y.$$

## B. Examples of Tropical Polynomial Curves and Surfaces

1) *Tropical Polynomial Curves*: Consider the analytic expressions for a Euclidean line and parabola

$$p_1(x) = ax + b, \quad p_2(x) = ax^2 + bx + c. \quad (4)$$

“Tropicalization,” that is, replacing sum with max and multiplication with addition, yields the corresponding max-plus tropical polynomials

$$\begin{aligned} p_1^{\max}(x) &= \max(a + x, b) \\ p_2^{\max}(x) &= \max(a + 2x, b + x, c). \end{aligned} \quad (5)$$

The equations for the min-plus case are identical as in (5) by replacing max with min. The graphs of all the above can be seen in Fig. 1.

2) *Tropical Polynomial Surfaces*: Consider the equations of the following tropical planes represented as 2-D max-plus and min-plus polynomials of degree 1:

$$\begin{aligned} f(x, y) &= \max(x, 2 + y, 7) \\ g(x, y) &= \min(5 + x, 7 + y, 9) \end{aligned} \quad (6)$$

whose graphs can be seen as surfaces in Fig. 2(a) and (b).

Next, to the general Euclidean conic polynomial

$$p_{\text{conic}}(x, y) = ax^2 + bxy + cy^2 + dx + ey + f \quad (7)$$

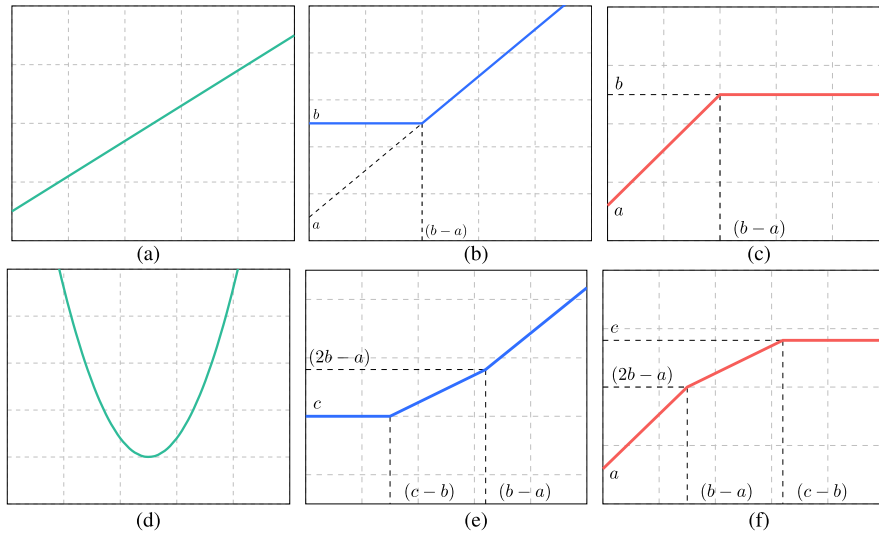
there corresponds the following two-variable max-plus tropical polynomial of degree 2:

$$p_{\text{conic}}^{\max}(x, y) = \max(a + 2x, b + x + y, c + 2y, d + x, e + y, f). \quad (8)$$

Its min-plus version is shown in Fig. 2(c).

## C. Tropicalization via Dequantization of Algebraic Geometry

The algebraic side of tropical geometry [68] results from a transformation of analytic Euclidean geometry where the traditional arithmetic of the real field  $(\mathbb{R}, +, \times)$  involved in



**Fig. 1.** Euclidean and tropical 1-D polynomial curves of first and second degrees. (a) Euclidean line. (b) Max-plus line. (c) Min-plus line. (d) Euclidean parabola. (e) Max-plus parabola. (f) Min-plus parabola.

the analytic expressions of geometric objects is replaced by the arithmetic of the max-plus or min-plus semiring. A geometric explanation and visualization of this transformation is obtained from Viro’s graphing of polynomial curves on log–log paper [111]. Consider the monomial curve  $v = cu^a$ ,  $c > 0$ , on the positive quadrant of the  $(u, v)$ -plane and consider the log–log transformation of both coordinates composed with a uniform scaling by  $\theta > 0$ :  $x = \theta \log u$  and  $y = \theta \log v$ . Then, on the  $(x, y)$ -plane, the curve becomes the line  $y = b/\theta + ax$ , where  $b = \log c$ . If we have a  $K$ -term polynomial curve  $v = P(u) = \sum_{k=1}^K c_k u^{a_k}$  with  $c_k = \exp(b_k) > 0$  and  $a_k \in \mathbb{R}$  (i.e., a posynomial [12]), then we convert it to

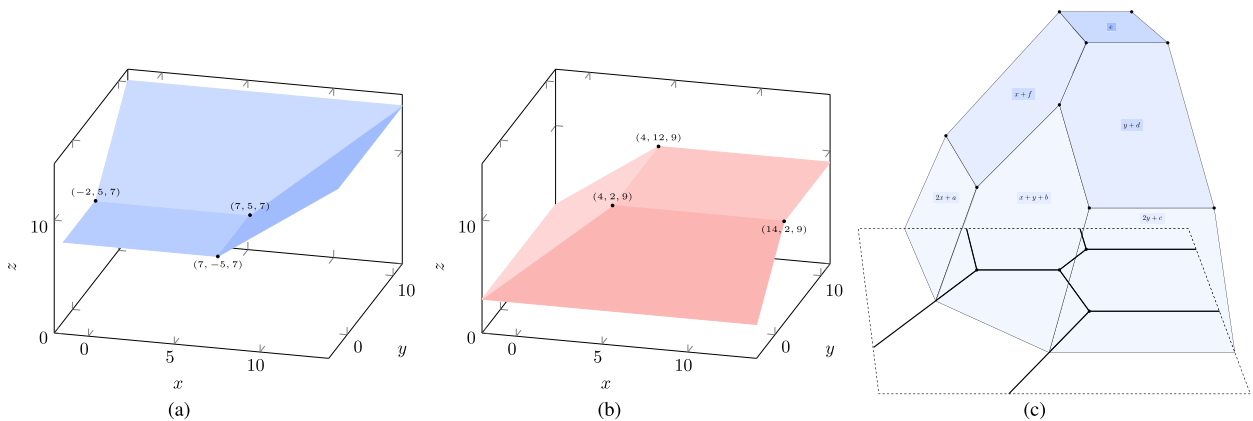
$$P_\theta(x) = \theta \log \left( \sum_{k=1}^K \exp(b_k/\theta) \exp(a_k x/\theta) \right). \quad (9)$$

As  $\theta \downarrow 0$ , this yields, via the Maslov dequantization, a  $K$ -term **1-D max-plus tropical polynomial**

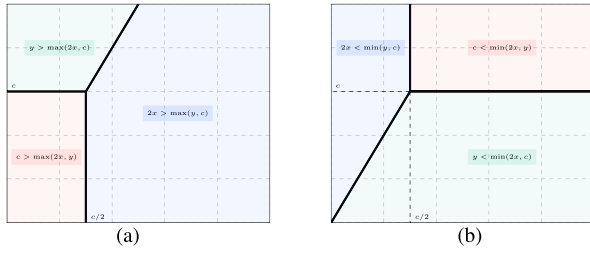
$$\lim_{\theta \downarrow 0} P_\theta(x) = p(x) = \max_{k=1}^K \{b_k + a_k x\}. \quad (10)$$

While each  $P_\theta(x)$  is a smooth function, their limit  $p(x)$  is a max-affine function and represents a PWL convex function. If we perform dequantization with negative exponents, we obtain a min-plus polynomial that is a PWL concave function.

The above procedure extends to multiple dimensions or higher degrees and shows us the way to tropicalize any classical  $d$ -variable polynomial (linear combination of power monomials)  $\sum_k c_k u_1^{a_{k1}}, \dots, u_d^{a_{kd}}$  defined over  $\mathbb{R}_{>0}^d$ , where  $c_k > 0$  and  $\mathbf{a}_k = [a_{k1}, \dots, a_{kd}]^T$  is traditionally



**Fig. 2.** (a) and (b) Surfaces (graphs) of the two tropical planes defined in (6). (a) Max-plus plane  $f$ . (b) Min-plus plane  $g$ . (c) Surface of the 2-D min-plus tropical polynomial function (tropic conic)  $p(x, y) = \min(a + 2x, b + x + y, c + 2y, d + x, e + y, f)$  and its tropical quadratic curve. (c) is inspired by [68, Fig. 1.3.2].



**Fig. 3.** Tropical curve of the max-polynomial  $p(x, y) = \max(2x, y, c)$  left and its dual min-polynomial  $p'(x, y) = \min(2x, y, c)$  right. (a) Max-plus curve. (b) Min-plus curve.

some nonnegative integer<sup>3</sup> vector, but, herein, we allow  $\mathbf{a}_k \in \mathbb{R}^d$ : replace the sum with max and log the individual monomials. Thus, a general  $d$ -variable max-plus polynomial  $p: \mathbb{R}^d \rightarrow \mathbb{R}$  has the expression

$$p(\mathbf{x}) = \bigvee_{k=1}^K \mathbf{a}_k^T \mathbf{x} + b_k, \quad \mathbf{x} = [x_1, \dots, x_d]^T. \quad (11)$$

We define the **rank** of a tropical polynomial  $p$  as the number of affine terms involved in the maximum; here,  $K = \text{rank}(p)$ . Its graph is a max of  $K$  hyperplanes with intercepts  $b_k = \log c_k \in \mathbb{R}$  and real slope vectors  $\mathbf{a}_k \in \mathbb{R}^d$ . The degree of  $p$  is  $|\mathbf{a}| = \max_k \|\mathbf{a}_k\|_1$ , where  $\|\mathbf{a}_k\|_1 = |a_{k1}| + \dots + |a_{kd}|$ . Thus, the curves or surfaces of real algebraic geometry become via dequantization the graphs of convex PWL functions represented by tropical polynomials.

#### D. Tropical Curves and Newton Polytopes

To the zero set of a classical polynomial, there corresponds the *tropical curve* or *hypersurface* of a max-plus tropical polynomial  $p: \mathbb{R}^d \rightarrow \mathbb{R}$ :

$$\mathcal{V}(p) := \left\{ \mathbf{x} \in \mathbb{R}^d : \begin{array}{l} \text{more than one terms of } p(\mathbf{x}) \\ \text{attain the max} \end{array} \right\}. \quad (12)$$

The above also defines the tropical curve of min-plus polynomials by replacing max with min. Thus,  $\mathcal{V}(p)$  consists of the singularity points (of nondifferentiability) of  $p(\mathbf{x})$ . Examples are shown in Fig. 3 for degree-1 tropical polynomials and in Fig. 2(c) for a degree-2 polynomial.

Another interesting geometric object related to a max-plus polynomial  $p$  is its *Newton polytope*, which is the convex hull [denoted by  $\text{conv}(\cdot)$ ] of the set of points corresponding to its slope coefficient vectors

$$\text{Newt}(p) := \text{conv}(\{\mathbf{a}_k : k = 1, \dots, \text{rank}(p)\}). \quad (13)$$

<sup>3</sup>Traditionally, “tropical polynomials” assume that the parameters  $a_{ki}$  are nonnegative integers. If we also allow negative integers, we get “Laurent tropical polynomials.” As in [15], we allow any real coefficients; this may be called “tropical posynomials” [16].

This satisfies several important properties [18]:

$$\text{Newt}(p_1 \vee p_2) = \text{conv}(\text{Newt}(p_1) \cup \text{Newt}(p_2)) \quad (14)$$

$$\text{Newt}(p_1 + p_2) = \text{Newt}(p_1) \oplus \text{Newt}(p_2) \quad (15)$$

where  $\oplus$  denotes Minkowski set addition, as defined in (21). Examples are shown in Fig. 4. Thus, the Newton polytope of the sum (resp. max) of two tropical polynomials is the Minkowski sum (resp. the convex hull of the union) of their individual polytopes.

#### E. Tropical Halfspaces and Polytopes

In pattern analysis problems on Euclidean spaces  $\mathbb{R}^d$ , we often use halfspaces  $\mathcal{H}(\mathbf{a}, b) := \{\mathbf{x} \in \mathbb{R}^d : \mathbf{a}^T \mathbf{x} \leq b\}$ , polyhedra (finite intersections of halfspaces), and polytopes (compact polyhedra formed as the convex hull of a finite set of points). Replacing linear inner products  $\mathbf{a}^T \mathbf{x}$  with max-plus versions yields *tropical halfspaces* [36], which are defined as the following subsets of  $\mathbb{R}_{\max}^d$  with parameters  $\mathbf{a} = [a_i], \mathbf{b} = [b_i] \in \mathbb{R}_{\max}^{d+1}$ :

$$\mathcal{T}(\mathbf{a}, \mathbf{b}) := \left\{ \mathbf{x} \in \mathbb{R}_{\max}^d : \begin{array}{l} \max\{a_1 + x_1, \dots, a_d + x_d, a_{d+1}\} \leq \\ \max\{b_1 + x_1, \dots, b_d + x_d, b_{d+1}\} \end{array} \right\} \quad (16)$$

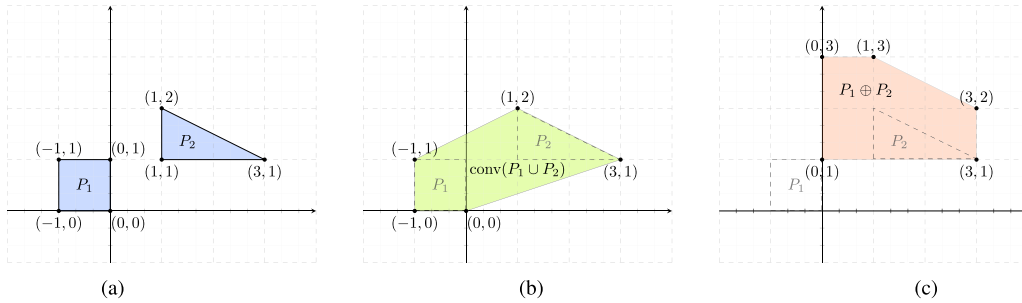
where  $\min(a_i, b_i) = -\infty$  for all  $i$ . Thus, for each  $i$ , only one coefficient is needed<sup>4</sup> either in the left or in the right side of inequality (16). Replacing max with min in (16) yields tropical halfspaces that are min-plus hyperplanes. Examples of tropical halfplanes are shown in Fig. 5, forming a planar polytope. It is obvious that their separating boundaries are tropical lines. Such regions in multiple dimensions were used in [18], [19], and [113] as morphological perceptrons.

As an example in the 3-D space, in Fig. 6 we can see the intersection of the tropical halfspaces corresponding to the two tropical polynomials in (6). This polytope is the polyhedral region formed by intersecting the half-space above the surface of the 2-D max-plus polynomial  $f$  with the half-space below the surface of the min-plus polynomial  $g$ .

We note from Figs. 5 and 6 that the number of tropical boundaries required to form polytopes, which could serve as decision regions in pattern classification problems, is smaller than the number of linear boundaries. See, for instance, the polytope  $R_P$  in Fig. 5(b). This observation remains valid in higher dimensions too, namely decision regions can be formed with fewer tropical lines or

<sup>4</sup>The general expression (16) of a tropical half-space includes as special cases expressions  $\{\mathbf{x} : \bigvee_i a_i + x_i \leq b\}$  which seem as a direct tropical analog of the expression  $\{\mathbf{x} : \sum_i a_i x_i \leq b\}$  for Euclidean halfspaces. For example, it is shown in [36] that  $\{\mathbf{x} : \max(a + x, c) \leq \max(b + x, d)\} = \{\mathbf{x} : \max(a + x, c) \leq d\}$  if  $a > b$ .





**Fig. 4.** Newton polytopes of (a) two max-polynomials  $p_1(x, y) = \max(x + y, 3x + y, x + 2y)$  and  $p_2(x, y) = \max(0, -x, y - x)$  (polytopes), (b) their max  $p_1 \vee p_2$  [Newton (max)], and (c) their sum  $p_1 + p_2$  [Newton (sum)].

hyperplanes than their Euclidean counterparts. Intuitively, the nonlinearity of a tropical half-space lets us to form more complex decision regions with possibly fewer parameters.

### III. ELEMENTS OF MAX-PLUS ALGEBRA, WEIGHTED LATTICES, AND MONOTONE OPERATORS

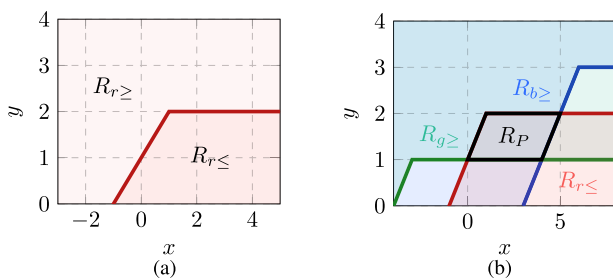
#### A. Lattices and Monotone Operators

Signals and vectors can be viewed as elements of complete lattices  $(\mathcal{L}, \vee, \wedge)$ , where  $\mathcal{L}$  is the set of lattice elements equipped with two binary operations,  $\vee$  and  $\wedge$ , which denote the lattice supremum and infimum, respectively. Each of these operations induces a partial ordering  $\leq$ , for example, for any  $X, Y \in \mathcal{L}$ ,  $X \leq Y \iff Y = X \vee Y$ . The lattice operations satisfy many properties, including associativity, commutativity, idempotence, and compatibility with the partial ordering. Completeness means that the supremum and infimum of any (even infinite) subset of  $\mathcal{L}$  exists and belongs to  $\mathcal{L}$ . Examples of complete lattices used in image processing include: 1) the lattice of Euclidean shapes, that is, subsets of  $\mathbb{R}^d$ , equipped with the set union and intersection, and 2) the lattice of functions  $f : E \rightarrow \mathbb{R}$  with (arbitrary) domain  $E$  and values in  $\mathbb{R}$ , equipped

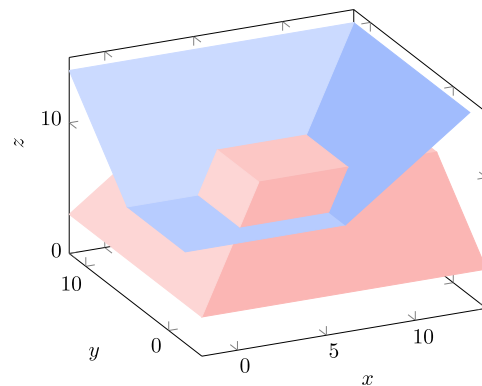
with the pointwise supremum and pointwise infimum of extended real-valued functions.

1) *Monotone Operators*: For data processing, we also consider operators  $\psi : \mathcal{L} \rightarrow \mathcal{M}$  between two complete lattices. A lattice operator  $\psi$  is called *increasing* if it is order preserving, that is, if, for any  $X, Y \in \mathcal{L}$ ,  $X \leq Y \implies \psi(X) \leq \psi(Y)$ . Given two operators  $\psi$  and  $\phi$ , we will write  $\psi \leq \phi \iff \psi(X) \leq \phi(X) \forall X$ . Examples of increasing operators are the lattice homomorphisms that preserve suprema and infima. If a lattice homomorphism is also a bijection, then it becomes an automorphism. Four fundamental types of increasing operators are: *dilations*  $\delta$  and *erosions*  $\varepsilon$  that satisfy, respectively,  $\delta(\bigvee_i X_i) = \bigvee_i \delta(X_i)$  and  $\varepsilon(\bigwedge_i X_i) = \bigwedge_i \varepsilon(X_i)$  over arbitrary (possibly infinite) collections; *openings*  $\alpha$  that are increasing, idempotent ( $\alpha^2 = \alpha$ ), and antiextensive ( $\alpha \leq \text{id}$ ), where  $\text{id}$  denotes the identity operator; and *closings*  $\beta$  that are increasing, idempotent, and extensive ( $\beta \geq \text{id}$ ).

A lattice operator  $\psi$  is called *decreasing* if it is order-inverting, that is,  $X \leq Y \implies \psi(X) \geq \psi(Y)$ . Dual homomorphisms interchange suprema with infima and, hence, are decreasing operators. For example, *antidilations*, denoted as  $\delta^a$ , satisfy  $\delta^a(\bigvee_i X_i) = \bigwedge_i \delta^a(X_i)$ . A lattice dual automorphism is a bijection that interchanges suprema



**Fig. 5.** Regions  $R_{c \geq}$  and  $R_{c \leq}$  formed by min-plus tropical halfspaces in  $\mathbb{R}^2$ , where  $c$  denotes the color of the tropical boundary and  $\geq$  (resp.  $\leq$ ) the set of points above (resp. below) the boundary. (a) Red boundary is the min-plus tropical line  $y = \min(1 + x, 2)$  (single region). (b) Green and blue boundaries are, respectively, the tropical lines  $y = \min(4 + x, 1)$  and  $y = \min(x - 3, 3)$  (multiple regions).  $R_P$  is the polytope formed by the intersection of three tropical halfplanes.



**Fig. 6.** Intersection of halfspaces of the 2-D max-plus and min-plus tropical polynomials in (6).

with infima. For example, a *negation*  $\nu$  is a dual automorphism that is also involutive, that is,  $\nu^2 = \mathbf{id}$ .

2) *Residuation and Adjunctions*: An increasing operator  $\psi : \mathcal{L} \rightarrow \mathcal{M}$  between two complete lattices is called *residuated* [8], [9] if there exists an increasing operator  $\psi^\# : \mathcal{M} \rightarrow \mathcal{L}$  such that

$$\psi\psi^\# \leq \mathbf{id} \leq \psi^\#\psi. \quad (17)$$

Here,  $\psi^\#$  is called the **residual** of  $\psi$ , is unique, and is the closest to being an inverse of  $\psi$ . Specifically, the residuation pair  $(\psi, \psi^\#)$  can solve inverse problems of the type  $\psi(X) = Y$  either exactly since  $\hat{X} = \psi^\#(Y)$  is the greatest solution of  $\psi(X) = Y$  if a solution exists, or approximately since  $\hat{X}$  is the *greatest subsolution* in the sense that

$$\hat{X} = \psi^\#(Y) = \bigvee \{X : \psi(X) \leq Y\}. \quad (18)$$

On complete lattices, an increasing operator  $\psi$  is residuated (resp. a residual  $\psi^\#$ ) if and only if it is a dilation (resp. erosion). The residuation theory has been used for solving inverse problems (mainly in matrix algebra) over the extended max-plus semiring  $(\overline{\mathbb{R}}, \vee, +)$  or other idempotent semirings as lattices are made complete [6], [23], [25], [26].

A pair  $(\delta, \varepsilon)$  of two operators  $\delta : \mathcal{L} \rightarrow \mathcal{M}$  and  $\varepsilon : \mathcal{M} \rightarrow \mathcal{L}$  between two complete lattices is called **adjunction** if

$$\delta(X) \leq Y \iff X \leq \varepsilon(Y) \quad \forall X \in \mathcal{L}, Y \in \mathcal{M}. \quad (19)$$

In any adjunction,  $\delta$  is a dilation and  $\varepsilon$  is an erosion. The double inequality (19) is equivalent to the inequality (17) satisfied by a residuation pair of increasing operators if we identify the residuated map  $\psi$  with  $\delta$  and its residual  $\psi^\#$  with  $\varepsilon$ . Furthermore, from (19) or (17), it follows that any adjunction  $(\delta, \varepsilon)$  automatically yields an opening  $\alpha = \delta\varepsilon$  and a closing  $\beta = \varepsilon\delta$ , where the composition of two operators is written as an operator product. Viewing  $(\delta, \varepsilon)$  as an adjunction instead of a residuation pair has the advantage of the additional geometrical intuition and visualization afforded by the dilation and erosion operators in image and shape analysis.

Given a dilation  $\delta$ , there is a unique erosion

$$\varepsilon(Y) = \delta^\#(Y) = \bigvee \{X \in \mathcal{L} : \delta(X) \leq Y\} \quad (20)$$

such that  $(\delta, \varepsilon)$  is an adjunction and conversely. Thus, dilations and erosions on complete lattices always come in pairs. In any adjunction  $(\delta, \varepsilon)$ ,  $\varepsilon$  is called the *adjoint erosion* of  $\delta$ , whereas  $\delta$  is the *adjoint dilation* of  $\varepsilon$ .

*Example 1*: Two adjunctions used in nonlinear image processing and shape analysis are the following:

- 1) A morphological set adjunction is the pair of Minkowski set addition  $\oplus$  and subtraction  $\ominus$ : for  $X, B \subseteq \mathbb{R}^d$ ,

$$\begin{aligned} \delta_B(X) = X \oplus B &:= \{x + b \in \mathbb{R}^d : x \in X, b \in B\} \\ \varepsilon_B(X) = X \ominus B &:= \{x - b \in \mathbb{R}^d : x \in X, b \in B\}. \end{aligned} \quad (21)$$

- 2) A signal adjunction is the supremal (max-plus) convolution  $f \oplus g$  of  $f : \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$  by  $g$  and the infimal convolution  $f \ominus g$  of  $f(x)$  by  $-g(-x)$  used in morphological image processing:

$$\begin{aligned} \delta_g(f)(x) = f \oplus g(x) &:= \sup_y \{f(y - x) + g(y)\} \\ \varepsilon_g(f)(x) = f \ominus g(x) &:= \inf_y \{f(x - y) - g(y)\}. \end{aligned} \quad (22)$$

## B. Max- $\star$ Algebra and Weighted Lattices

1) *Clodum—Extending Tropical Scalar Arithmetic*: A lattice  $\mathcal{M}$  is often endowed with an additional binary operation, called symbolically the “multiplication”  $\star$ , under which  $(\mathcal{M}, \star)$  is a semigroup, a monoid, or a group. Such ordered monoids have been studied in detail in [7], [39], and [116] and form the algebraic basis of the max-plus algebra.

Consider now an algebra  $(\mathcal{K}, \vee, \wedge, \star, \star')$  with four binary operations satisfying the following.

- C1:  $(\mathcal{K}, \vee, \wedge)$  is a complete distributive lattice. Thus, it contains its least  $\perp := \bigwedge \mathcal{K}$  and greatest element  $\top := \bigvee \mathcal{K}$ . The supremum  $\vee$  (resp. infimum  $\wedge$ ) plays the role of a generalized “addition” (resp. “dual addition”).
- C2:  $(\mathcal{K}, \star)$  is a monoid whose operation  $\star$  plays the role of a generalized “multiplication” with identity (“unit”) element  $e$  and is a dilation (i.e., distributes over  $\vee$ ).
- C3:  $(\mathcal{K}, \star')$  is a monoid with identity  $e'$  whose operation  $\star'$  plays the role of a generalized “dual multiplication” and is an erosion (i.e., distributes over  $\wedge$ ).

The least (greatest) element  $\perp$  ( $\top$ ) of  $\mathcal{K}$  is both the “zero” element for the “addition”  $\vee$  ( $\wedge$ ) and an absorbing null for the “multiplication”  $\star$  ( $\star'$ ).

We call the resulting algebra a *complete lattice-ordered double monoid*, shortly **clodum** [72], [74]. Previous studies on minimax or max-plus algebra have used alternative names for structures similar to the above definitions that emphasize semigroups and semirings instead of lattices [6], [25], [39] (see [74] for similarities and differences).

A clodum  $\mathcal{K}$  is called *self-conjugate* if it has a lattice negation  $a \mapsto a^*$  such that

$$\left( \bigvee_i a_i \right)^* = \bigwedge_i a_i^*, \quad \left( \bigwedge_i b_i \right)^* = \bigvee_i b_i^*, \quad (a \star b)^* = a^* \star' b^*. \quad (23)$$

Table 1 Examples of Scalar Clodums

Clodum	Set $\mathcal{K}$	'Add'	'Zero' ( $\perp$ )	'Dual Add'	'Dual Zero' ( $\top$ )	'Mult' ( $\star$ )	'Unit' ( $e$ )	'Dual Mult' ( $\star'$ )	'Dual Unit' ( $e'$ )	Conjugate ( $a^*$ )
Max-plus	$\mathbb{R} \cup \{\pm\infty\}$	$\vee$	$-\infty$	$\wedge$	$+\infty$	$+$	$0$	$+$	$0$	$-a$
Max-times	$[0, +\infty]$	$\vee$	$0$	$\wedge$	$+\infty$	$\times$	$1$	$\times'$	$1$	$a^{-1}$
Max-min	$[0, 1]$	$\vee$	$0$	$\wedge$	$1$	$\wedge$	$1$	$\vee$	$0$	$1 - a$
Max-softmin	$\mathbb{R} \cup \{\pm\infty\}$	$\vee$	$-\infty$	$\wedge$	$+\infty$	$\wedge_\theta$	$+\infty$	$\vee_\theta$	$-\infty$	$-a$

The suprema and infima in (23) may be over any collections.

Examples of scalar clodums are summarized in Table 1. The max-plus and max-times clodums have a richer structure. Specifically, if  $\star = \star'$  over  $G = \mathcal{K} \setminus \{\perp, \top\}$ , where  $(G, \star)$  is a group and  $(G, \vee, \wedge)$  is a conditionally complete lattice (i.e., every nonempty bounded subset has a supremum and infimum), then the clodum  $\mathcal{K}$  becomes a *complete lattice-ordered group*, shortly **clog**. Then, for each  $a \in G$ , there exists its “multiplicative inverse”  $a^{-1}$  such that  $a \star a^{-1} = e$ . Furthermore, the “multiplication”  $\star$  and its self-dual  $\star'$  can be extended over the whole  $\mathcal{K}$  by involving the nulls, and the clodum becomes self-conjugate by setting  $a^* = a^{-1}$  if  $\perp < a < \top$ ,  $\top^* = \perp$ , and  $\perp^* = \top$ . Thus, in a clog  $\mathcal{K}$ ,  $\star$  and  $\star'$  coincide in all cases with only one exception: the combination of the least and greatest elements.

All clodum examples of Table 1 have commutative “multiplications.” An example with *noncommutative* “multiplications” is the *matrix* max- $\star$  clodum  $(\mathcal{K}^{n \times n}, \vee, \wedge, \boxtimes, \boxtimes')$ , where  $\mathcal{K}^{n \times n}$  is the set of  $n \times n$  matrices with entries from a clodum  $\mathcal{K}$ ,  $\vee/\wedge$  denote here elementwise matrix sup/inf, and  $\boxtimes$  and  $\boxtimes'$  denote max- $\star$  and min- $\star'$  matrix “multiplications”

$$[A \boxtimes B]_{ij} = \bigvee_{k=1}^n a_{ik} \star b_{kj}, \quad [A \boxtimes' B]_{ij} = \bigwedge_{k=1}^n a_{ik} \star' b_{kj}. \quad (24)$$

For the max-plus clog  $(\overline{\mathbb{R}}, \vee, \wedge, +, +')$ , these matrix “multiplications” are denoted by  $\boxplus$  and  $\boxplus'$ , defined as

$$[A \boxplus B]_{ij} = \bigvee_{k=1}^n a_{ik} + b_{kj}, \quad [A \boxplus' B]_{ij} = \bigwedge_{k=1}^n a_{ik} + b_{kj}. \quad (25)$$

2) *Complete Weighted Lattices—Nonlinear Spaces*: Consider a nonempty collection  $\mathcal{W}$  of mathematical objects, which will be our space; examples of such objects include vectors in  $\overline{\mathbb{R}}^d$  or signals  $f : E \rightarrow \overline{\mathbb{R}}$ . Also, consider a clodum  $(\mathcal{K}, \vee, \wedge, \star, \star')$  of scalars with *commutative* operations  $\star, \star'$ , and  $\mathcal{K} \subseteq \overline{\mathbb{R}}$ . We define *two internal operations* among vectors/signals  $X, Y$  in  $\mathcal{W}$ : their supremum  $X \vee Y : \mathcal{W}^2 \rightarrow \mathcal{W}$  and their infimum  $X \wedge Y : \mathcal{W}^2 \rightarrow \mathcal{W}$ , which we denote using the same supremum symbol ( $\vee$ ) and infimum symbol ( $\wedge$ ) as in the clodum, hoping that the differences will be clear to the reader from the context. Furthermore,

we define *two external operations* among any vector/signal  $X$  in  $\mathcal{W}$  and any scalar  $c$  in  $\mathcal{K}$ : a “scalar multiplication”  $c \star X : (\mathcal{K}, \mathcal{W}) \rightarrow \mathcal{W}$  and a “scalar dual multiplication”  $c \star' X : (\mathcal{K}, \mathcal{W}) \rightarrow \mathcal{W}$ , again by using the same symbols as in the clodum. Now, we define  $\mathcal{W}$  to be a **weighted lattice** space over the clodum  $\mathcal{K}$  if it satisfies a set of axioms postulated in [74], which: 1) makes  $\mathcal{W}$  a distributive lattice with respect to its two internal vector operations  $\vee$  and  $\wedge$  and 2) endow the external operations  $\star$  and  $\star'$  between scalars and vectors with associativity and distributivity properties. These axioms bear a striking similarity with those of a linear space. One difference is that the vector/signal addition ( $+$ ) of linear spaces is now replaced by two dual superpositions: the lattice supremum ( $\vee$ ) and infimum ( $\wedge$ ); furthermore, the scalar multiplication ( $\times$ ) of linear spaces is now replaced by two operations  $\star$  and  $\star'$  that are dual to each other. Only one major property of linear spaces is missing from the weighted lattices: the existence of “additive inverses.” We define the space  $\mathcal{W}$  to be a **complete weighted lattice (CWL)** if: 1)  $\mathcal{W}$  is closed under any (possibly infinite) suprema and infima and 2) the distributivity laws between the scalar operations  $\star$  ( $\star'$ ) and the supremum (infimum) are of the infinite type.

3) *Vector and Signal Operators on Weighted Lattices*: We focus on CWLs whose underlying set is a *space*  $\mathcal{W}$  of functions  $f : E \rightarrow \mathcal{K}$  with values from a clodum  $(\mathcal{K}, \vee, \wedge, \star, \star')$  of scalars. Such functions include  $d$ -dimensional vectors if  $E = \{1, 2, \dots, d\}$  or  $d$ -dimensional signals of continuous ( $E = \mathbb{R}^d$ ) or discrete domain ( $E = \mathbb{Z}^d$ ). Then, we extend *pointwise* the supremum, infimum, and scalar multiplications of  $\mathcal{K}$  to functions, for example, for  $F, G \in \mathcal{W}$ ,  $a \in \mathcal{K}$ , and  $x \in E$ , we define  $(F \vee G)(x) := F(x) \vee G(x)$  and  $(a \star F)(x) := a \star F(x)$ . Furthermore, the scalar operations  $\star$  and  $\star'$ , extended pointwise to functions, distribute over any suprema and infima, respectively. If the clodum  $\mathcal{K}$  is self-conjugate, then we can extend the conjugation  $(\cdot)^*$  to functions  $F$  pointwise:  $F^*(x) := (F(x))^*$ .

Elementary increasing operators on  $\mathcal{W}$  are those that act as **vertical translations** (in short V-translations) of functions. Specifically, pointwise  $\star$  ( $\star'$ ) “multiplications” of functions in  $\mathcal{W}$  by scalars in  $\mathcal{K}$  yield the (dual) *V-translations*. A function operator  $\psi$  on  $\mathcal{W}$  is called **V-translation invariant** if it commutes with any V-translation  $\tau$ , that is,  $\psi\tau = \tau\psi$ ; similarly for dual translations.

More complex increasing operators are combinations of (dual) V-translations and dilations (erosions), called **dilation V-translation-invariant (DVI)** operators  $\delta$  or **erosion V-translation-invariant (EVI)** operators  $\varepsilon$ . Such



operators obey a sup- $\star$  or an inf- $\star'$  superposition

$$\begin{aligned} \delta \left( \bigvee_i c_i \star F_i \right) &= \bigvee_i c_i \star \delta(F_i) \\ \varepsilon \left( \bigwedge_i c_i \star' F_i \right) &= \bigwedge_i c_i \star' \varepsilon(F_i). \end{aligned} \quad (26)$$

On signal spaces, these properties create supremal and infimal *nonlinear convolutions*; details can be found in [74].

Next, we focus on finite-dimensional CWLs that are nonlinear vector spaces  $\mathcal{W} = \mathcal{K}^d$ , equipped with the pointwise partial ordering  $x \leq y$ , supremum  $x \vee y = [x_i \vee y_i]$ , and infimum  $x \wedge y = [x_i \wedge y_i]$  between any vectors  $x, y \in \mathcal{W}$ . Then,  $(\mathcal{W}, \vee, \wedge, \star, \star')$  is a CWL. Elementary increasing operators are the *vector V-translations*  $\tau_a(x) = a \star x = [a \star x_i]$  and their duals  $\tau'_a(x) = a \star' x$ , which “multiply” a scalar  $a$  with a vector  $x$  elementwise. A vector transformation on  $\mathcal{W}$  is called (dual) V-translation-invariant if it commutes with any vector (dual) V-translation. Each vector  $x = [x_1, \dots, x_d]^T$  can be expressed as the max of V-translated impulse vectors  $q_j = [q_j(i)]$ , where  $q_j(i) = e$  at  $i = j$  and  $\perp$  else, or as the min of dual V-translated impulses  $q'_j = [q'_j(i)]$ , where  $q'_j(i) = e'$  at  $i = j$  and  $\top$  else. Based on these vector representations, the following theorem establishes that all V-translation-invariant dilations and erosions of vectors are max- $\star$  and min- $\star'$  matrix-vector “products,” respectively.

*Theorem 1 [74]:* Consider mappings between two finite-dimensional CWLs.

- 1) Any vector transformation between two finite-dimensional CWLs, that is, from  $\mathcal{K}^n$  to  $\mathcal{K}^m$  is DVI iff it can be represented as a matrix-vector max- $\star$  product  $\delta_A(x) := A \boxtimes x$ , where  $A = [a_{ij}] \in \mathcal{K}^{m \times n}$  with  $a_{ij} = [\delta(q_j)]_i$ ,  $i = 1, \dots, m$  and  $j = 1, \dots, n$ .
- 2) Any vector transformation from  $\mathcal{K}^n$  to  $\mathcal{K}^m$  is EVI iff it can be represented as a matrix-vector min- $\star'$  product  $\varepsilon_A(x) := A \boxtimes' x$  where  $A = [a_{ij}]$  with  $a_{ij} = [\varepsilon(q'_j)]_i$ .

Given such a vector dilation  $\delta(x) = A \boxtimes x : \mathcal{K}^n \rightarrow \mathcal{K}^m$ , there corresponds a unique erosion  $\varepsilon : \mathcal{K}^m \rightarrow \mathcal{K}^n$  (equal to the residual operator  $\delta^\sharp$ ) so that  $(\delta, \varepsilon)$  is a *vector adjunction*, that is,  $\delta(x) \leq y \iff x \leq \varepsilon(y)$ . We can find the adjoint vector erosion by decomposing both vector operators based on *scalar operators*  $(\eta, \zeta)$  that form a *scalar adjunction* on  $\mathcal{K}$ :

$$\eta(a, v) \leq w \iff v \leq \zeta(a, w). \quad (27)$$

If we use as scalar “multiplication” a commutative binary operation  $\eta(a, v) = a \star v$  that is a dilation on  $\mathcal{K}$ , its scalar adjoint erosion becomes

$$\zeta(a, w) = \sup \{v \in \mathcal{K} : a \star v \leq w\} \quad (28)$$

which is a (possibly noncommutative) binary operation on  $\mathcal{K}$ . Then, the original vector dilation  $\delta(x) = A \boxtimes x$  is decomposed as

$$[\delta(x)]_i = \bigvee_{j=1}^n \eta(a_{ij}, x_j) = \bigvee_{j=1}^n a_{ij} \star x_j, \quad i = 1, \dots, m \quad (29)$$

whereas its adjoint vector erosion (i.e., the residual  $\delta^\sharp$  of  $\delta$ ) is decomposed as

$$[\delta^\sharp(y)]_j = [\varepsilon(y)]_j = \bigwedge_{i=1}^m \zeta(a_{ij}, y_i), \quad j = 1, \dots, n. \quad (30)$$

Furthermore, if  $\mathcal{K} = (\vee, \wedge, \star, \star')$  is a *clog*, then  $\zeta(a, w) = w \star' a^*$ , and hence

$$\varepsilon(y) = A^* \boxtimes' y, \quad [\varepsilon(y)]_j = \bigwedge_{i=1}^m y_i \star' a_{ij}^*, \quad j = 1, \dots, n \quad (31)$$

where  $A^* = [a_{ji}^*]$  is the *adjoint matrix* (i.e., conjugate transpose) of  $A = [a_{ij}]$ .

## IV. SOLVING MAX- $\star$ EQUATIONS AND OPTIMIZATION

### A. $\ell_\infty$ Optimal Solutions of Max-Plus Equations

Consider the max-plus clog  $(\overline{\mathbb{R}}, \vee, \wedge, +, +')$ , a matrix  $A \in \overline{\mathbb{R}}^{m \times n}$ , and a vector  $b \in \overline{\mathbb{R}}^m$ . The set of solutions of the max-plus equation

$$A \boxtimes x = b \quad (32)$$

over  $\overline{\mathbb{R}}$  is either empty or forms an idempotent semigroup under vector  $\vee$  because, if  $x_1$  and  $x_2$  are two solutions, then  $x_1 \vee x_2$  is also a solution. A related problem in applications of the max-plus algebra to scheduling is when a vector  $x$  represents start times, a vector  $b$  represents finish times, and the matrix  $A$  represents processing delays. Then, if (32) does not have an exact solution, it is possible to find the optimum  $x$  such that we minimize a norm of earliness subject to zero lateness

$$\text{Minimize } \|A \boxtimes x - b\|_p \quad \text{s.t. } A \boxtimes x \leq b \quad (33)$$

where  $\|\cdot\|_p$  denotes some  $\ell_p$ -norm. Both problem (32) and the constrained minimization problem (33) for  $p = 1$  or  $p = \infty$  have been solved by Cuninghame–Green [25].

*Theorem 2 [25]:* If (32) has a solution, then<sup>5</sup>

$$\hat{x} = A^* \boxtimes' b = \left[ \bigwedge_{i=1}^m b_i - a_{ij} \right] \quad (34)$$

<sup>5</sup>To cover all cases of combining finite and infinite scalar numbers in the max-plus clog  $(\overline{\mathbb{R}}, \vee, \wedge, +, +')$ , we should write the subtractions  $b_i - a_{ij}$  in (34) as  $b_i +' (-a_{ij})$  and use the rules (1).

is its greatest solution and the optimum solution to problem (33).

The proof results since  $\hat{x}$  is the greatest solution of  $A \boxplus x \leq b$ , as shown in [15] and [25]. It can also be directly seen from the adjunction  $(\delta, \varepsilon)$  where

$$A \boxplus x = \delta(x) \leq b \iff x \leq \varepsilon(b) = A^* \boxplus' b. \quad (35)$$

The solutions of (32) and of (33) for the  $\ell_\infty$  case have been further analyzed in [15] both algebraically and combinatorially. It is also possible to search and find *sparse solutions* of either the exact equation (32) or the approximate problem (33), as done in [109], where sparsity here means a large number of  $-\infty$  values in the solution vector.

Furthermore, there is actually a stronger result that is not biased to be a subsolution of (32) but provides the *unconstrained optimal solution* of the following problem:

$$\text{Minimize } \|A \boxplus x - b\|_\infty. \quad (36)$$

*Theorem 3 [15], [25]:* If  $2\mu = \|A \boxplus \hat{x} - b\|_\infty = \|A \boxplus (A^* \boxplus' b) - b\|_\infty$  is the  $\ell_\infty$  error corresponding to the greatest subsolution of  $A \boxplus x = b$ , then the unique solution of (36) is

$$\tilde{x} = \mu + \hat{x} = \mu + A^* \boxplus' b. \quad (37)$$

The computational complexity to find both optimal solutions  $\hat{x}$  and  $\tilde{x}$  is  $O(mn)$ .

## B. Projections on Weighted Lattices

The optimal subsolution of (33) can be viewed as a nonlinear “projection” of  $b$  onto the column space of  $A$  [26]. To understand this, note first that any adjunction  $(\delta, \varepsilon)$  automatically yields two lattice projections, an opening  $\alpha = \delta\varepsilon$  and a closing  $\beta = \varepsilon\delta$ , such that

$$\alpha^2 = \alpha \leq \mathbf{id} \leq \beta = \beta^2.$$

We call them “projections” because, in analogy to projection operators on linear spaces, they preserve the structure of the lattice space w.r.t. the partial ordering, and they are idempotent.

Projections on idempotent semimodules<sup>6</sup> have been studied in [23] for the general case and with more details in [2] and [26] for the max-plus case to which we focus

<sup>6</sup>Idempotent semimodules are like vector spaces with vector “addition”  $\vee$  whose vector and scalar arithmetic are defined over idempotent semirings. If, in our definition of a weighted lattice, one focuses only on one vector “addition,” say the supremum, and its corresponding scalar “multiplication,” then the weaker algebraic structure becomes an idempotent semimodule over an idempotent semiring  $(\mathcal{K}, \vee, \star)$ . This has been studied in [23], [39], and [64].

herein. Let  $\mathcal{X} = \overline{\mathbb{R}}^n$  be viewed as a complete idempotent semimodule over the complete max-plus semiring  $\mathbb{R}_{\max} \cup \{\infty\} = \overline{\mathbb{R}}$ , and let  $S$  be a complete subsemimodule of  $\mathcal{X}$ . Then, a *canonical projector* on  $S$  is defined as the nonlinear map [23]

$$P_S : \mathcal{X} \rightarrow S, \quad P_S(x) := \bigvee \{v \in S : v \leq x\}. \quad (38)$$

Its definition implies that  $P_S$  is a lattice opening, that is, increasing, antiextensive, and idempotent. Furthermore, there is a concept of “distance” on such semimodules, which allows to use a nonlinear projection theorem for best approximations. Specifically, let us consider a distance between two vectors  $x$  and  $y$  defined via the *range semi-metric* [25]

$$d_H(x, y) = \max_i (x_i - y_i) - \min_i (x_i - y_i), \quad x, y \in \mathbb{R}^n \quad (39)$$

also known, in a more general form, as the *Hilbert projective metric* [23]. Then, for any vector  $x \in \mathcal{X}$ ,  $P_S(x)$  is the best approximation (but not necessarily unique) of  $x$  by elements of  $S$  in the sense that  $P_S(x)$  is an element of  $S$  attaining the shortest distance from  $x$ , that is, [2], [23]

$$d_H(x, P_S(x)) = d_H(x, S) \quad (40)$$

where the distance between a vector  $x$  and the subspace  $S$  is defined by  $d_H(x, S) := \inf\{d_H(x, v) : v \in S\}$ . Note the analogy with Euclidean spaces  $\mathbb{R}^n$  where the linear projection of a point  $x \in \mathbb{R}^n$  to a linear subspace  $S$  is given by the unique point  $y \in S$  such that  $x - y$  is orthogonal to  $S$ .

Now, if we consider the optimization problem (33) and define the subsemimodule  $S$  in (38) as the max-plus span of the columns of matrix  $A$ , then the canonical projection of  $b$  onto it equals

$$P_S(b) = A \boxplus \hat{x} = A \boxplus (A^* \boxplus' b) \leq b \quad (41)$$

which is a lattice opening  $\delta(\varepsilon(b)) \leq b$  from (35).

## C. $\ell_p$ Optimal Solutions of Max- $\star$ Equations

Herein, we generalize the results of Section IV-A from max-plus to max- $\star$  algebra. Consider a scalar commutative clodum  $(\mathcal{K}, \vee, \wedge, \star, \star')$ , a matrix  $A \in \mathcal{K}^{m \times n}$ , and a vector  $b \in \mathcal{K}^m$ . We consider the set of solutions of both the exact max- $\star$  equation

$$A \boxtimes x = b \quad (42)$$

as well as its approximate solutions that are optimal solutions of the following constrained minimization problem:

$$\text{Minimize } \|A \boxtimes x - b\|_p \quad \text{s.t. } A \boxtimes x \leq b \quad (43)$$

where  $\|\cdot\|_p$  is any  $\ell_p$  norm with  $p = 1, 2, \dots, \infty$ . By using adjunctions, we provide next a more general result (than Theorem 2) for the general case when  $\mathcal{K}$  is a general clog or just a clodum (which may have no inverses for its “multiplication” operations).

*Theorem 4 [74]:* Consider the vector dilation  $\delta(x) = A \boxtimes x : \mathcal{K}^n \rightarrow \mathcal{K}^m$ , and let  $\varepsilon$  be its adjoint vector erosion.

- 1) If (42) has a solution, then

$$\hat{x} = \varepsilon(\mathbf{b}) = \left[ \bigwedge_{i=1}^m \zeta(a_{ij}, b_i) \right] \quad (44)$$

is its greatest solution, where  $\zeta$  is the scalar adjoint erosion of  $\star$  as in (28).

- 2) If  $\mathcal{K}$  is a clog, the solution (44) becomes

$$\hat{x} = A^* \boxtimes' \mathbf{b} = \left[ \bigwedge_{i=1}^m b_i \star' a_{ij} \right]. \quad (45)$$

- 3) The solution to the optimization problem (43) for any  $\ell_p$  norm  $\|\cdot\|_p$  is generally (44) or (45) in the case of a clog.

A main idea for solving (43) is to consider vectors  $x$  that are *subsolutions* in the sense that  $\delta(x) = A \boxtimes x \leq \mathbf{b}$  and find the greatest such subsolution  $\hat{x} = \varepsilon(\mathbf{b})$ , which yields either the greatest exact solution of (42) or an optimum subsolution in the sense of (43). To prove the latter, note that, since  $\mathbf{y} = \delta(\varepsilon(\mathbf{b}))$  is the greatest lower estimate (GLE) of  $\mathbf{b}$ ,  $b_i - y_i$  is nonnegative and minimum for all  $i$ , and hence, the norm  $\|\mathbf{b} - \mathbf{y}\|_p$  is minimum for any  $p = 1, 2, \dots, \infty$ .

Unfortunately, the type of unconstrained  $\ell_\infty$  optimal solution offered by Theorem 3 in the max-plus case does not generally carry over to a general clodum, as shown for the max–min clodum in [27].

## V. TROPICAL GEOMETRY OF NEURAL NETWORKS WITH PWL ACTIVATIONS

In this section, we present some applications of concepts and techniques from tropical geometry in studying neural networks with PWL activations. Early connections between tropical geometry and neural networks were sketched in [18] and later developed in greater detail in [19] and [114]. Tools from tropical geometry (in particular, the *Maslov dequantization*) have also been used to design neural networks that approximate convex and log–log convex data [16] and general continuous functions over convex sets [17]. For the remainder of the section, we primarily develop the tropical–geometric characterization of neural network layers following [19] and [114] and describe other applications near its end.

A central motivation for the use of tropical geometry in the study of neural networks is characterizing their *expressive power*. Tools used for this purpose range from the *Vapnik–Chervonenkis* (VC) dimension to the *activation*

*pattern* of a neural network. The seminal work of [84] and [88] proposed studying the expressive power of networks whose output is a PWL function via the number of its *inference regions* (also interchangeably called **linear regions**)—defined to be the *maximally connected partitions of the input space, in which the output of the network is a linear function*. Intuitively, networks with many linear regions can represent more complicated functions compared to networks with only a few regions of linearity. Upper and lower bounds on the number of linear regions of ReLU networks have been derived in, for example, [4], [84], [88], and [98], using arguments from combinatorics and/or polyhedral geometry. As tropical geometry is centered around the study of PWL curves and surfaces, it emerges as a natural tool for tackling this problem.

### A. Geometric Characterization of NN Layers

Motivated by the approach of [86], we seek a similar characterization of the geometry of a neural network in terms of the vertices of the appropriate Newton polytopes. With such a characterization at hand, we will then proceed to derive a simple geometric algorithm for enumerating the number of these vertices given a fixed network, to serve as a proxy for its expressive power. Our initial observation is that all PWL activation functions used in practice are tropical polynomials:

*Example 2 (ReLU/Leaky ReLU):* Given input  $v = \mathbf{w}^T \mathbf{x} + b$  with  $\mathbf{w}, \mathbf{x} \in \mathbb{R}^d$ , a rectifier linear unit computes

$$\text{ReLU}(v) = \max(0, v). \quad (46)$$

A commonly used variation is the Leaky ReLU [67], which computes (for some  $\alpha \in (0, 1)$ )

$$\text{LReLU}_\alpha(v) = \max(v, \alpha v). \quad (47)$$

Since the number of affine pieces in the expression of an ReLU/LReLU unit is just 2, we deduce that the latter are tropical polynomials of rank 2.

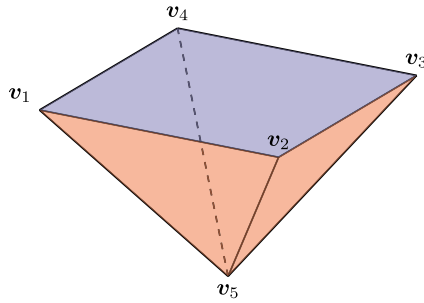
*Example 3 (Maxout):* Given  $\mathbf{W} \in \mathbb{R}^{d \times K}$  and  $\mathbf{b} \in \mathbb{R}^K$ ,  $\mathbf{x} \in \mathbb{R}^d$

$$\text{maxout}(\mathbf{x}) = \max_{j \in [K]} \left( \mathbf{W}_j^T \mathbf{x} + b_j \right) \quad (48)$$

where we denote  $\mathbf{W}_j$  for the  $j$ th column of  $\mathbf{W}$ . Thus, a maxout unit is a tropical polynomial of rank  $K$  since its expression is a maximum of  $K$  affine pieces.

These connections were observed in [18]. In particular, the authors showed that, for a single maxout unit, the number of linear regions it determines is equal to the number of vertices in the *upper hull* of its **extended Newton polytope**, defined as

$$\text{ENewt}(p) := \text{conv}(\{(b_j, \mathbf{a}_j) : j \in [K]\}) \quad (49)$$



**Fig. 7.**  $P := \text{conv}\{v_1, \dots, v_5\}$ . The upper hull,  $P^{\max}$ , is depicted in light blue.

where  $p(x)$  is given in the form of (11). For a polytope  $P$ , its **upper hull** is defined as

$$P^{\max} := \{(\lambda, \mathbf{x}) : \lambda = \sup\{t \in \mathbb{R} : (t, \mathbf{x}) \in P\}\}. \quad (50)$$

A simple example is shown in Fig. 7.

*Proposition 1* [18]: The linear regions defined by a PWL function  $p$  of the form (11) are in bijection with the number of vertices in  $\text{ENewt}^{\max}(p)$ .

*Proof Sketch:* The function computed by the tropical polynomial at each  $\mathbf{x}$  is the value of the following linear program:

$$p(\mathbf{x}) := \max\{b + \mathbf{a}^T \mathbf{x} : (b, \mathbf{a}) \in \text{ENewt}(p)\}. \quad (51)$$

It is straightforward to show that maximizers to (51) cannot exist outside  $\text{ENewt}(p)$ ; an appeal to the fundamental theorem of linear programming completes the proof.  $\square$

Proposition 1 is limited as it only characterizes a single PWL unit. However, it forms the basis for a geometric characterization of an entire NN layer. In particular, we can view each layer as a collection of tropical polynomials. Recall that the tropical hypersurface  $\mathcal{V}(p)$  of a tropical polynomial  $p$  is the set of points  $\mathbf{x}$  at which  $p(\mathbf{x})$  is nondifferentiable. Given a collection  $p_1, \dots, p_m$ , the union  $\bigcup_i \mathcal{V}(p_i)$  contains all the points  $\mathbf{x}$  for which at least one of the polynomials is nondifferentiable. Thus, each region of linearity of a neural network layer corresponds to an open cell induced by  $\bigcup_i \mathcal{V}(p_i)$ .

We may now appeal to a fundamental duality result from tropical geometry, restated in the language necessary for our application. For a proof, see [19, Proposition 1] and the discussion following [114, Definition 3.2].

*Proposition 2:* Let  $p_1, \dots, p_m : \mathbb{R}^d \rightarrow \mathbb{R}$  denote a collection of tropical polynomials. Moreover, let  $\mathcal{V}(p)$  denote the tropical hypersurface of a polynomial  $p$ . Then, the number of open cells induced by  $\bigcup_{i=1}^m \mathcal{V}(p_i)$  is equal to the number of vertices in  $\text{Newt}(p_1) \oplus \dots \oplus \text{Newt}(p_m)$ .

An illustration appears in Fig. 8. By Proposition 2 and the preceding discussion, we have reduced the problem of counting linear regions to that of counting the number of vertices of Minkowski sums of Newton polytopes. However,

as the tropical polynomials involved may also involve monomials with constant terms, we need to apply a “lifting” argument to treat the  $p_i$ ’s as functions on  $\mathbb{R}^{d+1}$  and apply Proposition 2. The resulting Minkowski sum is precisely

$$\text{ENewt}(p_1) \oplus \dots \oplus \text{ENewt}(p_m) = \text{ENewt}\left(\sum_{i=1}^m p_i\right). \quad (52)$$

Thus, it suffices to count the number of vertices in the upper hull of the Minkowski sum of (52). Based on this observation, one may appeal to standard results on the number of vertices of Minkowski sums.

*Theorem 5* [41]: Let  $P_1, \dots, P_k$  be polytopes in  $\mathbb{R}^d$ , and let  $m$  denote the number of their nonparallel edges. Then, the number of vertices of  $P_1 \oplus \dots \oplus P_k$  is bounded above by

$$2 \sum_{j=0}^{d-1} \binom{m-1}{j}. \quad (53)$$

Moreover, the bound of (53) is tight when  $2k > d$ .

When  $P_1, \dots, P_m$  are the extended Newton polytopes of ReLU units, the number of nonparallel edges is at most  $m$  since each polytope is a line segment. When  $P_1, \dots, P_m$  are generic maxout units of rank  $k$ , each polytope has at most  $K$  vertices; hence, the number of nonparallel edges is at most  $m \cdot \binom{k}{2} = m \cdot (k(k-1))/2$ . Since Theorem 5 gives upper bounds for the total number of vertices, it is not clear *a priori* how loose these bounds are for the number of vertices in the upper hull; when each  $P_i$  is the Minkowski sum of line segments, a symmetry argument can be invoked to yield bounds for the upper hull. The results are summarized in the following.

*Corollary 1:* The number  $\mathcal{N}_m^d$  of linear regions of a neural network with a single hidden layer of  $m$  neurons and  $d$  inputs is upper bounded as follows.

1) In the case of ReLU/LReLU activations, we have

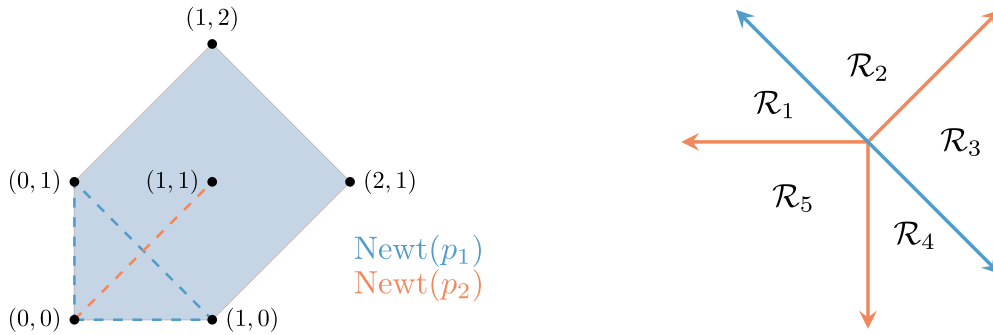
$$\mathcal{N}_m^d \leq \sum_{j=0}^d \binom{m}{j}. \quad (54)$$

The bound in (54) is tight if all the line segments generating the extended Newton polytopes, as well as their projections to the last  $d$  coordinates, are in *general position*.

2) In the case of Maxout activations of rank  $k$ , we have

$$\mathcal{N}_m^d \leq \min \left\{ k^m, 2 \sum_{j=0}^d \binom{m \cdot \frac{k(k-1)}{2}}{j} \right\}. \quad (55)$$

Similar upper bounds are straightforward to derive for convolutional layers [19]. Upper bounds for multilayer networks—which are most common in practice—are also available. However, for multilayer networks, these bounds



**Fig. 8.** Visualization of Proposition 2 for  $p_1(x, y) = \max(x, y, 0)$  and  $p_2(x, y) = \max(x + y, 0)$ . On the left,  $\text{Newt}(p_1) \oplus \text{Newt}(p_2)$  has five vertices, equal to the number of open cells formed by  $\mathcal{V}(p_1) \cup \mathcal{V}(p_2)$  (corresponding to linear regions of the polynomial  $p \equiv p_1 + p_2$ ), shown on the right.

are not a direct result of the equivalence between neural network layers and Newton polytopes. In particular, it is unclear how the composition of two neural network layers acts on their Newton polytopes and whether or not the resulting object admits a similar correspondence with linear regions of the output. The following upper bound applies to multilayer ReLU networks.

**Proposition 3** [114, Th. 6.3]: Consider a ReLU network with  $L$  layers of size  $n_1, \dots, n_L$  and inputs of dimension  $d$ . If  $n_\ell \geq d$ ,  $\ell = 1, \dots, L - 1$ , the number of linear regions of the network is upper bounded by

$$\prod_{\ell=1}^{L-1} \sum_{j=0}^d \binom{n_\ell}{j}. \quad (56)$$

**Example 4:** Suppose that we are given inputs of dimension  $d = 5$ . Consider the two following cases.

- 1) *One hidden layer with  $n_1 = 20$ :* Applying the formula from (54), the number of linear regions generated by this network is at most  $\sum_{j=0}^d \binom{20}{j} = 21\,700$ .
- 2) *Two hidden layers with  $n_1 = n_2 = 10$ :* Using (56), the number of linear regions generated by this network is at most

$$\left( \sum_{j=0}^d \binom{10}{j} \right)^2 = 407\,044.$$

We see that distributing  $m = 20$  hidden units over 2 layers instead of forming a “wide” layer increases the expressiveness of the resulting network by at least one order of magnitude.

Note that naively chaining  $L - 1$  applications of Corollary 1 gives a strictly worse bound than that of Proposition 3, as the size of the intermediate inputs for each layer can be arbitrarily larger than  $d$ . However, since the dimension of the input to the neural network is  $d$ , the “effective” dimension of each intermediate output can be at most  $d$  as well.

We conclude this section with a discussion of other research directions intimately related to PWL neural networks and their implications.

1) *Lower Bounds:* Earlier studies have provided almost-matching lower bounds for the number of linear regions of a multilayer network. In particular, Montufar et al. [84] showed—via a constructive proof—that the number of linear regions of a DNN is on the order of  $\Omega((d/W)^{(L-1)W} \cdot d^W)$  when each layer consists of  $W$  units.

The lower bound is rather existential in nature; it merely exhibits a function with a large number of linear regions representable by a DNN, instead of providing sufficient conditions (as a function of network parameters) for networks to attain this lower bound. Nevertheless, it is another argument in favor of choosing deep versus shallow architectures for learning, piling on a wealth of existing theoretical, and/or empirical evidence; for example, Telgarsky [105] followed a different approach, constructing a “hard” family of functions that are representable by networks of constant width and polynomial depth but cannot be approximated by shallow networks of subexponential width. Similarly, Arora et al. [4] exhibited a family of hard functions that require a superexponential number of hidden units when represented by a shallow ReLU DNN, as opposed to a polynomial number of units for a deep ReLU DNN, resulting in a lower bound exponentially larger than that of [105], as well as a continuum (instead of a countable family) of hard functions. Subsequent work [90] employed a different measure of expressive power called “trajectory length,” which measures changes in the output of a network as its input is varied on a 1-D path, to arrive at a similar conclusion.

2) *Generative Priors in Signal Recovery:* In addition to the depth versus width discourse, the number of linear regions of neural networks plays an important role in signal recovery with *generative priors*; a motivating example is that of compressed sensing, where one observes a set of measurements

$$\mathbf{y} = \mathbf{A}\mathbf{x}^* + \boldsymbol{\eta}$$

where  $\mathbf{x}^* \in \mathbb{R}^d$  is an unknown signal to be recovered,  $\boldsymbol{\eta}$  is observation noise, and  $\mathbf{A}$  is a known *design matrix*, typically consisting of standard Gaussian elements, with



far fewer rows than columns. The resulting problem is underdetermined and calls for further assumptions to be placed on  $\mathbf{x}^*$  (for a comprehensive review of compressed sensing, see [31]).

The most common assumption in the literature is that  $\mathbf{x}^*$  is *sparse*, in which case, the information-theoretic requirement for recovery is  $k \ll d$  measurements, where  $k$  is the number of nonzero entries of  $\mathbf{x}^*$ . However, known tractable algorithms exhibit a *computational-statistical gap* for certain problems (such as sparse phase retrieval or low-rank matrix recovery), in the sense that their sample complexity scales *quadratically*, instead of *linearly*, in  $k$ . To overcome this, researchers have proposed replacing sparsity with a less restrictive assumption; in particular, that  $\mathbf{x}^*$  **lies in the range of a ReLU network**  $G: \mathbb{R}^k \rightarrow \mathbb{R}^d$ ; in other words,  $\mathbf{x}^* = G(\mathbf{z}^*)$  for some latent vector  $\mathbf{z}^*$ . This assumption places a so-called *generative prior* on  $\mathbf{x}^*$  [10], [44].

Generative priors are known to “close” the statistical-computational gap in several applications of interest. Developing the theory behind this crucially relies on the fact that *the output of an ReLU network lies in the union of linear subspaces*, the number of which is sufficiently bounded for reasonable architectures. Tight upper bounds on the number of linear regions of ReLU networks enable precise statements about the sample complexity of recovery algorithms under a generative prior.

## B. Counting Linear Regions in Practice

In this section, we provide a geometric algorithm for approximating the number of linear regions of a neural network layer after a brief overview of existing approaches.

1) *Mixed-Integer Formulations*: A number of works have used mixed-integer programming (MIP) formulations to obtain empirical bounds on the number of linear regions; Serra *et al.* [98] showed that deep rectifier networks are mixed-integer representable when the input is restricted to a polytope. Their proof is constructive and crucially depends on a mixed-integer formulation, as summarized in the following.

Fix  $i$  and  $\ell$  to index a neuron  $i$  within a layer  $\ell$ , we denote as  $\mathbf{h}^\ell$  the vector containing the output of the  $\ell$ th layer, and let  $\mathbf{h}_0 = \mathbf{x}$  be the input to the neural network. The MIP from [98] enforces the following constraints for all  $i, \ell$ :

$$\left\{ \begin{array}{l} \mathbf{W}_i^\ell \mathbf{h}^{\ell-1} + b_i^\ell = h_i^\ell - \bar{h}_i^\ell \\ h_i^\ell \leq M z_i^\ell \\ \bar{h}_i^\ell \leq M (1 - z_i^\ell) \\ \mathbf{h}^\ell, \bar{\mathbf{h}}^\ell \geq 0 \\ z_i^\ell \in \{0, 1\}. \end{array} \right. \quad (57)$$

Let us parse the constraints in (57). First,  $z_i^\ell$  is an indicator that reveals whether neuron  $i$  in layer  $\ell$  is active or

not.  $M$  is an unspecified, sufficiently large constant that enforces  $h_i^\ell$  to be 0 when  $z_i = 0$ . If  $h_i^\ell$  denotes the output of the neuron,  $\bar{h}_i^\ell$  is a complementary “output” that satisfies  $\bar{h}_i^\ell = \max(0, -\mathbf{W}_i^\ell \mathbf{h}^{\ell-1} - b_i^\ell)$ . In [98, Th. 11], it is shown that, for a fixed  $\mathbf{x}$  and as long as  $|\mathbf{W}_i^\ell \mathbf{h}^{\ell-1} + b_i^\ell| \leq M$ , enforcing the constraints in (57) for every neuron returns a feasible solution, yielding the output of the rectifier network. Given that result, we can allow  $\mathbf{x}$  to vary over the input domain  $\mathcal{X}$  and enumerate the integer solutions  $\mathbf{z}$  of the following MIP:

$$\begin{aligned} & \text{Maximize } f \\ & \text{s.t. (57) holds } \quad \forall i, \ell \\ & \quad f \leq h_i^\ell + (1 - z_i^\ell) M \quad \forall i, \ell \\ & \quad \mathbf{x} \in \mathcal{X}. \end{aligned} \quad (58)$$

However, enumerating solutions of a mixed-integer program can be computationally intractable. To address this issue, a probabilistic algorithm was proposed in [97] to produce lower bounds to the number of possible solutions.

2) *Enumeration via Reverse Search*: The MIP-based approach above makes the simplifying assumption that the input domain is bounded, which helps determine a lower bound for  $M$  so that (58) is a valid formulation. Even though Serra *et al.* [98] discuss the issue of unbounded input domains, it tends to complicate algorithm design. In contrast, treating the enumeration problem from the scope of Newton polytope vertices applies to general domains. It is known that extreme points of Minkowski sums of polytopes are sums of extreme points of the individual polytopes; moreover, enumerating vertices of Minkowski sums of polytopes in vertex representation is possible via the so-called *reverse search* method [5], [35].

The resulting algorithm for vertex enumeration has runtime

$$\mathcal{O}(\delta \cdot \text{LP}(d, \delta) \cdot N), \quad \delta := \sum_{i=1}^m \delta_i$$

where  $N$  is the number of vertices,  $P_i \subset \mathbb{R}^d$ ,  $i = 1, \dots, m$ , are the polytopes in the Minkowski sum,  $\delta_i$  denotes the maximal degree of the vertex adjacency graph of  $P_i$ , and  $\text{LP}(d, \delta)$  denotes the complexity of solving a linear program (LP) in  $d$  variables and  $\delta$  constraints. The above implies straightforward bounds for exact counting of linear regions of ReLU/maxout layers. For ReLUs,  $\delta_i = 2$  for all  $i$ , so  $\delta = 2m$ . In the latter case, denoting  $k_i$  for the rank of the  $i$ -th unit,  $\delta = \sum_i k_i$ .

Unfortunately, reverse search requires solving a prohibitive number of LPs, rendering the above approach impractical. We attack this problem from a different angle, by considering the “dual” problem of counting vertices of convex polytopes by sampling.

### C. A Geometric Algorithm

We present a randomized method for “sampling” the extreme points of the upper hull of a polytope  $P = P_1 \oplus \dots \oplus P_m$ . We generate  $K$  standard normal vectors, that is,  $\mathbf{g}_k \stackrel{i.i.d.}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ , and compute  $(\mathbf{g}_k)^T \mathbf{v}_i$  for all extreme points  $\mathbf{v}_i$ . We record the minimizers/maximizers for each polytope  $P_j$  and repeat the trial. Denoting by  $\mathbf{V}_i \in \mathbb{R}^{k_i \times d}$  the matrix whose rows contain the coordinates of each vertex of  $P_i$ , the above procedure essentially counts the number of unique tuples giving the row indices of the extrema of  $\mathbf{V}_i \mathbf{g}_k$  for all  $i$ .

From our discussion motivating the use of the reverse search method, it is clear that the resulting number is a lower bound on the number of vertices in the Minkowski sum. The resulting Algorithm 1 leverages the techniques in [28]. This method and its specialization to upper hulls work for *general* polytopes, whereas the MIP-based methods in the literature are only presented for rectifier networks. On the other hand, it should be noted that MIP formulations can be used to enumerate the number of linear regions of **deep** neural networks; in contrast, it is unclear how to adapt our geometric algorithm or the reverse search method for neural networks with more than one layer. Finally, we note that adapting Algorithm 1 for counting vertices in upper hulls of Minkowski sums is described in [19, Sec. 4.1]. The idea for extending the technique to the upper hull is simple: to ensure that the maximizers of the linear forms lie on the upper hull, we restrict ourselves to samples with a positive first coordinate. The final guarantee is similar to the one given in Proposition 4, though stated in terms of a restricted normal cone.

---

#### Algorithm 1 Sampling Points in the Convex Hull

---

**Input:** polytopes  $P_1, \dots, P_m$  in vertex representation  
 $I_{\text{ext}} := \emptyset$ .  
**for**  $j = 1, \dots, K$  **do**  
     Sample  $\mathbf{g}_j \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$   
     Compute  $\mathbf{z}^i := \mathbf{V}_i \mathbf{g}_j, \forall i \in [m]$ .  
     Collect  $\begin{cases} \mathbf{z}_{\max} := (\text{argmax } \mathbf{z}^1, \dots, \text{argmax } \mathbf{z}^m) \\ \mathbf{z}_{\min} := (\text{argmin } \mathbf{z}^1, \dots, \text{argmin } \mathbf{z}^m) \end{cases}$   
      $I_{\text{ext}} := I_{\text{ext}} \cup \{\mathbf{z}_{\max}, \mathbf{z}_{\min}\}$   
**end for**

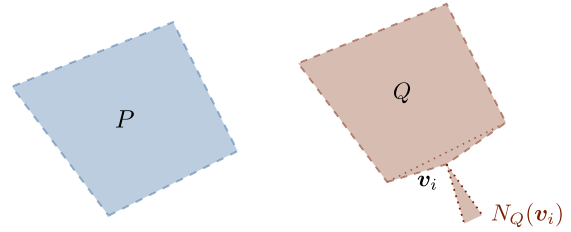
---

Algorithm 1 provides a nontrivial lower bound to the number of extreme points of the resulting Minkowski sum with high probability, as shown in Proposition 4.

*Proposition 4:* Let  $N$  denote the number of vertices of  $P = P_1 \oplus \dots \oplus P_m$ , a failure probability  $\delta$ , and define

$$\tilde{N} := \left( \log \left( \frac{1}{\max_i (1 - 2\omega(N_P(\mathbf{v}_i)))} \right) \right)^{-1}$$

where  $\omega(N_P(\mathbf{v}_i))$  is the *solid angle* of the normal cone  $N_P(\mathbf{v}_i)$  of the  $i$ th vertex. Then, for  $K \geq \tilde{N} \log(N/\delta)$  in Algorithm 1, the algorithm



**Fig. 9.** Polytopes  $P$  and  $Q$  and their solid angles. All the solid angles of  $P$  (left) are bounded away from zero. On the other hand, for  $Q$  (right), we have  $\omega(N_Q(\mathbf{v}_i)) \ll 1$ .

records all the vertices with probability at least  $1 - \delta$ .

*Proof Sketch:* The key idea in the proof is the following: extreme points of Minkowski sums are also extreme points of individual summands. Consequently, missing a “configuration” of minimizers across our trials is equivalent to missing an extreme point  $\mathbf{v}$  of the Minkowski sum.

Moreover, it is not hard to see (e.g., [19, Corollary 1]) that the solid angles of the normal cones of the vertices of a polytope  $P$  form a probability distribution, with

$$\omega(N_P(\mathbf{v}_i)) = \mathbb{P}_{\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)}(\mathbf{g} \in N_P(\mathbf{v}_i)) \quad (59)$$

the probability that  $\mathbf{g}$  is in the normal cone at  $\mathbf{v}_i$  and, consequently,  $\mathbf{v}_i$  being the minimizer of the linear function  $x \mapsto x^T \mathbf{g}$ . The rest follows from a coupon collector-style argument; a detailed proof is available in [19, Sec. 4].  $\square$

*Example 5:* Suppose that  $P$  has all-equal solid angles, that is,  $\omega(N_P(\mathbf{v}_i)) = (1/N)$ , for all  $i$ , in which case,  $\tilde{N} = \log(N/(N-2))^{-1}$ . Rewriting  $\log N/(N-2) = \log(1 + 2/(N-2))$  and combining with the inequality  $\log(1+x) \leq x$ , we see that

$$\tilde{N} \geq \frac{N-2}{2} \Rightarrow K \geq \left( \frac{N}{2} - 1 \right) \log(N/\delta)$$

is necessary to achieve probability failure at most  $\delta$ . Note that this shows that Algorithm 1 will require at least this number of samples for **any** polytope  $P$ ; indeed, it is easy to see that  $\min_i \omega(N_P(\mathbf{v}_i)) \leq (1/N)$  for any polytope with  $N$  vertices.

Our guarantee heavily depends on the cones  $N_P(\mathbf{v}_i)$ . If there are vertices that only slightly “extend” out of the polytope, our required sample size will be a large multiple of  $N$ . Fig. 9 illustrates (nonzonotopal) examples in  $\mathbb{R}^2$ ;  $Q$  has a vertex where the solid angle of the normal cone is close to 0, in contrast to  $P$  that is more “regular.” Nevertheless, the proposed algorithm can be easily parallelized, only relies on computing inner products, and crucially utilizes the geometric insights from the Newton polytope characterization of neural network layers.

### D. Other Connections Between Tropical Geometry and Neural Networks

1) *Tropical Polynomial Division and Network Simplification:* Another problem where tropical geometry can be of

use is neural network minimization; as neural networks increase in complexity, so do their needs in computing time and memory, limiting their use in time-sensitive applications. Therefore, we seek to reduce the size of a neural network while maintaining its accuracy. Several methods have attempted to solve this problem, by removing either connections between neurons [43] or neurons themselves [47], [66] from the network. The former is referred to as *weight* or *unstructured pruning* and the latter as *channel/neuron* or *structured pruning*. These studies show that minimal drops in accuracy (roughly 1% on the VGG-16 architecture) are possible, despite a significant decrease in network complexity. Note that neural network compression is distinct from the so-called dropout technique [102]; the latter is a technique applied during the training stage and aims to address the problem of overfitting by setting a random subset of the neurons to zero during each training epoch.

Tropical geometry can also provide novel methods for neural network simplification.

- 1) Given a fixed ReLU network, we can attempt to construct a smaller neural network whose Newton polytopes closely approximate the polytopes of the original network. The resulting algorithm is *constructive* and relies on the concept of *tropical polynomial division* [100], which approximates the dividend using the Newton polytopes of the divisor and quotient. Since this method constructs a network from scratch, it can be much faster than pruning methods in practice. It was originally applied to minimize the second-to-last layer of networks with a single output neuron in the context of binary classification problems, with less than 0.5% loss in accuracy even when only 1% of the hidden units are retained. Extensions to multiclass problems are considered in [101].
- 2) A complementary approach appeared in [3]. The authors first obtain a tropical geometric characterization of the decision boundaries of neural networks using their Newton polytopes; following that, they present a regularization method that balances a sparsity-inducing penalty with an objective that attempts to preserve the decision boundaries of the neural network. In contrast to the previous two approaches, this is a pruning method.

2) *Morphological Neural Networks*: Though feedforward networks with PWL activations have become the de-facto standard in neural computation, the paradigm of so-called *morphological computation* is also closely related to tropical geometry. In morphological computation, linear operations are replaced with their tropical versions; thus, the building blocks of a morphological neural network are replaced by *dilations* and *erosions* instead of linear operations. In its most elementary version, a morphological (max, +)-perceptron computes the function

$$x \mapsto w^T \boxplus x = \max_i \{w_i + x_i\} \quad (60)$$

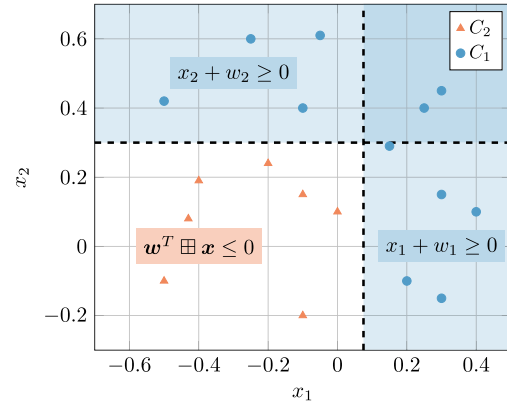


Fig. 10. Example of tropically separable patterns in  $\mathbb{R}^2$ , using the weight vector  $w = [0.075, 0.3]^T$ .

where  $w \in \mathbb{R}_{\max}^d$  is a set of trainable weights. In binary classification, the decision regions induced by a (max, +) perceptron are collections of so-called *tropical halfspaces*. A (max, +)-perceptron can separate two classes if and only if a certain *tropical polyhedron* is nonempty, a condition that can be checked efficiently for this particular case.

*Proposition 5 [18, Proposition 1]*: Consider  $N_1$  points from class  $C_1$  and  $N_2$  points from class  $C_2$ , forming the matrices  $X_1 \in \mathbb{R}^{N_1 \times d}$  and  $X_2 \in \mathbb{R}^{N_2 \times d}$ . Then, these points can be separated by a morphological perceptron of the form (60) if and only if

$$\begin{aligned} \{w \in \mathbb{R}_{\max}^d : X_1 \boxplus w \geq \mathbf{0}_{N_1}, \quad X_2 \boxplus w \leq \mathbf{0}_{N_2}\} \neq \emptyset \\ \Leftrightarrow X_1 \boxplus (X_2^* \boxplus' \mathbf{0}) \geq \mathbf{0}. \end{aligned} \quad (61)$$

An example of a tropically separable configuration of points is shown in Fig. 10. Even though the morphological paradigm dates back almost 3 decades [91], [92], [103], [113], a recent resurgence of interest has led to new developments; for example, it was recently shown [83], [115] that a morphological neural network with a hidden layer consisting of dilations and erosions followed by a linear layer is a universal approximator. In a more recent publication [34], the authors focus on deep learning for image processing, treating all nonlinear operations (e.g., max-pooling) as trainable morphological operators to complement trainable convolutional operations and achieve competitive results in tasks, such as boundary detection using considerably fewer parameters than other architectures.

## VI. TROPICAL GEOMETRY AND GRAPHICAL MODELS

### A. Hidden Markov Models

The use of tropical geometry within the framework of parametric statistics was pioneered by Pachter and Sturmfels [86]. Specifically, they consider graphical models that are formally represented by directed acyclic graphs

with two sets of vertices: the *hidden variables*  $\mathbf{X} = (X_1, \dots, X_m)$  and the *observed variables*  $\mathbf{Y} = (Y_1, \dots, Y_n)$ . Moreover, we use  $s_1, \dots, s_d$  to denote the *model parameters*. Given an observation  $\sigma = (\sigma_1, \dots, \sigma_n)$ , the observation probabilities are polynomials of degree  $E$  in the model parameters, where  $E$  is the number of edges of the aforementioned graph. We use  $f_\sigma(s_1, \dots, s_d) = \mathbb{P}_{s_1, \dots, s_d}(\mathbf{Y} = \sigma)$  to denote the observation probability. Pachter and Sturmfels [86] asked a fundamental question about this family of models:

*How do the solutions to inference problems depend on the model parameters?*

The authors fix the numbers  $d$  and  $n$  of model parameters and observations, and furthermore, assume that each of the observed variables can take  $\ell$  different values. Mathematically, this model is a polynomial map  $f : \mathbb{R}^d \rightarrow \mathbb{R}^{\ell^n}$ , each of the coordinates being one of the aforementioned polynomials of degree  $E$ . Let  $u_i := -\log(s_i)$  determine the associated logarithmic parameter space. Moreover, define

$$g_\sigma(u_1, \dots, u_d) := -\max_{\mathbf{h}} \log \mathbb{P}_{s_1, \dots, s_d}(\mathbf{X} = \mathbf{h} \mid \mathbf{Y} = \sigma). \quad (62)$$

Pachter and Sturmfels [86] show that  $g_\sigma$  is PWL and concave on the logarithmic parameter space, with the normal cones of  $\text{Newt}(f_\sigma)$  identifying its domains of linearity. As the parameters  $u_1, \dots, u_d$  vary, they define inference functions  $\sigma \mapsto \hat{\mathbf{h}}$ , where  $\hat{\mathbf{h}}$  is the most likely tuple of hidden variables given an observation  $\sigma$ . This leads to the following.

*Proposition 6 [86, Proposition 6]:* The inference functions  $\sigma \mapsto \hat{\mathbf{h}}$  of a graphical model  $f$  are in bijection with the vertices of the Newton polytope of the map  $f$ . The explanations  $\hat{\mathbf{h}}$  for a fixed observation  $\sigma$  in a graphical model are in bijection with the vertices of the Newton polytope of the polynomial  $f_\sigma$ .

This is the main ingredient in [86], which the authors employ to deduce upper bounds on the number of inference functions and explanations of graphical models, by leveraging known bounds on the number of vertices of Newton polytopes. Finally, they motivate theoretically the use of the so-called *polytope propagation* algorithm to enumerate the vertices of the aforementioned polytopes, including an application to inference for biological sequence analysis [87].

The authors of a later publication [24] study the *Restricted Boltzmann Machine (RBM)*, a graphical model that is the building block of deep belief networks [52], using techniques from algebraic and tropical geometry. Formally, RBMs are represented by a bipartite graph on hidden variables  $\mathbf{h} \in \{0, 1\}^k$  and observed variables  $\mathbf{v} \in \{0, 1\}^n$ , with “activation”

$$\psi(\mathbf{v}, \mathbf{h}) := \exp(\mathbf{h}^T \mathbf{W} \mathbf{v} + \mathbf{b}^T \mathbf{v} + \mathbf{c}^T \mathbf{h}) \quad (63)$$

which determines a probability distribution

$$p(\mathbf{v}) := \frac{1}{Z} \sum_{\mathbf{h} \in \{0, 1\}^k} \psi(\mathbf{v}, \mathbf{h}), \quad Z := \sum_{\mathbf{v}, \mathbf{h}} \psi(\mathbf{v}, \mathbf{h}) \quad (64)$$

where  $Z$  is the induced *log-partition function*. The authors then define the **tropical RBM** model by applying the Maslov dequantization principle to  $\log p(\mathbf{v})$ , leading to the PWL convex model in (65)

$$q(\mathbf{v}) := \max\{\mathbf{h}^T \mathbf{W} \mathbf{v} + \mathbf{b}^T \mathbf{v} + \mathbf{c}^T \mathbf{h} : \mathbf{h} \in \{0, 1\}^k\}. \quad (65)$$

Similar to [86], varying the parameters  $(\mathbf{b}, \mathbf{c}, \mathbf{W})$  determines a collection of inference functions. Cueto et al. [24] then obtained the following characterization of an RBM’s inference functions (recall that a *linear threshold function* is a function of the form  $f(\mathbf{x}) = \text{sign}(\boldsymbol{\alpha}^T \mathbf{x} + \beta)$ ):

*Proposition 7 [24, Proposition 5.1]:* The inference functions for the RBM in  $k$  hidden and  $n$  observed variables are precisely those Boolean functions  $\{0, 1\}^n \rightarrow \{0, 1\}^k$  for which each of the coordinate functions is a linear threshold function.

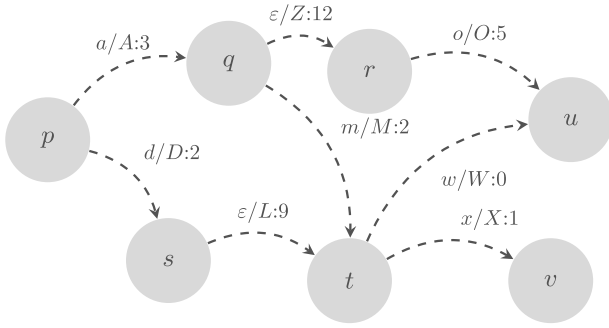
## B. Tropical Algorithms on WFSTs

*1) Introduction:* Weighted Finite State Transducers (WFSTs) introduce a computational framework that extends traditional automata, with applications in automatic speech recognition, natural language processing, computational biology, and more. The workhorse of the framework is the Viterbi algorithm, a decoding procedure that performs inference over graphs. The framework also includes a variety of algorithms aiming to reduce the computational footprint, which can be split into two categories: 1) algorithms that respect the initial topology of the network, refactoring the weights or removing extraneous transitions and 2) algorithms that fundamentally alter the structure of the network, via network minimization or composition. In any case, WFSTs, complete with their suite of diverse algorithms, present a formal mathematical framework whose properties have been analyzed for decades. A simple WSFT is shown in Fig. 11.

WFST algorithms historically employed tropical arithmetic [81], [82] for practical reasons. However, their formal modeling using tropical matrix algebra was only recently explored. A recent work [106] tropicalized the Viterbi algorithm<sup>7</sup> and its pruning variant, both seminal communications algorithms, by expressing the symbol observation probabilities as a tropical diagonal matrix. A following work [107] extended the tropicalization to other instrumental WFST algorithms, namely, *epsilon removal* and *weight pushing*, via the strong and weak transitive closures of the network.

<sup>7</sup>The framework of weighted lattices allows us to analyze the max-product form of the Viterbi algorithm as a nonlinear dynamical system in state space and extend it to more general forms that accept control inputs [74].





**Fig. 11. Toy WFST.** Transitions are of the form  $i/o:c$ , where  $i$  is the input symbol,  $o$  is the output symbol, and  $c$  is the cost of the transition. For example, the input sequence  $ax$  would be decoded to  $DLX$  with a total cost of 12. Here, the  $\varepsilon$ -transition denotes the lack of an input symbol.

In addition, a tropical analog to spectral graph theory can be found, which studies the existence and characterization of solutions to the tropical (sup-)eigenvalue problem. While mathematically analyzing WFSTs, certain elements from tropical spectral theory arise, which are introduced in the next section.

2) *Background:* A WFST is mainly characterized by the *transition matrix* of a network, which we denote  $\mathbf{A} \in \mathbb{R}_{\min}^{d \times d}$ , and where each entry  $a_{ij}$  corresponds to the cost of transitioning from state  $i$  to state  $j$ . The initial states are denoted by  $\pi \in \mathbb{R}_{\min}^d$ , where each initial state has a finite cost, and  $+\infty$  otherwise. Similarly, emitting (or final) states are denoted by  $\rho \in \mathbb{R}_{\min}^d$  and have also finite costs (and  $+\infty$  if they are not final states). We define the *weak transitive closure* of  $\mathbf{A}$  as

$$\Gamma(\mathbf{A}) := \mathbf{A} \wedge \mathbf{A}^2 \wedge \cdots \wedge \mathbf{A}^d \wedge \cdots \quad (66)$$

and the *strong transitive closure* as

$$\Delta(\mathbf{A}) := \mathbf{I} \wedge \mathbf{A} \wedge \mathbf{A}^2 \wedge \cdots \wedge \mathbf{A}^d \wedge \cdots \quad (67)$$

where  $\mathbf{A}^k = \overbrace{\mathbf{A} \boxplus \dots \boxplus \mathbf{A}}^{k \text{ times}}$ . The *minimum cycle mean* of  $\mathbf{A}$  is defined as

$$\lambda(\mathbf{A}) = \min_{c \in C(\mathbf{A})} \frac{\text{weight}(c)}{\text{length}(c)}$$

where  $C(\mathbf{A})$  is the set of cycles of the network, and  $\text{weight}(\cdot)$  and  $\text{length}(\cdot)$  denote the weight (sum of the costs along the cycle) and length of a cycle, respectively.

In *tropical spectral analysis*, the min-plus eigenproblem<sup>8</sup> of  $\mathbf{A}$  consists of finding the *eigenvalues*  $\lambda$  and *eigenvectors*  $\mathbf{v}$  such that

$$\mathbf{A} \boxplus \mathbf{v} = \lambda + \mathbf{v}. \quad (68)$$

<sup>8</sup>In [15], the discussion revolves around the max-plus eigenproblems; instead, in this section, the analysis will focus on the min-plus eigenproblems.

The minimum cycle mean  $\lambda(\mathbf{A})$  plays a fundamental role in the min-plus eigenproblem; indeed, it is the *smallest* eigenvalue and the only one whose corresponding eigenvectors may be *finite* [15]. For the spectral analysis component, we will heavily rely on the following theorem, which characterizes the *sup-eigenvectors* of  $\mathbf{A}$ , which are defined as the solutions to:

$$\mathbf{A} \boxplus \mathbf{v} \geq \lambda + \mathbf{v}. \quad (69)$$

*Theorem 6 [15, Dual of Theorem 1.6.18]:* Suppose that  $\mathbf{A}$  has at least one finite entry. If  $\lambda \leq \lambda(\mathbf{A})$  and  $\lambda < +\infty$ , then the following holds.

- 1)  $\mathbf{A} \boxplus \mathbf{v} \geq \lambda + \mathbf{v}$  has a finite solution.
- 2) The set of finite sup-eigenvectors is

$$V^*(\mathbf{A}, \lambda) = \left\{ \Delta(\mathbf{A} - \lambda) \boxplus \mathbf{u} : \mathbf{u} \in \mathbb{R}^d \right\}. \quad (70)$$

- 3)  $\mathbf{A} \boxplus \mathbf{v} \geq \lambda + \mathbf{v}$  only holds if  $\mathbf{v} = \Delta(\mathbf{A} - \lambda) \boxplus \mathbf{u}$ ,  $\mathbf{u} \in \mathbb{R}_{\min}^d$ .

The characterization of the eigenvectors is of significant importance in tropical settings. Indeed, in max-plus dynamical systems modeling manufacturing processes, it is desired that some systems whose dynamics are governed by  $\mathbf{x}(t) = \mathbf{A} \boxplus \mathbf{x}(t-1)$  eventually reach a steady state where the processing occurs at regular intervals, that is,  $\mathbf{x}(t) = \lambda + \mathbf{x}(t-1)$ . If  $\mathbf{x}(0)$  is an eigenvector, the steady state is immediately reached; therefore, the characterization is fundamental, as it provides well-behaved configurations, in terms of reachability of a steady state, for a dynamical system.

3) *Tropicalization of WFST Algorithms:* The *Viterbi algorithm*, stemming from the field of communications, attempts to decode the most probable series of latent states from a data sequence. At the heart of this algorithm is the following recursive computation: given a sequence of observations  $\{\sigma_t\}_{t=0}^T$ , observation probabilities  $b(\sigma_t)$ , and a transition matrix  $\mathbf{W}$ , the highest probability of a single partial state sequence ending at state  $i$  at time  $t$  and accounting for the first  $t+1$  observations is given by

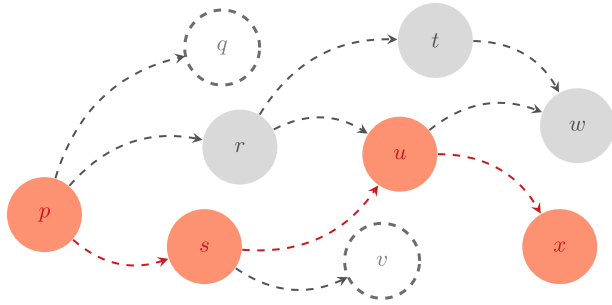
$$q_i(t) = \max_j b_i(\sigma_t) W_{ji} q_j(t-1). \quad (71)$$

We can formally tropicalize (71) and provide a recursive solution for the state vector  $\mathbf{x}(t)$

$$\mathbf{x}(t) = \mathbf{P}(\sigma_t) \boxplus \mathbf{A}^T \boxplus \mathbf{x}(t-1) \quad (72)$$

where  $\mathbf{x}(t) = -\log q(t)$ ,  $\mathbf{A} = -\log \mathbf{W}$ , and  $\mathbf{P}(\sigma_t) = \text{diag}(-\log b(\sigma_t))$ , with  $\text{diag}(\cdot)$  denoting a matrix with the argument in the diagonal and  $+\infty$  elsewhere.





**Fig. 12.** Illustration of the Viterbi pruning. The path of optimal states is denoted by red. States colored gray were examined by the algorithm for optimality, whereas the dashed states had high costs and were pruned.

Viterbi pruning is a practical technique that is frequently used in order to reduce the computational burden of the decoding. In essence, the optimal path is computed at each step, and only the paths whose cost is within a certain threshold are further expanded. An intuitive example is given in Fig. 12. Viterbi pruning can be thought as the problem

$$\mathbf{X}(t) \boxplus \mathbf{y} \geq \eta \quad (73)$$

where  $\mathbf{X}(t) = \text{diag}(x(t))$  and  $\eta$  is a vector with  $\eta_i = \frac{1}{2}(x(t)^T \boxplus x(t)) + \theta$ , where  $\theta$  is the pruning parameter. We can then interpret pruning as finding the smallest solution  $\bar{\mathbf{y}} \in \mathbb{R}_{\min}^d$ , satisfying the min-plus inequality (73), which can be done using the dual of Theorem 2

$$\bar{\mathbf{y}} = \mathbf{X}^*(t) \boxplus \eta \quad (74)$$

where  $\mathbf{X}^*(t) = -\mathbf{X}^T(t)$  and the negative entries of  $\bar{\mathbf{y}}$  indicate the indices to be pruned. A geometrical interpretation can be given to the Viterbi pruning; in particular, the set of feasible solutions at each step is a tropical polytope (see Fig. 13)

$$\mathcal{T}(\mathbf{x}(t), \eta) = \left\{ z \in \mathbb{R}_{\min}^d : z \geq \mathbf{x}(t), z \leq \eta \right\}. \quad (75)$$

*Example 6:* Let the state vector be

$$\mathbf{x}(t_0) = [1 \quad 7 \quad 4]^T$$

at some time  $t_0$ , and suppose that the pruning parameter is  $\theta = 5$ . Then,  $\eta_i = \frac{1}{2}(\mathbf{x}(t_0)^T \boxplus \mathbf{x}(t_0)) + \theta = 6$ . The optimal solution then is given by (74)

$$\bar{\mathbf{y}} = \begin{bmatrix} -1 & -\infty & -\infty \\ -\infty & -7 & -\infty \\ -\infty & -\infty & -4 \end{bmatrix} \boxplus \begin{bmatrix} 6 \\ 6 \\ 6 \end{bmatrix} = \begin{bmatrix} 5 \\ -1 \\ 2 \end{bmatrix}.$$

As  $\bar{y}_2$  is negative, it gets pruned, and the resulting vector is

$$\mathbf{x}_p(t_0) = [1 \quad \infty \quad 4]^T.$$

We emphasize that  $\eta$ ,  $\bar{\mathbf{y}}$ , and the resulting polytope are different for each time step  $t$ .

The weight pushing algorithm is an essential component of the WFST framework [82]. The algorithm improves the effectiveness of the Viterbi pruning by pushing weights toward earlier transitions and states, without altering the overall path statistics (i.e., the decoded sequences and their probabilities). After weight pushing, low-probability sequences can be identified and pruned early during decoding and increasing efficiency.

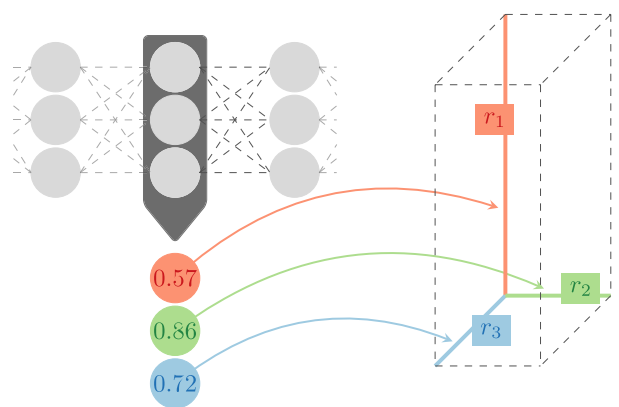
Integral to the weight pushing algorithm is the computation of a potential for each state of the graph. In short, the potential value is the weight amount that can be “pushed” to earlier states and can be computed via an iterative evaluation. A single iteration of the potential vector can be expressed as [107]

$$\mathbf{v}_{i+1} = \mathbf{v}_i \wedge \mathbf{A} \boxplus \mathbf{v}_i \quad (76)$$

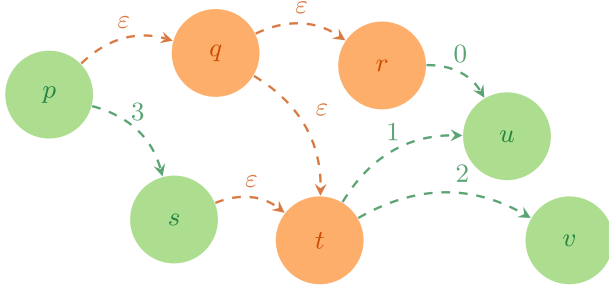
with  $\mathbf{v}_0 = \rho$  being the emission vector. Recursively iterating (76), we arrive at the final potential vector

$$\begin{aligned} \mathbf{v}_\infty &= \rho \wedge \mathbf{A} \boxplus \rho \wedge \mathbf{A}^2 \boxplus \rho \wedge \dots \wedge \mathbf{A}^n \boxplus \rho \wedge \dots \\ &= \Delta(\mathbf{A}) \boxplus \rho \end{aligned} \quad (77)$$

where the computation of  $\Delta(\mathbf{A})$  is finite under very mild assumptions; namely, that the graph does not contain cycles of negative weight, and thus,  $\lambda(\mathbf{A}) \geq 0$  (a standard assumption for WFSTs). Henceforth, we assume that these conditions hold. After the potential computation,



**Fig. 13.** At each decoding step,  $\mathbf{x}(t)$  and  $\eta$  of (75) define a polytope. The vector  $\mathbf{r} = \eta - \mathbf{x}(t)$  denotes the range of each dimension (negative ranges indicate that the index is pruned).



**Fig. 14. Illustration of epsilon transitions and states.** orange denotes states and transitions that will be removed by the epsilon removal algorithm, whereas surviving ones are denoted by green.

the network parameters can be updated via the rules

$$\pi' = \pi + v_\infty, \quad \rho' = \rho - v_\infty, \quad A' = V^- \boxplus A \boxplus V^+ \quad (78)$$

where  $V^+ = \text{diag}(v_\infty)$  and  $V^- = \text{diag}(-v_\infty)$ .

Another instrumental algorithm to the WFST framework is *epsilon removal* [82]. Similar to weight pushing, this algorithm facilitates decoding by decreasing the size of the network by removing extraneous transitions and states. Examples of these states and transitions can be seen in Fig. 14. The removal of extraneous transitions and states is achieved through the computation of the *epsilon closure* of every state, which encapsulates the states that are reachable using only epsilon transitions. To that end, the network matrix  $A$  can be decomposed [107] into two components

$$A = A_\varepsilon \wedge A_{\varepsilon^\perp} \quad (79)$$

where  $A_\varepsilon$  contains only the epsilon transitions and  $A_{\varepsilon^\perp}$  contains the nonepsilon transitions. The epsilon closure is then computed as the shortest distances of the network matrix  $A_\varepsilon$ , which, via the definition of the weak transitive closure, is given by  $\Gamma(A_\varepsilon)$ , where the computation is finite and equal to  $\Gamma(A_\varepsilon) = A_\varepsilon \wedge \dots \wedge A_\varepsilon^d$ . Having computed the epsilon closure, the updated network parameters take the form

$$\begin{aligned} A' &= A_{\varepsilon^\perp} \wedge (\Gamma(A_\varepsilon) \boxplus A_{\varepsilon^\perp}) = \Delta(A_\varepsilon) \boxplus A_{\varepsilon^\perp} \\ \rho' &= \rho \wedge (\Gamma(A_\varepsilon) \boxplus \rho) = \Delta(A_\varepsilon) \boxplus \rho. \end{aligned} \quad (80)$$

4) *Spectral Analysis of Tropical WFST Algorithms*: The representation that we developed in the previous sections offers a unified computational framework that enables a holistic analysis of WFSTs; in certain cases, it also enables a geometrical characterization of the algorithms via elements of algebraic geometry, such as polytopes. Herein, the computational framework of tropical algebra further enables the *spectral characterization* of the graph algorithms, that is, we are able to characterize the introduced algorithms via their *eigenvalues*. This characterization

introduces a new dimension to these algorithms, as we are now able to examine their properties for different eigenvalues.

We established that a mild (and realistic) assumption for the class of networks is that the cycles of the network have nonnegative weights, and therefore,  $\lambda(A) \geq 0$ . Therefore, we can view (77) in the scope of Theorem 6

$$v_\infty = \Delta(A) \boxplus \rho = \Delta(A - \lambda) \boxplus u \quad (81)$$

with  $u = \rho$  and  $\lambda = 0$ . Thus,  $v_\infty$  is a tropical sup-eigenvector of  $A$  for the tropical eigenvalue 0. Then, contextualizing (78) under the prism of Theorem 6,  $\pi'$  comprises two subsystems: the original system model  $\pi$  (which is required to maintain the system dynamics) and a new, well-behaved in terms of steady state, subsystem  $v_\infty$ . The rest of the updates in (78) ensures that the cost of each path remains unaffected.

Similarly, we can revisit (80) and express the updated network parameters as

$$\begin{aligned} \rho' &= \Delta(A_\varepsilon) \boxplus \rho = \Delta(A_\varepsilon - \lambda) \boxplus u \\ A' &= \Delta(A_\varepsilon) \boxplus A_{\varepsilon^\perp} = \Delta(A_\varepsilon - \lambda) \boxplus U \end{aligned} \quad (82)$$

where  $\lambda = 0$ ,  $u = \rho$ , and  $U = A_{\varepsilon^\perp}$ . Note that  $\rho' = \Delta(A_\varepsilon - \lambda) \boxplus u$  is similar to (81); it simply refers to the sup-eigenproblem of  $A_\varepsilon$ . The second equation of (82) consists of a *collection* of tropical sup-eigenvectors of  $A_\varepsilon$ . In this case, the immediate effects of Theorem 6 are less pronounced; while  $\rho'$  and  $A'$  are (collections of) eigenvectors, they are not employed to send the system to a steady state.

From this analysis, we make two remarks: first,  $A'$  of (78) is, by definition, *visualized* [15], meaning that it has a simpler structure than  $A$ , while still maintaining the same spectral properties. As a second remark, we note that tropical eigenvalue problems have *infinite* solutions. Indeed, it is a well-known fact in tropical algebra [15] that sup-eigenvectors exist for each eigenvalue  $\lambda \in [0, \lambda(A)]$ . Therefore, this creates a whole family of WFSTs that all solve some eigenvalue problem for all  $\lambda$  in the aforementioned range.

## VII. TROPICAL REGRESSION

Herein, we expand on our previous work [76] and apply tropical geometry and max-plus algebra to a fundamental regression problem of approximating the shape of curves and surfaces by fitting *piecewise linear (PWL)* functions, represented by tropical polynomials (11), to data possibly sampled from a functional form and in the presence of noise. We begin with a brief sampling of PWL models.

### A. PWL Function Representation and Data Fitting

PWL functions  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  are defined as follows: 1) their domain is divided into a finite number of

polyhedral regions separated by linear  $(d - 1)$ -dimensional boundaries that are hyperplanes or subsets of hyperplanes and 2) they are affine over each region and continuous on each boundary. In using them for regression, two major problems are *representation*, that is, finding better analytical expressions to represent them, and their *parameter estimation* for modeling a nonlinear system or fitting some data. Furthermore, while these problems are well-explored in the 1-D case, they remain relatively underdeveloped for multidimensional data.

The so-called *canonical representation* for continuous PWL functions, consisting of an affine function plus a weighted sum of absolute-value affine functions, has been extensively studied and applied to nonlinear circuit analysis and modeling [21], [42], [59], [62]. However, it is complete only for 1-D PWL functions. In higher dimensions it needs multilevel nestings of the absolute-value functions [58], [59]. The *lattice representation*, developed in [104], is a constructive way to generate min-max combinations of affine functions that provide a complete representation of continuous PWL functions in arbitrary dimensions. Combining the canonical with the lattice representations in [112] involved producing an equivalent representation as a difference of two convex max-affine functions.

A more recent approach is to focus on the class of *convex* PWL functions represented by a *maximum of affine functions*, which are essentially max-plus tropical polynomials as in (11), and use them for data fitting. Starting from early least-squares solutions [51], [54], some representative recent approaches to solve this *convex regression* problem include [45], [46], [53], [60], and [69]. In all these approaches, there is an iteration that alternates between partitioning the data domain and locally fitting affine functions (using least-squares or some linear optimization procedure) to update the local coefficients. For a known partition, the convex PWL function is formed as the max of the local affine fits. Then, a PWL function generates a new partition, which can be used to refit the affine functions and improve the estimate. As explained in [69], this iteration can be viewed as a Gauss-Newton algorithm to solve the above nonlinear least-squares problem. The rank  $K$  of the model can be increased until some error threshold is reached. Interesting and promising generalizations of the above max-affine representation for convex functions include works that use softmax instead of max, via the *log-sum-exp* models for convex and log-log convex data [16], [17], [53]. Other iterative approaches for convex PWL data fitting include [108]. Closer to our work is [55] that solves max-plus equations using least-squares and assumes that the slope parameters  $a_k$  in (11) are known. Reaching a local minimum of the  $\ell_2$  error norm for approximately solving max-plus equations was approached in [55] both via steepest descent (which was found computationally infeasible for large problems) and via Newton’s method with undershooting (which could not guarantee convergence to a local minimum). Very recently, it was

shown in [38] that, under certain assumptions, a carefully initialized alternating minimization algorithm converges linearly for max-affine regression. Finally, it was demonstrated in [20] how to efficiently solve large-scale convex regression—albeit with an unconstrained number of affine pieces. For additional references, we refer the reader to the bibliography in the above works.

Next, we focus on convex PWL regression via the max-affine model, which has a tropical interpretation, and propose a direct *noniterative* and *low-complexity* approach to estimate its parameters by using the optimal solutions of max-plus equations of Section IV-A.

## B. Optimal Fitting Max-Plus Tropical Lines and Planes

Given data  $(x_i, y_i) \in \mathbb{R}^2, i = 1, \dots, N$ , if we wish to fit a Euclidean line  $y = ax + b$  by minimizing the  $\ell_2$  error norm  $\|\mathbf{y} - a\mathbf{x} - b\|_2$ , where  $\mathbf{y} = [y_i]$  and  $\mathbf{x} = [x_i]$ , the optimal solution, that is, the *least-squares estimate (LSE)*, for the parameters  $a$  and  $b$  is

$$\hat{a}_{LS} = \frac{N \sum_i x_i y_i - (\sum_i x_i) (\sum_i y_i)}{N \sum_i (x_i)^2 - (\sum_i x_i)^2}$$

$$\hat{b}_{LS} = \frac{\sum_i (y_i - \hat{a}_{LS} x_i)}{N} \tag{83}$$

Suppose that, now, we wish to fit a max-plus tropical line  $p(x) = \max(a + x, b)$  by minimizing some  $\ell_p$  error norm. The equations to solve for finding the optimal parameter vector  $\mathbf{w} = [a, b]^T$  become

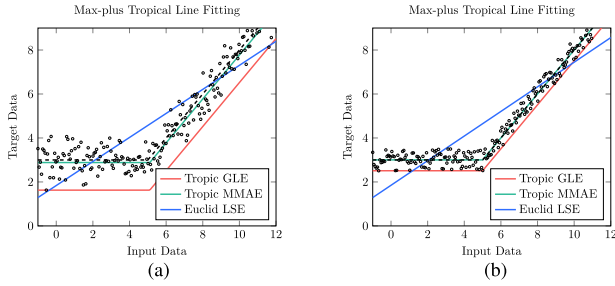
$$\underbrace{\begin{bmatrix} x_1 & 0 \\ \vdots & \vdots \\ x_N & 0 \end{bmatrix}}_{\mathbf{X}} \boxplus \underbrace{\begin{bmatrix} a \\ b \end{bmatrix}}_{\mathbf{w}} = \underbrace{\begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}}_{\mathbf{y}} \tag{84}$$

By Theorem 2, the optimal (min  $\ell_p$  error) subsolution is

$$\hat{\mathbf{w}} = \begin{bmatrix} \hat{a} \\ \hat{b} \end{bmatrix} = \mathbf{X}^* \boxplus \mathbf{y} = \begin{bmatrix} \bigwedge_i y_i - x_i \\ \bigwedge_i y_i \end{bmatrix} \tag{85}$$

where  $\mathbf{X}^* = -\mathbf{X}^T$  and  $\bigwedge_i = \bigwedge_{i=1}^N$ . This vector  $\hat{\mathbf{w}}$  yields (after max-plus “multiplication” with  $\mathbf{X}$ ) the *greatest lower estimate (GLE)* of the data  $\mathbf{y}$ . Thus, the above approach allows to optimally fit (w.r.t. any  $\ell_p$  error norm) max-plus tropical lines to arbitrary data from below. In addition, we can obtain the best (unconstrained) approximation with a tropical line that yields the smallest  $\ell_\infty$  error. This *minimum max absolute error (MMAE)* solution is, by Theorem 3

$$\tilde{\mathbf{w}} = \hat{\mathbf{w}} + \mu, \quad \mu = \frac{1}{2} \|\mathbf{X} \boxplus \hat{\mathbf{w}} - \mathbf{y}\|_\infty. \tag{86}$$



**Fig. 15.** (a) Optimal fitting via (85) or (86) of a max-plus tropical line  $y = \max(x - 2, 3)$  (shown in black dashed curve) to data from the line corrupted by additive i.i.d. Gaussian noise  $\sim \mathcal{N}(0, 0.25)$ . Blue line: Euclidean line fitting via least-squares. Red line: best subsolution (GLE). Green line: best unconstrained (MMAE) solution. (b) Same experiment as in (a) but with uniform noise  $\sim \text{Unif}[-0.5, 0.5]$ .

*Example 7:* Suppose that we have  $N = 200$  data observations  $(x_i, y_i)$  from the tropical line  $p(x) = \max(x - 2, 3)$ , where the 200 abscissae  $x_i$  were uniformly spaced in  $[-1, 12]$  and their corresponding values  $y_i = p(x_i) + \epsilon_i$  are contaminated with two different types of zero-mean noise i.i.d. random variables  $\epsilon_i$ , Gaussian noise  $\sim \mathcal{N}(0, 0.25)$ , and uniform noise  $\sim \text{Unif}[-0.5, 0.5]$ . Fig. 15 shows the two optimal solutions (85) and (86) for fitting a max-plus tropical line, superimposed with the least-squares Euclidean line fit. The parameter estimates and errors are shown in Table 2.

The above approach and tropical solution can also be extended to fitting planes. Specifically, we wish to fit a general max-plus tropical plane  $p(x, y)$

$$p(x, y) = \max(a + x, b + y, c) \quad (87)$$

to the given data  $(x_i, y_i, z_i) \in \mathbb{R}^3$ ,  $i = 1, \dots, N$ , where  $z_i = p(x_i, y_i) + \text{error}$ , by minimizing some  $\ell_p$  error norm. The equations to solve for finding the parameters  $w = [a, b, c]^T$  become

$$\underbrace{\begin{bmatrix} x_1 & y_1 & 0 \\ \vdots & \vdots & \vdots \\ x_N & y_N & 0 \end{bmatrix}}_{\mathbf{X}} \boxplus \underbrace{\begin{bmatrix} a \\ b \\ c \end{bmatrix}}_{\mathbf{w}} = \underbrace{\begin{bmatrix} z_1 \\ \vdots \\ z_N \end{bmatrix}}_{\mathbf{z}}. \quad (88)$$

By Theorem 2, the optimal subsolution, which yields approximations of  $z = [z_i]$  from below, is  $\hat{w} = \mathbf{X}^* \boxplus' \mathbf{z}$ .

**Table 2** Errors and Parameter Estimates for Optimal Fitting of a Max-Plus Tropical Line to Data Corrupted by Uniform Noise

Line fit Method	$\ \text{error}\ _{\text{RMS}}$	$\ \text{error}\ _{\infty}$	$\hat{a}$	$\hat{b}$
Tropical GLE	0.598	0.988	-2.492	2.509
Tropical MMAE	0.288	0.494	-1.998	3.003
Euclidean LSE	0.968	2.135	0.560	1.849

Hence

$$\underbrace{\begin{bmatrix} \hat{a} \\ \hat{b} \\ \hat{c} \end{bmatrix}}_{\hat{\mathbf{w}}} = \underbrace{\begin{bmatrix} -x_1 & \cdots & -x_N \\ -y_1 & \cdots & -y_N \\ 0 & \cdots & 0 \end{bmatrix}}_{\mathbf{X}^*} \boxplus' \underbrace{\begin{bmatrix} z_1 \\ \vdots \\ z_N \end{bmatrix}}_{\mathbf{z}} = \begin{bmatrix} \bigwedge_i z_i - x_i \\ \bigwedge_i z_i - y_i \\ \bigwedge_i z_i \end{bmatrix}. \quad (89)$$

Furthermore, the MMAE solution is given by  $\tilde{w} = \hat{w} + \mu$ , where  $\mu = \frac{1}{2} \|\mathbf{X} \boxplus \hat{w} - \mathbf{z}\|_{\infty}$ .

### C. Optimally Fitting Tropical Polynomial Curves and Surfaces

The above approach and solution can also be generalized to polynomial curves of higher degree and to multidimensional data  $x \in \mathbb{R}^d$ . We wish to fit a max-plus tropical polynomial

$$p(x) = \max\left(\mathbf{a}_1^T x + b_1, \dots, \mathbf{a}_K^T x + b_K\right) = \bigvee_{k=1}^K \mathbf{a}_k^T x + b_k \quad (90)$$

to the given data  $(x_i, y_i) \in \mathbb{R}^{d+1}$ ,  $i = 1, \dots, N$ , where  $y_i = p(x_i) + \text{error}$ , by minimizing some  $\ell_p$  error norm. The exact equations are

$$\underbrace{\begin{bmatrix} \mathbf{a}_1^T x_1 & \mathbf{a}_2^T x_1 & \cdots & \mathbf{a}_K^T x_1 \\ \mathbf{a}_1^T x_2 & \mathbf{a}_2^T x_2 & \cdots & \mathbf{a}_K^T x_2 \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{a}_1^T x_N & \mathbf{a}_2^T x_N & \cdots & \mathbf{a}_K^T x_N \end{bmatrix}}_{\mathbf{X}} \boxplus \underbrace{\begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_K \end{bmatrix}}_{\mathbf{w}} = \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}}_{\mathbf{y}}. \quad (91)$$

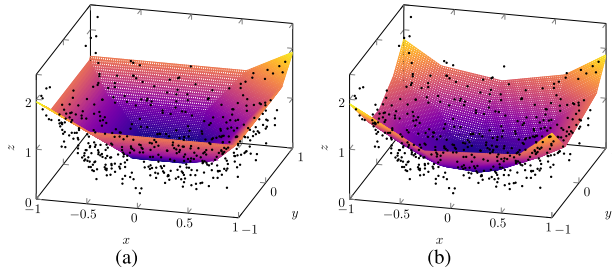
We assume that the slope vectors  $\mathbf{a}_k$  are given, and we optimize for the parameters  $\{b_k\}$ . By Theorem 2, the optimal subsolution for minimum  $\ell_p$  error is

$$\hat{w} = \begin{bmatrix} \hat{b}_1 \\ \vdots \\ \hat{b}_K \end{bmatrix} = \mathbf{X}^* \boxplus' \mathbf{y} = \begin{bmatrix} \bigwedge_{i=1}^N y_i - \mathbf{a}_1^T x_i \\ \vdots \\ \bigwedge_{i=1}^N y_i - \mathbf{a}_K^T x_i \end{bmatrix}. \quad (92)$$

Note that  $\mathbf{X} \boxplus \hat{w} \leq \mathbf{y}$ . Furthermore, by Theorem 3, the unconstrained solution that yields the minimum  $\ell_{\infty}$  error is

$$\tilde{w} = \mu + \hat{w}, \quad \mu = \frac{1}{2} \|\mathbf{X} \boxplus \hat{w} - \mathbf{y}\|_{\infty}. \quad (93)$$

There are two major categories of problems to which the above general tropical regression model can be applied:



**Fig. 16.** Two-dimensional tropical fitting using the optimal unconstrained (MMAE) approach to data from (94). (a) Tropic conic with known integer slopes [2-D conic ( $K = 11$ )]. (b) Slopes found via  $K$ -means on gradients ( $K = 25$ ).

First, if the slopes  $a_k$  are known for all  $K$  terms, then the above optimal solutions estimate the rest of the tropical model parameters (i.e., the intercepts  $b_k$ ) with a linear complexity  $O(dNK)$ . Second, in the case of *unknown slopes*, we can cluster the numerical gradients of the data using  $K$ -means, use the centroids of the  $K$  clusters as estimates of the slope vectors  $a_k$ , and then optimally solve for the intercepts  $b_k$ ; this approach was proposed in [76]. (An alternative heuristic approach is to discretize the range of the numerical gradients of the data and use as slopes all integer multiples of a slope step up to the desired accuracy.) In both approaches, if, for some  $k$ , we set or estimate the intercept  $b_k$  to be equal or close to  $-\infty$ , this essentially removes the corresponding line or hyperplane from the max-affine combination. Next, we illustrate both approaches via an example.

*Example 8:* Suppose that we are given  $N = 500$  data observations  $(x_i, y_i, z_i)$  as in Fig. 16 from the noisy paraboloid surface [45]

$$z = x^2 + y^2 + \epsilon \quad (94)$$

where  $\epsilon \sim \mathcal{N}(0, 0.25^2)$  is the zero-mean noise and the planar locations  $(x_i, y_i)$  of the data points were drawn as i.i.d. random variables  $\sim \text{Unif}[-1, 1]$ . First, as a tropical regression example with known slopes, let us fit to the above data a symmetric (with all positive and negative integer slopes in  $[-2, 2]$ ) max-plus tropical conic polynomial

$$p(x, y) = \bigvee_{0 \leq |k+\ell| \leq 2, k, \ell \geq 0} b_{k\ell} + kx + \ell y \quad (95)$$

where  $z_i = p(x_i, y_i) + \text{error}$ , by minimizing some  $\ell_p$  error norm. The equations to solve for finding the 11 parameters  $\mathbf{w} = [b_{0,-2}, \dots, b_{0,0}, b_{1,0}, b_{0,1}, b_{1,1}, b_{2,0}, b_{0,2}]^T$  become

$$\underbrace{\begin{bmatrix} -2y_1 & \cdots & 0 & x_1 & \cdots & 2x_1 & 2y_1 \\ \vdots & & \vdots & \vdots & & \vdots & \vdots \\ -2y_N & \cdots & 0 & x_N & \cdots & 2x_N & 2y_N \end{bmatrix}}_{\mathbf{X}} \boxplus \mathbf{w} = \mathbf{z}. \quad (96)$$

By Theorem 3, the optimal unconstrained solution for MMAE is  $\tilde{\mathbf{w}} = \mu + \hat{\mathbf{w}}$ , where  $\hat{\mathbf{w}} = \mathbf{X}^* \boxplus' \mathbf{z}$  and  $\mu$  is half the  $\ell_\infty$  error incurred by  $\hat{\mathbf{w}}$ . The resulting MMAE conic surface is shown in Fig. 16(a).

As a second approach, let us fit a tropical model of rank  $K$

$$p_K(x, y) = \max(a_1x + b_1y + c_1, \dots, a_Kx + b_Ky + c_K). \quad (97)$$

This consists of  $K$  planes of unknown slopes estimated by using  $K$ -means on the numerical gradients of the 2-D data, whereas the intercepts  $c_k$  are computed using the tropical fitting algorithm as in (93), which solves the unconstrained  $\ell_\infty$  problem. In this combined approach, the first step ( $K$ -means) is heuristic, yielding probably a local minimum for the slope estimation subproblem, whereas the second step (tropical regression for the intercepts) yields a global minimum optimally solving the unconstrained  $\ell_\infty$  problem. By varying  $K$ , we empirically find that even a small number of planes with adaptive slopes (e.g., see the case  $K = 25$  shown in Fig. 16) can yield both better approximations than the fixed slope case (but of course at a higher computational cost as discussed next) and generally good approximations, as seen by the errors in Table 3.

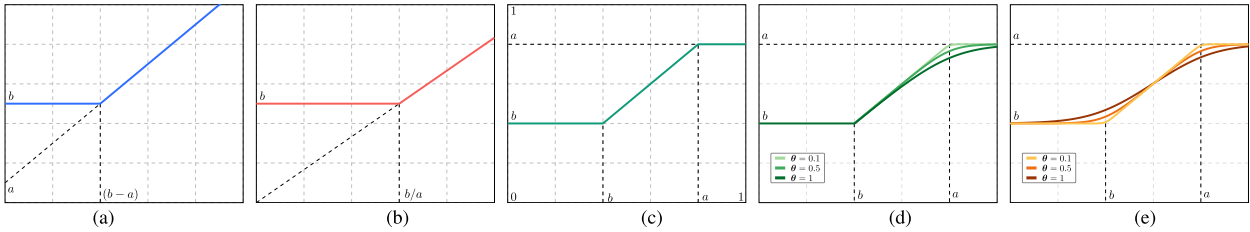
*Computational complexity:* Recent methods for convex PWL data fitting are commonly variations of iterative nonlinear least-squares algorithms. The standard least-squares estimator (LSE) [54] scales cubically in  $d$  and  $N$ , becoming intractable in the high-dimensional and/or large sample setting. The nonlinear least-squares problems in [53] and [69] are solved iteratively via a partitioning algorithm, with each iteration taking time  $O((d+1)^2N)$ ; however, these algorithms may not converge or fit the data poorly. This obstacle, largely due to nonconvexity, is empirically overcome by running multiple instances of the algorithm from random initializations. The convex adaptive partitioning (CAP) algorithm proposed in [45] solves a linear regression problem for each partition, leading to time complexity  $O(d(d+1)^2N \log(N) \log \log(N))$ .

In contrast, our algorithm for the case of unknown slopes has a complexity of  $O(dNKi_K)$ , where  $i_K$  is the number of  $K$ -means iterations. After computing the  $K$  centroids  $\mathbf{a}_k$ , it performs a single pass over the data to form (91) and solve for  $b_k$ , with total complexity  $O(dNK)$ . If the true slopes have some clustering structure,  $K$ -means will converge quickly, and the cost of our algorithm will

**Table 3** Errors for Optimal Tropical Fitting of the Function (94) Using 2-D Max-Plus Polynomials

$K$	GLE		MMAE	
	error <sub>RMS</sub>	error   <sub>∞</sub>	error <sub>RMS</sub>	error   <sub>∞</sub>
11 (conic)	0.6307	1.7049	0.4167	0.8524
10	0.6659	1.6022	0.3641	0.8011
25	0.5674	1.2779	0.3016	0.6389
50	0.5489	1.3068	0.3159	0.6534
100	0.5364	1.2828	0.3135	0.6414





**Fig. 17.** (a)–(d)  $\max$ - $\star$  tropical lines  $y = \max(a \star x, b)$ : (a) Max-plus line:  $y = \max(a + x, b)$ , (b) Max-times line:  $y = \max(a \cdot x, b)$ , (c) Max-min line:  $y = \max(a \wedge x, b)$ , and (d) Max-softmin line:  $y = \max(a \wedge_\theta x, b)$ . (e) Softmax-softmin line:  $s(x) = (a \wedge_\theta x) \nabla_\theta b$ . In (d) and (e), the parameter  $\theta$  varies.

be practically linear. As such, in nonpathological cases, we can assume that  $i_K$  may be treated as a small constant, thus improving on both the CAP algorithm and on the traditional LSE, resulting in a complexity  $O(dNK)$ . Finally, note that, for the case with known slopes, our algorithm has a complexity of  $\Theta(dNK)$ .

Tropical regression with unknown slopes  $a_k$  is equivalent to the problem of max-affine regression. When the number of terms is  $K \in [2, N/(d+1)]$ , a recent result [61] shows that the problem is, in fact, NP-hard for any choice of loss function  $\ell$  that satisfies  $\ell(x) = 0 \Leftrightarrow x = 0$ .

Some interesting connections arise when the number of pieces  $K = 2$ . On the one hand, if we fix  $a_1 = 0, b_1 = 0$ , we recover the problem of *ReLU regression*, which is known to be NP-hard [29], [70]. On the other hand, if we constrain  $-a_1 = a_2 \equiv a$ , and similarly,  $-b_1 = b_2 \equiv b$ , and furthermore denoting  $v = [a^T, b]^T$  and  $\tilde{x}_i = [x_i^T, 1]^T$ , tropical regression with squared  $\ell_2$  error becomes

$$\arg \min_{v \in \mathbb{R}^{d+1}} \sum_{i=1}^N (y_i - |v^T \tilde{x}_i|)^2$$

which is the problem of *phase retrieval*; the latter is also known to be NP-hard [33].

### VIII. TROPICAL ALGEBRA AND GEOMETRY ON WEIGHTED LATTICES

#### A. Generalized Tropical Lines and Planes

In the same way that weighted lattices generalize max-plus algebra and extend it to other types of clodum arithmetic, we can extend the basic objects of max-plus tropical geometry (i.e., tropical lines and planes) to other max- $\star$  geometric objects. For example, over a clodum  $(\mathcal{K}, \vee, \wedge, \star, \star')$ , we can generalize<sup>9</sup> max-plus tropical lines  $y = \max(a+x, b)$  as  $y = \max(a \star x, b)$ . Fig. 17(a)–(d) shows some generalized tropical lines where the  $\star$  operation is sum (+), product ( $\times$ ), min ( $\wedge$ ), and softmin ( $\wedge_\theta$ ). In the first three cases, the generalized tropical lines are PWL functions. However, in Fig. 17(d), a portion of the line is curving. To further illustrate this curving and create a symmetry between the max and min operations, we show

<sup>9</sup>The generalization in this section is done only w.r.t. to the “generalized multiplications” of the clodum, which becomes arbitrary scalar operations  $\star$  and  $\star'$  (instead of  $+$  and  $+$ ) that distribute over max and min, respectively. However, the “generalized additions” of the clodum remain the operations maximum and minimum.

in Fig. 17(e) a smooth function

$$s(x) = (a \wedge_\theta x) \nabla_\theta b = \theta \log \left[ \exp \left( -\log \left( e^{-a/\theta} + e^{-x/\theta} \right) + e^{b/\theta} \right) \right] \quad (98)$$

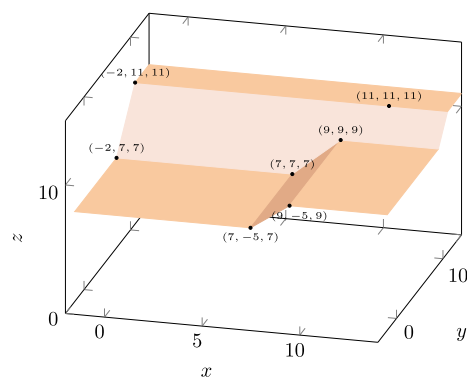
which goes beyond the max- $\star$  framework and is actually a softmax-softmin.

Similarly, we can generalize max-plus tropical planes  $z = \max(a+x, b+y, c)$  to max- $\star$  as  $z = \max(a \star x, b \star y, c)$ . Fig. 18 shows a max-min plane where  $\star = \min$ . This is an interesting geometrical polyhedral object that consists of portions of planes, either sloped or horizontal, at several levels.

Furthermore, we can generalize max-plus halfspaces (16) to max- $\star$  tropical halfspaces

$$\mathcal{T}(a, b) := \left\{ x \in \mathcal{K}^d : a^T \boxtimes \begin{bmatrix} x \\ e \end{bmatrix} \leq b^T \boxtimes \begin{bmatrix} x \\ e \end{bmatrix} \right\}. \quad (99)$$

Examples of max-plus tropical halfspaces are shown in Figs. 5 and 6. The slopes of their bounding line segments or faces are either zero or equal to 1. Max-product halfspaces can give boundaries that are PWL but have arbitrary slopes. Max-min halfspaces have PWL boundaries with more corner points or edges (see examples in Figs. 17(c) and 18). Finally, a totally different generalization results if we replace the “multiplication”  $\star$  in a generalized tropical line with the (log-sum-exp) softmin operation of (2), as shown in Fig. 17(d)–(e), in which case,



**Fig. 18.** Max-min plane  $z = \max(9 \wedge x, 11 \wedge y, 7)$ .

the line segments of a tropical line will become partially or totally smooth exponential curves.

### B. Generalized Tropical Regression

Suppose that we wish to fit a general max- $\star$  tropical plane

$$p(x, y) = \max(a \star x, b \star y, c) \tag{100}$$

to the given data  $(x_i, y_i, z_i) \in \mathbb{R}^3, i = 1, \dots, N$ , where  $z_i = p(x_i, y_i) + \text{error}$ , by minimizing some  $\ell_p$  error norm. The equations to solve for finding the optimal parameters  $w = [a, b, c]^T$  become

$$\underbrace{\begin{bmatrix} x_1 & y_1 & e \\ \vdots & \vdots & \vdots \\ x_N & y_N & e \end{bmatrix}}_{\mathbf{X}} \boxtimes \underbrace{\begin{bmatrix} a \\ b \\ c \end{bmatrix}}_w = \underbrace{\begin{bmatrix} z_1 \\ \vdots \\ z_N \end{bmatrix}}_z. \tag{101}$$

If we can accept subsolutions, which yield approximations of the given data from below, then, by Theorem 4, the optimal subsolution for any clodum arithmetic is

$$\hat{w} = \begin{bmatrix} \hat{a} \\ \hat{b} \\ \hat{c} \end{bmatrix} = \begin{bmatrix} \bigwedge_i \zeta(x_i, z_i) \\ \bigwedge_i \zeta(y_i, z_i) \\ \bigwedge_i \zeta(e, z_i) \end{bmatrix} \tag{102}$$

where  $\zeta$  is the scalar adjoint erosion (28) of  $\star$ . This vector  $\hat{w}$  yields (after max- $\star$  “multiplication” with  $\mathbf{X}$ ) the GLE of  $z$ . Next, we write in detail the solution for the three special cases where the scalar arithmetic is based either on the max-plus, or the max-times, or the max–min clodum:

$$\begin{bmatrix} \hat{a} \\ \hat{b} \\ \hat{c} \end{bmatrix}^T = \begin{cases} \left( \bigwedge_i z_i - x_i, \bigwedge_i z_i - y_i, \bigwedge_i z_i \right) \\ \left( \bigwedge_i z_i / x_i, \bigwedge_i z_i / y_i, \bigwedge_i z_i \right) \\ \left( \bigwedge_i (z_i \vee \mathbf{1}_{[z_i \geq x_i]}), \bigwedge_i (z_i \vee \mathbf{1}_{[z_i \geq y_i]}), \bigwedge_i z_i \right) \end{cases} \tag{103}$$

where  $\mathbf{1}_{[\cdot]}$  is 1 if the predicate  $[\cdot]$  is true and 0 otherwise. This approach allows to optimally fit (w.r.t. any  $\ell_p$  error norm) general max- $\star$  tropical planes to arbitrary data from below.

### IX. CONCLUSION AND FUTURE DIRECTIONS

Tropical geometry and max-plus algebra offer a rich collection of ideas and tools to model and solve problems in machine learning. In this work, we have surveyed the state

of the art and some recent progress in three areas: 1) DNNs with PWL activation functions; 2) probabilistic graphical models and algorithms for WFSTs; and 3) nonlinear regression with PWL functions. Furthermore, we have introduced extensions to general max algebras that allowed us to: 1) express the optimal solutions of several of the above problems as projections onto nonlinear vector spaces called weighted lattices and 2) generalize tropical geometrical objects. We conclude by outlining below some future research directions.

- 1) This work developed a Newton polytope representation of neural network layers, which was explored in the context of single-layer networks; due to the stability of Newton polytopes under addition and multiplication, one could try to derive similar representations for compositions of multiple layers. On the one hand, this may allow for refined empirical estimates on their complexity measured in terms of linear regions; on the other hand, further developing the aforementioned representation can pave the way for better network minimization methods, possibly combined with ideas from sparse regression.
- 2) It was briefly mentioned in Section II-E that tropical polytopes are more “economical” in their number of required parameters. This can introduce a whole new field of study, where there is an explicit characterization of Euclidean polytopes that can be exactly represented by a more efficient, tropical polytope while also providing quantifiable metrics for the relative gain.
- 3) The results of Section VI-B4 can be extended to more concrete benefits. While mainly algebraic, there is an extensive theory on the *reachability* and *robustness* [15] of the tropical matrices via their eigenvalue characterizations. Adapting these results to the WFST setting is nontrivial and a possible avenue for future study.
- 4) Regarding our work on tropical regression, we note that the max-affine representation is not limited to PWL functions only because we can represent any convex function as a supremum of a (possibly infinite) number of affine functions via the Fenchel–Legendre transform [32], [65], [94]. Closely related ideas are based on morphological slope transforms [30], [50], [71] that offer generalizations of this result to non-convex functions and approximate representations via adjunctions. ■

### Acknowledgment

The authors would like to thank the two anonymous reviewers for their constructive comments and G. Smyrnis for insightful discussions on tropical geometry and neural networks.

## REFERENCES

- [1] M. Akian, S. Gaubert, and A. Guterman, "Tropical polyhedra are equivalent to mean payoff games," *Int. J. Algebra Comput.*, vol. 22, no. 1, Feb. 2012, Art. no. 1250001.
- [2] M. Akian, S. Gaubert, V. Nićić, and I. Singer, "Best approximation in max-plus semimodules," *Linear Algebra Appl.*, vol. 435, no. 12, pp. 3261–3296, Dec. 2011.
- [3] M. Alfarra, A. Bibi, H. Hammoud, M. Gaafar, and B. Ghanem, "On the decision boundaries of neural networks: A tropical geometry perspective," 2020, *arXiv:2002.08838*. [Online]. Available: <http://arxiv.org/abs/2002.08838>
- [4] R. Arora, A. Basu, P. Mianjy, and A. Mukherjee, "Understanding deep neural networks with rectified linear units," in *Proc. Int. Conf. Learn. Represent.*, 2018.
- [5] D. Avis and K. Fukuda, "Reverse search for enumeration," *Discrete Appl. Math.*, vol. 65, nos. 1–3, pp. 21–46, Mar. 1996.
- [6] F. Baccelli, G. Cohen, G. J. Olsder, and J.-P. Quadrat, *Synchronization and Linearity: An Algebra for Discrete Event Systems*. Hoboken, NJ, USA: Wiley, 2001.
- [7] G. Birkhoff, *Lattice Theory*. Providence, RI, USA: AMS, 1967.
- [8] T. S. Blyth, *Lattices and Ordered Algebraic Structures*. New York, NY, USA: Springer-Verlag, 2005.
- [9] T. S. Blyth and M. F. Janowitz, *Residuation Theory*. New York, NY, USA: Pergamon, 1972.
- [10] A. Bora, A. Jalal, E. Price, and A. G. Dimakis, "Compressed sensing using generative models," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2017, pp. 537–546.
- [11] G. Borgefors, "Distance transformations in arbitrary dimensions," *Comput. Vis., Graph., Image Process.*, vol. 27, no. 3, pp. 321–345, Sep. 1984.
- [12] S. Boyd, S.-J. Kim, L. Vandenberghe, and A. Hassibi, "A tutorial on geometric programming," *Optim. Eng.*, vol. 8, no. 1, pp. 67–127, May 2007.
- [13] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge, U.K.: Cambridge Univ. Press, 2004.
- [14] R. W. Brockett and P. Maragos, "Evolution equations for continuous-scale morphological filtering," *IEEE Trans. Signal Process.*, vol. 42, no. 12, pp. 3377–3386, Dec. 1994.
- [15] B. Butković, *Max-Linear Systems: Theory and Algorithms*. London, U.K.: Springer-Verlag, 2010.
- [16] G. C. Calafiore, S. Gaubert, and C. Possieri, "Log-sum-exp neural networks and posynomial models for convex and log-log-convex data," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 30, no. 5, pp. 1–12, Mar. 2020.
- [17] G. C. Calafiore, S. Gaubert, and C. Possieri, "A universal approximation result for difference of log-sum-exp neural networks," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 31, no. 12, pp. 5603–5612, Dec. 2020.
- [18] V. Charisopoulos et al., "Morphological perceptrons: Geometry and training algorithms," in *Proc. Int. Symp. Math. Morphol.*, in Lecture Notes in Computer Science, vol. 10225, J. Angulo, Ed. Cham, Switzerland: Springer, 2017, pp. 3–15.
- [19] V. Charisopoulos and P. Maragos, "A tropical approach to neural networks with piecewise linear activations," 2018, *arXiv:1805.08749*. [Online]. Available: <http://arxiv.org/abs/1805.08749>
- [20] W. Chen and R. Mazumder, "Multivariate convex regression at scale," 2020, *arXiv:2005.11588*. [Online]. Available: <http://arxiv.org/abs/2005.11588>
- [21] L. O. Chua and A.-C. Deng, "Canonical piecewise-linear representation," *IEEE Trans. Circuits Syst.*, vol. 35, no. 1, pp. 101–111, Jan. 1988.
- [22] G. Cohen, D. Dubois, J. Quadrat, and M. Viot, "A linear-system-theoretic view of discrete-event processes and its use for performance evaluation in manufacturing," *IEEE Trans. Autom. Control*, vol. 30, no. 3, pp. 210–220, Mar. 1985.
- [23] G. Cohen, S. Gaubert, and J.-P. Quadrat, "Duality and separation theorems in idempotent semimodules," *Linear Algebra Appl.*, vol. 379, pp. 395–422, Mar. 2004.
- [24] M. A. Cueto, J. Morton, and B. Sturmfels, "Geometry of the restricted Boltzmann machine," in *Algebraic Methods in Statistics and Probability II (Contemporary Mathematics)*, vol. 516, M. A. G. Viana and H. P. Wynn, Eds. Providence, RI, USA: AMS, 2010, pp. 135–153.
- [25] R. Cuninghame-Green, *Minimax Algebra*. New York, NY, USA: Springer-Verlag, 1979.
- [26] R. A. Cuninghame-Green, "Projections in minimax algebra," *Math. Program.*, vol. 10, no. 1, pp. 111–123, Dec. 1976.
- [27] R. A. Cuninghame-Green and K. Cechlárová, "Residuation in fuzzy algebra and some applications," *Fuzzy Sets Syst.*, vol. 71, no. 2, pp. 227–239, Apr. 1995.
- [28] A. Damle and Y. Sun, "A geometric approach to archetypal analysis and nonnegative matrix factorization," *Technometrics*, vol. 59, no. 3, pp. 361–370, Jul. 2017.
- [29] S. S. Dey, G. Wang, and Y. Xie, "An approximation algorithm for training one-node ReLU neural network," 2018, *arXiv:1810.03592*. [Online]. Available: <http://arxiv.org/abs/1810.03592>
- [30] L. Dorst and R. Van den Boomgaard, "Morphological signal processing and the slope transform," *Signal Process.*, vol. 38, no. 1, pp. 79–98, Jul. 1994.
- [31] Y. Eldar and G. Kutyniok, *Compressed Sensing: Theory and Applications*. Cambridge, U.K.: Cambridge Univ. Press, 2012.
- [32] W. Fenichel, "On conjugate convex functions," *Can. J. Math.*, vol. 1, no. 1, pp. 73–77, Feb. 1949.
- [33] M. Fickus, D. G. Mixon, A. A. Nelson, and Y. Wang, "Phase retrieval from very few measurements," *Linear Algebra Appl.*, vol. 449, pp. 475–499, May 2014.
- [34] G. Franchi, A. Fehri, and A. Yao, "Deep morphological networks," *Pattern Recognit.*, vol. 102, Jun. 2020, Art. no. 107246.
- [35] K. Fukuda, "From the zonotope construction to the Minkowski addition of convex polytopes," *J. Symbolic Comput.*, vol. 38, no. 4, pp. 1261–1272, Oct. 2004.
- [36] S. Gaubert and R. D. Katz, "Minimal half-spaces and external representation of tropical polyhedra," *J. Algebr. Combinatorics*, vol. 33, no. 3, pp. 325–348, May 2011.
- [37] S. Gaubert and M. Plus, "Methods and applications of (max, +) linear algebra," in *Proc. Symp. Theor. Aspects Comput. Sci.*, 1997, pp. 261–282.
- [38] A. Ghosh, A. Pananjady, A. Guntuboyina, and K. Ramchandran, "Max-affine regression: Provable, tractable, and near-optimal statistical estimation," 2019, *arXiv:1906.09255*. [Online]. Available: <http://arxiv.org/abs/1906.09255>
- [39] M. Gondran and M. Minoux, *Graphs, Droids and Semirings: New Models and Algorithms*. New York, NY, USA: Springer, 2008.
- [40] I. Goodfellow, D. Warde-Farley, M. Mirza, A. Courville, and Y. Bengio, "Maxout networks," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2013, pp. 1319–1327.
- [41] P. Gritzmann and B. Sturmfels, "Minkowski addition of polytopes: Computational complexity and applications to Gröbner bases," *SIAM J. Discrete Math.*, vol. 6, no. 2, pp. 246–269, May 1993.
- [42] C. Güzelis and I. C. Göknaar, "A canonical representation for piecewise-affine maps and its applications to circuit analysis," *IEEE Trans. Circuits Syst.*, vol. 38, no. 11, pp. 1342–1354, Nov. 1991.
- [43] S. Han, J. Pool, J. Tran, and W. Dally, "Learning both weights and connections for efficient neural network," in *Proc. Adv. Neural Inf. Process. Syst.*, 2015.
- [44] P. Hand and V. Voroninski, "Global guarantees for enforcing deep generative priors by empirical risk," in *Proc. Conf. Learn. Theory*, 2018, pp. 970–978.
- [45] L. A. Hannah and D. B. Dunson, "Multivariate convex regression with adaptive partitioning," 2011, *arXiv:1105.1924*. [Online]. Available: <http://arxiv.org/abs/1105.1924>
- [46] L. A. Hannah and D. B. Dunson, "Ensemble methods for convex regression with applications to geometric programming based circuit design," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2012.
- [47] Y. He, J. Lin, Z. Liu, H. Wang, L.-J. Li, and S. Han, "AMC: AutoML for model compression and acceleration on mobile devices," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, 2018, pp. 784–800.
- [48] B. Heidergott, G. J. Olsder, and J. van der Woude, *Max Plus at Work: Modeling and Analysis of Synchronized Systems: A Course on Max-Plus Algebra and Its Applications*. Princeton, NJ, USA: Princeton Univ. Press, 2006.
- [49] H. Heijmans, *Morphological Image Operators*. New York, NY, USA: Academic, 1994.
- [50] H. J. A. M. Heijmans and P. Maragos, "Lattice calculus of the morphological slope transform," *Signal Process.*, vol. 59, no. 1, pp. 17–42, May 1997.
- [51] C. Hildreth, "Point estimates of ordinates of concave functions," *J. Amer. Stat. Assoc.*, vol. 49, no. 267, pp. 598–619, Sep. 1954.
- [52] G. E. Hinton and R. R. Salakhutdinov, "Reducing the dimensionality of data with neural networks," *Science*, vol. 313, no. 5786, pp. 504–507, Jul. 2006.
- [53] W. Hoberg, P. Kirschen, and P. Abbeel, "Data fitting with geometric-programming-compatible softmax functions," *Optim. Eng.*, vol. 17, no. 4, pp. 897–918, Dec. 2016.
- [54] C. A. Holloway, "On the estimation of convex functions," *Oper. Res.*, vol. 27, no. 2, pp. 401–407, 1979.
- [55] J. Hook, "Linear regression over the max-plus semiring: Algorithms and applications," 2017, *arXiv:1712.03499*. [Online]. Available: <http://arxiv.org/abs/1712.03499>
- [56] T. Hori and A. Nakamura, *Speech Recognition Algorithms Using Weighted Finite-State Transducers*. San Rafael, CA, USA: Morgan & Claypool, 2013.
- [57] J.-E. Perin, "Tropical semirings," in *Idempotency*, J. Gunawardena, Ed. Cambridge, U.K.: Cambridge Univ. Press, 1998, pp. 50–69.
- [58] P. Julian, "The complete canonical piecewise-linear representation: Functional form minimal degenerate intersections," *IEEE Trans. Circuits Syst. I, Fundam. Theory Appl.*, vol. 50, no. 3, pp. 387–396, Mar. 2003.
- [59] C. Kahlert and L. O. Chua, "The complete canonical piecewise-linear representation. I. The geometry of the domain space," *IEEE Trans. Circuits Syst. I, Fundam. Theory Appl.*, vol. 39, no. 3, pp. 222–236, Mar. 1992.
- [60] J. Kim, L. Vandenberghe, and C.-K.-K. Yang, "Convex piecewise-linear modeling method for circuit optimization via geometric programming," *IEEE Trans. Comput.-Aided Design Integr. Circuits Syst.*, vol. 29, no. 11, pp. 1823–1827, Nov. 2010.
- [61] F. Lauer, "On the complexity of piecewise affine system identification," *Automatica*, vol. 62, pp. 148–153, Dec. 2015.
- [62] J.-N. Lin, H.-Q. Xu, and R. Unbehauen, "A generalization of canonical piecewise-linear functions," *IEEE Trans. Circuits Syst. I, Fundam. Theory Appl.*, vol. 41, no. 4, pp. 345–347, Apr. 1994.
- [63] G. L. Litvinov, "Maslov dequantization, idempotent and tropical mathematics: A brief introduction," *J. Math. Sci.*, vol. 140, no. 3, pp. 426–444, Jan. 2007.
- [64] G. L. Litvinov, V. P. Maslov, and G. B. Shpiz, "Idempotent functional analysis: An algebraic approach," *Math. Notes*, vol. 69, no. 5, pp. 696–729, 2001.
- [65] Y. Lucet, "What shape is your conjugate? A survey of computational convex analysis and its applications," *SIAM J. Optim.*, vol. 20, no. 1,

- pp. 216–250, Jan. 2009.
- [66] J.-H. Luo, J. Wu, and W. Lin, “ThiNet: A filter level pruning method for deep neural network compression,” in *Proc. IEEE Int. Conf. Comput. Vis.*, Oct. 2017, pp. 5058–5066.
- [67] A. Maas, A. Hannun, and A. Ng, “Rectifier nonlinearities improve neural network acoustic models,” in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2013, p. 3.
- [68] D. Maclagan and B. Sturmfels, *Introduction to Tropical Geometry*. Providence, RI, USA: AMS, 2015.
- [69] A. Magnani and S. P. Boyd, “Convex piecewise-linear fitting,” *Optim. Eng.*, vol. 10, no. 1, pp. 1–17, Mar. 2009.
- [70] P. Manurangsi and D. Reichman, “The computational complexity of training ReLU(s),” 2018, [arXiv:1810.04207](https://arxiv.org/abs/1810.04207). [Online]. Available: <http://arxiv.org/abs/1810.04207>
- [71] P. Maragos, “Morphological systems: Slope transforms and max-min difference and differential equations,” *Signal Process.*, vol. 38, no. 1, pp. 57–77, Jul. 1994.
- [72] P. Maragos, “Lattice image processing: A unification of morphological and fuzzy algebraic systems,” *J. Math. Imag. Vis.*, vol. 22, nos. 2–3, pp. 333–353, May 2005.
- [73] P. Maragos, “Morphological filtering for image enhancement and feature detection,” in *Image and Video Processing Handbook*, A. Bovik, Ed., 2nd ed. Amsterdam, The Netherlands: Elsevier, 2005, pp. 135–156.
- [74] P. Maragos, “Dynamical systems on weighted lattices: General theory,” *Math. Control, Signals, Syst.*, vol. 29, no. 4, pp. 1–49, Dec. 2017.
- [75] P. Maragos, “Tropical geometry, mathematical morphology and weighted lattices,” in *Proc. Int. Symp. Math. Morphol.*, in Lecture Notes in Computer Science, vol. 11564, B. Burgeth, Ed. Cham, Switzerland: Springer, 2019, pp. 3–15.
- [76] P. Maragos and E. Theodosis, “Multivariate tropical regression and piecewise-linear surface fitting,” in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, May 2020, pp. 3822–3826.
- [77] V. P. Maslov, “On a new superposition principle for optimization problems,” *Russian Math. Surv.*, vol. 42, no. 3, pp. 39–48, 1987.
- [78] W. M. McEneaney, *Max-Plus Methods for Nonlinear Control and Estimation*. Cambridge, MA, USA: Birkhäuser, 2006.
- [79] F. Meyer, *Topographic Tools for Filtering and Segmentation 1 & 2*. Hoboken, NJ, USA: Wiley, 2019.
- [80] G. Mikhalkin, “Enumerative tropical algebraic geometry in  $\mathbb{R}^2$ ,” *J. Amer. Math. Soc.*, vol. 18, no. 2, pp. 313–377, 2005.
- [81] M. Mohri, “Weighted automata algorithms,” in *Handbook Weighted Automata*. Berlin, Germany: Springer, 2009, pp. 213–254.
- [82] M. Mohri, F. Pereira, and M. Riley, “Weighted finite-state transducers in speech recognition,” *Comput. Speech Lang.*, vol. 16, no. 1, pp. 69–88, Jan. 2002.
- [83] R. Mondal, S. Sundar Mukherjee, S. Santra, and B. Chanda, “Dense morphological network: An universal function approximator,” 2019, [arXiv:1901.00109](https://arxiv.org/abs/1901.00109). [Online]. Available: <http://arxiv.org/abs/1901.00109>
- [84] G. F. Montufar, R. Pascanu, K. Cho, and Y. Bengio, “On the number of linear regions of deep neural networks,” in *Proc. Adv. Neural Inf. Process. Syst.*, 2014.
- [85] J.-J. Moreau, “Inf-convolution, Sous-additivité, Convexité des Fonctions Numériques,” *J. Mathématiques Pures Appliquées*, vol. 49, pp. 109–154, 1970.
- [86] L. Pachter and B. Sturmfels, “Tropical geometry of statistical models,” *Proc. Nat. Acad. Sci. USA*, vol. 101, no. 46, pp. 16132–16137, Nov. 2004.
- [87] L. Pachter and B. Sturmfels, “Parametric inference for biological sequence analysis,” *Proc. Nat. Acad. Sci. USA*, vol. 101, no. 46, pp. 16138–16143, Nov. 2004.
- [88] R. Pascanu, G. Montufar, and Y. Bengio, “On the number of response regions of deep feed forward networks with piece-wise linear activations,” 2013, [arXiv:1312.6098](https://arxiv.org/abs/1312.6098). [Online]. Available: <http://arxiv.org/abs/1312.6098>
- [89] L. F. C. Pessoa and P. Maragos, “Neural networks with hybrid morphological/rank/linear nodes: A unifying framework with applications to handwritten character recognition,” *Pattern Recognit.*, vol. 33, no. 6, pp. 945–960, Jun. 2000.
- [90] M. Raghun, B. Poole, J. Kleinberg, S. Ganguli, and J. Sohl-Dickstein, “On the expressive power of deep neural networks,” in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2017, pp. 2847–2854.
- [91] G. X. Ritter and P. Sussner, “An introduction to morphological neural networks,” in *Proc. 13th Int. Conf. Pattern Recognit.*, 1996.
- [92] G. X. Ritter, P. Sussner, and J. L. Diza-de-Leon, “Morphological associative memories,” *IEEE Trans. Neural Netw.*, vol. 9, no. 2, pp. 281–293, Mar. 1998.
- [93] G. X. Ritter and G. Urcid, “Lattice algebra approach to single-neuron computation,” *IEEE Trans. Neural Netw.*, vol. 14, no. 2, pp. 282–295, Mar. 2003.
- [94] R. T. Rockafellar, *Convex Analysis*. Princeton, NJ, USA: Princeton Univ. Press, 1970.
- [95] J. Serra, *Image Analysis and Mathematical Morphology*. New York, NY, USA: Academic, 1982.
- [96] J. Serra, Ed., *Image Analysis and Mathematical Morphology*, vol. 2. New York, NY, USA: Academic, 1988.
- [97] T. Serra and S. Ramalingam, “Empirical bounds on linear regions of deep rectifier networks,” in *Proc. AAAI Conf. Artif. Intell.*, 2020, pp. 5628–5635.
- [98] T. Serra, C. Tjandraatmadja, and S. Ramalingam, “Bounding and counting linear regions of deep neural networks,” in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2018, pp. 4558–4566.
- [99] I. Simon, “On semigroups of matrices over the tropical semiring,” *RAIRO Theor. Inform. Appl.*, vol. 28, nos. 3–4, pp. 277–294, 1994.
- [100] G. Smyrnis, P. Maragos, and G. Retinas, “Maxpolynomial division with application to neural network simplification,” in *Proc. ICASSP IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, May 2020, pp. 4192–4196.
- [101] G. Smyrnis and P. Maragos, “Multiclass neural network minimization via tropical newton polytope approximation,” in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2020, pp. 9068–9077.
- [102] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, “Dropout: A simple way to prevent neural networks from overfitting,” *J. Mach. Learn. Res.*, vol. 1, no. 15, pp. 1929–1958, 2014.
- [103] P. Sussner and E. L. Esmi, “Morphological perceptrons with competitive learning: Lattice-theoretical framework and constructive learning algorithm,” *Inf. Sci.*, vol. 181, no. 10, pp. 1929–1950, May 2011.
- [104] J. M. Tarela and M. V. Martínez, “Region configurations for realizability of lattice piecewise-linear models,” *Math. Comput. Model.*, vol. 30, nos. 11–12, pp. 17–27, Dec. 1999.
- [105] M. Telgarsky, “Benefits of depth in neural networks,” in *Proc. Conf. Learn. Theory*, 2016, pp. 1517–1539.
- [106] E. Theodosis and P. Maragos, “Analysis of the viterbi algorithm using tropical algebra and geometry,” in *Proc. IEEE Int. Workshop Signal Process. Adv. Wireless Commun.*, Jun. 2018, pp. 1–5.
- [107] E. Theodosis and P. Maragos, “Tropical modeling of weighted transducer algorithms on graphs,” in *Proc. ICASSP IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, May 2019, pp. 8653–8657.
- [108] A. Toriello and J. P. Vielma, “Fitting piecewise linear continuous functions,” *Eur. J. Oper. Res.*, vol. 219, no. 1, pp. 86–95, May 2012.
- [109] A. Tsiamis and P. Maragos, “Sparsity in max-plus algebra,” *Discrete Events Dyn. Syst.*, vol. 29, no. 2, pp. 163–189, 2019.
- [110] T. van den Boom and B. D. Schutter, “Modeling and control of switching max-plus-linear systems with random and deterministic switching,” *Discrete Event Dyn. Syst.*, vol. 22, no. 2, pp. 293–332, 2012.
- [111] O. Viro, “Dequantization of real algebraic geometry on logarithmic paper,” in *Proc. Eur. Congr. Math.*, 2001, pp. 135–146.
- [112] S. Wang, “General constructive representations for continuous piecewise-linear functions,” *IEEE Trans. Circuits Syst. I, Reg. Papers*, vol. 51, no. 9, pp. 1889–1896, Sep. 2004.
- [113] P.-F. Yang and P. Maragos, “Min-max classifiers: Learnability, design and application,” *Pattern Recognit.*, vol. 28, no. 6, pp. 879–899, Jun. 1995.
- [114] L. Zhang, G. Naitzat, and L.-H. Lim, “Tropical geometry of deep neural networks,” in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2018, pp. 5824–5832.
- [115] Y. Zhang et al., “Max-plus operators applied to filter selection and model pruning in neural networks,” in *Proc. Int. Symp. Math. Morphol.*, in Lecture Notes in Computer Science, B. Burgeth, Eds., vol. 11564. Cham, Switzerland: Springer, 2019, pp. 310–322.
- [116] U. Zimmermann, *Linear and Combinatorial Optimization in Ordered Algebraic Structures*. Amsterdam, The Netherlands: North Holland, 1981.

## ABOUT THE AUTHORS

**Petros Maragos** (Fellow, IEEE) received the M.Eng. degree in electrical engineering from the National Technical University of Athens (NTUA), Athens, Greece, in 1980, and the M.Sc. and Ph.D. degrees from the Georgia Institute of Technology (Georgia Tech), Atlanta, GA, USA, in 1982 and 1985, respectively.

In 1985, he joined the Faculty of the Division of Applied Sciences, Harvard University, Boston, MA, USA, where he worked for eight years as a Professor of electrical engineering with the Harvard Robotics Laboratory. In 1993, he joined the Faculty of the School of Electrical and Computer Engineering (ECE), Georgia Tech, where he was with the Center for Signal and Image Processing. From 1996 to 1998, he had a joint appointment



as the Director of Research at the Institute of Language and Speech Processing, Athens. Since 1999, he has been a Professor with the School of ECE, NTUA, where he is currently the Director of the Intelligent Robotics and Automation Laboratory. He has held visiting positions at the Massachusetts Institute of Technology, Cambridge, MA, USA, in 2012, and the University of Pennsylvania, Philadelphia, PA, USA, in 2016. His research and teaching interests include signal processing, systems theory, machine learning, image processing and computer vision, speech & language processing, and robotics. In these areas, he has published numerous articles and book chapters and has also coedited three Springer research books: one on multimodal processing and two on shape analysis.

Dr. Maragos has served as a member of the Greek National Council for Research and Technology. He has also served as a member of three IEEE SPS technical committees and, recently, the IEEE SPS



Education Board. For his research contributions, he was elected as a Fellow of IEEE in 1995 and EURASIP in 2010. He was a recipient or co-recipient of several awards for his academic work, including the 1987–1992 US NSF Presidential Young Investigator Award, the 1988 IEEE ASSP Young Author Best Paper Award, the 1994 IEEE SPS Senior Best Paper Award, the 1995 IEEE W.R.G. Baker Prize for the most outstanding original paper, the 1996 Pattern Recognition Society's Honorable Mention Best Paper Award, the Best Paper Award from the CVPR-2011 Workshop on Gesture Recognition, and the 2007 EURASIP Technical Achievement Award. He has served as an Associate Editor for IEEE TRANSACTIONS ON ASSP and IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE. He has served as an editorial board member and a guest editor for several journals on signal processing, image analysis, and vision. He has served as a co-organizer of several conferences and workshops, including, recently, the 2017 European Signal Processing Conference (EUSIPCO) and the 2023 International Conference on Acoustics, Speech, and Signal Processing (ICASSP) as the General Chair. He was elected as the IEEE SPS Distinguished Lecturer for the term 2017–2018.

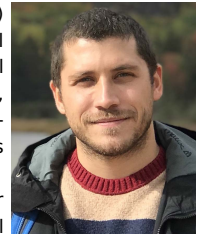
**Vasileios Charisopoulos** received the M.Eng. degree in electrical and computer engineering from the National Technical University of Athens, Athens, Greece, in 2017, completing his thesis under the supervision of Professor Petros Maragos. He is currently working toward the Ph.D. degree at the School of Operations Research & Information Engineering, Cornell University, Ithaca, NY, USA, where he is advised by Prof. Damek Davis.



From May to August 2017, he worked as a Researcher in INRIA Paris-Saclay, Palaiseau, France, hosted by Prof. Xavier Allamigeon. He is broadly interested in the mathematics of data science and large-scale numerical linear algebra.

Mr. Charisopoulos was a recipient of the Andreas G. Leventis Fellowship, the NSF-Dagstuhl Travel Award for junior researchers, and a Cornell University Fellowship for first-year Ph.D. studies.

**Emmanouil Theodosis** (Member, IEEE) received the M.Eng. degree in electrical and computer engineering from the National Technical University of Athens, Athens, Greece, in 2018, conducting his master thesis under the supervision of Professor Petros Maragos.



He worked as a Research Assistant for the following year at the National Technical University of Athens. In the fall of 2019, he joined the Computer Science Department, School of Engineering and Applied Sciences, Harvard University, Boston, MA, USA, as a Ph.D. Student, where he is advised by Prof. Demba Ba. His research interests are mainly theoretical, revolving around model-based deep learning, algebraic geometry, optimization on nonlinear surfaces, and the theory of machine learning. He has coauthored various publications in these fields, including a book chapter and papers in leading conferences in signal processing.

Mr. Theodosis was a Student Member of the Signal Processing Society. He received a Certificate of Distinction and Excellence in Teaching from Harvard University, the Gerondelis Foundation Scholarship, the Robert L. Wallace Prize Fellowship for his first-year Ph.D. studies, and a scholarship from Eurobank EFG for his undergraduate studies. He was an Invited Reviewer for EUSIPCO 2020 and 2021.