
ADVANCES IN THE TRAINING, PRUNING AND ENFORCEMENT OF SHAPE CONSTRAINTS OF MORPHOLOGICAL NEURAL NETWORKS USING TROPICAL ALGEBRA

Dimitriadis Nikolaos

École Polytechnique Fédérale de Lausanne
Lausanne Switzerland
nikolaos.dimitriadis@epfl.ch

Petros Maragos

National Technical University of Athens
Athens, Greece
maragos@cs.ntua.gr

ABSTRACT

In this paper we study an emerging class of neural networks based on the morphological operators of dilation and erosion. We explore these networks mathematically from a tropical geometry perspective as well as mathematical morphology. Our contributions are threefold. First, we examine the training of morphological networks via Difference-of-Convex programming methods and extend a binary morphological classifier to multiclass tasks. Second, we focus on the sparsity of dense morphological networks trained via gradient descent algorithms and compare their performance to their linear counterparts under heavy pruning, showing that the morphological networks cope far better and are characterized with superior compression capabilities. Our approach incorporates the effect of the training optimizer used and offers quantitative and qualitative explanations. Finally, we study how the architectural structure of a morphological network can affect shape constraints, focusing on monotonicity. Via Maslov Dequantization, we obtain a softened version of a known architecture and show how this approach can improve training convergence and performance.

Keywords Tropical Geometry · Mathematical Morphology · Morphological Neural Networks · Monotonicity · Sparsity · Pruning · Maslov Dequantization

1 Introduction

In the past decade, the field of neural networks has garnered research interest in the machine learning community, paving the way for the formation of a novel field called *Deep Learning*. The cell of the models is the neuron, introduced by Rosenblatt. The neuron mimics the transformations of data performed in biological organisms. In mathematical terms, it consists of a multiply-accumulate scheme that is succeeded by a nonlinearity, called *activation function*. An alternative lies in morphology-based models.

In morphological neural networks the operations of addition and multiplication of the aforementioned multiply-accumulate scheme are replaced by maximum (or minimum) and addition, respectively. This process is called *tropicalization* and yields a path towards tropical mathematics. An important aspect of the operator change is the lack of need for an activation function, since maximum (or minimum) are inherently nonlinear operations. Networks with this modified neuron have been studied in the context of neural networks [1, 2, 3, 4, 5]. Recently, the field of tropical geometry has been linked with this class of morphological networks [6, 7, 8, 9, 9]. Tropical geometry studies *piecewise linear (PWL)* surfaces whose arithmetic is governed by a tropical semiring, where ordinary addition is replaced by the maximum or minimum and ordinary multiplication is replaced by ordinary addition. We refer to these algebraic structures as $(\max, +)$ and $(\min, +)$ semirings, respectively. These two semirings are dual and linked via the isomorphism $\phi(x) = -x$.

* This work was performed when N.Dimitriadis was at the National Technical University of Athens.

In this paper, we study morphological networks via mathematical constructs stemming from mathematical morphology, lattice theory and tropical algebra. Our contributions are multifaceted and can be summarized as follows:

- An alternative training process for morphological networks for general multiclass classification tasks is presented. By formulating the classification problem as an instance of a Difference-of-Convex optimization problem, the training can be performed with analytical and robust methods *not* based on stochastic optimization.
- We explore the compression abilities of morphological networks and draw favorable conclusions in comparison with their linear counterparts. We show that morphological networks encode information with fewer parameters.
- We present a method to enforce monotonicity on the output of a morphological network using smooth operator approximations and show how the softened versions improve performance and training convergence.

The structure of the paper is the following. Section 2 reviews past work that motivate our findings. In Section 3 we introduce some mathematical background concepts that accompany the ideas of the next sections and highlight the connections of mathematical morphology and tropical geometry. In Section 4 we introduce an alternative method of morphological neural network training based on a heuristic called *Convex-Concave Procedure* and proposed by Charisopoulos and Maragos [6] and propose an extension to general (i.e. multiclass) classification problems. In Section 5 we study denser morphological networks and shift our focus from accuracy to sparsity. We use a pruning method to evaluate how the removal of parameters affects performance on both morphological networks and their linear counterparts concluding with the former’s superiority in that regard. Section 6 explores regression with monotonicity constraints and focuses on neural network methods, exploring the smoothening of the "hard" morphological operators min and max to achieve better performance and simultaneously alleviate some training issues. Section 7 includes concluding remarks.

2 Related Work

Davidson and Hummer introduced the morphological neuron [10], aiming at learning morphological elements such as erosion and dilation. These efforts were intensified and led to more general models, such as a simple network with a hidden layer tasked to solve binary classification problems [1]. These models use different operators than addition and multiplication and are rooted in Mathematical Morphology and Lattice Theory [11]. The model proposed by Ritter and Sussner [1] constructs a decision boundary parallel to the axes, a limitation that has been addressed in two major ways. First, the architecture was extended to include a second hidden layer [3]. This modification allows the network to learn multiple axis-parallel decision boundaries and, thus, tackle more general problems. Second, Barmoutis and Ritter resolved this limitation by rotating the hyperplanes [12].

Ritter et al. introduced the term *morphological networks* by replacing addition and multiplication with maximum and addition [2]. They focused on the storage and computing capabilities of associative memories based on the morphological neuron. Ritter and Urcid link the morphological neuron with biological processes, introduce its dendritic structure and show that morphological perceptrons with a single hidden layer can approximate any high dimensional compact region within an arbitrary degree of accuracy [13]. Furthermore, Sussner and Esmi study such networks in the prism of competitive learning, where the operator argmax is employed to create a winner-take-all strategy at the output layer [14].

Yang and Maragos introduced the class of *min-max* classifiers which consist of two layers, one with max terms and one with min terms. The models are trained within the context of Probably Approximately Correct (PAC) learning and can be considered as a generalization of Boolean functions based on Lattice Theory [5]. Pessoa and Maragos proposed a hybrid neuron that combines morphological and linear terms [4]. A common theme in the morphological network literature is the formulation of gradient descent variants for training to address the issue of non-differentiability of morphological operators.

Using the tropical mathematics framework, Charisopoulos and Maragos studied these networks and proposed a training algorithm not rooted in stochastic optimization but rather on Difference-of-Convex Programming [6] and provided lower bounds for the linear regions of such networks [7]. Zhang et al. study feedforward network with Rectified Linear Unit (ReLU) activations and showed their equivalence with tropical rational maps [15]. Also, Calafiore et al. study a similar class of networks with *Log-Sum-Exp* terms, which can be thought of as smooth approximations of the hard morphological operators [16, 17]. Tropical mathematics have also been linked with sparse systems and the problem of pruning models. [18] studies the sparsity of max-plus systems. Other works have focused on neural network pruning; Smyrnis et al. develop a geometric algorithm for tropical polynomial division aimed at neural network minimization [19] and extend this method to multiclass problems [9].

In recent years, the field of morphological networks has adopted the trend of Deep Learning with increasingly more complex architectures, both in terms of width as well as depth. Mondal et al. study dense morphological networks by choosing the morphological operators of dilation and erosion as the basic operations of a neuron and overcome their undifferentiability by considering smooth approximations [20]. Other works have focused on extending the morphological setting to convolutional layers. Franchi et al. proposed a joint operation that includes pooling and the morphological non-linearities [21], while Mellouli et al. employed the counter-harmonic mean to inject morphological operations to convolutional layers and created an interpretable morphological convolutional neural network aimed at digit recognition [22].

Morphological Neural Networks have also been proposed to address monotonicity constraints. First, in the context of neural networks, Archer and Wang proposed to a heuristic which updates the weights of a binary classification model so that samples do not violate the monotonicity constraints. The resulting network has only positive weights (for the case of increasing monotonicity) [23]. A different approach lies in constraining the weights to positive values in a neural network with a single hidden layer. This can be achieved through a monotonic nonlinear transformation with a positive images such as the sigmoid function [24]. However, Velikova et al. showed that this approach requires K hidden layers to approximate a K -dimensional monotonic surface [25]. On the other hand Sill proposes a network that guarantees monotonicity for the output via architectural design [26]. This method is extended by Daniels and Velikova to include cases of partially monotone functions [27]. Also, Sill’s work can also be considered a special case of the min-max classifiers proposed by Yang and Maragos [5].

Apart from morphological networks, other monotonic methods have been proposed in the machine learning field. An overview appears in [28]. Some efforts have been focused on interpolated look-up tables, where shape constraints including monotonicity can be imposed [29, 30] and these ideas have also been applied in the context of deep networks [31]. An alternative approach has been proposed by Wehenkel and Louppe, where monotonic functions are modeled via neural network techniques by enforcing a positive image on the derivative of the monotonic function [32].

3 BACKGROUND CONCEPTS

In machine learning literature, the perceptron uses a multiply-accumulate scheme that feeds into an activation function $\phi(\cdot)$, which performs a nonlinear transformation. A graphical representation appears in Fig. 1. Specifically, the perceptron consists of a weight vector $\mathbf{a} \in \mathbb{R}^n$ and a bias term $b \in \mathbb{R}$. These parameters are combined with the input $\mathbf{x} \in \mathbb{R}^n$ to compute the weighted sum $f(\mathbf{x}) = \mathbf{a}^\top \mathbf{x} + b = \sum_{i=1}^n a_i x_i + b$, which then passes from the aforementioned nonlinear filter. The tropical or morphological neuron has similar structure but performs different operations. The operations of addition and multiplication are replaced by \max (or \min) and addition. This process is known as *tropicalization*. A notable difference is the absence of a nonlinear activation function, since the operators \max (or \min) are inherently nonlinear. In the context of tropical algebra, we use the symbols \boxplus and \boxplus' to denote max-plus and min-plus matrix multiplication, respectively. These operations are defined as $(\mathbf{A} \boxplus \mathbf{B})_{ij} = \bigvee_{q=1}^k a_{iq} + b_{qj}$ and $(\mathbf{A} \boxplus' \mathbf{B})_{ij} = \bigwedge_{q=1}^k a_{iq} + b_{qj}$ for matrices of appropriate dimensions.

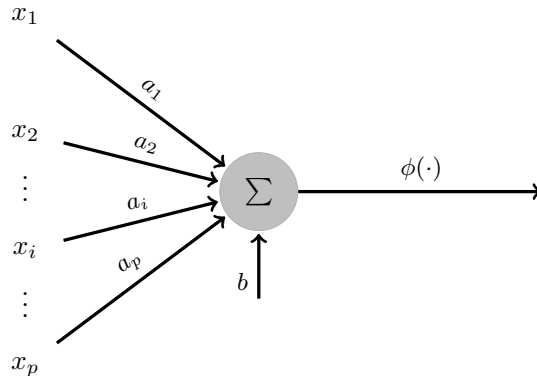


Figure 1: Perceptron

The operator change yields piecewise linear (PWL) surfaces. For $\mathbf{x} \in \mathbb{R}^n$, the tropical max-plus polynomial p_{\vee} is defined as the maximum of multiple affine terms:

$$p_{\vee}(\mathbf{x}) = \max_i \{\mathbf{a}_i^{\top} \mathbf{x} + b_i\} = \bigvee_i \mathbf{a}_i^{\top} \mathbf{x} + b_i, \quad (1)$$

whilst its min-plus equivalent arises from replacing max with min. Examples of max-plus and min-plus polynomials are presented in Fig. 2. The max-plus (min-plus) polynomial generates convex (concave) PWL surfaces. In the context of mathematical morphology, these operators are often referred to as dilation δ and erosion ε , respectively, and are defined as:

$$\delta_{\mathbf{w}}(\mathbf{x}) = w_0 \vee \left(\bigvee_i w_i + x_i \right) \quad (2)$$

$$\varepsilon_{\mathbf{m}}(\mathbf{x}) = m_0 \wedge \left(\bigwedge_i m_i + x_i \right) \quad (3)$$

where \mathbf{w} , \mathbf{m} correspond to the weights of the dilation and the erosion neurons, respectively. These expressions are special cases of tropical polynomials, since the term \mathbf{a}_i is specific and can be expressed as tropical matrix multiplications as $\delta_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^{\top} \boxplus \tilde{\mathbf{x}}$ and $\varepsilon_{\mathbf{m}}(\mathbf{x}) = \mathbf{m}^{\top} \boxminus \tilde{\mathbf{x}}$ where $\tilde{\mathbf{x}} = [1 \ \mathbf{x}]$, i.e. includes the bias term. Dilation and erosion are commonly used in image processing. Given a pixel neighborhood, dilation outputs the pixel with the largest value. Thus, light objects are enlarged, while dark regions are shrunk. The erosion operator performs the dual transformations.

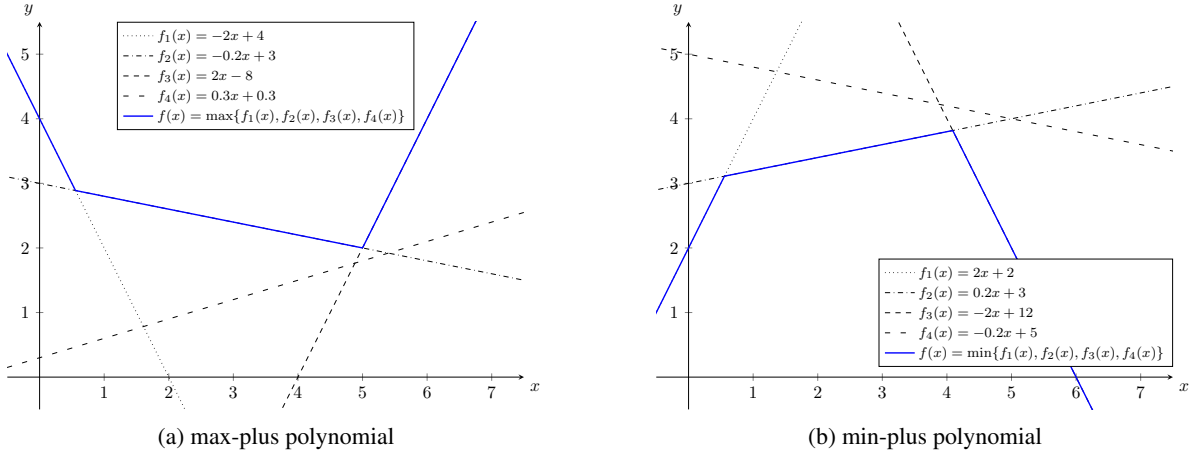


Figure 2: Constructing polynomials in the tropical semiring. In both cases, $f_4(x)$ has no bearing in the resulting surface. This term can be removed.

A generalization of the tropical polynomials lies in the notion of hinging hyperplanes by Wang and Sun [33], which are a collection of affine terms joined by points where multiple terms *dominate*, i.e. are maximizers or minimizers. These surfaces can be neither convex or concave and correspond to a lattice-theoretic view of these PWL surfaces, also found in the works of Tarela et al. [34, 35]. The maximizers or minimizers are points that are not differentiable and are responsible for the hard edges. In the case of strictly max-plus (or min-plus) PWL surfaces, this collection of points is known as the tropical hypersurface $\mathcal{T}(p)$. The hard surfaces can be softened by approximating the min and max operators, a process known as Maslov Dequantization:

Definition 1 (Maslov Dequantization [36]). Let $x, y \in \mathbb{R}$ and $h > 0$. The transformation $x \vee_h y = h \log(e^{x/h} + e^{y/h})$ defines the Maslov Dequantization of the max operator, yielding its soft approximation. Similarly, the Maslov Dequantization of the min operator is $x \wedge_h y = -h \log(e^{-x/h} + e^{-y/h})$. As $h \rightarrow 0$, the soft versions approach the hard ones: $\lim_{h \rightarrow 0} x \vee_h y = x \vee y$ and $\lim_{h \rightarrow 0} x \wedge_h y = x \wedge y$ (see proof in Appendix).

For small positive values of h , these approximations are part of the Log-Sum-Exp family, used in convex analysis and recently linked to tropical polynomials [16, 17]. We use the reciprocal of h , the hardness parameter $\beta = h^{-1}$.

4 Convex-Concave Procedure Training for General Classification Tasks

4.1 Binary Classification Problem formulation

In the neural network training context, the Gradient Descent algorithm and its variants [37] dominate the literature. Network parameters are optimized with stochastic algorithms, allowing parallelization of the process, thus faster

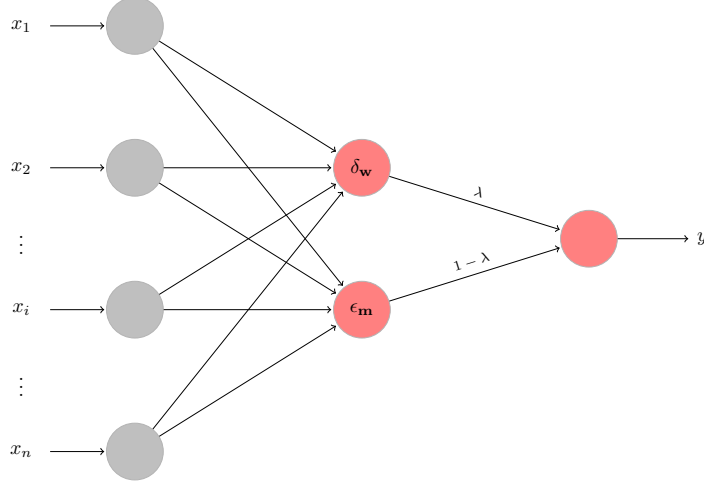


Figure 3: Dilation-Erosion Perceptron (DEP)

training, but sacrificing on robustness. Alternative methods have been proposed that eliminate the stochastic element and offer low variance in the resulting models. A salient example lies in the Support Vector Machines (SVMs) literature, where the training is formulated as a convex program [38], which allows the use of fast and robust algorithms. In the current section, we explore a similar training scheme, based on Difference-of-Convex programming and extend its applicability to multiclass problems.

Let us consider the task of classifying the pattern $\mathbf{x}_i \in \mathbb{R}^n, i = 1, 2, \dots, N$ to two distinct classes, \mathcal{N} for negative ($y_i = -1$) and \mathcal{P} for positive ($y_i = +1$). Given an input $\mathbf{x} \in \mathbb{R}^n$ and weight vectors $\mathbf{w}, \mathbf{m} \in \mathbb{R}^{n+1}$, a dilation term $\delta_{\mathbf{w}}$ and an erosion term $\epsilon_{\mathbf{m}}$ are combined to produce an output. The network is shown in Fig. 3 and carries the name *Dilation-Erosion Perceptron (DEP)*. A similar architecture lies in the maxout network [39], which uses the more general notion of tropical polynomials (with arbitrary slopes \mathbf{a}_i , see (1)). Charisopoulos and Maragos formulate the training of this classifier [6] as:

$$\begin{aligned} \min \quad & \sum_{i=1}^N v_i \max\{0, \xi_i\} \\ \text{s.t.} \quad & \lambda \delta_{\mathbf{w}}(\mathbf{x}_i) + (1 - \lambda) \epsilon_{\mathbf{m}}(\mathbf{x}_i) \geq -\xi_i \quad \forall \mathbf{x}_i \in \mathcal{P}, \\ & \lambda \delta_{\mathbf{w}}(\mathbf{x}_i) + (1 - \lambda) \epsilon_{\mathbf{m}}(\mathbf{x}_i) \leq +\xi_i \quad \forall \mathbf{x}_i \in \mathcal{N} \end{aligned} \quad (4)$$

where ξ_i are slack variables and v_i correspond to a weighting scheme. Specifically, the slack variables ξ_i ensure that only misclassified patterns are taken into account for the objective function, whereas the weighting scheme allows the attachment of higher importance to some misclassified patterns. Charisopoulos and Maragos propose penalizing patterns with greater chances of being outliers. Specifically, for classes $\mathcal{C}_0 = \mathcal{N}$ and $\mathcal{C}_1 = \mathcal{P}$:

$$v_i = \frac{\lambda_i}{\max_j \lambda_j}, \quad \lambda_j = \frac{1}{\|\mathbf{x}_j - \boldsymbol{\mu}_k\|_p}, \quad \boldsymbol{\mu}_k = \frac{1}{|\mathcal{C}_k|} \sum_{\mathbf{x}_i \in \mathcal{C}_k} \mathbf{x}_i$$

The above training formulation (4) corresponds to a Difference-of-Convex (DC) optimization problem, since the dilation term is convex whereas the erosion term is concave. This optimization class is vast and includes all C^2 functions [40]. Various methods have been proposed to tackle such problems with many focusing on the Fenchel Conjugate. Another approach lies in a heuristic called *Convex-Concave Procedure (CCP)*, proposed by Yuille and Rangarajan [41] and extended in [42, 43]. This heuristic is based on the observation that a DC program is convex iff the concave term, in this case the erosion term, is affine. Hence, by convexifying this term, i.e. calculating its linear approximation via a first-order Taylor series expansion, the problem becomes convex unlocking the ability of efficient and fast solutions. Repeating this process yields the CCP. Lipp and Boyd extended this procedure by allowing the initialization of the algorithm without a feasible point and improved convergence with the use of penalty terms [42]. CCP was further extended [43] by the introduction of a structured method of defining DC problems which automatically converts them in a form suitable for generic solvers [43]. We use such solvers in our experiments, see CVXPY [44].

Valle [45] proposes a modification of the formulation (4), where a greedy algorithm combines the dilation and the erosion perceptrons, which are trained separately beforehand. This method allows the inclusion of a regularization term

$R(\mathbf{u})$ in the objective function, since one set of parameters is optimized. Valle proposes $R(\mathbf{u}) = C\|\mathbf{u} - \mathbf{r}\|_1$, where $\mathbf{u} = \mathbf{w}$ or $\mathbf{u} = \mathbf{m}$ and \mathbf{r} is a reference term. The linear combination is calculated by minimizing the average hinge loss:

$$\lambda^* = \operatorname{argmin}_{0 \leq \lambda \leq 1} \sum_{i=1}^N \max\{0, -y_i [\lambda \delta_{\mathbf{w}}(\mathbf{x}_i) + (1 - \lambda) \varepsilon_{\mathbf{m}}(\mathbf{x}_i)]\} \quad (5)$$

As a lattice-based model, the Dilation-Erosion Perceptron suffers from a major flaw. Particularly, it presupposes a partial ordering both on the features and the classes. This results in counterintuitive behavior. By simply inverting the classes $\mathcal{N} \rightleftharpoons \mathcal{P}$, the problem intuitively does not change but mathematically it does, which is reflected on a severe drop in performance [45]. A method to placate this issue can be found on the use of reduced morphological operators based on a reduced ordering:

Definition 2 (reduced ordering). Let R be a nonempty set, \mathcal{L} be a complete lattice and $\rho : R \rightarrow \mathcal{L}$ be a surjective mapping. A reduced ordering, or r-ordering, is defined as:

$$\mathbf{x} \leq_{\rho} \mathbf{y} \Leftrightarrow \rho(\mathbf{x}) \leq \rho(\mathbf{y}), \forall \mathbf{x}, \mathbf{y} \in R. \quad (6)$$

Effectively, this means that a mapping $\rho : \mathbb{R}^n \rightarrow \mathbb{R}^m$ which transforms the input is used and the classification is performed on the resulting dataset. Let the original dataset be $\mathcal{D} = \{(\mathbf{x}_i, y_i) \in \mathbb{R}^n \times \{-1, +1\} : i \in [N]\}$, the mapping produces the following dataset $\mathcal{D}_{\text{new}} = \{(\rho(\mathbf{x}_i), y_i) \in \mathbb{R}^m \times \{-1, +1\} : i \in [N]\}$ and consists of m mappings $\rho(\mathbf{x}) = [\rho_1(\mathbf{x}), \rho_2(\mathbf{x}), \dots, \rho_m(\mathbf{x})]^T$ where $\rho_i(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}, i \in [m]$ are mappings from the original input space to the real numbers and are called *kernels*. Kernels are extensively used in Support Vector Machine modeling. Some characteristic examples are presented in Table 1.

Kernel	$k(\mathbf{x}, \mathbf{y})$
Linear	$\langle \mathbf{x}, \mathbf{y} \rangle$
Polynomial	$\langle 1 + \mathbf{x}, \mathbf{y} \rangle^d$
Gaussian	$e^{-\ \mathbf{x} - \mathbf{y}\ ^2 / (2\sigma^2)}$
Sigmoid	$\tanh(\gamma \langle \mathbf{x}, \mathbf{y} \rangle + r)$

Table 1: Kernels

4.2 Extension to multiclass problems

The Dilation-Erosion Perceptron is inherently a binary classifier. As is the case with Support Vector Machines, this classifier can be extended to tackle multiclass problems. Notably, there are two major approaches: *one-versus-the-rest* and *one-versus-one*. Let us consider a problem with $K > 2$ classes. On both categories, the idea lies in the construction of several models, all binary, and combine their outputs to produce a single prediction.

In the former approach, one classifier is trained for each class. In order to do so, the positive class consists of the elements of said class \mathcal{C}_k whereas the negative patterns consist of all the other classes \mathcal{C}_{-k} . A straightforward issue with this approach lies on the imbalance of the datasets. Let N be the total number of datapoints and consider a relatively uniform distribution of patterns among classes. Then, the positive class is much smaller than the negative: $|\mathcal{C}_k| \simeq \frac{N}{K} \ll |\mathcal{C}_{-k}| \simeq \frac{(K-1)N}{K}$. A way to address this issue is the use of weighting scheme: $+1$ for positive and $-\frac{1}{K}$ for negative patterns. In the latter approach, one classifier is trained for every pair of classes and the prediction is determined by the (hard) majority vote of the corpus of classifiers. The issue that emerges is the sheer number of classifiers that have to be trained. For K classes, there are $\frac{K(K-1)}{2}$ pairs. However, each pair is trained on only $\sim \frac{2}{K} \cdot N \ll N$ patterns.

In our experiments, we employ the *one-versus-one* approach with a greedy Dilation-Erosion Perceptron [45] on the MNIST [46] and FashionMNIST [47] datasets. Thus, $= \frac{10 \times 9}{2} = 45$ distinct classifiers must be trained. Each classifier uses a Bagging approach with n Radial Basis Function (RBF) kernels. A Bagging classifier uses the same kernel, in this case RBF, but the parameters are trained in different subsets of the original dataset. The results are presented in Table 2 and are comparable with gradient-descent methods of training neural networks, presented in the next section. Nevertheless, the networks presented in next sections are deeper and denser and, thus, have many more parameters. An important aspect of this non-stochastic training method is robustness. Repeating the experiments 10 times yielded very low variance in accuracy, a phenomenon that does not present when training with stochastic gradient descent variants. An important aspect of this method is the kernel selection. Also, there is a different approach other than Bagging, the Ensemble method, where different kernels are used but on the entire training set. Other parameterizations include the

	MNIST	FashionMNIST
$n = 5$	97.72 \pm 0.01	88.21 \pm 0.01
$n = 10$	97.72 \pm 0.01	88.07 \pm 0.01
$n = 15$	97.67 \pm 0.01	88.11 \pm 0.01
$n = 20$	97.64 \pm 0.01	88.12 \pm 0.01

Table 2: Results of Bagging *multiclass* reduced DEP with n RBF kernels.

number and type of kernels. The combinations are vast and, with further experimentation, the results of Table 2 can be improved.

5 Pruning Morphological Neural Nets

The Dilation-Erosion Perceptron has only two neurons in the hidden layer. A straightforward extension is the population of the hidden layer with more neurons and/or the use of multiple hidden layers. This family of networks is called Dense Morphological [20] and an instance is depicted in Fig. 4. The single dilation neuron of DEP is replaced by n_1 and n_2 terms in the first and second hidden layer, respectively. Similarly for the erosion neuron. By stacking dilation and erosion neurons on hidden layers, the network calculates activations similar to morphological openings and closings.

By increasing the number of parameters, both in depth and width, the network is able to search for effective hidden representations of the input data, removing the major problem of selecting a surjective mapping ρ , which plagues the multiclass DEP. The use of stochastic optimization methods for training allows for parallelization of the training process and the utilization of optimized deep learning libraries that take advantage of Graphical Process Units (GPUs), resulting in faster training. In this case, the models are trained with standard gradient descent methods. We study two variants: (mini-batch) Stochastic Gradient Descent (SGD) and Adaptive Momentum Estimation (Adam) [48]. Our focus lies *not* on achieving the highest possible accuracy, but on showing the compression ability of morphological networks compared to traditional ones. To this end, we apply pruning techniques to evaluate the ability of the various networks to retain information with a fraction of the original nodes.

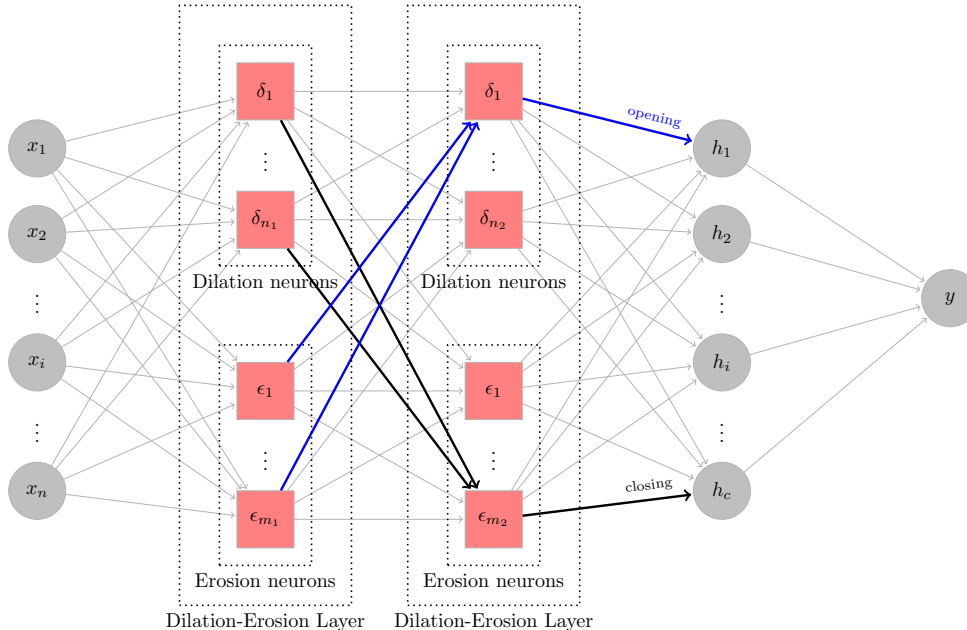


Figure 4: Dense Morphological Network with 2 hidden layers. Square nodes correspond to morphological operators, either min (erosion) or max (dilation). The output layer is fully connected.

A note on the notation of the models used. Regarding the hidden layer(s), the models studied include: only dilation neurons (denoted as δ), only erosion neurons (ϵ), a mixed network with both types of neurons ((δ, ϵ)) as well as a Feedforward neural networks with Rectified Linear Unit (ReLU) activations (FF-ReLU) for comparison. The smooth

morphological networks are denoted with the subscript β . The experiments focus on visual recognition tasks and the datasets selected remain the same as in Section 4, i.e. MNIST and FashionMNIST. The training lasts 50 epochs. After experimentation, we conclude that the best learning rates differ for each optimizer and correspond to the values $\eta = 0.001$ and $\eta = 0.09$ for Adaptive Momentum Estimation and Stochastic Gradient Descent, respectively. For the Stochastic Gradient Descent, we omit networks with 2 hidden layers and the models with smooth operators, since our experimental setup of 50 epochs does not yield competitive results, due to gradient propagation issues.

The results are presented in Table 3. The Adam algorithm yields higher accuracy for both datasets and both optimizers. Also, the increase in density of the hidden layer is beneficial to accuracy but there is a diminishing returns effect; after a certain amount of nodes in the hidden layer the performance boost is minimal. In both datasets, the linear model FF-ReLU slightly outperforms the morphological equivalents. It is important to note that Mondal et al. [20] used a different experimental setup (w.r.t. epochs¹ etc.) and showed that morphological networks can in fact achieve higher accuracy on these datasets compared to FF-ReLU. However, as stated above, our focus lies *not* on accuracy but on compression ability.

	#	Adaptive Momentum Estimation						Stochastic Gradient Descent			
		δ	ϵ	δ, ϵ	$2(\delta, \epsilon)$	$\delta_\beta, \epsilon_\beta$	FF-ReLU	δ	ϵ	δ, ϵ	FF-ReLU
MNIST	32	94.36	87.12	92.19	92.73	92.85	96.54	92.83	82.40	90.00	97.00
	64	96.12	91.74	94.92	95.15	95.03	97.53	94.65	87.98	93.00	97.56
	128	96.90	94.91	96.59	96.33	96.93	98.07	95.29	91.68	94.73	97.92
	256	97.30	94.98	97.42	97.08	96.95	98.13	96.10	91.81	95.74	98.02
	400	97.62	96.17	97.63	97.09	97.84	98.03	95.87	94.59	96.07	98.08
FashionMNIST	32	80.98	81.87	83.16	83.60	81.35	87.06	77.75	79.43	79.82	87.49
	64	84.09	84.38	85.20	85.34	83.36	87.83	80.41	82.44	83.10	88.29
	128	85.92	86.13	86.76	86.51	84.47	88.44	81.21	84.18	83.98	88.66
	256	86.22	87.15	87.11	87.12	85.81	89.09	78.65	85.77	85.25	87.69
	400	86.62	88.05	88.34	87.47	87.13	89.44	80.09	86.20	86.21	88.81

Table 3: Comparison of the accuracy of various network architectures on the MNIST and FashionMNIST datasets. Different sizes of the hidden layers, denoted as #, are evaluated. See above for columns description.

The compression ability of the various networks is evaluated by the effect on performance that the removal of a nontrivial amount of parameters has. We apply an ℓ_1 -norm pruning scheme [49], which retains salient features and removes unimportant ones. By specifying a percentage $p \in [0, 100]$, only the $p\%$ of the parameters with the highest ℓ_1 -norm are retained. We focus on the hidden layer, where the bulk of the parameters is concentrated. Various percentages of units are pruned and the performance of the remaining network is evaluated on the test set.

The results are presented in Table 4. With shades of red, the deterioration in terms of performance is highlighted. All variants of the morphological network significantly outperform their linear counterpart in their ability to retain information and the effectiveness of their constructed hidden representations. This effect becomes more evident in the more complicated dataset FashionMNIST. For example, in the case of the SGD optimizer, the erosion and mixed morphological networks have virtually zero decline in performance (the erosion net has actually a slight increase) even when $p = 1\%$ but the pruning has negatively affected the performance by $p = 25\%$, i.e. FF-ReLU requires ~ 25 times more parameters in the hidden layer. This difference can be explained by Fig. 5, where the hidden layer activations are depicted for a morphological and a linear network. Lighter colors correspond to higher values. The representations constructed by the morphological network are sparse; only a fraction of the parameters have high values. Intuitively, this means that the removal of the rest (depicted in blue) will have minimal impact in information loss and, thus, accuracy. However, this effect is not present in the linear net, since the weights have more uniform values.

Another insight lies in the optimizer selection. Even though models trained with Stochastic Gradient Descent have lower accuracy scores, their ability to retain information is superior than those trained via Adaptive Momentum Estimation. Quantitatively, the optimizer effect can be seen in Table 4 by comparing the first and last rows ($p = 100\%$ vs $p = 1\%$). The SGD models have no loss in accuracy, whereas the Adam models have a slight decrease. Quantitatively, this effect can be explained by Fig. 6, where a single neuron of networks with 400 neurons in the hidden layer is depicted for each model category. In the case of the Adam-trained morphological network, the activation is reminiscent of the the outline of a MNIST pattern and uses a nontrivial amount of parameters. On the other hand, Fig. 6b shows notable sparsity

¹[20] uses 400 epochs for training, but we only use 50 epochs.

in the SGD-trained model activations. The effect holds for the linear networks but is milder given the already dense activations.

In conclusion, the morphological networks are characterized by extreme economy in their constructed representations. This is due to the fact that a morphological neuron, dilation or erosion, allows only one element x_i of the input $\mathbf{x} \in \mathbb{R}^n$ to determine its output. With a small (to none) decrease in accuracy compared to their linear counterparts, the morphological networks are able to retain their performance with a small fraction of the original parameters and display higher compression ability.

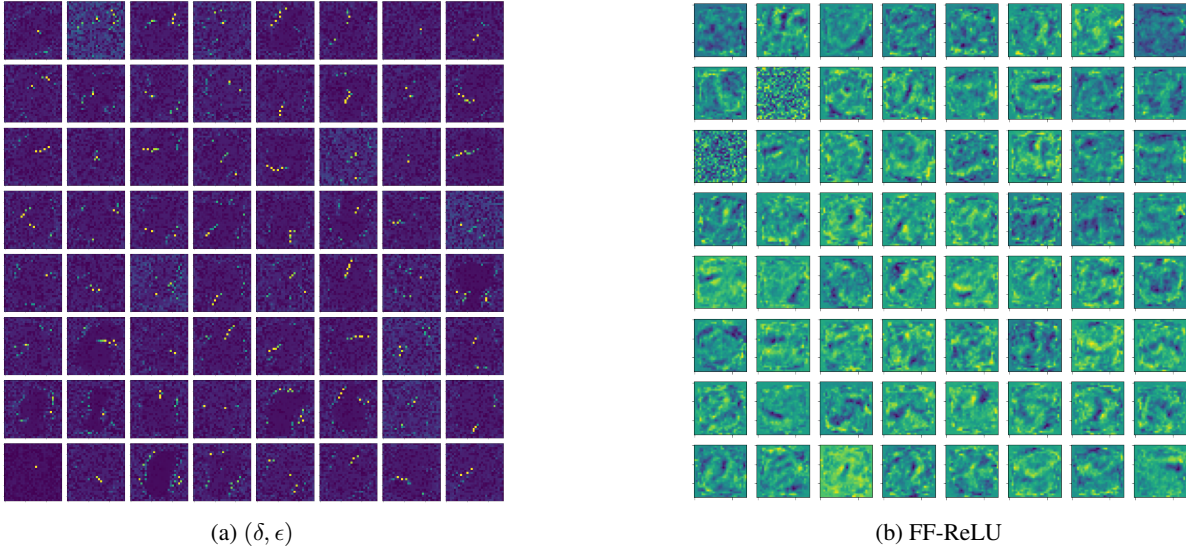


Figure 5: Examples of hidden layer activations for various models (MNIST dataset). (a) corresponds to a morphological network (dilation neurons) and (b) to a feedforward network with ReLU activations (FF-ReLU). Both models have 64 neurons in the hidden layer (forming an 8×8 grid). Each element of the grid has dimensions 28×28 pixels, i.e. equal to the input images of the MNIST dataset.

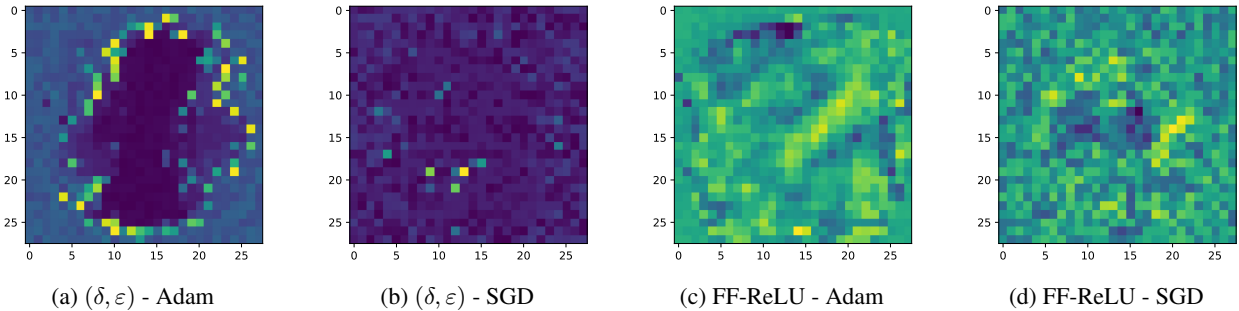


Figure 6: Examples of hidden layer activations for various models (MNIST dataset). (a) and (b) correspond to a morphological network (dilation neurons), whereas (c) and (d) correspond to a feedforward network with ReLU activations.

6 Monotonicity Constraints

In this section, we explore how architectural structure can be used to enforce shape constraints. We examine the case of monotonicity. The output of a network is called monotone if it does not decrease with the increase of the input² The network into consideration was introduced by Sill [26] and is depicted in Fig. 7. Max-affine terms, i.e. dilations or max-plus polynomials, are succeeded by a min-pooling layer, i.e. an erosion term or a min-plus polynomial. Sill names as a *group* each of the K max terms, each of which comprises of J_k affine terms or hyperplanes, $k = 1, 2, \dots, K$. For

²Without loss of generality, we consider only monotonically increasing functions.

		Adaptive Momentum Estimation				Stochastic Gradient Descent				
		p	δ	ϵ	(δ, ϵ)	FF-ReLU	δ	ϵ	(δ, ϵ)	FF-ReLU
MNIST	100%		97.62	96.17	97.95	98.13	94.86	93.36	96.07	98.16
	75%		97.62	96.18	97.93	98.15	94.86	93.36	96.07	98.12
	50%		97.62	96.22	97.90	98.17	94.86	93.37	96.07	98.08
	25%		97.62	96.09	97.87	97.51	94.86	93.40	96.06	98.01
	10%		97.62	95.78	97.74	93.38	94.86	93.38	96.09	96.67
	7.5%		97.62	95.42	97.76	90.17	94.86	93.38	96.10	95.56
	5%		97.62	94.51	97.66	83.39	94.86	93.40	96.10	92.96
	2.5%		97.62	93.43	97.37	68.93	94.86	93.39	96.09	80.48
	1%		97.62	91.17	97.08	44.22	94.86	93.38	96.08	58.07
FashionMNIST	100%		86.31	86.82	88.32	88.82	82.06	85.23	86.21	87.79
	75%		86.30	86.81	88.30	88.88	82.00	85.23	86.21	87.75
	50%		86.22	86.80	88.33	88.18	82.05	85.25	86.20	87.19
	25%		85.95	86.85	88.31	82.15	81.90	85.26	86.28	84.35
	10%		85.58	86.27	88.05	65.89	81.67	85.27	86.23	73.22
	7.5%		85.47	86.15	87.99	57.93	81.63	85.27	86.21	63.95
	5%		85.37	85.81	87.76	49.12	81.52	85.24	86.22	47.73
	2.5%		84.91	85.47	87.56	42.48	81.14	85.26	86.22	38.84
	1%		81.14	84.86	86.85	28.13	80.68	85.27	86.18	35.46

Table 4: Performance of pruned networks on the MNIST and FashionMNIST datasets for various model architectures. With shades of red, we show the rapid deterioration in performance. With green, we draw attention to the *absence* of loss in accuracy performance between the full (unpruned) network and the pruned network with only 1% of the parameters in the hidden layer.

simplicity, we consider the case where $J_k = J, \forall k = 1, 2, \dots, K$. Then, for input $\mathbf{x} \in \mathbb{R}^n$, the output is:

$$y = f(\mathbf{x}) = \bigwedge_{k \in [K]} \bigvee_{j \in [J]} \{\mathbf{w}_{k,j}^\top \mathbf{x} + b_{k,j}\} \quad (7)$$

Monotonicity constraints are enforced by limiting the weight vector \mathbf{w} to nonnegative values via a function f with a positive image $f: \mathbb{R}^n \rightarrow \mathbb{R}_{\geq 0}^n$. Various such functions have been proposed. Sill used an exponential transformation $w_i = e^{z_i}, z_i \in \mathbb{R}$, whilst Velikova et al. [25] propose the transformation $w_i = z_i^2$, which suffers no arithmetic issues and offers the possibility of flat surfaces when $\mathbf{z} = \mathbf{0}$. The network produces a PWL approximation of the input data. Max-plus terms construct convex surfaces, whereas min-plus terms construct concave surfaces. Their combination in this network structure accommodates for the more general nature of a monotonic surface, which combines both convex and concave parts. The resulting surface can approximate any monotone function to an arbitrary degree of accuracy [26, Theorem 3.1].

From a mathematical morphology viewpoint, the application of dilation succeeded by an erosion term results in a closing. For strictly positive weights, the transformation is reversible and corresponds to a morphological opening $x = f^{-1}(y) = \bigvee_{k \in [K]} \bigwedge_{j \in [J]} \{w_{k,j}^{-1}(y - b_{k,j})\}$ [50]. An opening $\alpha = \delta\epsilon$ is an increasing, idempotent and anti-extensive operator, while a closing $\beta = \epsilon\delta$ shares the first two properties but is extensive. Opening and closing form an adjunction pair [11].

The morphological terms are not differentiable, which has immediate effect in training via stochastic gradient descent variants, since the backpropagation step does not have an analytical form. Subgradient methods can be used in this case. The backpropagation equations for the dilation and erosion terms of Sill’s monotonic network are:

$$\frac{\partial \epsilon}{\partial \delta_k} = \begin{cases} 1 & \text{argmax}_h \{\delta_h\} = k \\ 0 & \text{otherwise} \end{cases} \quad (8)$$

$$\frac{\partial \delta_k}{\partial \mathbf{w}_{k,j}} = \begin{cases} \mathbf{x} & \text{argmax}_h \{\mathbf{w}_{k,h}^\top \mathbf{x}\} = j \\ 0 & \text{otherwise} \end{cases} \quad (9)$$

In other words, the morphological operators are selective in the sense that a single element x_i of the input vector \mathbf{x} is solely responsible for the output. Thus, in the backpropagation step, only x_i ’s parameters are updated. The effect is

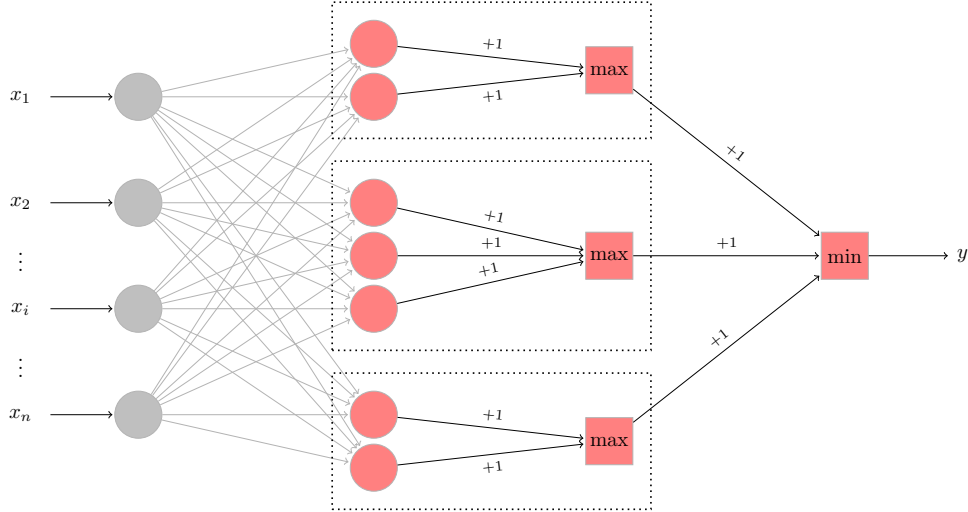


Figure 7: Monotonic network. The gray edges correspond to nonnegative weights.

intensified in Sill’s network where two morphological layers are stacked. This implies that only the active hyperplane’s weights are updated for a given pattern. Out of the $K \times J$ affine terms, only a single one gets updated. This results in slow convergence and, potentially, poor approximations. Given low initialized weight parameters, the updated hyperplanes have acquired higher values in order to approximate the data and, subsequently, dominate the groups, not allowing the remaining (hyperplanes) to update their weights.

Sill proposes a variant of the gradient descent algorithm where the gradient for each hyperplane is computed by the error of the patterns corresponding to the active hyperplane at each iteration. Another method to address this issue is the use of a gain parameter G which magnifies the initialized weights. This way, the weights of the dominating hyperplane get diminished during the backpropagation step allowing the other hyperplanes to dominate in the next epochs. We propose a method to circumvent this issue completely, which lies in the softening of morphological operators \max and \min via Maslov Dequantization, alleviating the undifferentiability of their hard counterparts.

We illustrate this method via a simple example. We consider the (strictly) increasing function $f(x) = x^3 + x + \sin x$, $x \in [-4, 4]$ and scale both domain and image to $[-1, 1]$. Glorot uniform initialization [51] is used for all network weights. 100 observations are sampled uniformly and corrupted with additive i.i.d zero-mean Gaussian noise $\epsilon \sim \mathcal{N}(0, \sigma^2)$. The training lasts 1000 epochs using Adaptive Momentum Estimation with learning rate $\eta = 0.01$. For comparison purposes, we use isotonic regression [52], a method based on quadratic programming that yields monotone piecewise constant surfaces and has the following formulation in \mathbb{R} :

$$\begin{aligned} \min \quad & \sum_i w_i (y_i - \hat{y}_i)^2 \\ \text{s.t.} \quad & \hat{y}_i < \hat{y}_j \quad x_i < x_j \end{aligned}$$

For the monotonic net, we select gain parameter $G = 20$ for initializing the weights. For the smooth monotonic net, we select hardness parameter $\beta = 5$. The results are presented in Table 5 for $K = J = 5$ for various noise levels. Particularly, for $\sigma = 0.15$ the surfaces that each method produces are depicted in Fig. 8. From Table 5, we conclude that the smooth monotonic net outperforms the other methods for all noise levels σ . Furthermore, its training procedure needs not the selection of an arbitrary gain parameter G and is overall simpler and more intuitive.

σ	0.05	0.1	0.15	0.2
Linear Reg.	0.0236	0.03077	0.04827	0.0505
Isotonic Reg.	0.0042	0.01112	0.02557	0.0417
Sill Net	0.00305	0.01107	0.02401	0.0390
Smooth Sill Net	0.00294	0.00938	0.02302	0.0386

Table 5: RMS error of monotonic regression methods with noise $\epsilon \sim \mathcal{N}(0, \sigma^2)$

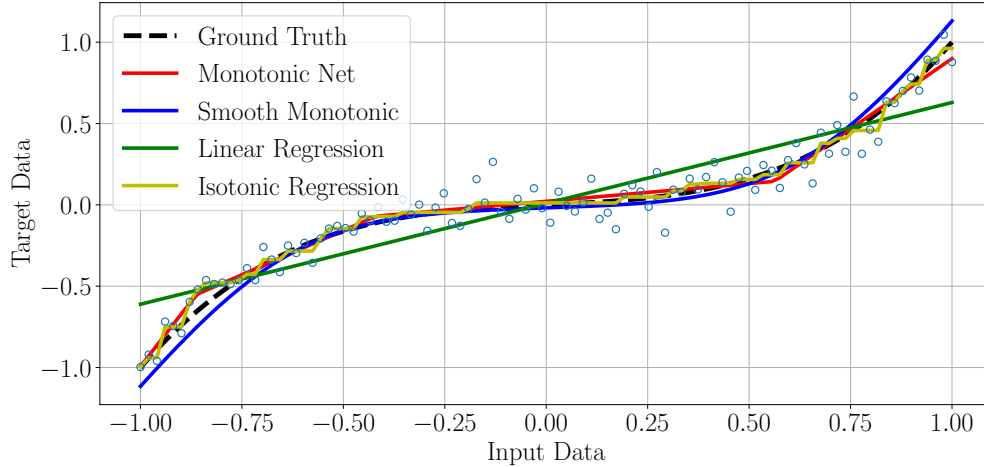


Figure 8: Comparison of monotonic regression methods

7 CONCLUSION

In this paper, we have proposed extensions of morphological neurons regarding their training, enforcing shape constraints such as monotonicity via architectural choices and studied their compression ability. We extended the binary classifier *Dilation-Erosion Perceptron* to general classification tasks using the *one-versus-one* approach and employed a Bagging method with Radial Basis Functions kernels to alleviate the partial ordering of lattice-based models and developed robust classifiers with accuracies comparable to more general morphological networks. We studied the accuracy of dense morphological networks under heavy pruning and compared their ability to construct efficient hidden representations to their linear counterparts, drawing favorable conclusions for the morphological models. Finally, we proposed the use of smooth operators in a monotonic network which improves not only performance but convergence in training as well.

References

- [1] Ritter, G. X. and Sussner, P. “An introduction to morphological neural networks”. In: *Proceedings of 13th International Conference on Pattern Recognition*. Vol. 4. IEEE. 1996, pp. 709–717.
- [2] Ritter, G. X., Sussner, P., and Diza-de-Leon, J. “Morphological associative memories”. In: *IEEE Transactions on neural networks* 9.2 (1998), pp. 281–293.
- [3] Sussner, P. “Morphological perceptron learning”. In: *Proceedings of the 1998 IEEE International Symposium on Intelligent Control (ISIC) held jointly with IEEE International Symposium on Computational Intelligence in Robotics and Automation (CIRA) Intell.* IEEE. 1998, pp. 477–482.
- [4] Pessoa, L. F. and Maragos, P. “Neural networks with hybrid morphological/rank/linear nodes: a unifying framework with applications to handwritten character recognition”. In: *Pattern Recognition* 33.6 (2000), pp. 945–960.
- [5] Yang, P.-F. and Maragos, P. “Min-max classifiers: Learnability, design and application”. In: *Pattern Recognition* 28.6 (1995), pp. 879–899.
- [6] Charisopoulos, V. and Maragos, P. “Morphological Perceptrons: Geometry and Training Algorithms”. In: *Mathematical Morphology and Its Applications to Signal and Image Processing*. Cham: Springer International Publishing, 2017. ISBN: 978-3-319-57240-6.
- [7] Charisopoulos, V. and Maragos, P. “A Tropical Approach to Neural Networks with Piecewise Linear Activations”. In: *arXiv* (May 2018). URL: <http://arxiv.org/abs/1805.08749>.
- [8] Zhang, Y. et al. “Max-plus Operators Applied to Filter Selection and Model Pruning in Neural Networks”. In: *International Symposium on Mathematical Morphology and Its Applications to Signal and Image Processing*. Springer. 2019, pp. 310–322.
- [9] Smyrnis, G. and Maragos, P. “Multiclass Neural Network Minimization via Tropical Newton Polytope Approximation”. In: *Proceedings of the 37th International Conference on Machine Learning (ICML)*. PMLR, July 2020.

- [10] Davidson, J. L. and Hummer, F. “Morphology neural networks: An introduction with applications”. In: *Circuits, Systems and Signal Processing* 12.2 (1993), pp. 177–210.
- [11] Maragos, P. “Dynamical systems on weighted lattices: general theory”. In: *Mathematics of Control, Signals, and Systems* 29.4 (Dec. 2017). ISSN: 0932-4194, 1435-568X. DOI: [10.1007/s00498-017-0207-8](https://doi.org/10.1007/s00498-017-0207-8).
- [12] Barmpoutis, A. and Ritter, G. X. “Orthonormal basis lattice neural networks”. In: *Computational Intelligence Based on Lattice Theory*. Springer, 2007, pp. 45–58.
- [13] Ritter, G. X. and Urcid, G. “Lattice algebra approach to single-neuron computation”. In: *IEEE Transactions on Neural Networks* 14.2 (2003), pp. 282–295.
- [14] Sussner, P. and Esmi, E. L. “Morphological perceptrons with competitive learning: Lattice-theoretical framework and constructive learning algorithm”. In: *Information Sciences* 181.10 (2011), pp. 1929–1950.
- [15] Zhang, L., Naitzat, G., and Lim, L.-H. “Tropical Geometry of Deep Neural Networks”. In: *Proc. ICML*. 2018.
- [16] Calafiore, G. C., Gaubert, S., and Possieri, C. “Log-sum-exp neural networks and posynomial models for convex and log-log-convex data”. In: *arXiv:1806.07850 [cs]* (2018).
- [17] Calafiore, G. C. et al. “A Universal Approximation Result for Difference of log-sum-exp Neural Networks”. In: *arXiv:1905.08503 [cs]* (2019).
- [18] Tsiamis, A. and Maragos, P. “Sparsity in max-plus algebra and systems”. In: *Discrete Event Dynamic Systems* (May 2019). ISSN: 0924-6703, 1573-7594. DOI: [10.1007/s10626-019-00281-1](https://doi.org/10.1007/s10626-019-00281-1).
- [19] Smyrnis, G., Maragos, P., and Retsinas, G. “Maxpolynomial Division with Application To Neural Network Simplification”. In: *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 4192–4196.
- [20] Mondal, R., Santra, S., and Chanda, B. “Dense Morphological Network: An Universal Function Approximator”. In: *arXiv* (2019). URL: <http://arxiv.org/abs/1901.00109>.
- [21] Franchi, G., Fehri, A., and Yao, A. “Deep morphological networks”. In: *Pattern Recognition* 102 (2020), p. 107246.
- [22] Mellouli, D. et al. “Morphological convolutional neural network architecture for digit recognition”. In: *IEEE transactions on neural networks and learning systems* 30.9 (2019), pp. 2876–2885.
- [23] Archer, N. P. and Wang, S. “Application of the Back Propagation Neural Network Algorithm with Monotonicity Constraints for Two-Group Classification Problems*”. In: *Decision Sciences* 24.1 (1993), pp. 60–75. ISSN: 1540-5915. DOI: [10.1111/j.1540-5915.1993.tb00462.x](https://doi.org/10.1111/j.1540-5915.1993.tb00462.x).
- [24] Kay, H. and Ungar, L. H. “Estimating monotonic functions and their bounds”. In: *AIChE Journal* 46.12 (2000), pp. 2426–2434.
- [25] Velikova, M., Daniels, H., and Feelders, A. “Solving Partially Monotone Problems with Neural Networks”. In: 12 (2006), p. 6.
- [26] Sill, J. “Monotonic networks”. In: *Advances in neural information processing systems*. 1998.
- [27] Daniels, H. and Velikova, M. “Monotone and Partially Monotone Neural Networks”. In: *IEEE Transactions on Neural Networks* 21.6 (2010), pp. 906–917. ISSN: 1941-0093. DOI: [10.1109/TNN.2010.2044803](https://doi.org/10.1109/TNN.2010.2044803).
- [28] Gupta, M. et al. “Monotonic calibrated interpolated look-up tables”. In: *The Journal of Machine Learning Research* 17.1 (2016), pp. 3790–3836.
- [29] Gupta, M. et al. “Diminishing Returns Shape Constraints for Interpretability and Regularization”. In: *Advances in Neural Information Processing Systems* 31. Ed. by S. Bengio et al. Curran Associates, Inc., 2018, pp. 6834–6844.
- [30] Milani Fard, M. et al. “Fast and Flexible Monotonic Functions with Ensembles of Lattices”. In: *Advances in Neural Information Processing Systems* 29. Ed. by D. D. Lee et al. Curran Associates, Inc., 2016, pp. 2919–2927.
- [31] You, S. et al. “Deep Lattice Networks and Partial Monotonic Functions”. In: *arXiv:1709.06680 [cs, stat]* (2017).
- [32] Wehenkel, A. and Louppe, G. “Unconstrained Monotonic Neural Networks”. In: *Advances in Neural Information Processing Systems* 32. Ed. by H. Wallach et al. Curran Associates, Inc., 2019, pp. 1545–1555.
- [33] Wang, S. and Sun, X. “Generalization of hinging hyperplanes”. In: *IEEE Transactions on Information Theory* 51.12 (2005), pp. 4425–4431. ISSN: 1557-9654. DOI: [10.1109/TIT.2005.859246](https://doi.org/10.1109/TIT.2005.859246).
- [34] Tarela, J., Alonso, E., and Martinez, M. “A representation method for PWL functions oriented to parallel processing”. In: *Mathematical and Computer Modelling* 13.10 (1990), pp. 75–83.
- [35] Tarela, J. and Martinez, M. “Region configurations for realizability of lattice piecewise-linear models”. In: *Mathematical and Computer Modelling* 30.11-12 (1999), pp. 17–27.
- [36] Litvinov, G. L. “Maslov dequantization, idempotent and tropical mathematics: A brief introduction”. In: *J. Math. Sciences* 140.3 (2007). ISSN: 1072-3374, 1573-8795.
- [37] Ruder, S. “An overview of gradient descent optimization algorithms”. In: *arXiv* (2017). URL: <http://arxiv.org/abs/1609.04747>.

-
- [38] Boyd, S. P. and Vandenberghe, L. *Convex optimization*. Cambridge University Press, 2004. ISBN: 978-0-521-83378-3.
- [39] Goodfellow, I. et al. “Maxout Networks”. In: vol. 28. PMLR, 2013, pp. 1319–1327. URL: <http://proceedings.mlr.press/v28/goodfellow13.html>.
- [40] Hartman, P. “On functions representable as a difference of convex functions”. In: *Pacific Journal of Mathematics* 9.3 (Sept. 1959), pp. 707–713. ISSN: 0030-8730, 0030-8730. DOI: [10.2140/pjm.1959.9.707](https://doi.org/10.2140/pjm.1959.9.707).
- [41] Yuille, A. L. and Rangarajan, A. “The Concave-Convex Procedure”. In: *Neural Computation* 15.4 (Apr. 2003), pp. 915–936. ISSN: 0899-7667, 1530-888X. DOI: [10.1162/08997660360581958](https://doi.org/10.1162/08997660360581958).
- [42] Lipp, T. and Boyd, S. “Variations and extension of the convex-concave procedure”. In: *Optimization and Engineering* 17.2 (June 2016), pp. 263–287. ISSN: 1389-4420, 1573-2924. DOI: [10.1007/s11081-015-9294-x](https://doi.org/10.1007/s11081-015-9294-x).
- [43] Shen, X. et al. “Disciplined convex-concave programming”. In: *2016 IEEE 55th Conference on Decision and Control (CDC)*. IEEE, Dec. 2016. ISBN: 978-1-5090-1837-6. URL: <http://ieeexplore.ieee.org/document/7798400/>.
- [44] Diamond, S. and Boyd, S. “CVXPY: A Python-embedded modeling language for convex optimization”. In: *Journal of Machine Learning Research* 17.83 (2016), pp. 1–5.
- [45] Valle, M. E. “Reduced Dilation-Erosion Perceptron for Binary Classification”. In: *Mathematics* 8.4 (2020), p. 512. ISSN: 2227-7390. DOI: [10.3390/math8040512](https://doi.org/10.3390/math8040512).
- [46] LeCun, Y. “The MNIST database of handwritten digits”. In: <http://yann.lecun.com/exdb/mnist/> (1998).
- [47] Xiao, H., Rasul, K., and Vollgraf, R. “Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms”. In: *arXiv* (Aug. 2017).
- [48] Kingma, D. P. and Ba, J. “Adam: A Method for Stochastic Optimization”. In: *arXiv:1412.6980 [cs]* (2017).
- [49] Li, H. et al. “Pruning filters for efficient convnets”. In: *arXiv* (2016).
- [50] Duetting, P. et al. “Optimal Auctions through Deep Learning”. In: *ICML*. 2019. URL: <http://proceedings.mlr.press/v97/duetting19a.html>.
- [51] Glorot, X. and Bengio, Y. “Understanding the difficulty of training deep feedforward neural networks”. In: *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. 2010, pp. 249–256.
- [52] Barlow, R. E. and Brunk, H. D. “The isotonic regression problem and its dual”. In: *J. Amer. Stat. Assoc.* 337 (1972).

A Appendix

Proof. We will prove that $\delta_\beta(\mathbf{x}) = \frac{1}{\beta} \log \left(\sum_k e^{\beta x_k} \right) \rightarrow \max_k \{x_k\} = \delta(\mathbf{x})$ as $\beta \rightarrow \infty$.

$$\begin{aligned}
\lim_{\beta \rightarrow \infty} \delta_\beta(\mathbf{x}) &= \lim_{\beta \rightarrow \infty} \frac{1}{\beta} \log \left(\sum_k e^{\beta y_k} \right) \\
&= \lim_{\beta \rightarrow \infty} \frac{\log \left(\sum_k e^{\beta y_k} \right)}{\beta} &> \text{DLH} \\
&= \lim_{\beta \rightarrow \infty} \frac{\sum_k y_k e^{\beta y_k}}{\sum_k e^{\beta y_k}} \\
&= \lim_{\beta \rightarrow \infty} \sum_j \frac{y_j e^{\beta y_j}}{\sum_k e^{\beta y_k}} \\
&= \lim_{\beta \rightarrow \infty} \sum_j \frac{y_j}{1 + \sum_{k \neq j} e^{\beta(y_k - y_j)}} &> \text{divide by } e^{\beta y_j} \\
&= \sum_j \underbrace{\lim_{\beta \rightarrow \infty} \frac{y_j}{1 + \sum_{k \neq j} e^{\beta(y_k - y_j)}}}_{A_j} &(10)
\end{aligned}$$

where with "DLH" we note the use of the *De L'Hôpital* rule. Let K be the number of elements that are equal to the maximum term. Then, we have the following cases for the term A_j :

- $j \neq j^* = \operatorname{argmax}_k y_k$. Then, there exist three sets $I^>, I^<, I^=$ such that $I^> = \{i : y_i > y_j\}$, $I^< = \{i : y_i < y_j\}$ and $I^= = \{i : y_i = y_j\}$. But $I^> \neq \emptyset$. Hence,

$$\begin{aligned}
A_j &= \lim_{\beta \rightarrow \infty} \frac{y_j}{1 + \sum_{k \neq j} e^{\beta(y_k - y_j)}} \\
&= \lim_{\beta \rightarrow \infty} \frac{y_j}{1 + \sum_{k \in I^<} e^{\beta(y_k - y_j)} + \sum_{k \in I^=} e^{\beta(y_k - y_j)} + \sum_{k \in I^>} e^{\beta(y_k - y_j)}} \\
&= \frac{y_j}{1 + \sum_{k \in I^<} e^{-\infty} + \sum_{k \in I^=} e^0 + \sum_{k \in I^>} e^\infty} = 0
\end{aligned}$$

- $j = j^* = \operatorname{argmax}_k y_k$. Then, $I^> = \emptyset, |I^|= K - 1$ from our hypothesis and each term $y_k - y_j$ of the exponential at the denominator is negative for all $k \in I^<$. Thus:

$$\begin{aligned}
A_j &= \lim_{\beta \rightarrow \infty} \frac{y_j}{1 + \sum_{k \neq j} e^{\beta(y_k - y_j)}} \\
&= \frac{y_j}{1 + \sum_{k \in I^=} e^0 + \sum_{k \in I^<} e^{-\infty}} \\
&= \frac{y_j}{1 + (K - 1) + 0} \\
&= \frac{y_j}{K} = \frac{y_{\max}}{K}
\end{aligned}$$

Subsequently, the quantity A_j is zero for every y_j that is not equal to the maximum y_{\max} and is equal to $\frac{y_{\max}}{K}$ for each of the K terms that are equal to the maximum. Thus, using 10 concludes the proof. \square