

Neural Spectrahedra and Semidefinite Lifts: Global Convex Optimization of Degree-Two Polynomial Activation Neural Networks in Polynomial-Time

Burak Bartan · Mert Pilanci

Received: date / Accepted: date

Abstract The training of two-layer neural networks with nonlinear activation functions is an important non-convex optimization problem with numerous applications and promising performance in layerwise deep learning. In this paper, we develop exact convex optimization formulations for two-layer neural networks with second degree polynomial activations based on dual relaxations and semidefinite programming. Remarkably, we show that our semidefinite relaxations are always tight. Therefore, the computational complexity for global optimization is polynomial in the input dimension and sample size for all input data. The developed convex formulations are proven to achieve the same globally optimal solution set as their non-convex counterparts. Specifically, globally optimal two-layer neural networks with degree-two polynomial activations can be found by solving a semidefinite program (SDP) and decomposing the solution using a procedure we call Neural Decomposition. Moreover, the choice of regularizers plays a crucial role in the computational tractability of neural network training. We show that the standard weight decay regularization formulation is NP-hard, whereas other simple convex penalties render the problem tractable in polynomial time via convex programming. The techniques go beyond the fully connected architecture to encompass various neural network structures including those with vector outputs and convolutional architectures. We provide extensive numerical simulations showing that the standard backpropagation approach often fails to achieve the global optimum of the training loss. The proposed approach is significantly faster to obtain better test accuracy compared to the standard backpropagation procedure. ¹

B. Bartan
Department of Electrical Engineering, Stanford University, CA, USA
E-mail: bbartan@stanford.edu

M. Pilanci
Department of Electrical Engineering, Stanford University, CA, USA
E-mail: pilanci@stanford.edu

¹ The supplementary material for this work can be accessed by following the link in [4].

Keywords Neural networks Convex optimization Semidefinite programming Polynomial activation

Mathematics Subject Classification (2020) 90C22 90C25 90C26 68T07

1 Introduction

We study neural networks from the optimization perspective by deriving equivalent convex optimization formulations with identical global optimal solution sets. The derived convex problems have important theoretical and practical implications for the computational complexity of optimal neural network training. Moreover, the convex optimization perspective provides a more concise parameterization of neural networks and reveal novel insights.

In non-convex optimization, the choice of optimization method and its internal hyperparameters, such as initialization, mini-batching and step sizes, have a considerable effect on the quality of the learned model. This is in sharp contrast to convex optimization problems, where locally optimal solutions are globally optimal and optimizer parameters have no influence on the solution and therefore the model. Moreover, the solutions of convex optimization problems can be obtained in a very robust, efficient and reproducible manner thanks to the elegant and extensively studied structure of convex programs. Therefore, our convex optimization based globally optimal training procedure enables the study of the neural network model and the optimization procedure in a *decoupled* way. For instance, step sizes employed in the optimization can be considered hyperparameters of non-convex models, which affect the model quality and may require extensive tuning. For a classification task, in our convex optimization formulation, step sizes as well as the choice of the optimizers are no longer hyperparameters to obtain better classification accuracy. Any convex optimization solver can be leveraged to solve the convex problem to obtain a globally optimal model.

Various nonlinearities have been proposed in the literature as activation functions. Among the most widely adopted ones is the ReLU (rectified linear unit) activation given by $\sigma(u) = \max(0, u)$. A recently proposed alternative is the swish activation $\sigma(u) = u(1 + e^{-u})^{-1}$, which performs comparably well [40]. Another important class is the polynomial activation where the activation function is a scalar polynomial of a fixed degree. We focus on second degree polynomial activation functions, i.e., $\sigma(u) = au^2 + bu + c$. Although polynomial coefficients $a; b; c$ can be regarded as hyperparameters, it is often sufficient to select them to approximate a target nonlinear activation function such as the ReLU or swish activation. ReLU and swish activations are plotted in Figure 1 along with their second degree polynomial approximations.

Our derivation of the convex program for degree-two polynomial activations leverages convex duality and the S-procedure, and results in a simple semidefinite program (SDP). We refer the reader to [39] for a survey of the S-procedure and applications in SDPs.

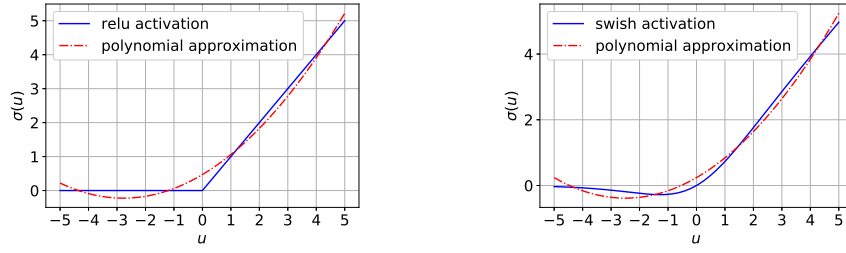


Fig. 1 ReLU (left) and swish (right) activation functions and their second degree polynomial approximations. ReLU activation: $\sigma(u) = \max(0; u)$ and its polynomial approximation: $\sigma(u) = 0.09u^2 + 0.5u + 0.47$. Swish activation: $\sigma(u) = u(1 + e^{-u})^{-1}$ and its polynomial approximation: $\sigma(u) = 0.1u^2 + 0.5u + 0.24$.

Main aspects of our work that differ from others in the literature that study the *optimization landscape* of two-layer neural networks (e.g. see prior work section) are the following: Our results (1) provide global optimal solutions in polynomial time, (2) uncover an important role of the regularizer in computational tractability, (3) hold for arbitrary convex loss functions and other network architectures such as vector output, convolutional and pooling, (4) are independent of the choice of the numerical optimizer and its parameters.

1.1 Overview of Our Contributions

- We show that the standard optimization formulation for training neural networks $f(x) = \sum_{j=1}^m (x^T u_j)_j$ with trainable parameters $u = (u_1; \dots; u_{m-1}; \dots; u_m)$, degree-two polynomial activations $\sigma(u) = au^2 + bu + c$, training data $X = [x_1; \dots; x_n]^T \in \mathbb{R}^{n \times d}$, $y \in \mathbb{R}^n$ and weight decay (WD) regularization given by

$$(P_{WD}) \quad \min_{\substack{u_j \in \mathbb{R}^d; \\ j \in \{1, \dots, m\}}} \sum_{j=1}^n (X u_j)_j; y + \sum_{j=1}^n (k u_j k_2^2 + \frac{1}{2} j^2) \quad (1)$$

is computationally intractable via reduction from the NP-hard subset sum problem for all values of m . The loss function $\ell(\hat{y}; y)$ is a function of the network output $\hat{y} \in \mathbb{R}^n$. The hardness result is valid for any loss function that satisfies $\arg \min_{\hat{y} \in \mathbb{R}^n} \ell(\hat{y}; y) = y$.

- Surprisingly, for square activation, i.e., $\sigma(u) = u^2$, we show that modifying the quadratic weight decay regularization to *cubic regularization*

$$(P_{SQ}) \quad \min_{\substack{u_j \in \mathbb{R}^d; \\ j \in \{1, \dots, m\}}} \sum_{j=1}^n (X u_j)_j; y + \sum_{j=1}^n (k u_j k_2^3 + j j^3) \quad (2)$$

enables global optimization in polynomial time via convex semidefinite programming when the number of neurons m is above a threshold m^* .

The threshold m is proportional to the rank of the solution of the convex semidefinite program. An exact expression is provided in the sequel. The computational complexity is polynomial in all problem parameters $(n; d; m)$. The result holds for any convex and closed function $\ell(\hat{y}; y)$; a standard example is the squared loss $\ell(\hat{y}; y) = k\hat{y} - yk_2^2$.

- Furthermore, for any degree-two polynomial activation σ , the non-convex neural network training problem

$$(P) \quad \min_{u_j \in \mathbb{R}^d; j \in [2R]; \delta_j \in [m]} \sum_{j=1}^m \ell(Xu_j; y) + \sum_{j=1}^m \lambda_j \|u_j\|_2^2$$

s.t. $\|u_j\|_2 = 1; \delta_j \in [m]$ (3)

can be equivalently stated as a convex semidefinite problem and solved in polynomial time. In fact, the cubic regularization strategy in (2) is a special case of this convex program. The result holds universally for all input data without any conditions and also holds when $\lambda \neq 0$.

- In deriving the convex formulations, we identify a concise re-parameterization of the neural network parameters that enables exact convexification by removing the redundancy in their classical formulation. This is similar in spirit to the semidefinite lifting procedure in relaxations of combinatorial optimization problems. In contrast to these relaxations, we show that our lifting is always exact as soon as the network width exceeds a critical threshold which can be efficiently determined.
- We develop a matrix decomposition procedure called Neural Decomposition to extract the optimal network parameters from the solution of our convex program, which is guaranteed to produce an optimal neural network. Neural Decomposition transforms the convex re-parameterization to the overparameterized, i.e., redundant, formulation in a similar spirit to (a non-orthogonal version of) Eigenvalue Decomposition.
- In addition to the fully connected architecture, we derive equivalent convex programs for other architectures such as convolutional, pooling and vector output architectures (see the supplementary material [4] for details).
- We provide extensive numerical simulations showing that the standard backpropagation approach fails to achieve the global optimum of the training loss. Moreover, the test accuracy of the proposed convex optimization is considerably higher in standard datasets as well as random planted models. Our convex optimization solver is significantly faster in total computation time to achieve similar or better test accuracy.

1.2 Prior Work

A considerable fraction of recent works on the analysis of optimization landscape of neural networks focuses on explaining why gradient descent performs well. The works [12, 44] consider the optimization landscape of a restricted class of neural networks with square activation and quadratic regularization

where the second layer weights are $\mathcal{X} \in \mathbb{R}^{d \times d}$. They show that when the neural network is overparameterized, i.e., $m \geq d$, the non-convex loss function has benign properties: all local minima are global and all saddle points have a direction of negative curvature. However, in this paper we show that training both the first and second layer weights with quadratic regularization in fact makes global optimization NP-hard for any value of m . In contrast, we provide a different formulation to obtain the global optimal solution via convex optimization in the more general case when the second layer weights are also optimized, the activation function is any arbitrary degree-two polynomial, and global optimum is achieved for all values of m . The work in [30] similarly studies two-layer neural networks with square activation function and squared loss and states results on both optimization and generalization properties. The authors in [17] focus on square activation networks from the perspectives of optimization landscape and generalization performance, where the setting is based on a planted model with a full rank weight matrix. In [25,28] it was shown that sufficiently wide ReLU networks have a benign landscape when each layer is sufficiently wide, satisfying $m \geq n + 1$.

Another recent work analyzing the training of neural networks with quadratic-like activations for deeper architectures is [2]. Authors in [2] consider degree-two polynomial activation functions and investigate layerwise training and compare with end-to-end training of layers. It is demonstrated in [2] that the degree-two polynomial activation function performs comparably to ReLU activation in deep networks. More specifically, it is reported in [2] that for deep neural networks, ReLU activation achieves a classification accuracy of 0.96 and a degree-two polynomial activation yields an accuracy of 0.95 on the Cifar-10 dataset. Similarly for the Cifar-100 dataset, they obtain an accuracy of 0.81 for ReLU activation and 0.76 for the degree-two polynomial activation. These numerical results are obtained for the activation $\sigma(u) = u + 0.1u^2$, which the authors prefer over the square activation $\sigma(u) = u^2$ to make the neural network training stable. Moreover, the performance of layerwise learning with such activation functions is considerably high, although there is a gap between end-to-end trained models. In addition, neural networks with polynomial activations have immediate applications in encrypted computing [22,18,31,34]. In encrypted computing, it is desirable to have a low degree polynomial as the activation function. For instance, homomorphic encryption can only support additions and multiplications in a straightforward way, which necessitates low degree polynomials as activations. In [18], degree-two polynomial approximations were shown to be effective for accurate neural network predictions with encryption.

In a recent series of papers, the authors derived convex formulations for training ReLU neural networks to global optimality [38,14,15,13,41,42]. Our work takes a similar convex duality approach in deriving the convex equivalents of non-convex neural network training problems. In particular, the previous work in this area deals with ReLU activations while in this work we focus on degree-two polynomial activations. Hence, the techniques involved in deriving the convex programs and the resulting convex programs are substan-

tially different. The convex program derived for ReLU activation in [38] has polynomial time trainability for fixed rank data matrices, whereas the convex programs developed in this work are polynomial time trainable with respect to all problem dimensions. More specifically, the convex program in [38] is

$$\begin{aligned} \min_{v_i, w_i \in \mathbb{R}^d; \delta_i \in [H]} \quad & \sum_{i=1}^H D_i X(v_i \ w_i); y + \sum_{i=1}^H (k v_i k_2 + k w_i k_2) \\ \text{s.t.} \quad & (2D_i \ I_n) X v_i \geq 0; (2D_i \ I_n) X w_i \geq 0; \delta_i \in [H]; \end{aligned} \quad (4)$$

where the neural network weights are constructed from $v_i \in \mathbb{R}^d$ and $w_i \in \mathbb{R}^d$, $i = 1, \dots, H$. The matrices D_i are diagonal matrices whose diagonal entries consist of $1_{x_1^T u \geq 0}; 1_{x_2^T u \geq 0}; \dots; 1_{x_n^T u \geq 0}$ for all possible $u \in \mathbb{R}^d$. The number of distinct D_i matrices, denoted by H is the number of hyperplane arrangements corresponding to the data matrix X . It is known that H is bounded by $2r \frac{e(n-1)}{r}$ where $r = \text{rank}(X)$ (see [38] for the details). In particular, convolutional neural networks have a fixed value of r , for instance m filters of size 3×3 yield $r = 9$. This is an exponential improvement over previously known methods that train optimal ReLU networks which are exponential in the number of neurons m and/or the number of samples n [3, 19, 6].

The work in [7] presents formulations for convex factorization machines with nuclear norm regularization, which is known to obtain low rank solutions. Vector output extension for factorization machines and *polynomial networks*, which are different from *polynomial activation networks*, is developed in [8]. Polynomial networks are equivalent to square activation networks with an addition of a linear neuron. In [8], the authors consider learning an infinitely wide square activation layer by a greedy algorithm. The proposed method in [8] does not provide optimal finite width networks. Furthermore, [29] presents a greedy algorithm for training polynomial networks. The algorithm provided in [29] is based on gradually adding neurons to the neural network to reduce the loss. More recently, [43] considers applying lifting for square activation neural networks and presents non-convex algorithms for low rank matrix estimation for two-layer neural network training.

1.3 Notation

Throughout the text, $\sigma: \mathbb{R} \rightarrow \mathbb{R}$ denotes the activation function of the hidden layer. We refer to the function $\sigma(u) = u^2$ as square activation and $\sigma(u) = au^2 + bu + c$ where $a, b, c \in \mathbb{R}$ as degree-two polynomial activation. We use $X \in \mathbb{R}^{n \times d}$ to denote the data matrix, where its rows $x_i \in \mathbb{R}^d$ correspond to data samples and columns are the features. In the text, whenever we have a function mapping from \mathbb{R} to \mathbb{R} with a vector argument (e.g., $\sigma(v)$ or v^2 where v is a vector), this means the elementwise application of that function to all the components of the vector v . We denote a column vector of ones by $\mathbf{1}$ and its dimension can be understood from the context. $\text{vec}(\cdot)$ denotes the vectorized

version of its argument. In writing optimization problems, we use \min and \max to refer to “minimize” and “maximize”. We use the notations $[m]$ and $1; \dots; m$ interchangeably.

We use $\ell(y; y)$ for convex loss functions throughout the text. The notation $\ell^*(v) = \sup_z (v^T z - \ell(z; y))$ denotes the Fenchel conjugate of the function $\ell(\cdot; y)$. Furthermore, we assume $\ell^* = \ell$ which holds when ℓ is a convex and closed function [9]. We use $Z \succeq 0$ for positive semidefinite matrices (PSD). $S^{d \times d}$ refers to the set of $d \times d$ dimensional symmetric matrices. tr refers to matrix trace. \otimes is used for outer product. The operator **conv** stands for the convex hull of a set.

1.4 Preliminaries on Semidefinite Lifting

We defer the discussion of semidefinite lifting for two-layer neural networks to next section. We now briefly discuss a class of problems where SDP relaxations lead to exact convex optimization solutions of the original non-convex problem and also instances where they fail to be exact. Let us consider a quadratic objective problem with a single quadratic constraint:

$$\begin{aligned} \min_{u \in \mathbb{R}^d} \quad & u^T Q_1 u + b_1^T u + c_1 \\ \text{s.t.} \quad & u^T Q_2 u + b_2^T u + c_2 \leq 0 \end{aligned} \quad (5)$$

where $Q_1, Q_2 \in \mathbb{R}^{d \times d}$ are indefinite, i.e., not assumed to be positive semidefinite and $b_1, b_2 \in \mathbb{R}^d, c_1, c_2 \in \mathbb{R}$. Due to the indefinite quadratics, this is a non-convex optimization problem. By introducing a matrix variable $U = uu^T$, one can equivalently state this problem as

$$\begin{aligned} \min_{\substack{U \in \mathbb{R}^{d \times d} \\ u \in \mathbb{R}^d}} \quad & \text{tr}(Q_1 U) + b_1^T u + c_1 \\ \text{s.t.} \quad & \text{tr}(Q_2 U) + b_2^T u + c_2 \leq 0 \\ & U = uu^T : \end{aligned} \quad (6)$$

This problem can be relaxed by replacing the equality by the matrix inequality $U \succeq uu^T$. Re-writing the expression $U \succeq uu^T$ as a linear matrix inequality via the Schur complement formula yields the following SDP

$$\begin{aligned} \min_{\substack{U \in \mathbb{R}^{d \times d} \\ u \in \mathbb{R}^d}} \quad & \text{tr}(Q_1 U) + b_1^T u + c_1 \\ \text{s.t.} \quad & \text{tr}(Q_2 U) + b_2^T u + c_2 \leq 0 \\ & \begin{bmatrix} 1 & u^T \\ u & U \end{bmatrix} \succeq 0 : \end{aligned} \quad (7)$$

Remarkably, it can be shown that the original non-convex problem in (5) can be solved exactly by solving the convex SDP in (7) via duality, under the

mild assumption that the original problem is strictly feasible (see [9]). This shows that the SDP relaxation is exact in this problem, returning a globally optimal solution when one exists. We note that there are alternative numerical procedures to compute the global optimum of quadratic programs with one quadratic constraint [9].

As an alternative to the lifting based primal approach, it is possible to obtain the same convex SDP formulation via Lagrangian duality. First, we note that the dual of the original non-convex problem in (5) is the SDP

$$\begin{aligned} \max_{\substack{Q_1, Q_2 \\ b_1, b_2}} \quad & \frac{1}{4} + c_2 + c_1 \\ \text{s.t.} \quad & \begin{pmatrix} Q_1 + Q_2 & b_1 + b_2 \\ (b_1 + b_2)^T & 0 \end{pmatrix} \preceq 0 \end{aligned} \quad (8)$$

Taking the dual of this SDP, we obtain the bidual of (5), which gives us the same SDP relaxation in (7).

We also note that the lifting approach $U = uu^T$ followed by the relaxation $U \succeq uu^T$ is not tight for quadratic programs with more than two quadratic constraints [33, 10]. A notable case with multiple constraints is the NP-hard Max-Cut problem and its SDP relaxation [20]

$$\max_{\substack{u \in \mathbb{R}^d \\ u_i^2 = 1; \forall i \in [d]}} u^T Qu = \max_{\substack{u \in \mathbb{R}^d \\ u_i^2 = 1; \forall i \in [d]}} \text{tr}(Quu^T) = \max_{\substack{U \in \mathbb{R}^{d \times d}, U \succeq 0 \\ U_{ii} = 1; \forall i \in [d]}} \text{tr}(QU) \quad (9)$$

The SDP relaxation of Max-Cut is not tight in general since its feasible set contains the cut polytope

$$\text{conv} \{ uu^T : u_i \in \{-1, 1\}; \forall i \in [d] \}$$

and other non-integral extreme points [26]. Nevertheless, an approximation ratio of 0.878 can be obtained via the Goemans-Williamson randomized rounding procedure under certain restrictions on Q [20]. It is conjectured that this is the best approximation ratio for Max-Cut [23], whereas it can be formally proven to be NP-hard to approximate within a factor of $\frac{16}{17}$ [21, 45]. Hence, in general we cannot expect to obtain exact solutions to problems of combinatorial nature, such as Max-Cut and variants using computationally efficient SDP relaxations.

It is instructive to note that a naive application of the SDP lifting strategy is not immediately tractable for two-layer neural networks. For simplicity, consider a scalar output degree-two polynomial activation network $f(x) = \sum_{j=1}^m (x^T u_j)^2$ where $(u) = \{u^2 + u\}$, and $f(u_j; \{g_{j=1}^m\})$ are trainable parameters. The corresponding training problem for a given loss function $\ell(\cdot; y)$ and

its SDP relaxation are as follows

$$\begin{aligned} & \min_{\substack{u_j \in \mathbb{R}^d; \\ j \in \mathbb{R}; \delta_j \in [m]}} \times \prod_{j=1}^n \mathcal{X}^n ((x^T u_j)^2 + x^T u_j) \quad j; y \\ & \min_{\substack{U_j \in \mathbb{R}^d; \\ U_j \in \mathbb{R}^d; \delta_j \in [m]}} \times \prod_{j=1}^n \mathcal{X}^n x^T U_j x + x^T u_j \quad j; y : \end{aligned} \quad (10)$$

The above problem is *non-convex* due to the bilinear terms $f U_j \quad j \in \mathbb{R}^m$. Moreover, a variable change $\hat{U}_j = U_j \quad j$ does not respect semidefinite constraints $U_j \quad u_j u_j^T$ when $j \in \mathbb{R}$. Another limitation is the prohibitively high number of variables in the lifted space, which is $d^2 m + dm + m$ as opposed to $dm + m$ in the original problem. Therefore, a different convex analytic formulation is needed to address all these problems.

Although SDP relaxations are extensively studied for various non-convex problems (see e.g. [46] for a survey of applications), instances with exact SDP relaxations are exceptionally rare. As will be discussed in the sequel, our main result for two-layer neural networks is another instance of an SDP relaxation leading to *exact* formulations where the semidefinite relaxation is tight.

In convex geometry, a *spectrahedron* is a convex body that can be represented as a linear matrix inequality which are the feasible sets of semidefinite programs. An example is the *elliptope* defined as the feasible set of the Max-Cut relaxation given by $U \succeq 0; U_{ii} = 1 \quad \delta_i$, which is a subset of $n \times n$ symmetric positive-definite matrices. Due to the existence of efficient projection operators and barrier functions of linear matrix inequalities, optimizing convex objectives over spectrahedra can be efficiently implemented, which renders SDPs tractable. We will show that degree-two polynomial activation networks can be represented via a class of linear matrix inequalities, dubbed *neural spectrahedra* (see Figure 2 for an example), which enables global optimization in polynomial time and elucidates their parameterization in convex analytic terms.

1.5 Paper Organization

The paper is structured as follows. Section 2 gives an overview of the theory developed in this work. Section 3 describes the convex optimization formulation via duality and S-procedure for degree-two polynomial activation neural networks. Section 4 establishes via the neural decomposition method that the convex problem developed in Section 3 can be used to train two-layer neural networks to global optimality. The hardness result for the weight decay regularization is studied in Section 5. We discuss the implementation details for solving the convex programs and give experimental results in Section 6. Omitted proofs can be found in the appendix. The supplementary material [4] has sections on the special case of square activation neural networks with cubic regularization, and on vector output and convolutional architectures, as well as additional numerical results.

Table 1 List of the neural network architectures considered in this work and the corresponding convex programs. Abbreviations are as follows. Poly (scalar): Degree-two polynomial activation scalar output, Poly (vector): Degree-two polynomial activation vector output, Convolutional: Convolutional neural network (CNN) with degree-two polynomial activation, Pooling: CNN with degree-two polynomial activation and average pooling, Quad (scalar, cubic reg): Square activation scalar output with cubic regularization, Quad (scalar, quad reg): Square activation scalar output with quadratic regularization. K is the number of patches and f is the filter size for the convolutional architecture. C is the output dimension for the vector output case. P is the pool size for average pooling. (u) is defined as u^2 for square activation, and $au^2 + bu + c$ for degree-two polynomial activation.

	Non-convex objective	Convex formulation	Upper bound on critical width m
Poly (scalar)	$\prod_{j=1}^m (Xu_j)_j; y + \prod_{j=1}^m j j$ s.t. $ku_j k = 1$	(21)	$2(d+1)$
Poly (vector)	$\prod_{j=1}^m (Xu_j)_j; Y + \prod_{j=1}^m k j k_1$ s.t. $ku_j k = 1$	Supp. Mat. [4]	$2(d+1)C$
CNN	$\prod_{j=1}^m \prod_{k=1}^K (X_k u_j)_{jk}; y + \prod_{j=1}^m k j k_1$ s.t. $ku_j k = 1$	Supp. Mat. [4]	$2(f+1)K^2$
Pooling	$\prod_{j=1}^m \prod_{k=1}^{K/P} \prod_{l=1}^P (X_{(k-1)P+l} u_j)_{jk}; y + \prod_{j=1}^m k j k_1$ s.t. $ku_j k = 1$	Supp. Mat. [4]	$2(f+1)\frac{K^2}{P^2}$
Quad (scalar, cubic reg)	$\prod_{j=1}^m (Xu_j)_j; y + \prod_{j=1}^m j j$ s.t. $ku_j k = 1$, or $\prod_{j=1}^m (Xu_j)_j; y + \frac{1}{\varepsilon} \prod_{j=1}^m (j j^3 + ku_j k_2^3)$	Supp. Mat. [4]	d
Quad (scalar, quad reg)	$\prod_{j=1}^m (Xu_j)_j; y + \prod_{j=1}^m j j^{2-3}$ s.t. $ku_j k = 1$, or $\prod_{j=1}^m (Xu_j)_j; y + \frac{1}{\varepsilon} \prod_{j=1}^m (j j^2 + ku_j k_2^2)$	NP-hard (intractable)	-

We summarize the types of neural network architectures considered in this work and the corresponding convex problems in Table 1. The fourth column of Table 1 shows the upper bounds for critical width m , i.e., the optimal number of neurons that one needs for global optimization of any problems with number of neurons m .

2 Lifted Representations of Networks with Degree-Two Polynomial Activations

Consider the network $f(x) = \prod_{j=1}^m (x^T u_j)_j$ where the activation function is the degree-two polynomial $(u) = au^2 + bu + c$. First, we note that the neural network output can be written as

$$\begin{aligned}
 f(x) &= \prod_{j=1}^m (a(x^T u_j)^2 + bx^T u_j + c)_j = \prod_{j=1}^m (haxx^T; u_j u_j^T i + hb x; u_j i + c)_j \\
 &= \frac{1}{4} \text{vec}(axx^T) \begin{matrix} 3 \\ 2 \end{matrix} \frac{1}{5} \text{vec} \left(\prod_{j=1}^m u_j u_j^T \right) \begin{matrix} 3 \\ 4 \end{matrix} + \frac{1}{5} \prod_{j=1}^m u_j \\
 &= h(x); (\hat{f} u_j; j g_{j=1}^m) i; \tag{11}
 \end{aligned}$$

where $h: \mathbb{R}^d \rightarrow \mathbb{R}^{d^2+d+1}$ and $i: \mathbb{R}^{m(d+1)} \rightarrow \mathbb{R}^{d^2+d+1}$ are formally defined in the sequel. The above identity shows that the nonlinear neural network output is linear over the *lifted features*

$$h(x) := \text{vec}(axx^T); bx; c \in \mathbb{R}^{d^2+d+1};$$



Fig. 2 (Left) The Neural Cone C_2^1 described by $(u^2; u; 1) \in \mathbb{R}^3$ where $u_j \in \mathbb{R}; \|u\|_2 = 1$. (Right) Neural Spectrahedron $\mathcal{M}(1)$ described by $(Z_{11}; Z_{12}; Z_{22}) \in \mathbb{R}^3$ where $Z = \begin{bmatrix} Z_{11} & Z_{12} & Z_{13} \\ Z_{12} & Z_{22} & Z_{23} \\ Z_{13} & Z_{23} & Z_{33} \end{bmatrix} \succeq 0; Z_{11} + Z_{22} = Z_{33} = 1$ (constrained to the slice $Z_{22} = Z_{11}$ and $Z^0 = 0$ in (14)). Note that we display a slice of the spectrahedron, which is higher dimensional.

In turn, the nonlinear model $f(x)$ is completely characterized by the *lifted parameters* which we define as the following matrix-vector-scalar triplet

$$(f; u_j; g_{j=1}^m) := \begin{pmatrix} \mathcal{X}^n \\ u_j u_j^T \\ g_{j=1}^m \end{pmatrix} \begin{matrix} j; \\ j; \\ j; \end{matrix} \begin{pmatrix} \mathcal{X}^n \\ u_j \\ g_{j=1}^m \end{pmatrix} \begin{matrix} j; \\ j; \\ j; \end{matrix} \in \mathbb{R}^{d^2+d+1};$$

Optimizing over the lifted parameter space initially appears as hard as the original non-convex neural network training problem. This is due to the cubic and quadratic terms involving the weights of the hidden and output layer in the lifted parameters. Furthermore, norms of the network weights are nonlinear in the lifted parameters, which complicates regularization terms, e.g., $\sum_{j=1}^m \|u_j\|_2^2$ typically included in training. Nevertheless, one of our main results shows that the lifted parameters can be exactly described using linear matrix inequalities.

We begin by characterizing the lifted parameter space as a non-convex cone.

Definition 1 (Neural Cone of degree two) We define the non-convex cone $C_2^m \subset \mathbb{R}^{d^2+d+1}$ as

$$C_2^m := \begin{pmatrix} \mathcal{X}^n \\ u_j u_j^T \\ g_{j=1}^m \end{pmatrix} \begin{matrix} j; \\ j; \\ j; \end{matrix} \begin{pmatrix} \mathcal{X}^n \\ u_j \\ g_{j=1}^m \end{pmatrix} \begin{matrix} j; \\ j; \\ j; \end{matrix} : u_j \in \mathbb{R}^d; \|u_j\|_2 = 1; j \in \{1, \dots, m\} \quad (12)$$

See Figure 2 (left) for a depiction of $C_2^1 \subset \mathbb{R}^3$ corresponding to the case $m = 1; d = 1$.

Surprisingly, we will show that the original non-convex neural network problem is solved exactly to *global optimality* when the optimization is performed over a convex set which we define as the *Neural Spectrahedron*, given by the convex hull of the cone C_2 . In other words, every element of the convex hull can be associated with a neural network of the form $f(x) = \sum_{j=1}^m (x^T u_j) g_j$ through

a special matrix decomposition procedure which we introduce in the Neural Decomposition section. Moreover, a Neural Spectrahedron can be described by a simple linear matrix inequality. Consequently, these two results enable global optimization of neural networks with polynomial activations of degree two in polynomial time with respect to all problem parameters: dimension d , number of samples n and number of neurons m . To the best of our knowledge, this is the first instance of a method that globally optimizes a standard neural network architecture with computational complexity polynomial in all problem dimensions. We refer the reader to the recent work [38] for a convex optimization formulation of networks with ReLU activation, where the worst case computational complexity is $O((\frac{n}{r})^r)$ with $r = \text{rank}(X)$.

It is equally important that our results characterize neural networks as constrained linear learning methods $h(x)$ in the lifted feature space (x) , where the constraints on the lifted parameters are precisely described by a Neural Spectrahedron via linear matrix inequalities. These constraints can be easily tackled with convex semidefinite programming or closed-form projections onto these sets in iterative first-order algorithms.

Next, we describe a compact set that we call *neural spectrahedron* which describes the lifted parameter space of networks with a constraint on the ℓ_1 norm of output layer weights.

Definition 2 A neural spectrahedron $S_2^m(t) \subset \mathbb{R}^{d^2+d+1}$ is defined as the compact convex set

$$S_2^m(t) := \text{conv} \left(\begin{array}{l} \mathcal{X}^n \\ u_j u_j^T \quad j=1 \dots m \\ \mathcal{X}^n \\ u_j \quad j=1 \dots m \\ \mathcal{X}^n \\ j \quad j=1 \dots m \\ : k u_j k_2 = 1; \\ \mathcal{X}^n \\ j \in \mathbb{R}; \delta_j \in [m]; \quad j, j, t : \end{array} \right) \quad (13)$$

We will show that a neural spectrahedron can be equivalently described as a linear matrix inequality via defining $S_2^m(t) = \mathcal{M}_1(t); \mathcal{M}_2(t); \mathcal{M}_4(t)$ with $\mathcal{M}_1(t) \in \mathbb{R}^{d \times d}$, $\mathcal{M}_2(t) \in \mathbb{R}^{d \times 1}$ and $\mathcal{M}_4(t) \in \mathbb{R}$ for all $m \leq m$ where

$$\mathcal{M}(t) = \left(\begin{array}{l} Z \quad Z^0 : Z = \begin{array}{cc} Z_1 & Z_2 \\ Z_2^T & Z_4 \end{array} \quad 0; Z^0 = \begin{array}{cc} Z_1^0 & Z_2^0 \\ Z_2^{0T} & Z_4^0 \end{array} \quad 0; \\ \text{tr}(Z_1) = Z_4; \text{tr}(Z_1^0) = Z_4^0; Z_4 + Z_4^0 \leq t \end{array} \right) \quad (14)$$

$Z; Z^0 \in \mathbb{S}^{(d+1) \times (d+1)}$, $Z_1; Z_1^0 \in \mathbb{S}^{d \times d}$, $Z_2; Z_2^0 \in \mathbb{R}^{d \times 1}$ and $Z_4; Z_4^0 \in \mathbb{R}_+$, and $m = m(t)$ is a critical number of neurons that satisfies $m(0) = 0$ and $m(t) \leq 2(d+1)\delta t$, which will be explicitly defined in the sequel. Therefore, an efficient description of the set $\mathcal{M}(t)$ in terms of linear matrix inequalities enables efficient convex optimization methods in polynomial time.

2.1 A Direct Derivation of SDP Formulations of the Neural Spectrahedron

In this section, we show how to obtain SDP based formulations of the neural network problem using a direct approach. In Section 3, we will show that this direct formulation is identical to the one obtained through Lagrangian duality, mirroring the relaxations of non-convex quadratic programs shown in Section 1.4.

First, we separate positive and negative components of the parameters $f_j g_{j=1}^m$, and parameterize $S_2^m(t)$ as follows

$$S_2^m(t) = {}^+S_2^m(t_+) + {}^-S_2^m(t_-) : t_+; t_- \geq 0; t_+ + t_- = t; \quad (15)$$

where $+$ denotes the Minkowski sum and we defined the one-sided Neural Spectrahedron ${}^+S_2^m(t_+)$ as

$${}^+S_2^m(t) := \mathbf{conv} \left(\begin{array}{l} \mathcal{X}^n \\ u_j u_j^T \\ j=1 \end{array} ; \begin{array}{l} \mathcal{X}^n \\ u_j \\ j=1 \end{array} ; \begin{array}{l} \mathcal{X}^n \\ u_j \\ j=1 \end{array} : \|u_j\|_2 = 1; \right. \\ \left. \begin{array}{l} j \in \mathbb{R}_+; \delta_j \in [m]; \\ j=1 \end{array} \right) \quad (16)$$

We observe that ${}^+S_2^m(t)$ is identical to the set ${}^+M_1; {}^+M_2; {}^+M_4 \subseteq \mathbb{R}^{d^2+d+1}$ where

$${}^+M(t) := \mathbf{conv} \left(\begin{array}{l} \mathcal{X}^n \\ u_j \\ j=1 \end{array} ; \begin{array}{l} u_j \\ 1 \end{array} ; \begin{array}{l} u_j \\ 1 \end{array} : u_j \in \mathbb{R}^d; \|u_j\|_2 = 1; \right. \\ \left. \begin{array}{l} j \in \mathbb{R}_+; \delta_j \in [m]; \\ j=1 \end{array} \right) \quad (17)$$

which is partitioned as ${}^+M(t) = \begin{array}{l} {}^+M_1 \\ {}^+M_2^T \\ {}^+M_4 \end{array}$ and ${}^+M_1 \subseteq \mathbb{S}^{d \times d}; {}^+M_2 \subseteq \mathbb{R}^{d \times 1}$ and ${}^+M_4 \subseteq \mathbb{R}_+$.

Next, we note that as soon as the network width² satisfies $m \geq d+1$, we have

$${}^+M(t) := \mathbf{conv} \left(\left(\begin{array}{l} u \\ 1 \end{array} ; \begin{array}{l} u \\ 1 \end{array} : \|u\|_2 = 1 \right) \right) \quad (18)$$

where $\mathbf{0}$ is the zero matrix, since $\sum_{j=1}^m \begin{array}{l} u_j \\ 1 \end{array} ; \begin{array}{l} u_j \\ 1 \end{array} \in \mathbb{S}^{(d+1) \times (d+1)}$ is a positive semidefinite matrix, and hence can be factorized³ as a convex combination of at most $d+1$ rank-one matrices of the form $\begin{array}{l} u \\ 1 \end{array} ; \begin{array}{l} u \\ 1 \end{array}$. Note that

² This assumption on the width is relaxed in Theorem 1 (see Section 3).

³ We describe the details of this factorization in the Neural Decomposition section.

the zero matrix is included to account for the inequality $\prod_{j=1}^m u_j \geq 1$ in (17). This important observation enables us to represent the convex hull of the non-convex Neural Cone (an example is shown in Figure 2), via the simple convex body ${}^+\mathcal{M}(t)$ given in (18).

Most importantly, the positive Neural Spectrahedron set ${}^+\mathcal{M}(t)$ provides a representation of the non-convex Neural Cone \mathcal{C}_2^m via its extreme points. ${}^+\mathcal{M}(t)$ has a simple description as a linear matrix inequality provided in the following lemma (the proof can be found in the appendix).

Lemma 1 *For $m = d + 1$, it holds that*

$${}^+\mathcal{M}(t) = \{Z : Z = \begin{bmatrix} Z_1 & Z_2 \\ Z_2^T & Z_4 \end{bmatrix} \succeq 0; \text{tr}(Z_1) = Z_4 = t\} \quad (19)$$

Therefore the positive Neural Spectrahedron can be represented as the intersection of the positive semidefinite cone and linear inequalities. Moreover, every element of ${}^+\mathcal{M}(t)$ can be factorized as $\prod_{j=1}^m \begin{bmatrix} u_j & u_j^T \\ u_j^T & u_j \end{bmatrix}$ for some $\|u_j\|_2 = 1; u_j \succeq 0; \delta_j \in [m]; \prod_{j=1}^m u_j = t$, which can be identified as an element of the non-convex Neural Cone \mathcal{C}_2^m and a neural network in the lifted parameter space as shown in (11).

The assumption $m = d + 1$ is only used here for the simplicity of the illustration. We improve this condition on m , the width of the network, for the same SDP formulation in Section 3. In the following sections, we only require $m \geq m_0$, where m_0 can be determined via the convex SDP. Furthermore, the regularization parameter directly controls the number of neurons m . We illustrate the effect of the regularization parameter on m in the numerical experiments section, and show that m_0 can be made arbitrarily small.

3 Convex Optimization and Duality for Degree-Two Polynomial Activation Networks

We consider the non-convex training of a two-layer fully connected neural network with degree-two polynomial activations $f(x) = \prod_{j=1}^m (x^T u_j)$ and derive a convex dual optimization problem. Here, \cdot is the degree-two polynomial $\phi(u) = au^2 + bu + c$. This neural network has m neurons with the first layer weights $u_j \in \mathbb{R}^d$ and second layer weights $v_j \in \mathbb{R}$. We refer to this case where $f : \mathbb{R}^d \rightarrow \mathbb{R}$ as the scalar output case. The results are extended to the vector output case in the supplementary material [4].

It is relatively easy to obtain a weak dual that provides a lower-bound via Lagrangian duality. However, in non-convex problems, a duality gap may exist since strong duality does not hold in general. Remarkably, we show that strong duality holds as soon as the network width exceeds a threshold which can be easily determined.

We will assume ℓ_1 norm regularization on the second layer weights as regularization and include constraints that the first layer weights are unit norm.

We note that ℓ_1 norm regularization on the second layer weights results in a special dual problem and hence is crucial in the derivations. We show in Square Activation Networks section that this formulation is equivalent to cubic regularization when the activation function is the square activation. For the squared ℓ_2 norm regularization, i.e., weight decay, we will in fact show that the problem is NP-hard.

The training of a network under this setting requires solving the non-convex optimization problem given by

$$(P) \quad \min_{u_j \in \mathbb{R}^d; \delta_j \in \mathbb{R}; \delta_j \geq 0} \sum_{j=1}^m \sum_{i=1}^n (X u_j)_i^2 + \sum_{j=1}^m \delta_j \|u_j\|_2^2$$

$$\text{s.t.} \quad \|u_j\|_2 = 1; \delta_j \geq 0$$
(20)

Here, $\ell(\cdot; y)$ is a convex loss function, e.g., squared loss or logistic loss. We will use the notation $\rho := \min(P)$, i.e., the minimum value of the optimization problem (P) . Theorem 1 states the main result for degree-two polynomial activation neural networks that the non-convex optimization problem in (20) can be solved globally optimally via a convex problem.

Theorem 1 (Globally optimal convex program for degree-2 polynomial activation networks) *The solution of the convex problem*

$$(SDRP) \quad \min_{Z; Z^0 \in \mathcal{S}^{(d+1)}} \ell(\hat{y}; y) + \text{tr}(Z + Z^0)$$

$$\text{s.t.} \quad \hat{y}_i = a x_i^T (Z_1 - Z_1^0) x_i + b x_i^T (Z_2 - Z_2^0) + c (Z_4 - Z_4^0);$$

$$i = 1; \dots; n$$

$$\text{tr}(Z_1) = Z_4; \text{tr}(Z_1^0) = Z_4^0$$

$$Z = \begin{bmatrix} Z_1 & Z_2 \\ Z_2^T & Z_4 \end{bmatrix} \succeq 0; Z^0 = \begin{bmatrix} Z_1^0 & Z_2^0 \\ Z_2^{0T} & Z_4^0 \end{bmatrix} \succeq 0$$
(21)

provides a global optimal solution for the non-convex problem (P) in (20) when the number of neurons satisfies $m \geq m^*$ where

$$m^* = \text{rank}(Z) + \text{rank}(Z^0) \quad (22)$$

Z and Z^0 denote the optimal solution of the convex program in (21). The optimal network weights can be extracted from Z and Z^0 using the Neural Decomposition procedure. It follows that the optimal number of neurons is upper bounded by $m^* \leq 2(d+1)$.

To the best of our knowledge, Theorem 1 provides the first method that globally optimizes a standard non-trivial neural network architecture with computational complexity polynomial in all problem dimensions.

The convex SDP in (21) and the other SDP formulations presented in this paper can be solved using interior point algorithms. Since the Slater's condition holds for these SDPs, the interior point algorithm for these problems

are polynomial time in the real number model, where the representations of the numbers are ignored and only the number of iterations is bounded by a polynomial ([27]).

The proof of Theorem 1 consists of two parts. In the first part, we leverage convex duality to find a lower-bounding convex SDP in Section 3.1. In the second part, we show a matching upper bound via the proposed Neural Decomposition algorithm in Section 4.

3.1 Lower Bound Proof for Theorem 1

In this section we show that the solution of the convex program (21) provides a lower bound for the solution of the non-convex problem (20). In the Neural Decomposition section, we prove, via the method of neural decomposition, that the solution of the convex problem provides also an upper bound, which concludes the proof of Theorem 1.

In proving the lower bound, we leverage duality. Minimizing over first y_j 's and then u_j 's, we can restate the problem (P) given in (20) as

$$\begin{aligned} p = \min_{u_j \in S^d; \delta_j \in [m]} \min_{y_j \in \mathbb{R}^d; \delta_j \in [m]; y \in \mathbb{R}^n} & \langle \hat{y}; y \rangle + \sum_{j=1}^n \lambda_j \|y_j\|^2 \\ \text{s.t.} \quad & \hat{y} = \sum_{j=1}^n (X u_j) \delta_j; \end{aligned} \quad (23)$$

where the set S^d is defined as $S^d := \{u \in \mathbb{R}^d : \|u\|_2 = 1\}$. Here, we use the notation $\rho := \min(P)$. The dual problem for the inner minimization problem is given by

$$\max_{v \in \mathbb{R}^n} \langle v; \hat{y} \rangle \quad \text{s.t.} \quad jv^T (X u_j) \delta_j \leq \lambda_j; \delta_j \in [m]; \quad (24)$$

Next, let us call the optimal solution of the following problem d

$$\begin{aligned} d = \min_{u_j \in S^d; \delta_j \in [m]} \max_{v \in \mathbb{R}^n} & \langle v; \hat{y} \rangle \\ \text{s.t.} \quad & jv^T (X u_j) \delta_j \leq \lambda_j; \delta_j \in [m]; \end{aligned} \quad (25)$$

By changing the order of the minimization and maximization operations, we obtain the following bound

$$\begin{aligned} d & \leq \max_{v \in \mathbb{R}^n} \langle v; \hat{y} \rangle \\ \text{s.t.} \quad & jv^T (X u) \delta_j \leq \lambda_j; u \in S^d; \end{aligned} \quad (26)$$

$\mathbb{R}^{1 \times 1}$. We note that because of the symmetry of Z and Z^θ , we have $Z_2^T = Z_3$ and $Z_2^{\theta T} = Z_3^\theta$. The Lagrangian for the problem in (28) is

$$L(v; \lambda; \mu; Z; Z^\theta) = \lambda (v) + \mu_1 \text{tr}(Z_1) + \mu_2 \text{tr}(Z_1^\theta) - \sum_{i=1}^n a_i v_i x_i^T (Z_1 - Z_1^\theta) x_i + b v^T X (Z_2 - Z_2^\theta) + (\lambda_1) Z_4 + (\lambda_2) Z_4^\theta - \sum_{i=1}^n c_i v_i (Z_4 - Z_4^\theta). \quad (30)$$

Maximizing the Lagrangian with respect to v , λ , μ , we obtain the problem (SDRP) given in (21), which concludes the lower bound part of the proof. In Neural Decomposition section, we introduce a method for decomposing the solution of this convex program (i.e. Z and Z^θ) into feasible neural network weights to prove the upper bound.

4 Neural Decomposition

A lower bound on the optimal value of the non-convex problem (P) in (20) is obtained via the solution of the convex program (SDRP) in (21) that we have derived using Lagrangian duality. Now we show that this lower bound is in fact identical to the optimal value of the non-convex problem, thus proving strong duality. Our approach is based on proving an upper bound by constructing neural network weights from the solution of the convex problem such that the convex objective achieves the same objective as the non-convex objective. Suppose that $(Z; Z^\theta)$ is a solution to (21). Let us denote the rank of Z by r and the rank of Z^θ by r^θ . We will discuss the decomposition for Z and then complete the picture by considering the same decomposition for Z^θ . We begin by noting that Z satisfies the constraints of (21), i.e.,

$$Z \succeq 0; \text{tr}(Z_1) = Z_4; \text{or } \text{tr} \begin{bmatrix} I_d & 0 \\ 0 & \mathbf{1} \end{bmatrix} Z = 0. \quad (31)$$

Suppose that we have a decomposition of Z as a sum of rank-1 matrices such that $Z = \sum_{j=1}^r p_j p_j^T$ where $p_j \in \mathbb{R}^{d+1}$ and $\text{tr}(p_j p_j^T G) = p_j^T G p_j = 0$ for $j = 1; \dots; r$. We show how this can always be done in the following subsection by introducing a new matrix decomposition method, dubbed the *neural decomposition* procedure.

Letting $p_j := c_j^T d_j$ with $c_j \in \mathbb{R}^d$ and $d_j \in \mathbb{R}$, we note that $p_j^T G p_j = 0$ implies $k c_j k_2^2 = d_j^2$. We may assume $p_j \neq 0$; d_j in the decomposition (otherwise we can simply remove zero components), implying $k c_j k_2^2 > 0$; d_j . Furthermore, this expression for p_j 's allows us to establish that

$$\sum_{j=1}^r p_j p_j^T = \sum_{j=1}^r \begin{bmatrix} c_j \\ d_j \end{bmatrix} \begin{bmatrix} c_j^T & d_j \end{bmatrix} = \sum_{j=1}^r \begin{bmatrix} c_j c_j^T & c_j d_j \\ d_j c_j^T & d_j^2 \end{bmatrix} = \begin{bmatrix} Z_1 & Z_2 \\ Z_3 & Z_4 \end{bmatrix}. \quad (32)$$

As a result, we have the following decompositions:

$$Z_1 = \prod_{j=1}^r c_j c_j^T = \prod_{j=1}^r u_j u_j^T k c_j k_2^2 = \prod_{j=1}^r u_j u_j^T d_j^{\rho} \quad (33)$$

$$Z_2 = \prod_{j=1}^r c_j d_j = \prod_{j=1}^r u_j d_j k c_j k_2 = \prod_{j=1}^r u_j d_j j d_j j \quad (34)$$

$$Z_4 = \prod_{j=1}^r d_j^{\rho}; \quad (35)$$

where we have introduced the normalized weights $u_j = \frac{c_j}{k c_j k_2}$, $j = 1; \dots; r$. If $d_j = 0$ for some j , we redefine the corresponding ρ_j as ρ_j , which does not modify the decomposition $\prod_{j=1}^r \rho_j \rho_j^T$ and the equality $\rho_j^T G \rho_j = 0$. Hence, without loss of generality, we can assume that $d_j = 0$ for all $j = 1; \dots; r$, which leads to

$$Z_1 = \prod_{j=1}^r u_j u_j^T d_j^{\rho}; \quad Z_2 = \prod_{j=1}^r u_j d_j^{\rho}; \quad Z_4 = \prod_{j=1}^r d_j^{\rho}; \quad (36)$$

Similarly for Z^0 , we will form the following decompositions:

$$Z_1^0 = \prod_{j=1}^{r^0} u_j^{\rho} u_j^{\rho T} d_j^{\rho^2}; \quad Z_2^0 = \prod_{j=1}^{r^0} u_j^{\rho} d_j^{\rho^2}; \quad Z_4^0 = \prod_{j=1}^{r^0} d_j^{\rho^2}; \quad (37)$$

Considering the decompositions for both Z and Z^0 , finally we obtain a neural network with first layer weights as $[u_1; \dots; u_r; u_1^{\rho}; \dots; u_{r^0}^{\rho}]g$, and second layer weights as $[d_1^{\rho}; \dots; d_r^{\rho}; d_1^{\rho^2}; \dots; d_{r^0}^{\rho^2}]g$. We note that this corresponds to a neural network with $r + r^0$ neurons. If both Z and Z^0 are full rank, then we will have $2(d+1)$ neurons, which is the maximum.

To see why we can use the decompositions of Z and Z^0 to construct neural network weights, we plug the expressions (36) and (37) in the objective of the convex program in (21):

$$\begin{aligned} & \sum_{j=1}^r j d_j^{\rho} j + \sum_{j=1}^{r^0} j d_j^{\rho^2} j; \quad \text{where} \\ \hat{y}_i = a x_i^T & \prod_{j=1}^r u_j u_j^T d_j^{\rho} + \prod_{j=1}^{r^0} u_j^{\rho} u_j^{\rho T} (d_j^{\rho^2}) x_i + b x_i^T \prod_{j=1}^r u_j d_j^{\rho} + \prod_{j=1}^{r^0} u_j^{\rho} (d_j^{\rho^2}) \\ & + c \prod_{j=1}^r d_j^{\rho} + \prod_{j=1}^{r^0} (d_j^{\rho^2}); \quad i = 1; \dots; n; \end{aligned} \quad (38)$$

We note that this expression exactly matches the optimal value of the non-convex objective in (20) for a neural network with $r + r^0$ neurons. Also, the unit norm constraints on the first layer weights are satisfied (hence feasible) since

u_j 's and u_j^l 's are normalized. This establishes that the neural network weights obtained from the solution of the convex program provide an upper bound for the minimum value of the original non-convex problem. Consequently, we have shown that the optimal solution of the convex problem (SDRP) in (21) provides a global optimal solution to the non-convex problem (P) in (20) and this concludes the proof of Theorem 1.

4.1 Neural Decomposition Procedure

Here we describe the procedure for computing the decomposition given by $Z = \sum_{j=1}^r p_j p_j^T = 0$ such that $p_j^T G p_j = 0, j = 1; \dots; r$ where $G = \begin{bmatrix} I_d & 0 \\ 0 & 1 \end{bmatrix}$. The computational complexity of this procedure is at most $O(d^3)$, which is dominated by the cost of one Eigenvalue Decomposition of Z in the initial step. This algorithm is inspired by the constructive proof of the S-procedure given in Lemma 2.4 of [39] with modifications to account for the equalities $p_j^T G p_j = 0$.

Neural Decomposition for PSD Matrices:

0. **Compute a rank-1 decomposition** $Z = \sum_{j=1}^r p_j p_j^T$.
This can be done with the Cholesky decomposition, or the eigenvalue decomposition $Z = \sum_{j=1}^r q_j q_j^T$. Since $Z \succeq 0$, we have $\lambda_j > 0$, for $j = 1; \dots; r$. Then we can obtain the desired rank-1 decomposition $Z = \sum_{j=1}^r p_j p_j^T$ by defining $p_j = \frac{1}{\sqrt{\lambda_j}} q_j, j = 1; \dots; r$.
1. **If $p_1^T G p_1 = 0$, return $y = p_1$. If not, find a $j \in \{2; \dots; r\}$ such that $(p_1^T G p_1)(p_j^T G p_j) < 0$.**
We know such j exists since $\text{tr}(Z G) = \sum_{j=1}^r p_j^T G p_j = 0$ (this is true since it is one of the constraints of the convex program), and $p_1^T G p_1 \neq 0$. Hence, for at least one $j \in \{2; \dots; r\}$, $p_j^T G p_j$ must have the opposite sign as $p_1^T G p_1$.
2. **Return $y = \frac{p_1 + p_j}{1 + \frac{p_1^T G p_1}{p_j^T G p_j}}$ where $\alpha \in \mathbb{R}$ satisfies $(p_1 + \alpha p_j)^T G (p_1 + \alpha p_j) = 0$.**
We know that such α exists since the quadratic equation

$$(p_1 + \alpha p_j)^T G (p_1 + \alpha p_j) = \alpha^2 p_j^T G p_j + 2 \alpha p_1^T p_j + p_1^T G p_1 = 0 \quad (39)$$

has real solutions as the discriminant $4(p_1^T p_j)^2 - 4(p_1^T G p_1)(p_j^T G p_j)$ is positive due to step 1 where we picked j such that $(p_1^T G p_1)(p_j^T G p_j) < 0$. To find α , we simply solve the quadratic equation for α .

3. **Update $r = r - 1$, and then the vectors $p_1; \dots; p_r$ as follows:**
Remove p_1 and p_j and insert $u = \frac{p_1 + p_j}{1 + \frac{p_1^T G p_1}{p_j^T G p_j}}$. Consequently, we will be dealing with the updated matrix $Z = Z - yy^T$ in the next iteration, which is of rank $r - 1$:

$$Z - yy^T = uu^T + \sum_{i=2; i \neq j}^r p_i p_i^T \quad (40)$$

Fig. 3 Illustration of the neural decomposition procedure for $d = 2$ (i.e. $Z \in \mathbb{R}^{3 \times 3}$). The dashed red arrows correspond to the eigenvectors of Z ($q_1; q_2; q_3$) and the solid blue arrows show the decomposed vectors p_1 and p_2 . In this example, the rank of Z is 2 where q_1 and q_2 are its two principal eigenvectors. The eigenvalue corresponding to the eigenvector q_3 is zero. The light blue colored surface shows the Lorentz cones $z = \sqrt{x^2 + y^2}$ and $z = -\sqrt{x^2 + y^2}$. We observe that the decomposed vectors p_1 and p_2 lie on the boundary of Lorentz cones.

Note that Step 0 is carried out only once and then steps 1 through 3 are repeated $r - 1$ times. At the end of $r - 1$ iterations, we are left with the rank-1 matrix $p_1 p_1^T$ which satisfies $p_1^T G p_1 = 0$ since initial Z satisfies $\text{tr}(Z G) = 0$ and the following $r - 1$ updates are of the form $y y^T$ which satisfies $y^T G y = 0$. If we denote the returned y vectors as y_i for the iteration i and y_r is the last one we are left with, then y_i 's satisfy the desired decomposition that $Z = \sum_{i=1}^r y_i y_i^T$ and $y_i^T G y_i = 0, i = 1; \dots; r$.

Figure 3 is an illustration of the neural decomposition procedure for a toy example with $d = 2$ where the eigenvectors of Z and the vectors p_j are plotted together. Due to the constraints $p_j^T G p_j = 0, j = 1; 2$, the vectors p_j have to lie on the boundary of Lorentz cones $z = \sqrt{x^2 + y^2}$ and $z = -\sqrt{x^2 + y^2}$. Decomposing the solution of the convex problem Z and Z^0 onto these cones, i.e., neural decomposition, enables the construction of neural network weights from Z and Z^0 .

5 Standard Weight Decay Formulation is NP-Hard

In the previous sections, we have studied the training of two-layer neural networks and derived a convex program whose solution globally optimizes the non-convex problem. In contrast, Theorem 2 states that if we have quadratic regularization (i.e. weight decay), the resulting optimization problem given in

⁴ In special relativity, Lorentz cones describe the path that a flash of light, emanating from a single event traveling in all directions takes through spacetime (see Figure 1.3.1 in [32]).

(41) is an NP-hard problem:

$$(P_{WD}) \quad \min_{\substack{u_j \in \mathbb{R}^d; \\ j \in \{1, \dots, m\}}} \sum_{j=1}^n (Xu_j)^2 - y + \sum_{j=1}^n (\|u_j\|_2^2 + ku_j k_2^2): \quad (41)$$

Theorem 2 The two-layer neural network optimization problem with square activation and standard squared ℓ_2 norm regularization, i.e., weight decay, in (41) is NP-hard for any value of m as $n \rightarrow \infty$.

The proof of Theorem 2 is deferred to the appendix. The proof relies on the polynomial time reduction of the problem P_{WD} from the NP-hard problem of phase retrieval. This result illustrates the type of regularization function plays a significant role in the computational tractability of the optimal training problem.

Remark 1 It follows from Theorem 2 that the two-layer neural network training problem with degree-two polynomial activation and unit norm first layer weights and $\sum_{j=1}^n \|u_j\|_2^p$ as the regularization term with $p < 1$ is also NP-hard for $n \rightarrow \infty$ since it reduces to the square activation case for the polynomial coefficients $a = 1; b = 0; c = 0$.

6 Numerical Results

In this section, we present numerical results that verify the presented theory of the convex formulations along with experiments comparing the test set performance of the derived formulations. All experiments have been run on a MacBook Pro with 16GB RAM.

We have used CVXPY [11, 1] for solving the convex SDPs. In particular, we have used the open source solver SCS (splitting conic solver) [35, 36] in CVXPY, which is a scalable first order solver for convex cone problems.

Furthermore, we have solved the non-convex problems via backpropagation for which we have used PyTorch [37]. We have used the SGD algorithm for the non-convex models. For all the experiments involving SGD in this section, we show only the results corresponding to the best learning rate that we select via an offline hyperparameter search. The momentum parameter is 0.9. In the plots, the non-convex models are either labeled as 'Backpropagation (GD)' or 'Backpropagation (SGD)'. The first one, short for gradient descent, means that the batch size is equal to the number of samples, and the second one, short for stochastic gradient descent, means that the batch size is not and the exact batch size is explicitly stated in the figure captions.

Figure 4 compares the solution of the non-convex problem via backpropagation and the solution of the corresponding convex program via a convex solver. The training and test costs are shown for a regression task with randomly generated data for the two-layer square activation neural network. We observe that convex SDP takes a much shorter time to optimize and obtains a

Fig. 4 Cost against wall-clock time on the training (left) and test (right) sets for stochastic gradient descent (SGD) and the convex SDP for square activation networks. The solid lines show the training curve of the non-convex model with SGD (with learning rate tuned optimally via extensive grid search) and each different colored solid curve corresponds to an independent trial. The dotted horizontal curve shows the cost for the convex SDP and the cross indicates the time that it takes to solve the convex SDP. The dataset X is synthetically generated by sampling from the i.i.d. Gaussian distribution and has dimensions $n = 100; d = 10$. Labels y are generated by a teacher network with 10 planted neurons. The regularization coefficient is $\lambda = 10^{-6}$ and the batch size for SGD is 10.

globally optimal solution while the SGD algorithm converges to local minima in some of the trials where the initialization is different.

Figure 5 compares the accuracy of the non-convex degree-two polynomial activation model when it is trained with different optimizers (SGD and Adam [24]) for a range of step sizes. The plots in Figure 5 show results for the convolutional neural network model. We observe that the convex formulations outperform the non-convex solution via SGD and Adam. The extension of the main result to convolutional neural networks is discussed in the supplementary material [4].

6.1 Comparison with ReLU Networks

We compare the classification accuracies for degree-two polynomial activation and ReLU activation in Figure 6 on three different binary classification UCI datasets. The regularization coefficient has been picked separately for degree-two polynomial activation and ReLU activation networks to maximize the accuracy. Figure 6 demonstrates that the convex SDP shows competitive accuracy performance and faster run times compared to ReLU networks.

6.2 Regularization Parameter

Figure 7 shows how the accuracy changes as a function of the regularization coefficient λ for the convex problem for two-layer degree-two polynomial activation networks. Figure 7 highlights that the choice of the regularization coefficient is critical in the accuracy performance. In plot a, we see that the value of λ that maximizes the test set accuracy is $\lambda = 10$ for which the optimal number of neurons m is near 20. We note that for the dataset in plot a, the

(a) CNN, MNIST, training

(b) CNN, MNIST, test

(c) CNN, CIFAR, training

(d) CNN, CIFAR, test

(e) FC, oocytes, training

(f) FC, oocytes, test

Fig. 5 Classification accuracy for various learning rates and optimizers are plotted on the same figure. SGD and Adam are used in solving the non-convex optimization problem. The solid blue lines each correspond to a different learning rate for SGD and each dashed green curve corresponds to a different learning rate for the Adam algorithm. Plots a, b: CNN with degree-two polynomial activations and global average pooling for binary classification on the first two classes of the MNIST dataset (12000 training samples). Plots c, d: The same architecture as plots a, b and the dataset is the first two classes of the CIFAR-10 dataset (10000 training samples). Plots e, f: Fully connected architecture for binary classification on the dataset oocytes-merluccius-nucleus-4d.

optimal number of neurons is upper bounded by $m \leq 2(d+1) = 32$. Similarly for plot b, the best choice for the regularization coefficient is $\lambda = 1$ and the optimal number of neurons for $\lambda = 1$ is near 40. Furthermore, we observe that a higher value for λ tends to translate to a lower optimal number of neurons m^* (plotted on the right vertical axis). Even though the convex optimization problem in (21) has a fixed number of variables (in this case, $2(d+1)^2$) for a given dataset, a low number of neurons is still preferable for many reasons such as inference speed. We observe that the number of neurons can be controlled via the regularization coefficient λ .

(a) DS1, training accuracy

(b) DS1, test accuracy

(c) DS2, training accuracy

(d) DS2, test accuracy

(e) DS3, training accuracy

(f) DS3, test accuracy

Fig. 6 Comparison of classification accuracies for neural networks with ReLU activation, degree-two polynomial activation ($a = 0.09; b = 0.5; c = 0.47$), and the convex SDP. DS1: dataset 1 is the oocytes-merluccius-nucleus-4d ($n = 817; d = 41$), DS2: dataset 2 is the credit approval dataset ($n = 552; d = 15$), DS3: dataset 3 is the breast cancer dataset ($n = 228; d = 9$).

(a) credit approval

(b) ionosphere

Fig. 7 Accuracy (left vertical axis) and optimal number of neurons (right vertical axis) against the regularization coefficient on binary classification datasets. These results have been obtained using the convex program in (21). The dimensions of the datasets are $n = 552; d = 15$ for plot a and $n = 280; d = 33$ for plot b.

7 Discussion

In this paper, we have studied the optimization of two-layer neural networks with degree-two polynomial activations. We have shown that regularization plays an important role in the tractability of the problems associated with neural network training. We have developed convex programs for the cases where the regularization leads to tractable formulations. Convex formulations are useful since they have many well-known advantages over non-convex optimization such as having to optimize fewer hyperparameters and no risk of getting stuck at local minima.

The methods presented in this work optimize the neural network parameters in a higher dimensional space in which the problem becomes convex. For fully connected neural networks with square activation, the standard non-convex problem requires optimizing n neurons (i.e. n -dimensional first layer weight and a 1-dimensional second layer weight per neuron). The convex program for this neural network finds the optimal network parameters in the lifted space \mathbb{S}^d . For degree-two polynomial activations, convex optimization takes place for Z and Z^0 in $\mathbb{S}^{(d+1)}$. We note that the dimensions of the convex programs are polynomial with respect to all problem dimensions. In contrast, the convex program of [38] has $\mathcal{O}(H)$ variables where H grows exponentially with respect to the rank of the data matrix.

We have used the SCS solver with CVXPY for solving the convex problems in the numerical experiments. It is important to note that there is room for future work in terms of which solvers to use. Solvers specifically designed for the presented convex programs could enjoy faster run times.

The scope of this work is limited to two-layer neural networks. We note that it is a promising direction to consider the use of our convex programs for two-layer neural networks as building blocks in learning deep neural networks. Many recent works such as [2] and [5] investigate layerwise learning algorithms for deep neural networks. The training of individual layers in layerwise learning could be improved by the presented convex programs since the convex programs can be efficiently solved and eliminate much of the hyperparameter tuning involved in standard neural network training.

Declarations

Funding This work was partially supported by the National Science Foundation under grants IIS-1838179, ECCS-2037304, Facebook Research and the Army Research Office.

Conflicts of interest/Competing interests The authors have no conflicts of interest to declare that are relevant to the content of this article.

Availability of data and material (data transparency) The data used in this work are publicly available.

Code availability (software application or custom code) Code will be made available at an online repository.

Ethics approval Not applicable.

Consent to participate No participants, not applicable.

Consent for publication No participants, not applicable.

References

1. Agrawal, A., Verschueren, R., Diamond, S., Boyd, S.: A rewriting system for convex optimization problems. *Journal of Control and Decision* 5(1), 42{60 (2018)
2. Allen-Zhu, Z., Li, Y.: Backward feature correction: How deep learning performs deep learning. *arXiv preprint arXiv:2001.04413* (2020)
3. Arora, R., Basu, A., Mianjy, P., Mukherjee, A.: Understanding deep neural networks with rectified linear units. In: 6th International Conference on Learning Representations, ICLR 2018 (2018)
4. Bartan, B., Pilanci, M.: Neural spectrahedra and semidefinite lifts: Global convex optimization of polynomial activation neural networks in fully polynomial-time. *CoRR abs/2101.02429* (2021). URL <https://arxiv.org/abs/2101.02429>
5. Belilovsky, E., Eickenberg, M., Oyallon, E.: Greedy layerwise learning can scale to imagenet. *CoRR abs/1812.11446* (2018). URL <http://arxiv.org/abs/1812.11446>
6. Bienstock, D., Muñoz, G., Pokutta, S.: Principled deep neural network training through linear programming (2018)
7. Blondel, M., Fujino, A., Ueda, N.: Convex factorization machines. *European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML PKDD)* (2015)
8. Blondel, M., Niculae, V., Otsuka, T., Ueda, N.: Multi-output polynomial networks and factorization machines. In: *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS'17*, p. 3351{3361. Curran Associates Inc., Red Hook, NY, USA (2017)
9. Boyd, S., Vandenberghe, L.: *Convex optimization*. Cambridge university press (2004)
10. Burer, S.: Copositive programming. In: *Handbook on semidefinite, conic and polynomial optimization*, pp. 201{218. Springer (2012)
11. Diamond, S., Boyd, S.: CVXPY: A Python-embedded modeling language for convex optimization. *Journal of Machine Learning Research* 17(83), 1{5 (2016)
12. Du, S., Lee, J.: On the power of over-parametrization in neural networks with quadratic activation. In: J. Dy, A. Krause (eds.) *Proceedings of the 35th International Conference on Machine Learning, Proceedings of Machine Learning Research*, vol. 80, pp. 1329{1338. PMLR, Stockholmstrassan, Stockholm Sweden (2018). URL <http://proceedings.mlr.press/v80/du18a.html>
13. Ergen, T., Pilanci, M.: Convex geometry of two-layer relu networks: Implicit autoencoding and interpretable models. In: *International Conference on Artificial Intelligence and Statistics*, pp. 4024{4033. PMLR (2020)
14. Ergen, T., Pilanci, M.: Implicit convex regularizers of cnn architectures: Convex optimization of two- and three-layer networks in polynomial time. *arXiv preprint arXiv:2006.14798* (2020)
15. Ergen, T., Pilanci, M.: Revealing the structure of deep neural networks via convex duality. *arXiv preprint arXiv:2002.09773* (2020)
16. Fickus, M., Mixon, D.G., Nelson, A.A., Wang, Y.: Phase retrieval from very few measurements. *arXiv preprint arXiv:1307.7176* (2013)
17. Gamarnik, D., Kızıldağ, E.C., Zadik, I.: Stationary points of shallow neural networks with quadratic activation function. *arXiv preprint arXiv:1912.01599* (2020)
18. Gilad-Bachrach, R., Dowlin, N., Laine, K., Lauter, K., Naehrig, M., Wernsing, J.: Cryptonets: Applying neural networks to encrypted data with high throughput and accuracy. In: *International Conference on Machine Learning*, pp. 201{210 (2016)
19. Goel, S., Kanade, V., Klivans, A., Thaler, J.: Reliably learning the relu in polynomial time. In: S. Kale, O. Shamir (eds.) *Proceedings of the 2017 Conference on Learning Theory, Proceedings of Machine Learning Research*, vol. 65, pp. 1004{1042. PMLR, Amsterdam, Netherlands (2017)

20. Goemans, M.X., Williamson, D.P.: Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming. *Journal of the ACM* 42, 1115{1145 (1995)
21. Hastad, J.: Some optimal inapproximability results. *Journal of the ACM (JACM)* 48(4), 798{859 (2001)
22. Hesamifard, E., Takabi, H., Ghasemi, M.: Cryptodl: towards deep learning over encrypted data. In: *Annual Computer Security Applications Conference (ACSAC 2016)*, Los Angeles, California, USA, vol. 11 (2016)
23. Khot, S., Kindler, G., Mossel, E., O'Donnell, R.: Optimal inapproximability results for max-cut and other 2-variable csps? *SIAM Journal on Computing* 37(1), 319{357 (2007)
24. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings (2015)*. URL <http://arxiv.org/abs/1412.6980>
25. Lacotte, J., Pilanci, M.: All local minima are global for two-layer relu neural networks: The hidden convex optimization landscape. *arXiv preprint arXiv:2006.05900* (2020)
26. Laurent, M., Poljak, S.: On a positive semidefinite relaxation of the cut polytope. *Linear Algebra and its Applications* 223(224), 439{461 (1995)
27. Laurent, M., Rendl, F.: Semidefinite programming and integer programming. In: K. Aardal, G. Nemhauser, R. Weismantel (eds.) *Discrete Optimization, Handbooks in Operations Research and Management Science*, vol. 12, pp. 393{514. Elsevier (2005). DOI [https://doi.org/10.1016/S0927-0507\(05\)12008-8](https://doi.org/10.1016/S0927-0507(05)12008-8). URL <https://www.sciencedirect.com/science/article/pii/S0927050705120088>
28. Lederer, J.: No spurious local minima: on the optimization landscapes of wide and deep neural networks (2020)
29. Livni, R., Shalev-Shwartz, S., Shamir, O.: On the computational efficiency of training neural networks. *NIPS'14* p. 855{863 (2014)
30. Mannelli, S.S., Vanden-Eijnden, E., Zdeborová, L.: Optimization and generalization of shallow neural networks with quadratic activation functions. *arXiv preprint arXiv:2006.15459* (2020)
31. Mishra, P., Lehmkuhl, R., Srinivasan, A., Zheng, W., Popa, R.A.: Delphi: A cryptographic inference system for neural networks. *PPMLP'20* p. 27{30 (2020). DOI 10.1145/3411501.3419418. URL <https://doi.org/10.1145/3411501.3419418>
32. Naber, G.L.: *The geometry of Minkowski spacetime: An introduction to the mathematics of the special theory of relativity*, vol. 92. Springer Science & Business Media (2012)
33. Nesterov, Y., Wolkowicz, H., Ye, Y.: Semidefinite programming relaxations of nonconvex quadratic optimization. In: *Handbook of semidefinite programming*, pp. 361{419. Springer (2000)
34. Obla, S., Gong, X., Alou, A., Hu, P., Takabi, D.: Effective activation functions for homomorphic evaluation of deep neural networks. *IEEE Access* 8, 153098{153112 (2020)
35. O'Donoghue, B., Chu, E., Parikh, N., Boyd, S.: Conic optimization via operator splitting and homogeneous self-dual embedding. *Journal of Optimization Theory and Applications* 169(3), 1042{1068 (2016). URL <http://stanford.edu/~boyd/papers/scs.html>
36. O'Donoghue, B., Chu, E., Parikh, N., Boyd, S.: SCS: Splitting conic solver, version 2.1.2. <https://github.com/cvxgrp/scs> (2019)
37. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., Chintala, S.: Pytorch: An imperative style, high-performance deep learning library. *Advances in Neural Information Processing Systems* 32 pp. 8024{8035 (2019). URL <http://papers.nips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>
38. Pilanci, M., Ergen, T.: Neural networks are convex regularizers: Exact polynomial-time convex optimization formulations for two-layer networks. *Proceedings of the International Conference on Machine Learning (ICML 2020)* (2020)
39. Polik, I., Terlaky, T.: A survey of the s-lemma. *SIAM Review* 49(3), 371{418 (2007). DOI 10.1137/S003614450444614X. URL <https://doi.org/10.1137/S003614450444614X>

40. Ramachandran, P., Zoph, B., Le, Q.: Searching for activation functions. arXiv preprint arXiv:1710.05941 (2018). URL <https://arxiv.org/pdf/1710.05941.pdf>
41. Sahiner, A., Ergen, T., Pauly, J., Pilanci, M.: Vector-output relu neural network problems are copositive programs: Convex analysis of two layer networks and polynomial-time algorithms. arXiv preprint arXiv:2012.13329 (2020)
42. Sahiner, A., Mardani, M., Ozturkler, B., Pilanci, M., Pauly, J.: Convex regularization behind neural reconstruction. arXiv preprint arXiv:2012.05169 (2020)
43. Soltani, M., Hegde, C.: Fast and provable algorithms for learning two-layer polynomial neural networks. *IEEE Transactions on Signal Processing* **67**(13), 3361–3371 (2019). DOI 10.1109/TSP.2019.2916743
44. Soltanolkotabi, M., Javanmard, A., Lee, J.D.: Theoretical insights into the optimization landscape of over-parameterized shallow neural networks. *IEEE Trans. Inf. Theor.* **65**(2), 742–769 (2019). DOI 10.1109/TIT.2018.2854560. URL <https://doi.org/10.1109/TIT.2018.2854560>
45. Trevisan, L., Sorkin, G.B., Sudan, M., Williamson, D.P.: Gadgets, approximation, and linear programming. *SIAM Journal on Computing* **29**(6), 2074–2097 (2000)
46. Wolkowicz, H., Saigal, R., Vandenberghe, L.: *Handbook of semidefinite programming: theory, algorithms, and applications*, vol. 27. Springer Science & Business Media (2012)

Appendix

This appendix contains proofs not included in the main body of the paper.

A Proof of Theorem 2

This section breaks down the proof of Theorem 2. At the core of the proof is the polynomial time reduction of the problem from the NP-hard problem of phase retrieval.

A.1 Equivalent Problem Representation

The optimization problem for training a two-layer fully connected neural network with square activation and quadratic regularization can be stated as

$$\rho := \min_{\substack{u_j \in \mathbb{R}^d; \\ j \in \{1, \dots, m\}}} \left[\sum_{j=1}^m (Xu_j)^2 - j; y + \frac{1}{c} \sum_{j=1}^m (j^2 + ku_j k_2^2) \right]; \quad (42)$$

where the scaling factor c is the same as before (i.e. $c = 2^{\frac{1}{3}} + 2^{-\frac{2}{3}} \approx 1.88988$). Rescaling $u_j \leftarrow u_j t_j^{1=2}$ and $j \leftarrow j=t_j$ for $t_j > 0, j = 1; \dots; m$, we obtain the following equivalent optimization problem

$$\rho = \min_{\substack{u_j \in \mathbb{R}^d; \\ j \in \{1, \dots, m\}}} \left[\sum_{j=1}^m (Xu_j)^2 - j; y + \frac{1}{c} \sum_{j=1}^m (j^2=t_j^2 + ku_j k_2^2 t_j) \right]; \quad (43)$$

Note that the regularization term is convex in t_j for $t_j > 0$. Optimizing the regularization term with respect to t_j leads to $t_j = 2^{1=3} \frac{j j j}{ku_j k_2}^{2=3}$ and plugging this in yields

$$\rho = \min_{\substack{u_j \in \mathbb{R}^d; \\ j \in \{1, \dots, m\}}} \left[\sum_{j=1}^m (Xu_j)^2 - j; y + \sum_{j=1}^m j j^{2=3} ku_j k_2^{4=3} \right]; \quad (44)$$

Defining scaled weights $j^{\theta} = j ku_j k_2^2$ and $u_j^{\theta} = u_j = ku_j k_2$, we obtain the equivalent problem

$$\rho = \min_{\substack{u_j^{\theta} \in \mathbb{R}^d; j^{\theta} \in \mathbb{R} \\ \text{s.t. } ku_j^{\theta} k_2 = 1; j \in \{1, \dots, m\}}} \left[\sum_{j=1}^m (Xu_j^{\theta})^2 - j^{\theta}; y + \sum_{j=1}^m j^{\theta} j^{\theta 2=3} \right]; \quad (45)$$

This shows that solving the standard weight decay formulation is equivalent to having the power of $2=3$ on the second layer weights as regularization and unit norm first layer weights as constraints.

A.2 Hardness Result

We design a data matrix such that the solution coincides with solving the phase retrieval problem which is NP-hard (see [16]). We consider the equality constrained version of (45), i.e., $\beta = 0$, which is given by

$$\begin{aligned} & \min_{u_j \in \mathbb{R}^d; j \in [m]} \sum_{j=1}^m \|u_j\|^2 \\ & \text{s.t. } \|u_j\|_2 = 1; j \in [m] \\ & \text{s.t. } \sum_{j=1}^m (Xu_j)^2 = y \end{aligned} \quad (46)$$

A.2.1 Addition of a Simplex Constraint

Let the first d rows of the data matrix X be e_1^T, \dots, e_d^T and let the first d entries of y be $1/d$. Then, the constraint $\sum_{j=1}^m (Xu_j)^2 = y$ implies

$$\sum_{j=1}^m u_{jk}^2 = 1/d \text{ for } k = 1; \dots; d \quad (47)$$

Summing the above for all $k = 1; \dots; d$, and noting that $\sum_{k=1}^d u_{jk}^2 = 1$ lead to the constraint $\sum_{j=1}^m 1 = 1$.

A.2.2 Reduction from the NP-Hard Phase Retrieval and Subset Sum Problem

We let $X = [I; \tilde{X}]$ and $y = [\frac{1}{d} \mathbf{1}; \tilde{y}]$ to obtain the simplex constraint $\sum_{j=1}^m 1 = 1$ as shown in the previous subsection. In this case, the optimization problem reduces to

$$\begin{aligned} & \min_{u_j \in \mathbb{R}^d; j \in [m]} \sum_{j=1}^m \|u_j\|^2 \\ & \text{s.t. } \|u_j\|_2 = 1; j \in [m] \\ & \text{s.t. } \sum_{j=1}^m (\tilde{X}u_j)^2 = \tilde{y} \\ & \text{s.t. } \sum_{j=1}^m u_{jk}^2 = 1/d; \quad k = 1; \dots; d \\ & \text{s.t. } \sum_{j=1}^m 1 = 1 \end{aligned} \quad (48)$$

Suppose that there exists a feasible solution $f_j; u_j; g_{j=1}^m$, which satisfies $\|u_j\|_2 = 1$, where $u_1 = 1$ and $u_1^T u_1 = 1$ with only one nonzero neuron. Then, it follows from Lemma 3 that this solution is strictly optimal. Consequently, the

problem in (48) is equivalent to

$$\begin{aligned} & \text{find}_{u_1 \in \mathbb{R}^d} u_1 \\ & \text{s.t. } (\tilde{x}_i^T u_1)^2 = \tilde{y}_i; \quad i = 1; \dots; (n-d) \\ & \quad u_{1k}^2 = 1-d; \quad k = 1; \dots; d; \end{aligned} \quad (49)$$

Lemma 3 (ℓ_p minimization recovers 1-sparse solutions when $0 < p < 1$)

Consider the optimization problem

$$\begin{aligned} & \min_{j \in \{1, \dots, m\}^n} \sum_{i=1}^n |j_i|^p \\ & \text{s.t. } j \in C; \end{aligned} \quad (50)$$

where C is a convex set and $p \in (0, 1)$. Suppose that there exists a feasible solution $j \in C$ and $\sum_{i=1}^n |j_i| = 1$ such that $k_0 = 1$. Then, j is strictly optimal with objective value 1. More precisely, any solution with cardinality strictly greater than 1 has objective value strictly larger than 1.

A.2.3 NP-hardness Proof

Subset sum problem given in Definition 3 is a decision problem known to be NP-complete (e.g. [16]). The decision version of the problem in (49) can be stated as follows: Does there exist a feasible u_1 ? We show that this decision problem is NP-hard via a polynomial time reduction from the subset sum problem.

Definition 3 (Subset sum problem) Given a set of integers A , does there exist a subset A_S whose elements sum to z ?

Lemma 4 establishes the reduction of the decision version of (49) from the subset sum problem. The proof is given in the sequel and follows the same approach used in the proof for the NP-hardness of phase retrieval in [16], with the main difference being the additional constraints $u_{1k}^2 = 1-d, k = 1; \dots; d$ in (49). Finally, Lemma 4 concludes the proof of Theorem 2.

Lemma 4 Consider the problem in (49). Let the first d samples of $\tilde{X} \in \mathbb{R}^{(d+1) \times d}$, denoted $\tilde{X}_D \in \mathbb{R}^{d \times d}$, be any diagonal matrix with -1 's and $+1$'s on its diagonal, and let the $(d+1)$ 'st sample be $\tilde{x}_{d+1} = \sqrt{d} [a_1 \dots a_d]^T$. Then, the decision version of the resulting problem returns 'yes' if and only if the answer for the subset sum problem with $A = \{a_1, \dots, a_d\}$ is 'yes'.

B Proof of Lemma 1

Proof We will denote the set in (18) as S_1 and the set in (19) as S_2 to simplify the notation. We will prove $S_1 = S_2$ by showing $S_1 \subseteq S_2$ and $S_2 \subseteq S_1$.

We first show $S_1 \subseteq S_2$. Let us take a point $S \in S_1$. This implies that S is a matrix of the form

$$t \sum_{j=1}^m \begin{bmatrix} u_j & \\ & 1 \end{bmatrix} \begin{bmatrix} u_j^T \\ 1 \end{bmatrix} = t \sum_{j=1}^m \begin{bmatrix} u_j u_j^T & \\ & 1 \end{bmatrix} = t \sum_{j=1}^m \begin{bmatrix} u_j u_j^T & \\ & 1 \end{bmatrix} \begin{bmatrix} u_j \\ 1 \end{bmatrix} \begin{bmatrix} u_j^T \\ 1 \end{bmatrix} \quad (51)$$

with $\|u_j\|_2 = 1$ and $\|u_j\|_2 = 1$ for all j . We note that $\text{tr}(t \sum_{j=1}^m u_j u_j^T) = t \sum_{j=1}^m \text{tr}(u_j u_j^T) = t \sum_{j=1}^m 1 = t$. This shows that S satisfies the equality condition in the definition (19). Now, we show that S is

a PSD matrix. Note that each of the rank-1 matrices $\begin{bmatrix} u_j & \\ & 1 \end{bmatrix} \begin{bmatrix} u_j^T \\ 1 \end{bmatrix}$ is a PSD matrix and since the coefficients t 's and 1 's are nonnegative, it follows that S is PSD. This proves that $S \in S_2$.

We next show $S_2 \subseteq S_1$. Let us take a point $S \in S_2$. This implies that S is PSD and $\text{tr}(S) = S_{44} = t_0 = t$. We show in the Neural Decomposition section that it is possible to decompose S via the neural decomposition procedure to obtain the expressions given in (36). It follows that we can write S in the following form

$$S = \sum_{j=1}^m \frac{t_0}{d_j^2} \begin{bmatrix} u_j u_j^T d_j^2 \\ u_j^T d_j^2 \end{bmatrix} \begin{bmatrix} u_j d_j^2 \\ d_j^2 \end{bmatrix}; \quad (52)$$

where the scaling factor $\frac{t_0}{\sum_{j=1}^m d_j^2}$ is to ensure that $\text{tr}(S) = S_{44} = t_0 = t$. It is obvious to see that S is in S_1 when $t_0 = t$ by the definition of S_1 given in (18). When $t_0 < t$, we still have that S is in S_1 which can be seen by noting that S_1 is defined as the convex hull of rank-1 matrices and the zero matrix. We can scale all the rank-1 matrices in the convex combination with $\frac{t_0}{t}$ and change the weight of the zero matrix accordingly.