

Global Optimality in Neural Network Training

Benjamin D. Haeffele and René Vidal

Johns Hopkins University, Center for Imaging Science, Baltimore, MD 21218, USA.

bhaeffele@jhu.edu rvidal@cis.jhu.edu

Abstract

The past few years have seen a dramatic increase in the performance of recognition systems thanks to the introduction of deep networks for representation learning. However, the mathematical reasons for this success remain elusive. A key issue is that the neural network training problem is nonconvex, hence optimization algorithms may not return a global minima. This paper provides sufficient conditions to guarantee that local minima are globally optimal and that a local descent strategy can reach a global minima from any initialization. Our conditions require both the network output and the regularization to be positively homogeneous functions of the network parameters, with the regularization being designed to control the network size. Our results apply to networks with one hidden layer, where size is measured by the number of neurons in the hidden layer, and multiple deep subnetworks connected in parallel, where size is measured by the number of subnetworks.

1. Introduction

As a broad definition, feed-forward deep networks are a collection of feature extraction layers, where each layer applies some form of linear transformation (e.g., convolution, dot-product), followed by a nonlinearity (e.g., rectification, max-pooling) to the output of the preceding layer. Modern networks are similar to classical neural networks, except that they involve many more layers and typically replace classical sigmoid nonlinearities by linear rectification (i.e., ReLU units). These modifications, together with the availability of massive amounts of data for training, have led to dramatic improvements in classification performance for various applications in computer vision, speech and natural language processing. However, the mathematical reasons for this success remain elusive.

An important mathematical challenge is that the problem of learning the parameters of a deep network is non-convex, hence optimization algorithms can get stuck in non-global minima. In contrast, local minimizers of a convex optimization problem are also global minimizers, so convex formu-

lations of learning problems are often preferable as they facilitate the analysis of the properties of the learning algorithm. This is one of the reasons for the popularity of classical learning algorithms such as linear regression, support vector machines, ℓ_1 minimization, and nuclear norm minimization, all of which involve solving a *convex optimization problem* of the form:

$$\min_X \ell(Y, \Phi(X, S)) + \lambda \Theta(X). \quad (1)$$

For classification problems, $\ell(Y, \Phi(X, S))$ is a *loss function* that measures the agreement between the true labels, Y , and the predicted labels, $\Phi(X, S)$, where S is the input data to be classified and X represents the classifier parameters, while $\Theta(X)$ is a *regularization function* designed to prevent overfitting. Convex formulations require both the loss function and the regularization function to be convex on X , e.g., $\ell(Y, \Phi(X, S)) = \|Y - S^T X\|_F^2$ and $\Theta(X) = \|X\|_F^2$.

Unfortunately, in practice many learning algorithms – and particularly those that seek to learn an appropriate representation of features directly from the data, such as principal component analysis (PCA), low-rank matrix completion, nonnegative matrix factorization, sparse dictionary learning, tensor factorization and deep learning – involve solving a *non-convex optimization problem* such as:

$$\min_{\{W^i\}_{i=1}^K} \ell(Y, \Phi(W^1, \dots, W^K)) + \lambda \Theta(W^1, \dots, W^K), \quad (2)$$

where Φ is an arbitrary, convexity destroying mapping.¹ For example, in deep neural network training, the output of the network is typically generated by applying an alternating series of linear and non-linear functions to the input data:

$$\Phi(W^1, \dots, W^K) = \psi_K(\psi_{K-1}(\dots \psi_2(\psi_1(S^T W^1)W^2) \dots W^{K-1})W^K), \quad (3)$$

where each W^i is an appropriately sized matrix that contains the connection weights between layers $i - 1$ and i of the network, and the $\psi_i(\cdot)$ functions apply some form of

¹For the sake of notational simplicity, we will omit the dependency of Φ on the data, S , from now on.

non-linearity after each matrix multiplication, e.g., a sigmoid function, rectification, max-pooling.²

For a very small number of non-convex problems, e.g., PCA, one is fortunate, and a global minimizer can be found in closed form. For other problems, e.g., ℓ_0 minimization, rank minimization, and low-rank matrix completion, one can replace the non-convex objective by a convex surrogate and show that under certain conditions the solutions to both problems are the same [9, 5]. In most cases, however, the optimal solutions cannot be computed in closed form, and a good convex surrogate may not be easy to find. This presents significant challenges to existing optimization algorithms – including (but certainly not limited to) alternating minimization, gradient descent, stochastic gradient descent, block coordinate descent, back-propagation, and quasi-Newton methods – which are typically only guaranteed to converge to a critical point of the objective function [16, 20, 25, 26]. However, for non-convex problems, the set of critical points includes not only global minima, but also local minima, local maxima, saddle points and saddle plateaus, as illustrated in Figure 1. As a result, the non-convexity of the problem leaves the model somewhat ill-posed in the sense that it is not just the model formulation that is important but also implementation details, such as how the model is initialized and particulars of the optimization algorithm, which can have a significant impact on the performance of the model.

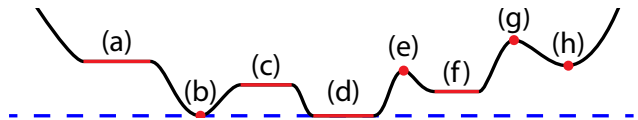


Figure 1. Example critical points of a non-convex function (shown in red). (a,c) Plateaus. (b,d) Global minima. (e,g) Local maxima. (f,h) Local minima.

To address the issue of non-convexity, a common strategy used in deep learning is to initialize the network weights, $\{W^k\}$, at random, update these weights using local descent, check if the training error decreases sufficiently fast, and if not, choose another initialization. In practice, it has been observed that if the size of the network is large enough, this strategy can lead to markedly different solutions for the network weights, which give nearly the same objective values and classification performance [6]. It has also been observed that when the size of the network is large enough and the nonlinearity is chosen to be a Rectified Linear Unit (ReLU), $\psi^+(x) = \max(x, 0)$, in lieu of a sigmoid function, many weights are zero, a phenomenon known as *dead neurons*, and the classification performance significantly improves [7, 15, 14, 27]. While this empir-

²Here we have shown the linear operations to be simple matrix multiplications to simplify notation, but this easily generalizes to other linear operators (e.g., convolution).

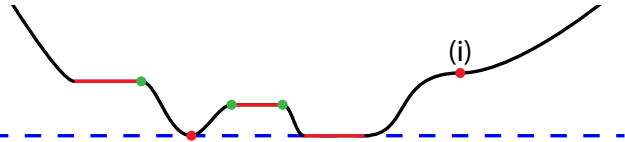


Figure 2. Guaranteed properties of the proposed framework. Starting from any initialization, a non-increasing path exists to a global minimizer. Starting from points on a plateau, a simple “sliding” method exists to find the edge of the plateau (green points).

ically suggests that when the size of the network is large enough and ReLU nonlinearities are used *all local minima could be global*, there is currently no rigorous theory that provides a precise mathematical explanation for these experimentally observed phenomena.

Paper Contributions. In this paper, we study conditions under which the optimization landscape for the non-convex optimization problem in (2) is such that *all critical points are either global minimizers or saddle points/plateaus*, as shown in Figure 2. We show that if the network size is large enough and the functions Φ and Θ are *sums of positively homogeneous functions of the same degree*, any local minimizer such that *some of its entries are zero* is also a global minimizer. Interestingly, ReLU and max-pooling nonlinearities are positively homogeneous, while sigmoids are not, which could provide a possible explanation for the improved performance of ReLU and max pooling. Furthermore, many state-of-the-art networks are not trained with *classical regularization*, such as an ℓ_1 or ℓ_2 norm penalty on the weight parameters but instead rely on techniques such as dropout. Our results also provide strong guidance on the design of network regularization to ensure the non-existence of spurious local minima, showing that traditional weight decay is not appropriate for deep networks. However, more recently proposed forms of regularization such as Path-SGD [18] or batch normalization [12] can be easily incorporated into our analysis framework, and stochastic regularization methods, such as a dropout [23], also have strong similarities to our framework.

Related Work. Prior work on global optimality of neural network training [3] showed that for neural networks with a single hidden layer, if the number of neurons in the hidden layer is not fixed, but instead fit to the data through a sparsity inducing regularization, then the process of training a globally optimal neural network is analogous to selecting a finite number of hidden units from a potentially infinite dimensional space of all possible hidden units. A weighted sum of the selected hidden units is then taken to produce the output. The specific optimization problem is of the form

$$\min_x \ell(Y, \sum_i h_i(S)x_i) + \lambda \|x\|_1, \quad (4)$$

where $h_i(S)$ represents one possible hidden unit activation

in response to the training data S from an infinite dimensional space $h_i(S) \in \mathcal{H}$ of all possible hidden unit activations. Clearly (4) is a convex optimization problem on x (assuming $\ell(Y, X)$ is convex on X) and straightforward to solve for a finite set of $h_i(S)$ activations. However, since \mathcal{H} is an infinite dimensional space the primary difficulty lies in how to select the appropriate hidden unit activations. Nonetheless, by using arguments from gradient boosting, it is possible to show that problem (4) can be globally optimized by sequentially adding hidden units to the network until one can no longer find a hidden unit whose addition will decrease the objective function [3, 10, 17].

Several recent works have also explored the error surface of multilayer neural networks using tools derived from random matrix theory and statistical physics. Applying ideas from random matrix theory to high-dimensional non-convex optimization, the authors of [8] argue that, under certain assumptions, for high-dimensional optimization problems if one is given a particular critical point, it is vastly more likely that the critical point will be a saddle point rather than a local minimizer and thus avoiding saddle points is the key difficulty in high-dimensional, non-convex optimization. Using arguments from statistical physics, the authors of [6] show that, under certain assumed distributions of the training data and the network weight parameters, as the number of hidden units in a network increases, the distribution of local minima becomes increasingly concentrated in a small band of objective function values near the global optimum (and thus all local minima become increasingly close to being global minima).

Additional recent work has analyzed the problem of training neural networks with a single hidden layer by estimating high order statistical moments of the network mapping using tensor decomposition methods and show that, with sufficient assumptions on the loss and data distribution, polynomial-time training is possible [13]. Further, the authors of [21] study the problem of when a given initialization of a neural network is likely to be within the basin of attraction of a global minimizer and provide conditions that ensure a random initialization will be within the basin of a global minimizer with high probability.

Our results will largely echo ideas from the above work, but we take a markedly different approach. Specifically, we will analyze the problem using a purely deterministic approach which does not make any assumptions regarding the distribution of the input data, the network weight parameter statistics, or the network initialization. With this approach, we will show that saddle points and plateaus are the *only* critical points that one needs to be concerned with due to the fact that for networks of sufficient size, local minima that require one to climb the objective surface to escape from, such as (f) and (h) in Figure 1, are guaranteed not to exist, as illustrated in Figure 2.

2. Problem Formulation

In this paper, we will study the non-convex problem:

$$\min_{r \in \mathbb{N}^+} \min_{\{W^i\}_{i=1}^K} \ell(Y, \Phi_r(W^1, \dots, W^K)) + \lambda \Theta_r(W^1, \dots, W^K). \quad (5)$$

For classification problems, Y typically contains the labels of the training examples, but in general this could be any arbitrary set of desired network outputs. K describes the number of weight layers in the network; the $\{W^k\}$ variables describe the parameters of different layers; and Φ_r defines the output of the network as a function of the network parameters. As an example, $K = 2$ is a network with one hidden layer where W^1 defines the weights from the input to the hidden layer, and W^2 defines the weights from the hidden layer to the output (Figure 3, left). The integer r defines the *size* of the network. For example, r could be the number of neurons in the hidden layer of a two-layer neural network (Figure 3, left), or the number of parallel sub-networks in a deep network (Figure 3, right). Note that a key feature of our formulation is that we optimize over the size of the network r , so as a result we need a means to control the network size and prevent overfitting. This is accomplished through a *regularization function* $\Theta_r(W^1, \dots, W^K)$, while the *loss function* ℓ measures how well Y is approximated by the network output, $\Phi_r(W^1, \dots, W^K)$, and $\lambda > 0$ balances the trade-off between regularization and loss.

3. Motivation: Matrix Factorization

Before considering neural networks, as a motivating example consider the following matrix factorization problem: Given a matrix $Y \in \mathbb{R}^{d_1 \times d_2}$, find factors $W^1 \in \mathbb{R}^{d_1 \times r}$ and $W^2 \in \mathbb{R}^{d_2 \times r}$ of small size, r , and small Frobenius norm that approximate Y as $W^1 W^{2\top}$, i.e.:

$$\min_{r \in \mathbb{N}^+} \min_{W^1, W^2} \frac{1}{2} \|Y - W^1 W^{2\top}\|_F^2 + \frac{\lambda}{2} (\|W^1\|_F^2 + \|W^2\|_F^2). \quad (6)$$

Notice that this problem is a particular case of (5) where $K = 2$; the loss function is chosen as the squared loss $\ell(Y, X) = \frac{1}{2} \|Y - X\|_F^2$; the factorization map Φ_r in (3) reduces to matrix multiplication - i.e., $\Phi_r(W^1, W^2) = W^1 W^{2\top}$, $S = I$, $\psi_i(Z) = Z$; and the regularization function is chosen as ℓ_2 (Tykhonov) regularization, i.e., $\Theta(W^1, W^2) = \frac{1}{2} (\|W^1\|_F^2 + \|W^2\|_F^2)$.

For the matrix factorization problem above, the regularizer is convex on (W^1, W^2) , but the overall objective is not due to the product $W^1 W^{2\top}$. While this makes the optimization problem in (6) non-convex, we can still analyze it by recalling the variational form of the nuclear norm $\|X\|_*$ (the sum of the singular values of the matrix X), which is

given by [22]:

$$\|X\|_* = \min_r \min_{W^1, W^2: W^1 W^{2\top} = X} \frac{1}{2} (\|W^1\|_F^2 + \|W^2\|_F^2). \quad (7)$$

The strong similarity between (7) and the regularizer in (6) suggests looking at the *convex* problem:

$$\min_X \frac{1}{2} \|Y - X\|_F^2 + \lambda \|X\|_*, \quad (8)$$

which is a classical nuclear norm minimization problem, whose solution can be found in closed form from the SVD of Y . Additionally, well-known results from positive semidefinite optimization have shown that although (6) is a non-convex optimization problem, all local minima of (6) will be globally optimal provided r is initialized to be sufficiently large [4, 2, 19, 11]. However, while for the particular case of (6) the problem can be easily recast as a semi-definite optimization problem [19], for alternative choices of regularization functions, Θ , this quickly becomes a non-trivial problem. Further, results from semi-definite optimization apply to problems of the form

$$\min_{r \in \mathbb{N}^+} \min_{W^1, W^2} \ell(Y, W^1 W^{2\top}) + \lambda \bar{\Theta}(W^1 W^{2\top}) \quad (9)$$

which are subtly (but critically) different from the problem we consider here, as we consider regularization on the factors directly, $\Theta(W^1, W^2)$, instead of on the product of the factors, $\bar{\Theta}(W^1 W^{2\top})$, resulting in problems that are typically considerably more challenging [1].

Nevertheless, we build on these ideas from matrix factorization to analyze a wide range of non-convex optimization problems, including neural network training and significant generalizations of the matrix factorization problem given in (6), and present a framework where a convex optimization problem, e.g., (8), provides an achievable lower bound of the non-convex problem of interest, e.g., (6). From this strong coupling between the convex and non-convex problems we then derive sufficient conditions to verify if a local-minimizer is a global-minimizer and show that if the size of the variables in the non-convex problem is initialized to be sufficiently large then from *any* initialization it is possible to reach a global-minimizer using purely local descent.

4. Neural Networks with One Hidden Layer

The above discussion on matrix factorization can be extended to neural networks with one hidden layer by properly adjusting the definitions of the maps Φ_r and Θ_r . In matrix factorization, Φ_r can be re-written as $\Phi_r(W^1, W^2) = W^1 W^{2\top} = \sum_{i=1}^r W_i^1 W_i^{2\top}$, where W_i^1 and W_i^2 are the i th columns of W^1 and W^2 , respectively. Likewise, Θ_r can be re-written as $\Theta_r(W^1, W^2) = \frac{1}{2} (\|W^1\|_F^2 + \|W^2\|_F^2) = \sum_{i=1}^r \frac{1}{2} (\|W_i^1\|_2^2 + \|W_i^2\|_2^2)$. This motivates the following

more general definitions for Φ_r and Θ_r :

$$\begin{aligned} \Phi_r(W^1, W^2) &= \sum_{i=1}^r \phi(W_i^1, W_i^2) \quad \text{and} \\ \Theta_r(W^1, W^2) &= \sum_{i=1}^r \theta(W_i^1, W_i^2), \end{aligned} \quad (10)$$

where ϕ and θ are positively homogeneous of degree 2, i.e., $\phi(\alpha w^1, \alpha w^2) = \alpha^2 \phi(w^1, w^2)$ for all $\alpha \geq 0$. Clearly, $\phi(w^1, w^2) = w^1 w^{2\top}$ and $\theta(w^1, w^2) = \|w^1\|_2^2 + \|w^2\|_2^2$ satisfy this property. But notice that it is also satisfied, for example, by the map $\phi(w^1, w^2) = \psi^+(S^\top w^1) w^{2\top}$, where recall $\psi^+(x, 0) = \max(x, 0)$ is a ReLU applied to each entry of $S^\top w^1$. The fundamental observation is that both linear transformations and ReLU nonlinearities³ are positively homogeneous functions of degree one, and so the output of a two-layer network is positively homogeneous of degree two.

With these definitions, it is easy to see that the output of a two-layer neural network with nonlinearity ψ^+ on the hidden units, such as the one illustrated in the left panel of Figure 3, can be expressed by the map Φ_r in (10), where r now represents the number of neurons in the hidden layer. Therefore, we can write the training problem for a two-layer network as:

$$\min_{r \in \mathbb{N}^+} \min_{W^1, W^2} \ell(Y, \Phi_r(W^1, W^2)) + \lambda \Theta_r(W^1, W^2). \quad (11)$$

This problem is non-convex due to the mapping Φ_r , so to analyze this non-convex problem, we define a generalization of the nuclear norm in (7) for two-layer neural networks as:

$$\Omega_{\phi, \theta}(X) = \min_{r \in \mathbb{N}^+} \min_{W^1, W^2: \Phi_r(W^1, W^2) = X} \Theta_r(W^1, W^2). \quad (12)$$

The intuition behind the above problem is that, given an output X generated by the network for some input S , we wish to find the network size r and weights (W^1, W^2) that produced the output X . Among all possible sizes and weights, we prefer those that minimize $\Theta_r(W^1, W^2)$.

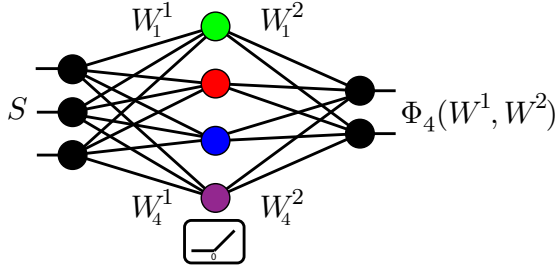
While $\Omega_{\phi, \theta}$ is no longer necessarily a norm, we will show in Proposition 1 that, under some additional conditions on θ , $\Omega_{\phi, \theta}$ is still convex. Therefore, if the loss ℓ is convex on X , so is the problem

$$\min_X \ell(Y, X) + \Omega_{\phi, \theta}(X). \quad (13)$$

As shown in Theorem 1, the convex problem (13) gives an achievable lower bound to the non-convex problem (11), and a local minimizer of the non-convex problem such that one of the columns of W^1 and W^2 is equal to zero gives a global minimizer for both the convex and non-convex problems.

³Notice that many other neural network operators such as max-pooling and convolution are also positively homogeneous.

ReLU Network with One Hidden Layer



Multilayer Parallel Network

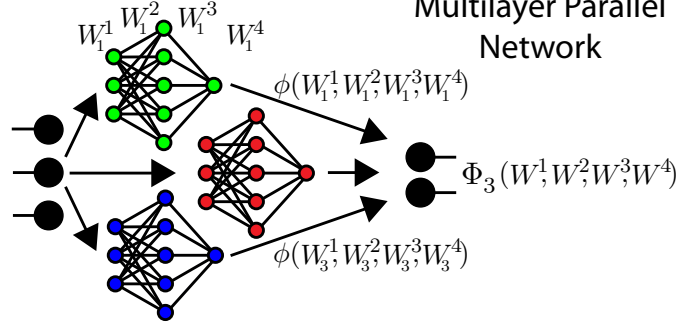


Figure 3. Example networks. (Left panel) ReLU network with a single hidden layer with the mapping Φ_r described by the equation in (10) with $(r = 4)$. Each color corresponds to one element of the elemental mapping $\phi(W_i^1, W_i^2)$. The colored hidden units have rectifying nonlinearities, while the black units are linear. (Right panel) Multilayer ReLU network with 3 fully connected parallel subnetworks $(r = 3)$, where each color corresponds to the subnetwork described the elemental mapping $\phi(W_i^1, W_i^2, W_i^3, W_i^4)$.

5. Deep Networks with Parallel Structure

In this section, we extend our results to networks formed by the addition of r parallel sub-networks, with each sub-network having the same architecture (see Figure 3, right). After introducing some specialized notation for the network weights and dimensions, we generalize the maps Φ and Θ to the context of deep networks. We then introduce a regularizer on the network weights and show that it induces a convex regularizer on the output. Finally, we state our main results to give sufficient conditions to guarantee that local minima are global minima and that for sufficiently large networks all local minima are guaranteed to be global minima.

Notation. We will use capital letters as a shorthand for a set of dimensions, and individual dimensions will be denoted with lower case letters. For example, $X \in \mathbb{R}^{d_1 \times \dots \times d_N} \equiv X \in \mathbb{R}^D$ for $D = d_1 \times \dots \times d_N$; we also denote the cardinality of $X \in \mathbb{R}^D$ as $\text{card}(X) = \prod_{i=1}^N d_i$. Similarly, $X \in \mathbb{R}^{D \times R} \equiv X \in \mathbb{R}^{d_1 \times \dots \times d_N \times r_1 \times \dots \times r_M}$ for $D = d_1 \times \dots \times d_N$ and $R = r_1 \times \dots \times r_M$. Given an element from a tensor space, we will use a subscript to denote a slice of the tensor along the last dimension. For example, given a matrix $W \in \mathbb{R}^{d_1 \times r}$, then $W_i \in \mathbb{R}^{d_1}, i \in \{1, \dots, r\}$, denotes the i 'th column of W . Similarly, given a cube $W \in \mathbb{R}^{d_1 \times d_2 \times r}$ then $W_i \in \mathbb{R}^{d_1 \times d_2}, i \in \{1, \dots, r\}$, denotes the i 'th slice along the third dimension. Further, given two tensors with matching dimensions except for the last dimension, $W \in \mathbb{R}^{D \times r_w}$ and $Q \in \mathbb{R}^{D \times r_q}$, we will use $[W Q] \in \mathbb{R}^{D \times (r_w + r_q)}$ to denote the concatenation of the two tensors along the last dimension. We'll also need the following definitions:

1. A size- r set of K factors $(W^1, \dots, W^K)_r \in \mathbb{R}^{(D^1 \times r)} \times \dots \times \mathbb{R}^{(D^K \times r)}$ is defined to be a set of K tensors where the final dimension of each tensor is equal to r .
2. A function $\theta : \mathbb{R}^{D^1} \times \dots \times \mathbb{R}^{D^K} \rightarrow \mathbb{R}^D$ is positively homogeneous with degree p if, $\forall \alpha \geq 0$,

$\theta(\alpha w^1, \dots, \alpha w^K) = \alpha^p \theta(w^1, \dots, w^K)$. Note that this implies that $\theta(0, \dots, 0) = 0$ for $p \neq 0$.

3. A function $\theta : \mathbb{R}^{D^1} \times \dots \times \mathbb{R}^{D^K} \rightarrow \mathbb{R}_+$ is positive semidefinite if $\theta(0, \dots, 0) = 0$ and $\theta(w^1, \dots, w^K) \geq 0, \forall (w^1, \dots, w^K)$.

Factorization and regularization maps. The maps Φ_r and Θ_r are defined as sums of positively homogeneous elemental mappings ϕ and θ , i.e.:

$$\begin{aligned} \Phi_r(W^1, \dots, W^K) &= \sum_{i=1}^r \phi(W_i^1, \dots, W_i^K) \quad \text{and} \\ \Theta_r(W^1, \dots, W^K) &= \sum_{i=1}^r \theta(W_i^1, \dots, W_i^K). \end{aligned} \quad (14)$$

Note that for matrix factorization, the elemental mapping $\phi : \mathbb{R}^{d_1} \times \mathbb{R}^{d_2} \rightarrow \mathbb{R}^{d_1 \times d_2}$ is defined as $\phi(w^1, w^2) = w^1 w^{2T}$, which is positively homogeneous of degree 2, and the factorization map $\Phi_r(W^1, W^2) = \sum_{i=1}^r W_i^1 W_i^{2T} = W_1 W_2^T$ is simply matrix multiplication for matrices with r columns. For neural networks, a typical elemental mapping ϕ would be defined as $\phi(w^1, \dots, w^K) = \psi_K(\dots \psi_2(\psi_1(S^T w^1) w^2) \dots w^K)$, which denotes the application of a linear transformation (dot-product, convolution) with parameters w^1 to the data S followed by a positively homogeneous nonlinearity ψ_i (ReLU, max-pooling), and so on for K weight layers. Therefore, the map Φ_r in (14) corresponds to the addition of r deep sub-networks in parallel, each one with the same number of layers and the same number of neurons per layer (see Figure 3, right). The well-known AlexNet network from [14], which consists of a series of convolutional layers, linear-rectification, max-pooling layers, response normalization layers, and fully connected layers, can be described by taking $r = 1$ and defining ϕ to be the entire transformation of the network (with slight modification of the response normalization layers, which are not positively homogenous, see supplement).

Note, however, that our results will rely on r potentially changing or being initialized to be sufficiently large, which limits the applicability of our results to current state-of-the-art network architectures (see discussion).

The *elemental regularization function* $\theta : \mathbb{R}^{D^1} \times \dots \times \mathbb{R}^{D^K} \rightarrow \mathbb{R}_+ \cup \infty$ takes as input the parameters of one sub-network and returns a non-negative number. The requirements we place on θ are that it must be positively homogeneous and positive semidefinite.

A regularizer on the parameters of the network that induces a convex regularizer on its output. To define our regularization function on the output of the network, $X = \Phi_r(W^1, \dots, W^K)$, it will be necessary that the elemental regularization function, θ , and the elemental mapping, ϕ , satisfy a few properties to be considered *compatible* for the definition of our regularization function. Specifically, we say that (ϕ, θ) are a *nondegenerate pair* if: 1) θ and ϕ are both positively homogeneous with degree p , for some $p > 0$ and 2) $\theta(z^1, \dots, z^K) > 0, \forall (z^1, \dots, z^K)$ such that $\phi(z^1, \dots, z^K) \neq 0$ and for all sequences $(z^1[n], \dots, z^K[n]), n = 1, \dots, \infty$, if $\|\phi(z^1[n], \dots, z^K[n])\| \rightarrow \infty$ then $\theta(z^1[n], \dots, z^K[n]) \rightarrow \infty$.⁴

Notice that any norm $\|w\|$ is positively homogeneous with degree 1, so by taking products of norms or sums of norms raised to an appropriate power we can match the degree of positive homogeneity between the mapping and regularizer. A typical regularizer for a mapping of degree K could be $\theta(w^1, \dots, w^K) = \prod_{i=1}^K \|w^i\|$ or $\theta(w^1, \dots, w^K) = \sum_{i=1}^K \|w^i\|^K$, where the choice of norms is arbitrary. However, as we place very few requirements on θ , our framework can include an extremely wide variety of potential regularization functions. For example, indicator functions on conic sets, such as requiring the factors to be non-negative, are also positively homogeneous and can be incorporated by our framework.

Given a nondegenerate pair (ϕ, θ) of an elemental mapping ϕ and an elemental regularization function θ , we define the *factorization regularization function*, $\Omega_{\phi, \theta}(X) : \mathbb{R}^D \rightarrow \mathbb{R}_+ \cup \infty$ to be

$$\Omega_{\phi, \theta}(X) \equiv \inf_{r \in \mathbb{N}^+} \inf_{(W^1, \dots, W^K)_r} \sum_{i=1}^r \theta(W_i^1, \dots, W_i^K) \quad (15)$$

s.t. $\Phi_r(W^1, \dots, W^K) = X$

with the additional condition that $\Omega_{\phi, \theta}(X) = \infty$ if $X \notin \bigcup_r \text{Im}(\Phi_r)$.

The following proposition shows that $\Omega_{\phi, \theta}$ is a convex function of X and that in general the infimum in (15) can

⁴Property 1 from the definition of a nondegenerate pair will be critical to our formulation. Property 2 is typically satisfied for most ‘interesting’ choices of (ϕ, θ) and is designed to avoid ‘pathological’ $\Omega_{\phi, \theta}$ functions (such as $\Omega_{\phi, \theta}(X) = 0 \forall X$).

be achieved with a finitely sized network (i.e., r does not need to approach ∞)⁵.

Proposition 1 *The function $\Omega_{\phi, \theta} : \mathbb{R}^D \rightarrow \mathbb{R} \cup \infty$ as defined in (15) has the following properties*

1. $\Omega_{\phi, \theta}$ is positive definite, i.e., $\Omega_{\phi, \theta}(0) = 0$ and $\Omega_{\phi, \theta}(X) > 0 \forall X \neq 0$.
2. $\Omega_{\phi, \theta}$ is positively homogeneous with degree 1.
3. $\Omega_{\phi, \theta}(X + Z) \leq \Omega_{\phi, \theta}(X) + \Omega_{\phi, \theta}(Z) \forall (X, Z)$
4. $\Omega_{\phi, \theta}(X)$ is convex w.r.t. $X \in \mathbb{R}^D$.
5. The infimum in (15) can be achieved with $r \leq \text{card}(X) \forall X$ s.t. $\Omega_{\phi, \theta}(X) < \infty$.

Global optimality from local minima. While $\Omega_{\phi, \theta}(X)$ is convex, unlike the nuclear norm $\|X\|_*$, it typically cannot be evaluated in polynomial time due to its complicated definition. Nonetheless, its convexity allows us to use $\Omega_{\phi, \theta}$ as an analysis tool to derive results for neural network training formulations. In particular, it allows us to consider the convex (but typically non-tractable) problem, given by

$$\min_X F(X) \equiv \ell(Y, X) + \lambda \Omega_{\phi, \theta}(X). \quad (16)$$

Here $X \in \mathbb{R}^D$ is the output of the factorization mapping $X = \Phi_r(W^1, \dots, W^K)$, $\ell(Y, X)$ is a loss function that is assumed to be once differentiable and convex in X , $\Omega_{\phi, \theta}(X)$ is as defined by (15) where (ϕ, θ) is assumed to be a nondegenerate pair, and $\lambda > 0$. Given these assumptions, we are now ready to state our main results.

Theorem 1 *Any local minimizer of the non-convex optimization problem*

$$\min_{(W^1, \dots, W^K)_r} f_r(W^1, \dots, W^K) \equiv \ell(Y, \Phi_r(W^1, \dots, W^K)) + \lambda \sum_{i=1}^r \theta(W_i^1, \dots, W_i^K) \quad (17)$$

such that $(W_{i_0}^1, \dots, W_{i_0}^K) = (0, \dots, 0)$ for some $i_0 \in \{1, \dots, r\}$ is a global minimizer of (17). Moreover, $X = \Phi_r(W^1, \dots, W^K)$ is a global minimizer of (16).

Proof Sketch. The proofs of all our results are available in the supplement, but here we outline the sketch of the argument. First, note that from the definition of

⁵In particular, the largest r needs to be is $\text{card}(X)$, and we note that $\text{card}(X)$ is a worst case upper bound on the size of the factorization. In certain cases the bound can be shown to be lower. As an example, $\Omega_{\phi, \theta}(X) = \|X\|_*$ when $\phi(u, v) = uv^T$ and $\theta(u, v) = \|u\|_2 \|v\|_2$. In this case the infimum can be achieved with $r \leq \text{rank}(X) \leq \min\{\text{card}(u), \text{card}(v)\}$.

$\Omega_{\phi, \theta}$ the convex optimization problem (16) globally lower bounds the non-convex factorization problem (17) for any $X = \Phi_r(W^1, \dots, W^K)$, and because (16) is a convex function of X , the conditions for global optimality are easily derived. The result is completed by showing that a local minimum of (17) which satisfies the statement of the Theorem also satisfies the conditions to be global minimum of (16) at $X = \Phi_r(W^1, \dots, W^K)$, and due to the fact that (16) globally lower bounds (17) this implies that the local minimum of (17) is a global minimum. ■

From this result, we can then test the global optimality of any local minimum from the immediate corollary:

Corollary 1 *Given a function $f_r(W^1, \dots, W^K)$ of the form given in (17), any local minimizer of the optimization problem*

$$\min_{(W^1, \dots, W^K)_r} f_r(W^1, \dots, W^K) \quad (18)$$

is a global minimizer if $f_{r+1}([W^1 \ 0], \dots, [W^K \ 0])$ is a local minimizer of f_{r+1} .

Global minima can be found by local descent. From the results of Theorem 1, we are now also able to show that if we let the number of subnetworks (r) become large enough, then from any initialization we can always find a global minimizer of $f_r(W^1, \dots, W^K)$ using a purely local descent strategy. Specifically, we have the following result, whose proof gives a meta-algorithm for solving the optimization problem.

Theorem 2 *Given a function $f_r(W^1, \dots, W^K)$ as defined by (17), if $r > \text{card}(X)$ then from any point (Z^1, \dots, Z^K) such that $f_r(Z^1, \dots, Z^K) < \infty$ there must exist a non-increasing path from (Z^1, \dots, Z^K) to a global minimizer of $f_r(W^1, \dots, W^K)$.*

Proof Sketch. The proof is done in a constructive manner and defines a meta-algorithm that can be combined with any local-descent algorithm to reach a global minimum.

1. Perform local descent until arriving at a local minimum.
2. If one of the parallel networks is all 0 - i.e., $\exists i_0 \in \{1, \dots, r\}$ such that $(W_{i_0}^1, \dots, W_{i_0}^K) = (0, \dots, 0)$ - then we are at a global minimum due to Theorem 1.
3. Else if there exists a nonzero $\beta \in \mathbb{R}^r$ such that $\sum_{i=1}^r \beta_i \phi(W_i^1, \dots, W_i^K) = 0$ then scale β so that $\min_i \beta_i = -1$ and set $W_i^k \leftarrow (1 + \beta_i)^{1/p} W_i^k$ for $k = 1, \dots, K$. Such a β is guaranteed to exist if $r > \text{card}(X)$, and the operation of scaling the variables by the $(1 + \beta_i)^{1/p}$ terms is shown to traverse a flat surface of the objective function until arriving at a point where one of the parallel networks is all 0. From there, if a

local descent direction exists continue local-descent. If no local descent direction exists then we are again at a global minimum due to Theorem 1.

4. Otherwise, if $r \leq \text{card}(X)$ and no β exists in the prior step then increment r by appending a subnetwork in parallel initialized as all-zeros. If a local-descent direction exists continue local-descent. Otherwise, we are at a global minimum due to Corollary 1.

6. Discussion: Limitations and Implications

While the framework described so far is very general, and provides guarantees of global optimality for various forms of neural network problems, we pause to note a few practical limitations and then discuss implications of our results in the design of neural networks.

First, note that the size of the network is controlled by a single parameter: r . While a single parameter is sufficient for matrix factorization, where r is the size of the factors, and two-layer neural networks, where r is the number of neurons in the hidden layer, this is insufficient to model deep networks where we want to control also the number of layers K as well as the number of neurons in each layer. In other words, while one could naively assume that the results so far apply to current deep networks by setting $r = 1$ (i.e., the network is not assumed to have a parallel structure), this is *not* the case because the assumption that the network needs to be “large enough” (i.e., r must be sufficiently large) is essential to prove that local minima are global. Therefore, to analyze current deep networks without this parallel structure, it is essential that we extend the framework to optimize over additional network size parameters, such as the number of layers and the numbers of neurons per layer.

Additionally, the maximum upper-bound size for r in Theorem 2 is typically much too large to be of practical use. Note, however, that the bounds we have shown here are for the most general case of mapping and regularizer, (ϕ, θ) , and that a very interesting line of future work is to explore sufficient conditions on these functions that allows the size of r to be greatly reduced. For example, in the nuclear norm case of matrix factorization, it is well known that the largest r will need to be is equal to the rank of the final solution. As an example from neural networks, if the architecture of a parallel sub-network (as defined by ϕ) is sufficiently rich to span the output space (i.e., $\text{Im}(\phi) = \mathbb{R}^D$), then if θ satisfies the triangle inequality it is easily seen from the definition of $\Omega_{\phi, \theta}$ that the infimum in (15) can always be achieved with a single parallel network (i.e., $r = 1$).

A final limitation to note is that the meta-algorithm used to construct the proof of Theorem 2 relies on using local descent, which is to be contrasted with typical optimization

algorithms such as gradient descent. While gradient descent is certainly a form of local descent, we remind the reader that in general finding a local descent direction of a non-convex function can be a NP-hard problem in general (for example, at a saddle-point), so our results do not necessarily imply the existence of polynomial time algorithms that can solve all of the potential formulations captured within our framework. Again, however, we emphasize that our analysis is a worst-case analysis (i.e., choose *any* possible initialization) for any potential mapping and regularizer, and significant potential exists to strengthen our results by considering specific families of mapping and regularizers, initialization strategies, or statistical distributions of the training data that allow for polynomial time guarantees.

Implications for Neural Networks. Despite the limitations discussed above, our analysis suggests several significant guiding principles regarding the training of neural networks which can facilitate more efficient optimization. The first is that balancing the degree of positive homogeneity between the regularization function and the mapping function is crucial. In fact, it can be shown (see supplement, Section 8) that if the degrees of positive homogeneity do not match between the mapping and the regularization function, then it either becomes impossible to make guarantees regarding the global optimality of a local minimum, or it becomes possible that the regularization function will do nothing to limit the size of the network, so the degrees of freedom in the model are largely determined by the user defined choice of r . In practice, this issue often arises in the context of weight decay, where the regularization function is typically chosen as $\Theta_r(W^1, \dots, W^K) = \sum_i \|W^i\|_F^2$ or $\Theta_r(W^1, \dots, W^K) = \sum_i \|W^i\|_1$. Since these functions are only positively homogeneous of degree 2 and 1, respectively, the mapping of a deep network will typically not be balanced with the regularization and one is guaranteed that non-optimal local minima will exist regardless of the size of the network (supplement, Proposition 2). We note that this is a potential explanation of the noted inferior performance of using weight decay versus dropout regularization [23, 14, 24] and that several more recently proposed successful regularization strategies are compatible with balanced degrees of positive homogeneity. For example, the Path-SGD regularizer proposed in [18] takes a product of weights along a path through the network and then calculates a norm of all possible paths through the network. Note that the product of weights along a path through the network will typically have a degree of positive homogeneity equal to K and thus will be balanced with the degree of positive homogeneity of the network output for typical network architectures. Likewise, Batch Normalization proposed in [12] essentially adds a whitening operation to the input of a layer which is a positively homogeneous transformation similar to contrast normalization but across training

examples (see Section 9 of the supplement). Taken together, our results suggest at several key properties one should account for in the design of network regularization and provide many interesting opportunities for experimental exploration in future work.

A final implication of our analysis is that neural networks which generate the output by taking the sum of multiple parallel subnetworks are highly conducive to efficient optimization. This idea, of linearly combining the outputs of multiple subnetworks, has clear analogies to ensemble methods like boosting and bagging and was a large motivation in the development of techniques such as dropout, which stochastically approximates the average output of an exponential number of subnetworks [23]. The framework we present here is not an exact analogy to dropout, as dropout enforces equality in network weights across parallel networks with a shared parameterization (i.e., if a neuron is present in a given subnetwork all of its input and output weights must be equal to the same neuron in the other subnetworks), but many interesting questions for future work can be asked regarding the concept of summing multiple subnetworks combined with considering more general forms of network mappings which allow for common parametrization of the subnetworks.

7. Conclusions

Here we have presented a general framework which allows for a wide variety of non-convex optimization problems, including certain forms of neural network training, to be analyzed with tools from convex analysis and induces a convex regularizer on the output of the non-convex mapping. In particular, we have shown sufficient conditions to guarantee that any local minimum is a global minimum of the non-convex factorization problem and that if the non-convex factorization problem is done with factors of sufficient size, then from any feasible initialization it is always possible to find a global minimizer using a purely local descent algorithm. Additionally, our results suggest that balancing the degrees of positive homogeneity between the network mapping and the regularization function is critical for preventing non-optimal local minima in the loss surface of modern neural network architectures and offer guidance for the design of network architectures and regularizers.

Acknowledgments. This work was supported by NSF grants 1447822, 1618485 and 1618637. We thank Laurent Younes for insightful discussions.

References

- [1] F. Bach. Convex relaxations of structured matrix factorizations. *arXiv:1309.3117v1*, 2013. 4
- [2] F. Bach, J. Mairal, and J. Ponce. Convex sparse matrix factorizations. *arXiv:0812.1869v1*, 2008. 4

- [3] Y. Bengio, N. L. Roux, P. Vincent, O. Delalleau, and P. Marcotte. Convex neural networks. In *Neural Information Processing Systems*, pages 123–130, 2005. 2, 3
- [4] S. Burer and R. D. C. Monteiro. Local minima and convergence in low-rank semidefinite programming. *Mathematical Programming, Series A*(103):427–444, 2005. 4
- [5] E. Candès and T. Tao. The power of convex relaxation: Near-optimal matrix completion. *IEEE Transactions on Information Theory*, 56(5):2053–2080, 2010. 2
- [6] A. Choromanska, M. Henaff, M. Mathieu, G. Ben Arous, and Y. LeCun. The loss surfaces of multilayer networks. In *International Conference on Artificial Intelligence and Statistics*, pages 192–204, 2015. 2, 3
- [7] G. E. Dahl, T. N. Sainath, and G. E. Hinton. Improving deep neural networks for lvcsr using rectified linear units and dropout. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 8609–8613, 2013. 2
- [8] Y. N. Dauphin, R. Pascanu, C. Gulcehre, K. Cho, S. Ganguli, and Y. Bengio. Identifying and attacking the saddle point problem in high-dimensional non-convex optimization. In *Neural Information Processing Systems*, pages 2933–2941, 2014. 3
- [9] D. L. Donoho. For most large underdetermined systems of linear equations the minimal ℓ^1 -norm solution is also the sparsest solution. *Communications on Pure and Applied Mathematics*, 59(6):797–829, 2006. 2
- [10] J. H. Friedman. Greedy function approximation: a gradient boosting machine. *Annals of Statistics*, pages 1189–1232, 2001. 3
- [11] B. Haeffele, E. Young, and R. Vidal. Structured low-rank matrix factorization: Optimality, algorithm, and applications to image processing. In *International Conference on Machine Learning*, pages 2007–2015, 2014. 4
- [12] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning*, pages 448–456, 2015. 2, 8
- [13] M. Janzamin, H. Sedghi, and A. Anandkumar. Beating the perils of non-convexity: Guaranteed training of neural networks using tensor methods. *arXiv preprint arXiv:1506.08473*, 2015. 3
- [14] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Neural Information Processing Systems*, pages 1097–1105, 2012. 2, 5, 8
- [15] A. L. Maas, A. Y. Hannun, and A. Y. Ng. Rectifier nonlinearities improve neural network acoustic models. In *International Conference on Machine Learning*, volume 30, 2013. 2
- [16] J. Mairal, F. Bach, J. Ponce, and G. Sapiro. Online learning for matrix factorization and sparse coding. *The Journal of Machine Learning Research*, 11:19–60, 2010. 2
- [17] L. Mason, J. Baxter, P. L. Bartlett, and M. R. Frean. Boosting algorithms as gradient descent. In *Neural Information Processing Systems*, pages 512–518, 2000. 3
- [18] B. Neyshabur, R. R. Salakhutdinov, and N. Srebro. Path-sgd: Path-normalized optimization in deep neural networks. In *Neural Information Processing Systems*, pages 2422–2430, 2015. 2, 8
- [19] B. Recht, M. Fazel, and P. Parrilo. Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM Review*, 52(3):471–501, 2010. 4
- [20] D. E. Rumelhart, G. E. Hinton, and R. J. Williams. Learning representations by back-propagating errors. *Cognitive Modeling*, 5, 1988. 2
- [21] I. Safran and O. Shamir. On the quality of the initial basin in overspecified neural networks. In *International Conference on Machine Learning*, pages 774–782, 2016. 3
- [22] N. Srebro, J. D. Rennie, and T. S. Jaakkola. Maximum-margin matrix factorization. In *Neural Information Processing Systems*, volume 17, pages 1329–1336, 2004. 4
- [23] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958, 2014. 2, 8
- [24] L. Wan, M. Zeiler, S. Zhang, Y. L. Cun, and R. Fergus. Regularization of neural networks using dropconnect. In *International Conference on Machine Learning*, pages 1058–1066, 2013. 8
- [25] S. J. Wright and J. Nocedal. *Numerical Optimization*, volume 2. Springer New York, 1999. 2
- [26] Y. Xu and W. Yin. A block coordinate descent method for regularized multiconvex optimization with applications to nonnegative tensor factorization and completion. *SIAM Journal on Imaging Sciences*, 6(3):1758–1789, 2013. 2
- [27] M. D. Zeiler, M. Ranzato, R. Monga, M. Mao, K. Yang, Q. V. Le, P. Nguyen, A. Senior, V. Vanhoucke, J. Dean, and G. E. Hinton. On rectified linear units for speech processing. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3517–3521, 2013. 2