

# Do Neural Networks Really Find the Best Parameters? Understanding the Global Landscape of Neural Networks

Bernardo Bianco Prado

## Abstract

Neural networks are omnipresent in contemporary machine learning and are promising in dealing with the immense amount of data that exists in today's world. Because networks are so useful, it's essential to understand their behavior and their outputs. Especially because their outputs influence people's lives and experiences. This paper surveys a series of results that seek to understand when neural networks converge to a good output.

## 1 Introduction: Formulating Our Discussion Mathematically

Neural networks have become one of the most promising tools in contemporary machine learning problems. Unfortunately, these networks are as mysterious as they are useful. Because of their complexity (they are built from many function compositions), it's hard to predict what kind of parameter fit (good or bad) will be returned by the neural network. Exploring what can be said about these outputs of neural networks has been subject of great research in recent years and many partial discoveries have been publicized. Unfortunately, a definitive result is still lacking. In this article we review the exposition found in [1] which surveys results that seek to understand how good are the outputs of neural networks. All of the facts presented in this article are taken from [1] unless specified otherwise.

Let us first formulate our discussion in more mathematical terms. We are interested in whether the output of a neural network is good, and the measurement of "goodness" is usually performed in terms of some loss function  $F(\theta)$  which is minimized by the best parameter  $\theta$ . In other words, we are searching for

$$\theta^* = \arg \min_{\theta \in \mathbb{R}^D} F(\theta).$$

The most common methods of finding the optimal  $\theta^*$  are essentially improved variants of gradient descent. The issues arise because, since  $F(\theta)$  is not convex, the gradient descent method can (in theory) get stuck in local minima or saddle points of the function  $F(\theta)$ . The article [1] surveys different results that try to understand the output of neural networks from two different angles. First, it seeks to understand whether it's possible that all local minima of  $F(\theta)$  are in fact global minima. From a different angle, the article tries to understand whether the gradient descent method can avoid the sub-optimal local minima of  $F(\theta)$  even if these bad local minima are present in the loss function.

### 1.1 Specifying the Objects of Study

For this discussion, we will need to talk about neural networks rigorously so it will be helpful introduce them in a mathematical language. Given an collection of input data  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^{d_x}$  and output data  $\mathbf{y}_1, \dots, \mathbf{y}_n \in \mathbb{R}^{d_y}$ , a fully connected neural network is a function  $f_{\theta} : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d_y}$  given by

$$f_{\theta}(\mathbf{x}) = \mathbf{W}_L \phi(\mathbf{W}_{L-1} \dots \phi(\mathbf{W}_2 \phi(\mathbf{W}_1 \mathbf{x} + b_1) + b_2) \dots + b_{L-1})$$

where  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  is called an activation function that is applied to vectors in an element-wise fashion, each  $\mathbf{W}_i \in \mathbb{R}^{d_i \times d_{i-1}}$  is a matrix and  $\boldsymbol{\theta} = (\mathbf{W}_1, b_1, \dots, \mathbf{W}_{L-1}, b_{L-1}, b_L)$ . By definition, we have specified  $d_0 = d_x, d_L = d_y$ . Our hope is that the neural network  $f_{\boldsymbol{\theta}}$  will predict the output  $\mathbf{y}$  from the input  $\mathbf{x}$  based on the available data inputs  $\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^n$ . To ensure that the neural network manages succeed in this prediction, we are trying to find the best parameters  $\boldsymbol{\theta}$  that will fit the input data. Mathematically, this is fomulated as a minimization problem

$$\min_{\boldsymbol{\theta}} F(\boldsymbol{\theta}) = \min_{\boldsymbol{\theta}} \sum_{i=1}^n \ell(\mathbf{y}_i, f_{\boldsymbol{\theta}}(\mathbf{x}_i))$$

where  $\ell(\mathbf{y}, \mathbf{z})$  is a function that measures the “distance” between  $\mathbf{y}$  and  $\mathbf{z}$ . Some examples are the Euclidean distance  $\ell(\mathbf{y}, \mathbf{z}) = \|\mathbf{y} - \mathbf{z}\|_2^2 = \sum_{j=1}^{d_y} (y_j - z_j)^2$  for regression problems and  $\ell(\mathbf{y}, \mathbf{z}) = \log(1 + \exp(-\mathbf{y}^T \mathbf{z}))$  for binary classification problems.

## 2 A Promising Example: Linear Neural Networks

Let us start by discussing a simpler problem in which the map  $\phi$  is just the identity map and so

$$f_{\boldsymbol{\theta}}(\mathbf{x}) = \mathbf{W}_L(\mathbf{W}_{L-1} \dots (\mathbf{W}_2(\mathbf{W}_1 \mathbf{x} + b_1) + b_2) \dots + b_{L-1}).$$

For simplicity, let us also assume  $b_1, \dots, b_{L-1} = 0$  so that

$$f_{\boldsymbol{\theta}}(\mathbf{x}) = \mathbf{W}_L \mathbf{W}_{L-1} \dots \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}.$$

We are also interested in the regression problem where  $\ell(\mathbf{y}, \mathbf{z}) = \|\mathbf{y} - \mathbf{z}\|_2^2$ . The linear neural network minimization problem then becomes

$$\min_{\boldsymbol{\theta}} \sum_{i=1}^n \|\mathbf{y}_i - f_{\boldsymbol{\theta}}(\mathbf{x}_i)\|_2^2 = \min_{\boldsymbol{\theta}} \|\mathbf{Y} - \mathbf{W}_L \mathbf{W}_{L-1} \dots \mathbf{W}_2 \mathbf{W}_1 \mathbf{X}\|_F^2 \quad (1)$$

where  $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_n]$ ,  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$  are concatenation matrices of the data and  $\|\cdot\|_F$  is the Frobenius norm.

The simplest example of this problem would be

$$\min_{u, v \in \mathbb{R}} (uv - 1)^2.$$

One can see that this has a collection of global minima when  $uv = 1$ , these are the only local minima of this function, but it does have other critical points such as a saddle point at  $(u, v) = (0, 0)$ . This is an example of a network for which every local minimum is a global minimum. More generally, we have the following result.

**Theorem 2.1** *If  $X, Y$  in the minimization problem (1) have full rank then every local minimum of the objective function is a global minimum.*  $\square$

In the example above, we have  $X = 1$  and  $Y = 1$  which are nonzero scalars and thus have full rank. This theorem generalizes the example for higher dimensional problems. Note that the result does not characterize critical points. For that we have a different theorem, which we write informally to get the intuition across

**Theorem 2.2 (Characterization of Critical Points for Overparametrized Linear Networks)** *Assume  $n \geq d_x, d_y$  (there are more data points than the dimension of the input and output),  $\mathbf{X}\mathbf{X}^T$  and  $\mathbf{X}\mathbf{Y}^T$  are full rank, and  $\mathbf{Y}\mathbf{X}^T(\mathbf{X}\mathbf{X}^T)^{-1}\mathbf{X}$  has distinct singular values. Assume further that the output and input layer have the lowest dimension out of all the layers of the network. Then for the problem (1), all critical points with  $\boldsymbol{\theta} = (\mathbf{W}_1, \dots, \mathbf{W}_L)$  so that  $\mathbf{W}_L \dots \mathbf{W}_2 \mathbf{W}_1$  being full rank is a global minimum and every critical point with  $\mathbf{W}_L \dots \mathbf{W}_2 \mathbf{W}_1$  being singular is a saddle point.*  $\square$

Observe that this theorem agrees with the example above. When  $uv = 1$ , then  $1$  is full rank and we have a global minimum. On the other hand, the critical point  $u = v = 0$  has  $uv = 0$  being singular and corresponds to a saddle point. This allowed us to completely classify critical points for certain overparametrized linear neural networks. Now we look at more general (nonlinear) overparametrized networks.

### 3 Moving Onward: Overparametrized Nonlinear Networks

As in the previous section, we start with a simple, one dimensional example of a nonlinear network. Consider

$$\min_{u,v \in \mathbb{R}} (u\phi(v) - 1)^2 \quad (2)$$

Consider, for example, when the activation function is the ReLU function  $\phi(t) = \max\{0, t\}$ . This has a sub-optimal local minimum at  $u = 0, v = -1$ ! This is explained in the following theorem.

**Theorem 3.1** *Let  $x, y$  be nonzero real numbers and consider the minimization problem*

$$\min_{u,v \in \mathbb{R}} (y - u\phi(vx))^2.$$

*Then the following are equivalent:*

(I) *The problem has no sub-optimal local minima.*

(II) *If  $\phi(t) = 0$  then  $t$  is not a local maximum or minimum of  $\phi$ .* □

This is a unsatisfying result since ReLU is a popular choice of activation function. There is in fact a result that, generally, for sigmoid activation functions, there is no guarantee that there will not be sub-optimal local minima for the single layer, one dimensional network minimization problem. This is formalized below.

**Proposition 3.2** *For  $n \geq 3$  with the data points  $x_1, \dots, x_n \in \mathbb{R}$  all distinct and for an sigmoid activation function  $\phi$ , there exist output data  $y_1, \dots, y_n \in \mathbb{R}$  so that the minimization problem*

$$\min_{u,v \in \mathbb{R}} \sum_{i=1}^n (y_i - u\phi(vx_i))^2$$

*has a sub-optimal local minimum.* □

In other words, for such sigmoid functions it's always possible that real world data will yield a neural network problem that has sub-optimal local minima. This proposition brings us to a halt in the discussion of sub-optimal local minima. We must then move away from this discussion of suboptimal local minima and towards something else. One useful notion of study is that of valleys and basins.

#### 3.1 A Less Local Approach: Valleys and Basins

The idea behind this discussion is to look at the landscape of the objective function and move away from the pointwise perspective. We are interested in what the neighborhood of local minima looks like. A useful notion is that of a *spurious valley*. A spurious valley is a connected component of the sub-level set

$$\{\theta : F(\theta) \leq c\}$$

that contains no strict global minimum of the objective function.

If a function has no spurious valley, even if a sub-optimal local minimum exists on a sufficiently small sub-level set of the objective function, there is a non-increasing path from that local minimum to the global minimum. This is useful because it might prevent the gradient descent method from reaching and getting stuck at sub-optimal local minima. Therefore, we are interested in cases when there are no spurious valleys. To this end, we have the following theorem.

**Theorem 3.3** Suppose that a one dimensional, arbitrarily deep neural network  $f_{\theta}$  satisfies the following:

1. The activation function  $\phi$  is strictly monotonic and  $\phi(\mathbb{R}) = \mathbb{R}$ .
2. The activation function is linearly independent from any shifted versions of itself. That is, for any  $m \geq 2$  and all combination of coefficients  $a_1, \dots, a_m$  and  $c_0, c_1, \dots, c_m$  with the  $c_i$  distinct, if  $c_0\phi(x) + c_1\phi(x - a_1) + \dots + c_m\phi(x - a_m) = 0$ , then we must have that  $c_0 = c_1 = \dots = c_m = 0$ .
3.  $d_{L-1} \geq d$ .
4. All the training data points are distinct.

Then the empirical loss objective function has no spurious valleys.

This would cover, for example  $\phi(t) = t$  but not the ReLU activation function. In any case, this is a promising result and moving away from the local perspective of local minima might be useful. Following this philosophy, it is of interest to talk about set-wise local minima. In particular, we are interested in *set-wise strict local minima*. For a function  $f : \mathbb{R}^d \rightarrow \mathbb{R}^l$ , the set  $S$  is a strict local minimum of  $f$  in the sense of sets if there is some  $\epsilon > 0$  so that for all  $\mathbf{x} \in S$  and all  $\mathbf{y} \in \mathbb{R}^d \setminus S$  satisfying  $\|\mathbf{x} - \mathbf{y}\|_2 \leq \epsilon$ , we have  $f(\mathbf{x}) < f(\mathbf{y})$ . In other words, a set-wise local minimum of  $f$  is a set so that the value of  $f$  all its points close to its boundary are smaller than all the nearby points outside of its boundary.

The notion of a set-wise strict local minimum gives rise to the definition of a *sub-optimal basin*, which is a set-wise local minimum that contains no global minima. This notion is essentially the set version of a sub-optimal local minimum. There is a result analogous to the theorem above to sub-optimal basins.

**Theorem 3.4** Suppose that a one dimensional, arbitrarily deep neural network  $f_{\theta}$  satisfies the following:

1. There is a  $k$  so that the  $k$ -th entry of all the input data  $(\mathbf{x}_i)_k$  are distinct.
2.  $d_{L-1} \geq d$ .
3. The activation function  $\phi$  is continuous.

Suppose, in addition that the loss function  $\ell(\mathbf{u}, \mathbf{v})$  is convex with respect to  $\mathbf{v}$ . Then the empirical loss objective function has no sub-optimal basins.

These results show a promising direction in understanding the behavior of descent methods for neural networks. These are the main theoretical results covered in the paper regarding the geometric landscape of neural networks. We now discuss some attempts at eliminating sub-optimal local minima from the network minimization problem.

## 4 Getting Rid of the Issue: Removing Bad Local Minima from Nonlinear Neural Networks

We are interested in identifying the regions where there are no sub-optimal local minima for the objective function of the neural network problem. For this, we are mostly interested in the last layer of the neural networks. In particular, we rewrite the neural network as

$$f_{\theta}(\mathbf{x}) = \mathbf{W}_L \mathbf{Z}_{\theta}(\mathbf{x}).$$

Where  $\mathbf{Z}_{\theta}$  essentially encompasses all the other layers of the neural network. Rewriting the problem like this is useful because we can impose conditions on  $\mathbf{Z}_{\theta}$  so that the minimization problem has no bad local minima as long as  $\mathbf{x}$  is restricted to a subset of  $\mathbb{R}^{d_x}$ . An example presented below.

**Claim 4.1** If  $\theta^*$  is a local minimum of  $F(\theta^*) = \sum_{i=1}^n \|\mathbf{y}_i - f_{\theta^*}(\mathbf{x}_i)\|_2^2$  so that the matrix

$$\mathbf{Z}_{\theta^*}(\mathbf{X}) = (\mathbf{Z}_{\theta^*}(\mathbf{x}_1), \dots, \mathbf{Z}_{\theta^*}(\mathbf{x}_1))$$

is full rank, then  $\theta^*$  is a global minimum.

Further, if the derivatives of the activation function  $\phi^{(k)}(0)$  for  $k = 1, \dots, n - 1$  are nonzero, then the set where  $\{\theta : \mathbf{Z}_{\theta}(\mathbf{X})$  has full rank $\}$  is dense.  $\square$

This result allows us to avoid local sub-optimal local minima with almost sure probability whenever the activation function satisfies the derivative condition in the claim. Unfortunately, this does not cover, for example, non-smooth activation functions such as ReLU.

An alternative is to modify the network and the objective function for the minimization problem so that the problem yields better results. Suppose we modify the neural networks by creating a new function  $\tilde{f}$  so that

$$\tilde{f}_{\tilde{\theta}}(\mathbf{x}) = f_{\theta}(\mathbf{x}) + a \exp(\mathbf{w}^T \mathbf{x} + b)$$

where  $\tilde{\theta} = (\theta, a, \mathbf{w}, b)$ . Suppose that the loss function  $\ell$  is the logistic loss function

$$\ell(\mathbf{y}, \mathbf{z}) = \log(1 + \exp(-\mathbf{y}^T \mathbf{z}))$$

then we consider the minimization problem

$$\min_{\tilde{\theta}} \tilde{F}(\tilde{\theta}) = \min_{\tilde{\theta}} \sum_{i=1}^n \ell(\mathbf{y}_i, \tilde{f}_{\tilde{\theta}}(\mathbf{x}_i)) + \frac{\lambda a^2}{2}.$$

Then we have the following result.

**Theorem 4.2** For the minimization problem above, we have

1.  $\tilde{F}(\tilde{\theta})$  has at least one local minimum.
2. At every local minimum,  $a = 0$ .
3. If  $\tilde{\theta}^* = (\theta^*, a^*, \mathbf{w}^*, b^*)$  is a local minimum of  $\tilde{F}(\tilde{\theta})$ , then it is also a global minimum of  $\tilde{F}(\tilde{\theta})$  and  $\theta^*$  is a global minimum of

$$F(\theta) = \sum_{i=1}^n \ell(\mathbf{Y}_i, \tilde{f}_{\theta}(\mathbf{x}_i)).$$

$\square$

This result seems like it solves all of our problems! At least for the logistic regression case. There is one issue, however, and that is that even though there are no sub-optimal local minima of this new minimization problem, there may be a decreasing path of  $\tilde{\theta}(s)$  for  $s > 0$  so that  $\|\tilde{\theta}(s)\|_2 \rightarrow \infty$  as  $s \rightarrow \infty$  and  $\tilde{F}(\tilde{\theta}(s))$  is strictly decreasing. This means that gradient descent may never converge to a local minimum and instead it may just send the  $\tilde{\theta}$  to infinity.

There is a way to avoid that for certain neural networks. For example, consider the one layer, one dimensional ReQU neural network defined by

$$f_{\theta}(x) = \sum_{j=1}^m a_j (\max\{\mathbf{w}_j x + b_j\})^2 = \sum_{j=1}^m a_j (\text{ReQU}(\mathbf{w}_j x + b_j))^2,$$

where  $m \geq n + 1$  with  $n$  being the number of input data points. As above, suppose the loss function is logistic regression and suppose we are minimizing the objective function

$$F(\boldsymbol{\theta}) = \sum_{i=1}^n \ell(y_i, f_{\boldsymbol{\theta}}(x_i)) + \frac{1}{3} \sum_{j=1}^m \lambda_j \left( |a_j|^3 + 2(\|\mathbf{w}_j\|_2^2 + b_j^2)^{3/2} \right)$$

where  $\boldsymbol{\theta} = (a_1, \mathbf{w}_1, b_1, \dots, a_m, \mathbf{w}_m, b_m)$ , and where  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_m)$  is a parameter vector that needs to be tuned in the problem. For the right  $\boldsymbol{\lambda}$ , one can prove that every local min achieves zero training error and thus every local minimum would be a global minimum. In particular, we have the following

**Theorem 4.3** *If  $m \geq n + 1$  and  $F(\boldsymbol{\theta})$  is defined as above, then for each collection of input and output data  $x_1, \dots, x_n, y_1, \dots, y_n \in \mathbb{R}$ , there is a  $\lambda_0 > 0$  and a set  $\mathcal{C} \subset \mathbb{R}$  of measure zero so that for any  $\boldsymbol{\lambda} \in (0, \lambda_0)^m \setminus \mathcal{C}$ , we have*

1. *The objective function  $F(\boldsymbol{\theta})$  is coercive, that is*

$$\lim_{\|\boldsymbol{\theta}\| \rightarrow \infty} F(\boldsymbol{\theta}) = +\infty.$$

2. *Every local minimum  $\boldsymbol{\theta}^*$  of the function  $F(\boldsymbol{\theta})$  is a global minimum of  $F(\boldsymbol{\theta})$  and achieves zero training error. □*

This says that basically whatever  $\boldsymbol{\lambda}$  one chooses, as long as its entries are small enough, gradient descent will likely always converge to a global minimum of  $F(\boldsymbol{\theta})$ .

We have covered useful properties of the minimization problem that will guarantee a desirable landscape of the objective function for the minimization problem. We have not, however discussed whether those properties translate to convergence of gradient descent to global minima. This is what we focus on below.

## 5 Was It All Worth It? The Behavior of Gradient Descent on the Discussed Problems

The hope is that for objective functions with the properties discussed above then for most initializations, gradient descent will converge to a global minimum. For example, if a function has sub-optimal basins the hope is that the sub-optimal basins have small measure and thus, gradient descent will converge to a global minimum with high probability. The goal result is outlined below.

**Principle (Blueprint for Gradient Descent)** *If  $F$  is an objective function represents the neural network minimization problem. Let  $\boldsymbol{\theta}^*$  be the point of convergence of gradient descent or stochastic gradient descent and  $F^*$  be the minimum value of  $F(\boldsymbol{\theta})$ . Then there are small constants  $\delta, \epsilon$  so that*

$$\mathbb{P}(F(\boldsymbol{\theta}^*) < F^* + \epsilon) > 1 - \delta.$$

Intuitively, this means that the gradient descent method will converge to a almost optimal value for  $F$  with high probability. Unfortunately, there is still not a lot of results that confirm this principle theoretically. Because we the conjecture deals with probability measures, stochastic gradient descent has proved more fitting to this discussion as it favors a probabilistic perspective.

There are some results, however, that guide us in the direction of this principle. The two strategies to achieve this are first showing that for almost all initial conditions avoid the bad regions of the objective function landscape. The other strategy is to show that whenever gradient descent reaches a bad region, it will likely leave it and move on towards a global minimum.

Some results are available for ultra-wide neural networks. For an example, consider the matrix  $\mathbf{G}(\boldsymbol{\theta}) = (\text{grad}_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}}(\mathbf{x}_1), \dots, \text{grad}_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}}(\mathbf{x}_n))$  and  $\mathbf{K}(\boldsymbol{\theta}) = \mathbf{G}(\boldsymbol{\theta})^T \mathbf{G}(\boldsymbol{\theta})$  and suppose  $d_y = 1$ . Suppose the loss function  $\ell$  is the Euclidean distance squared. Then we have the following result.

**Claim 5.1** *We have the following:*

1. If  $\boldsymbol{\theta}^*$  is a critical point of the objective function  $F(\boldsymbol{\theta})$  and  $\mathbf{K}(\boldsymbol{\theta}^*)$  has full rank, then  $\boldsymbol{\theta}^*$  is a global minimum and  $F(\boldsymbol{\theta}^*) = 0$ .
2. If  $\boldsymbol{\theta}_k$  for  $k = 0, 1, \dots$  are the iterates of gradient descent and if there is some constant  $c > 0$  so that

$$\lambda_{\min}(\mathbf{G}(\boldsymbol{\theta}_k)) \geq c$$

for all  $k = 0, 1, \dots$ , then we have that and limit point  $\boldsymbol{\theta}^*$  of the sequence  $\{\boldsymbol{\theta}_k\}_{k=0}^{\infty}$  which is a critical point of  $F(\boldsymbol{\theta})$  is a global minimum of  $F(\boldsymbol{\theta})$  with  $F(\boldsymbol{\theta}^*) = 0$ .

□

These results correspond to ultra wide neural networks and are very strong. The unfortunate fact is that many neural networks of interest are not ultra-wide and thus this does not necessarily apply to many cases of interest. This summarizes some results for gradient descent behaving well with certain neural network problems.

## 6 Discussion and Conclusion

This paper covers a comprehensive collection of results regarding neural networks. One question I had was whether some of the nonsmoothness of ReLU can be fixed by looking at the function

$$r(t) = \begin{cases} 0 & \text{if } t \leq 0 \\ (1 + t^2)^{1/2} - 1 & \text{if } t > 0 \end{cases}$$

which is a commonly used function in mathematical analysis that asymptotically behaves like  $|t|$  as  $t \rightarrow \infty$  and is smooth at  $t = 0$ .

In the future, I hope to implement neural networks and observe the behavior of gradient descent with different activation functions. I also hope to read some of the proofs for the results more deeply to understand ideas behind the behavior of neural networks. I have learned a lot from this paper and now I have a deeper knowledge of how neural networks behave and I understand a bit more about the theory of these essential objects in contemporary machine learning. I will surely use some of these results to better choose the setup whenever I'm using neural networks in the future.

## References

- [1] Ruoyu Sun, Dawei Li, Shiyu Liang, Tian Ding and R Srikant. *The Global Landscape of Neural Networks: An Overview*. 2020. <https://arxiv.org/abs/2007.01429>