

Complexity of Training ReLU Neural Network

Digvijay Boob, Santanu S. Dey, Guanghui Lan

Industrial and Systems Engineering, Georgia Institute of Technology

Abstract

In this paper, we explore some basic questions on the complexity of training neural networks with ReLU activation function. We show that it is NP-hard to train a two-hidden layer feedforward ReLU neural network. If dimension of the input data and the network topology is fixed, then we show that there exists a polynomial time algorithm for the same training problem. We also show that if sufficient over-parameterization is provided in the first hidden layer of ReLU neural network, then there is a polynomial time algorithm which finds weights such that output of the over-parameterized ReLU neural network matches with the output of the given data.

Keywords: NP-hardness, Neural Network, ReLU Activation Function, 2-hyperplane Separability

1. Introduction

Deep neural networks (DNNs) are functions computed on a graph parameterized by its edge weights. More formally, the graph corresponding to a DNN is defined by input and output dimensions $w_0, w_k \in \mathbb{Z}_+$, the number of hidden layers $k \in \mathbb{Z}_+$, and a sequence of k natural numbers w_1, w_2, \dots, w_k representing the number of nodes in each of the hidden k -layers. The function computed on a DNN graph is:

$$f := \tau \circ a_k \circ \dots \circ a_2 \circ \tau \circ a_1,$$

where \circ is function composition, τ is a nonlinear function (applied componentwise) called as the activation function, and $a_i : \mathbb{R}^{w_{i-1}} \rightarrow \mathbb{R}^{w_i}$ are affine functions.

Given the input and corresponding output data, the problem of training a deep neural network can be thought of as determining the edge weights of the directed layered graph that determine the affine functions a_i 's for which the output of the neural network matches the output data as closely as possible. Formally, given a set of input and output data $\{(x^i, y^i)\}_{i=1}^N$ where $(x^i, y^i) \in \mathbb{R}^{w_0} \times \mathbb{R}^{w_k}$, and a loss function $l : \mathbb{R}^{w_k} \times \mathbb{R}^{w_k} \rightarrow \mathbb{R}_{\geq 0}$ (for example, l can be the square loss function), the task is to determine the weights that define the affine function a_i 's

such that

$$\sum_{i=1}^N l(f(x^i), y^i) \tag{1}$$

is minimized.

Some commonly studied activation functions are: threshold function, sigmoid function and ReLU function. ReLU is one of the most important activation functions used widely in applications. Despite its wide use, the question of computational complexity of training multi-layer fully-connected ReLU neural network has remained open. This paper makes contributions in this direction. Before formally stating our results, we first present the current understanding of the computational complexity in the literature.

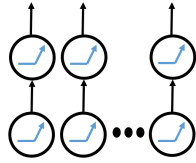
Complexity of training DNNs with threshold activation function. The threshold (sign) function is given by

$$\text{sgn}(x) := \begin{cases} 1 & \text{if } x > 0 \\ -1 & \text{if } x < 0 \end{cases}.$$

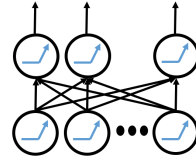
It was shown by Blum et al. [1] that the problem of training a simple two layer neural network with two nodes in the first layer and one node in the second layer while using threshold activation function at all the nodes is NP-complete. The training problem turns out to be equivalent to separation by two hyperplanes problem which was shown to be NP-complete by Megiddo [2]. There are other hardness results such as crypto hardness for intersection of k-hyperplanes which apply to neural networks with threshold activation function.

DNNs with rectified linear unit (ReLU) activation function. Theoretical worst case results presented above, along with limited empirical successes led to DNN's going out of favor by late 1990s. However, in recent times, DNNs became popular again due to the success of first-order gradient based heuristic algorithms for training. This success started with the work of [3] which presented empirical evidence that if DNNs are initialized properly, then we can find good solutions in reasonable runtime. This work was soon followed by series of early successes of deep learning in natural language processing [4], speech recognition [5] and visual object classification [6]. It was empirically shown in [7] that a sufficiently over-parameterized neural network can be trained to global optimality.

These gradient-based heuristics are not useful for neural networks with threshold activation function as there is no gradient information. Even networks with sigmoid activation function fell out of favor because gradient information is not valuable when input values are large [8]. The popular neural network architecture uses *ReLU activations* on which the gradient based methods are useful. Formally, the ReLU function is given by: $[x]_+ := \max(x, 0)$.



(a) ReLU network studied in [11, 12]



(b) Fully connected ReLU

Figure 1: Difference between ReLU model studied in [11, 12] and typical fully connected counterpart

Related literature. As discussed before, most hardness results so far are for neural networks with threshold activation function [1, 9, 10]. There are also limited results for ReLU that we discuss next: Recently, [11] examined ReLU activation from the point of view that if two connected ReLU nodes are appropriately designed, then it yields an approximation to the threshold function. Hence training problem for such a class of ReLU network should be as hard as training a neural network with threshold activation function. Similar results are shown by [12]. In both these papers, in order to approximate the threshold activation function, the neural network studied is not a fully connected network. More specifically, in the underlying graph of such a neural network, each node in the second hidden layer is connected to exactly one distinct node in the first hidden layer, weight of the connecting edge is set to -1 with the addition of some positive bias term. Figure 1 shows the difference between ReLU network studied by [11, 12] and fully connected ReLU network. The architecture artificially restricts the form of the affine functions in order to prove NP-hardness. In particular, it requires connecting hidden layer matrix to be a square diagonal matrix. Due to this restriction, it was unclear whether allowing non-diagonal entries of the matrix to be non-zero would make problem easier (more parameters hence higher power to neural network function) or more challenging (more parameters so more things to decide).

Another line of research on the hardness of training ReLU neural networks assumes that the data is coming from some distribution. More recent works in this direction include [13] which shows a smooth family of functions for which the gradient of squared error function is not informative while training neural network over Gaussian input distribution. Another study in this line of work considers Statistical Query (SQ) framework [14] (which contains SGD algorithms) and shows that there exists a class of special functions generated by single hidden layer neural network for which learning will require exponential number of queries (i.e. sample gradient evaluations) for the data coming from the product measure of the real valued log-concave distribution. These are interesting studies in their own right and generally consider hardness with respect to the algorithms that use stochastic gradient queries and require that such algorithm must perform minimization of the (expectation) objective functions. In comparison, we consider the framework of NP-hardness which takes into account

the complete class of the polynomial time algorithms, generally assumes that the data is given and requires an optimal solution to the corresponding empirical objective.

Recently, [15] showed that a single hidden layer ReLU network can be trained in polynomial time when dimension of input, w_0 , is constant.

Based on the above discussion, we see that the status of the complexity of training the multi-layer fully-connected ReLU neural network remains open. Given the importance of the ReLU NN, this is an important question. In this paper, we take the first steps in resolving this question.

Main Contributions.

- NP-hardness: We show that the training problem for a simple two hidden layer fully-connected NN which has two nodes in the first layer, one node in the second layer and ReLU activation function at all nodes is NP-hard (Theorem 3.1). Underlying graph of this network is exactly the same as that in Blum et al. [1] but all activation functions are ReLU instead of threshold function. Techniques used in the proof are different from earlier work in the literature because there is no combinatorial interpretation to ReLU as opposed to the threshold function.
- Polynomial-time solvable cases: We present two cases where the training problem with ReLU activation function can be solved in polynomial-time. The first case is when the dimension of the input is fixed (Theorem 3.3). This result generalizes the result from [15] and uses the hyperplane arrangement theorem for its proof.

We also observe that when the number of nodes in the first layer of the network is equal to the number of input data points (Proposition 3.4) then there exists a polynomial time algorithm. The proof of this fact follows from a simple observation that reduces the problem to fitting a single hidden layer neural network and then applying the polynomial time algorithm result for single hidden layer neural network in the work of [7]. This is the highly over-parameterized neural network setting. This result leads to some interesting open questions that we discuss later.

2. Notation and Definitions

We use the following standard set notation $[n] := \{1, \dots, n\}$. The letter d generally denotes the dimension of input data, N denotes the number of data-points and unless explicitly specified, the output data is one dimensional.

The main training problem of interest for the paper corresponds to a neural network with 3 nodes. The underlying graph is a layered directed graph with two layers. The first layer

contains two nodes and the second layer contains one node. The network is fully connected feedforward network. One can write the function corresponding to this neural network as follows:

$$F(x) = [w_0 + w_1[a_1(x)]_+ + w_2[a_2(x)]_+]_+, \quad (2)$$

where $a_i : \mathbb{R}^d \rightarrow \mathbb{R}$ for $i \in \{1, 2\}$ are real valued affine functions, and $w_0, w_1, w_2 \in \mathbb{R}$. The

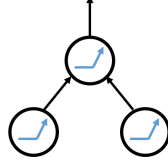


Figure 2: (2,1)-ReLU Neural Network. Also called 2-ReLU NN after dropping ‘1’. Here ReLU function is presented in each node to specify the type of activation function at the output of each node.

output of the two affine maps a_1, a_2 are the inputs to the two ReLU nodes in first hidden layer of network. The weights $\{w_0, w_1, w_2\}$ denote affine map for ReLU node in second layer. We refer to the network defined in (2) as (2,1)-ReLU Neural Network(NN). As its name suggests, it has 2 ReLU nodes in first layer and 1 ReLU node in second layer. We will refer to (k, j) -ReLU NN as a generalization of (2, 1)-ReLU NN where there are k ReLU nodes in first layer and j ReLU nodes in second layer. Note that the output of (k, j) -ReLU NN lies in \mathbb{R}^j . If there is only one node in the second layer, we will often drop the “1” and refer it as a 2-ReLU NN or k -ReLU NN depending on whether there are 2 or k nodes in the first layer, respectively. Figure 2 shows 2-ReLU NN.

Observation 1 *Note that*

$$w[ax + b]_+ \equiv \text{sgn}(w)[|w|(ax + b)]_+ = \text{sgn}(w)[\tilde{a}x + \tilde{b}],$$

so without loss of generality we will assume $w_1, w_2 \in \{-1, 1\}$ in (2).

Now we formally state the definition of the decision version of the training problem.

Definition 2.1 (Decision-version of the training problem) *Given a set of training data $(x^i, y^i) \in \mathbb{R}^d \times \{1, 0\}$ for $i \in S$, do there exist edge weights so that the resulting function F satisfies $F(x^i) = y^i$ for all $i \in S$.*

The decision version of the training problem in Definition 2.1 is asking if it is possible to find edge weights to obtain zero loss function value in the expression (1), assuming l is a norm, i.e., $l(a, b) = 0$ iff $a = b$.

3. Main Results

Theorem 3.1 *It is NP-hard to solve the training problem for 2-ReLU NN.*

An immediate corollary of Theorem 3.1 is the following:

Corollary 3.2 *Training problem of (2,j)-ReLU NN is NP hard, for all $j \geq 1$.*

The proof of Theorem 3.1 is obtained by reducing the 2-Hyperplane Separability Problem to the training problem of 2-ReLU NN. Details of this reduction and the proof of Theorem 3.1 and Corollary 3.2 are presented in Section 4.

While this paper was in submission, two more studies [16, 17] considered the computational complexity of training a single ReLU node and proved that it is a NP-hard problem. [16] also showed that it is NP-hard to train one hidden layer neural network with two nodes and ReLU activation at each node. This network basically removes the second layer ReLU activation and affine constant w_0 in (2) so that neural network function of their case can be rewritten as $F(x) = w_1[a_1(x)]_+ + w_2[a_2(x)]_+$. These are different network architectures and hence hardness of training any one of them does not necessarily imply hardness of training for remaining neural networks.

Megiddo [2] shows that the separability with fixed number of hyperplanes (generalization of 2-hyperplane separability problem) can be solved in polynomial-time in fixed dimension. Therefore 2-hyperplane separability problem can be solved in polynomial time given dimension is constant. Based on the reduction used to prove Theorem 3.1, a natural question to ask is “Can one solve the training problem of 2-ReLU NN problem in polynomial time under the same assumption?”. We answer this question in the affirmative.

Theorem 3.3 *Under the assumption that the dimension of input, d and the number of nodes in the first layer, k , are constant, then there exists a $\text{poly}(N)$ -time solution to the training problem of k -ReLU neural network, where N is the number of data-points.*

The high-level idea of the proof is the following: each data point “passes through” the three ReLU nodes and the activation function in these nodes is “turned on” or “turned off” (i.e., the output is 0 or not). We will enumerate all possible combinations of the data points being turned on or not, which we show is $\text{poly}(N)$ assuming d and k is fixed (by use of the Hyperplane Arrangement Theorem). Then we show that for each of these combinations and for each possible sign pattern of the weights defining the affine function applied at the second layer, corresponding optimal affine functions can be calculated via solving one convex program of poly size. Finally, we select the best optimal affine function which minimizes the loss function. Technique of Hyperplane Arrangement Theorem to enumerate partition

was used in [15] for proving poly(N)-time algorithms for single hidden layer neural networks. We extend this result for k -ReLU neural network which is a two hidden layer network. The complication due to second layer ReLU node are handled by solving a convex program of poly size. We present the formal proof of Theorem 3.3 in Appendix A.1.

We also study this problem under over-parameterization. Structural understanding of 2-ReLU NN yields an easy algorithm to solve training problem for N-ReLU neural network over N data points. In fact, the problem can be easily reduced to a single hidden layer NN.

Proposition 3.4 *Given data, $\{x^i, y^i\}_{i \in [N]}$ (where we assume that x^i s are distinct), then the training problem for N-ReLU NN has a poly(N,d)-time randomized algorithm, where N is the number of data-points and d is the dimension of input.*

The proof of this proposition is by first reducing the problem to that of training a single hidden layer network with N nodes on a dataset of size N. Then a polynomial time algorithm from [7] is applied for interpolating the data. The precise details are presented in Appendix A.2.

4. Training 2-ReLU NN is NP-hard

In this section, we give details about the NP-hardness reduction for the training problem of 2-ReLU NN. We begin with the formal definition of 2-Hyperplane Separability Problem.

Definition 4.1 (2-Hyperplane Separability Problem) *Given a set of points $\{x^i\}_{i \in [N]} \in \mathbb{R}^d$ and a partition of $[N]$ into two sets: S_1, S_0 , (i.e., $S_1 \cap S_0 = \emptyset$, $S_1 \cup S_0 = [N]$) decide whether there exists two hyperplanes $H_1 = \{x : \alpha_1^T x + \beta_1 = 0\}$ and $H_2 = \{x : \alpha_2^T x + \beta_2 = 0\}$ where $\alpha_1, \alpha_2 \in \mathbb{R}^d$ and $\beta_1, \beta_2 \in \mathbb{R}$ that separate the set of points in the following fashion:*

- i For each point x^i such that $i \in S_1$, both $\alpha_1^T x^i + \beta_1 > 0$ and $\alpha_2^T x^i + \beta_2 > 0$.*
- ii For each point x^i such that $i \in S_0$, $\alpha_1^T x^i + \beta_1 < 0$ or $\alpha_2^T x^i + \beta_2 < 0$.*

The 2-hyperplane separability problem is NP-complete [2]. Note the difference between conditions i and ii above. First one is an “AND” statement and second is an “OR” statement. Geometrically, solving 2-hyperplane separability problem means that finding two affine hyperplanes $\{\alpha_1, \beta_1\}$ and $\{\alpha_2, \beta_2\}$ such that all points in set S_1 lie in one quadrant formed by two hyperplanes and all points in set S_0 lie outside that quadrant. Due to this geometric intuition, the problem is called 2-hyperplane separability. We will construct a polynomial reduction from this NP-complete problem to training 2-ReLU NN, which will prove that training 2-ReLU NN is NP-hard.

Remark 4.1 (Variants of 2-hyperplane separability) *Note here that some sources also define 2-hyperplane separability problem with a minor difference. This difference is that strict inequalities, ‘>’, in Definition 4.1.i are diluted to inequalities, ‘≥’. In fact, these two problems are equivalent in the sense that there is a solution for the first problem if and only if there is a solution for the second problem. Solution for the first problem implies solution for the second problem trivially. Suppose there is a solution for the second problem, that implies there exist $\{\alpha_1, \beta_1\}$ and $\{\alpha_2, \beta_2\}$ such that for all $i \in S_0$ we have either $\alpha_1^T x^i + \beta_1 < 0$ or $\alpha_2^T x^i + \beta_2 < 0$. This implies $\epsilon := \min_{i \in S_0} \max\{-\alpha_1 x^i - \beta_1, -\alpha_2 x^i - \beta_2\} > 0$. So if we shift both planes by $\frac{1}{2}\epsilon$ i.e. $\beta_i \leftarrow \beta_i + \frac{1}{2}\epsilon$ then this is a solution to the first problem.*

Assumption: $\mathbf{0} \in S_1$. (Here $\mathbf{0} \in \mathbb{R}^d$ is a vector of zeros.) Suppose we are given a generic instance of 2-hyperplane separability problem with data-points $\{x^i\}_{i \in [N]}$ from \mathbb{R}^d and partition S_1 and S_0 of the set $[N]$. Since the answer of 2-hyperplane separability instance is invariant under coordinate translation, we can shift the origin to any x^i for $i \in S_1$, and therefore assume that the origin belongs to S_1 henceforth.

4.1. Reduction

Now we create a particular instance for 2-ReLU NN problem from a general instance of 2-hyperplane separability. We add two new dimensions to each data-point x^i . We also create a label, y^i , for each data-point. Moreover, we add a constant number of extra points to the training problem. Exact details are as follows:

Consider training set $\{(x^i, 0, 0), y^i\}_{i \in [N]}$ where $y^i = \begin{cases} 1 & \text{if } i \in S_1 \\ 0 & \text{if } i \in S_0 \end{cases}$.

Add additional 18 data points to the above training set as follows:

$$\begin{aligned} p_1 &\equiv \{(\mathbf{0}, 1, 1), 1\}, p_2 \equiv \{(\mathbf{0}, 1.5, 0.75), 1\}, p_3 \equiv \{(\mathbf{0}, 2, 1), 1\}, p_4 \equiv \{(\mathbf{0}, 2.25, 1.5), 1\}, \\ p_5 &\equiv \{(\mathbf{0}, 2, 2), 1\}, p_6 \equiv \{(\mathbf{0}, 1.5, 2.25), 1\}, p_7 \equiv \{(\mathbf{0}, 1, 2), 1\}, p_8 \equiv \{(\mathbf{0}, 0.75, 1.5), 1\}, \\ p_9 &\equiv \{(\mathbf{0}, 0, -1), 0\}, p_{10} \equiv \{(\mathbf{0}, 1, -1), 0\}, p_{11} \equiv \{(\mathbf{0}, 2, -1), 0\}, \\ p_{12} &\equiv \{(\mathbf{0}, 3, -1), 0\}, p_{13} \equiv \{(\mathbf{0}, 5, -1), 0\}, \\ p_{14} &\equiv \{(\mathbf{0}, -1, 0), 0\}, p_{15} \equiv \{(\mathbf{0}, -1, 1), 0\}, p_{16} \equiv \{(\mathbf{0}, -1, 2), 0\}, \\ p_{17} &\equiv \{(\mathbf{0}, -1, 3), 0\}, p_{18} \equiv \{(\mathbf{0}, -1, 5), 0\}. \end{aligned}$$

We call the set of additional data points with label 1 as T_1 and the set of additional data points with label 0 as T_0 . These additional data points (we refer to these points as the “gadget points”) are of fixed size. So this is a polynomial time reduction.

Figure 3 shows the gadget points. Note that origin is added to the gadget because there exists $i \in S_1$ such that $x^i = \mathbf{0}$. Hence training set has the data-point $\{(\mathbf{0}, 0, 0), 1\}$.

Henceforth we refer to the training problem of fitting 2-ReLU NN to this data as (\mathbf{P}) . In the context of the training problem (\mathbf{P}) , we abuse the notation and call the set of points $(x^i, 0, 0)$

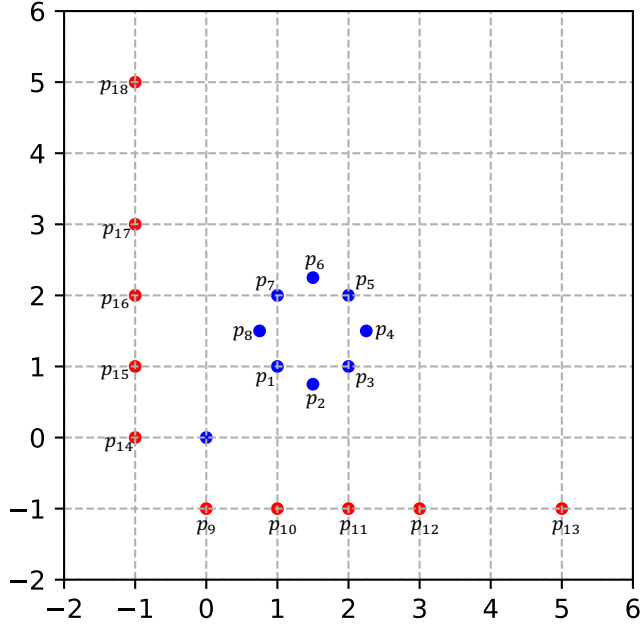


Figure 3: Gadget: Blue points represent set T_1 and red points represent set T_0 .

with label 1 as S_1 and the set of points $(x^i, 0, 0)$ with label 0 as S_0 . In particular, there is a direct correspondence between the sets S_1, S_0 defined in 2-hyperplane separability problem and sets S_1, S_0 defined for 2-ReLU NN training problem (\mathbf{P}). Use of our notation is generally clear from the context.

Now what remains is to show that the general instance of 2-hyperplane separability has a solution if and only if the constructed instance of 2-ReLU NN has a solution. In order to understand our approach better, we introduce the notion of “hard-sorting”. Hard-sorting is formally defined below, and its significance is stated in Lemma 4.5.

Definition 4.2 (Hard-sorting) *We say that a set of points $\{\pi^i\}_{i \in S}$, partitioned into two sets Π_0, Π_1 can be hard-sorted with respect to Π_1 if there exist two affine transformations l_1, l_2 and scalars w_1, w_2, c such that the following condition is satisfied:*

$$w_1[l_1(\pi)]_+ + w_2[l_2(\pi)]_+ \begin{cases} = c & \text{for all } \pi \in \Pi_1 \\ < c & \text{for all } \pi \in \Pi_0 \end{cases} \quad (3)$$

Being able to hard-sort implies that after passing the data through two nodes of the first hidden layer, the scalar input to the second hidden layer node must have a separation of the data-points in Π_1 and the data-points in Π_0 . Moreover, scalar input corresponding to all data points in Π_1 must be equal.

Remark 4.2 *If there exists scalars w_1, w_2, c and affine transformations l_1, l_2 such that*

$$w_1[l_1(\pi)]_+ + w_2[l_2(\pi)]_+ \begin{cases} = c & \text{for all } \pi \in \Pi_1; \\ > c & \text{for all } \pi \in \Pi_0, \end{cases}$$

then $-w_1, -w_2, -c, l_1, l_2$ satisfy condition (3) of hard-sorting.

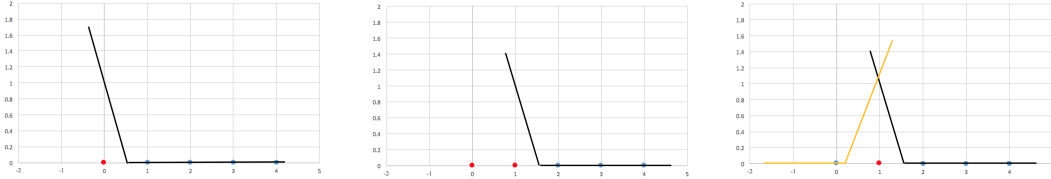
Remark 4.3 *Let $\bar{\Pi}_0 \subset \Pi_0$ and $\bar{\Pi}_1 \subset \Pi_1$. Then hard-sorting of $\Pi_0 \cup \Pi_1$ with respect to $\Pi_1 \Rightarrow$ hard-sorting of $\bar{\Pi}_0 \cup \bar{\Pi}_1$ with respect to $\bar{\Pi}_1$.*

Remark 4.4 *Without loss of generality, we may assume that $w_1, w_2 \in \{-1, 1\}$.*

It is not difficult to see that hard-sorting implies (\mathbf{P}) has a solution. We show that hard-sorting is also required for solving training problem. This is formally stated in lemma below.

Lemma 4.5 *The 2-ReLU NN training problem (\mathbf{P}) has a solution if and only if data-points $S_1 \cup T_1 \cup S_0 \cup T_0$ are hard-sorted with respect to $S_1 \cup T_1$.*

The proof of Lemma 4.5 can be found in Appendix A.4 . Figure 4 presents a geometric interpretation of Lemma 4.5



(a) Input is hard-sorted. This can give a perfect fit. (b) Since there are two red points so input is not hard-sorted. This cannot give a perfect fit. (c) Since blue points lies on different side of red points so input is not hard-sorted. This cannot give a perfect fit.

Figure 4: X-axis in figures above is output of the first layer of 2-ReLU NN i.e. $w_1[l_1(\pi)]_+ + w_2[l_2(\pi)]_+$. Y-axis is the output of second hidden layer node. Since output of first hidden layer goes to input of second hidden layer, we are essentially trying to fit ReLU node of second hidden layer. In particular, red and blue dots represent output of first hidden layer on data points with label 1 and 0 respectively. In Fig. (a) we see that hard-sorted input can be classified as 0/1 by a ReLU function. In Fig. (b) and Fig. (c) we see that input which is not hard-sorted cannot be classified exactly as 0/1 by a ReLU function.

We use the hard-sorting characterization of the solution of the training problem (\mathbf{P}) extensively. We first show the forward direction of the reduction in the lemma below. This is also the easier direction.

Lemma 4.6 *If 2-hyperplane separability problem has a solution, then problem (\mathbf{P}) has a solution.*

The proof of Lemma 4.6 can be found in Appendix A.3.

To prove the reverse direction we need to show that if a set of weights solve the training problem (\mathbf{P}) then we can generate a solution to the 2-hyperplane separability problem. In the rest of the proof we will argue that the only way to solve the training problem (\mathbf{P}) for 2-ReLU NN or equivalently hard-sort data-points is to find two affine functions $a_1, a_2 : \mathbb{R}^{d+2} \rightarrow \mathbb{R}$ such that i) $a_1(x) \leq 0$ and $a_2(x) \leq 0$ for all $x \in S_1 \cup T_1$ and ii) $a_1(x) > 0$ or $a_2(x) > 0$ for all $x \in S_0 \cup T_0$. If such a solution exists then there exists a solution to 2-hyperplane separability problem after dropping coefficients of last two dimensions of affine functions $-a_1$ and $-a_2$. Note that changing ' $<$ ' to ' \leq ' in 2-affine separability problem is valid in view of Remark 4.1. We will first show that we can hard-sort the gadget points only under the properties of a_1 and a_2 mentioned above. This implies that a solution to (\mathbf{P}) which hard-sorts all points (including the gadget points) must have same properties of a_1 and a_2 . This follows from counter-positive of Remark 4.3, i.e., if a subset of data-points cannot be hard-sorted then all data-points cannot be hard-sorted. Henceforth, we will focus on the gadget data-points (or the last two dimensions of the data).

4.1.1. Gadget Points and Hard-Sorting

In the following lemma, we show a necessary condition on a_1, a_2 satisfying hard-sorting of gadget data points $T_1 \cup T_0 \cup \{\mathbf{0}\}$ with respect to $T_1 \cup \{\mathbf{0}\}$.

Lemma 4.7 *If affine functions $a_1, a_2 : \mathbb{R}^{d+2} \rightarrow \mathbb{R}$ and scalars w_1, w_2, c satisfy hard-sorting of the data-points $T_1 \cup T_0 \cup \{\mathbf{0}\}$ with respect to $T_1 \cup \{\mathbf{0}\}$, then all points in T_1 must satisfy $a_1(x) \leq 0, a_2(x) \leq 0$. Moreover, we must have $w_1 = w_2 = -1$ and $c = 0$.*

Note that in view of Lemma 4.7 and counter-positive of Remark 4.3, we have that affine function $a_1, a_2 : \mathbb{R}^{d+2} \rightarrow \mathbb{R}$ and scalars w_1, w_2, c satisfying hard-sorting of $S_1 \cup T_1 \cup S_0 \cup T_0$ with respect to $S_1 \cup T_1$ must satisfy

$$-\left[a_1(x) \right]_+ - \left[a_2(x) \right]_+ \begin{cases} = 0 & \text{if } x \in S_1; \\ < 0 & \text{if } x \in S_0 \end{cases}.$$

The above condition is equivalent to the requirement that $a_1(x) \leq 0, a_2(x) \leq 0$ for all $x \in S_1$ and $a_1(x) > 0$ or $a_2(x) > 0$ for $x \in S_0$. After dropping the last two dimensions of $-a_1$ and $-a_2$, we obtain the solution for 2-affine separability problem. Now that we have reduced the problem to the key lemma above, the main purpose of this section is to prove Lemma 4.7.

Note that for each data point in the gadget $T_1 \cup T_0 \cup \{\mathbf{0}\}$, the first d elements are always 0. So for the sake of gadget, we may assume that $a_1, a_2 : \mathbb{R}^2 \rightarrow \mathbb{R}$ and the gadgets lies in \mathbb{R}^2 . They can be thought of as the projection of the original $a_i : \mathbb{R}^{d+2} \rightarrow \mathbb{R}$ and $\mathbf{0} \in \mathbb{R}^{d+2}$ to last two dimension which are relevant for gadget data points $T_1 \cup T_0$. Due to this observation, we

assume that $a_1, a_2 : \mathbb{R}^2 \rightarrow \mathbb{R}$ henceforth for this subsection, and provide a proof of Lemma 4.7 under this assumption.

The proof of Lemma 4.7 is divided into the following sequence of results.

Proposition 4.8 *Suppose that a_1, a_2 satisfy hard-sorting of $T_1 \cup T_0$ with respect to T_1 . Then there exists $x \in T_1$ such that $a_1(x) \leq 0, a_2(x) \leq 0$.*

Proof of Proposition 4.8 can be found in Appendix A.5.

Next we show one more simple proposition which is critical in proving the final result. The proof of this proposition can be found in Appendix A.7.

Proposition 4.9 *If affine functions a_1, a_2 and weights w_1, w_2 satisfy hard-sorting of $T_1 \cup T_0 \cup \{0\}$ with respect to $T_1 \cup \{0\}$, then w_1, w_2 **must** satisfy $w_1 = w_2 = -1$.*

We are now ready to present the prove Lemma 4.7.

Proof of Lemma 4.7. Since a_1, a_2 satisfy hard-sorting of the data points $T_1 \cup T_0 \cup \{0\}$ with respect to $T_1 \cup \{0\}$ then, in view of Proposition 4.8 and Proposition 4.9, we have

1. $\exists x \in T_1$ such that $a_1(x) \leq 0, a_2(x) \leq 0$.
2. $w_1 = w_2 = -1$.

Then we have that $-[a_1(x)]_+ - [a_2(x)]_+ = 0$ for all $x \in T_1$, due to condition (3) of hard-sorting. This implies $a_1(x) \leq 0, a_2(x) \leq 0$ for all $x \in T_1$. So we conclude the proof. \square

In the next section, we show that this result on the gadget data-points gives us the solution to the original 2-hyperplane separability problem.

4.1.2. From Gadget Data to Complete Data

Lemma 4.10 *If there is a solution to the problem (P), then there is a solution to corresponding 2-hyperplane separability problem.*

Proof. Note that if there is a solution to problem (P), then by Lemma 4.5, we must have $a_1, a_2 : \mathbb{R}^{d+2} \rightarrow \mathbb{R}$ and w_1, w_2, c hard-sorting $S_1 \cup T_1 \cup S_0 \cup T_0$ with respect to $S_1 \cup T_1$. In view of Lemma 4.7 and counter-positive of Remark 4.3, we have

1. $w_1 = w_2 = -1$.
2. $w_1[a_1(x)]_+ + w_2[a_2(x)]_+ = 0$ for all $x \in S_1 \cup T_1$ due to requirement (3) of hard-sorting.

Since $w_1 = w_2 = -1$, so 2 above implies $a_1(x) \leq 0$ and $a_2(x) \leq 0$ for all $x \in S_1 \cup T_1$. Moreover, we require $a_1(x) > 0$ or $a_2(x) > 0$ for all $x \in S_0 \cup T_0$ due to condition (3) of hard-sorting. Now as discussed earlier, $-a_1, -a_2$ after ignoring coefficients of last two dimensions will yield

solution to 2-hyperplane separability problem. Hence we conclude the proof. \square

Now we are ready to prove the main NP-hardness theorem.

Proof of Theorem 3.1. Using Lemma 4.6 and Lemma 4.10, we conclude the proof. \square

Below we state an immediate corollary of Theorem 3.1 whose proof can be found in Appendix A.8.

Corollary 4.11 *Training problem of (2,j)-ReLU NN is NP hard.*

5. Discussion

We showed that the problem of training 2-ReLU NN is NP-hard. Given the importance of ReLU activation function in neural networks, in our opinion, this result resolves a significant gap in understanding complexity class of the problem at hand. On the other hand, we show that the problem of training N -ReLU NN is in P. So a natural research direction is to understand the complexity status when input layer has more than 2 nodes and strictly less than N nodes. A particularly interesting question in that direction is to generalize the gadget we used for 2-ReLU NN to the case of k -ReLU NN.

Acknowledgments

We would like to thank anonymous reviewers whose comments helped in simplifying several results in this paper.

Appendices

A. Proofs of Auxiliary Results

In this appendix, we provide proof of all auxiliary results.

A.1. Proof of Theorem 3.3

Suppose we partition the set $[N]$ into sets Q_j and \bar{Q}_j such that all points in Q_j satisfy $a_j(x) \geq 0$ and all points in \bar{Q}_j satisfy $a_j(x) < 0$ for all $j \in [k]$. Given a set $S \subseteq [k]$, we define $T(S) := \left(\bigcap_{j \in S} Q_j\right) \cap \left(\bigcap_{j \in \bar{S}} \bar{Q}_j\right)$ where $\bar{S} = [k] \setminus S$. Let $z = (a_1, \dots, a_k, w_0, w_1, \dots, w_k)$. Then the objective function can be written as

$$f(z) = \sum_{S \subseteq [k]} \sum_{i \in T(S)} \left([w_0 + \sum_{j \in S} w_j a_j(x^i)]_+ - y_i \right)^2.$$

Now we can partition $T(S)$ into sets $T(S)_1$ and $T(S)_2$ for each $S \subseteq [k], S \neq \phi$. For $T(S)_1$, the ReLU term in the objective, $w_0 + \sum_{j \in S} w_j a_j(x)$ (note that this is an affine function), is

constrained to be non-negative and for $T(S)_2$ the ReLU terms is constrained to be non-positive. We need not enumerate partitions of $T(\phi)$ since ReLU terms for $T(\phi)$ do not depend on data-points. The key observation is that the partition of $T(S)$ into sets $T(S)_1$ and $T(S)_2$ is a partition due to a hyperplane.

Number of combinations: According to the Hyperplane Arrangement Theorem, given a set of points $\{x^i\}_{i \in N}$ in \mathbb{R}^d , the number of distinct partitions created by linear separators is $O(N^d)$. Moreover, due to [18], we can enumerate all possible partitions created by linear separators in $O(N^d)$ time. Therefore, there are a total of $O(N^{kd})$ possible combinations of $Q_j, j \in [k]$. For each such $Q_j, j \in [k]$, there are 2^k non-empty subsets $T(S) \subseteq [N]$. For each $T(S), S \neq \phi$, there are $O(|T(S)|^d) = O(N^d)$ possible ways to partition $T(S)$ into $T(S)_1$ and $T(S)_2$. So number of product combinations is $O(N^{(2^k-1)d})$. Hence there are a total of $O(N^{(kd+(2^k-1)d)})$ combinations.

Number of convex programs: By Observation 1 it suffices to check for $w_1, \dots, w_k = \pm 1$. We will divide the optimization problem in two cases $w_0 \geq 0$ and $w_0 \leq 0$. So there are a total of 2^{k+1} convex programs for each possible combination of $Q_j, T(S)_1$ for all $S \subseteq [k]$ of the following form:

$$\min \sum_{\substack{S \subseteq [k], \\ S \neq \phi}} \left\{ \sum_{i \in T(S)_1} \left(\left(w_0 + \sum_{j \in S} w_j a_j(x^i) \right) - y_i \right)^2 + \sum_{i \in T(S)_2} (0 - y_i)^2 \right\} + \sum_{i \in T_\phi} \left([w_0]_+ - y_i \right)^2$$

subject to constraints

$$\begin{aligned} a_j(x^i) &\geq 0, & \forall j, i \in Q_j \\ a_j(x^i) &\leq 0, & \forall j, i \in \overline{Q}_j \end{aligned} \tag{4}$$

$$\begin{aligned} w_0 + \sum_{j \in S} w_j a_j(x^i) &\geq 0, & \forall S \subseteq [k], S \neq \phi, i \in T(S)_1 \\ w_0 + \sum_{j \in S} w_j a_j(x^i) &\leq 0, & \forall S \subseteq [k], S \neq \phi, i \in T(S)_2 \end{aligned} \tag{5}$$

Moreover, we add constraint $w_0 \geq 0$ or $w_0 \leq 0$ and change the $[w_0]_+$ term in objective with w_0 or 0 respectively. Every program has $k(d+1) + 1$ variables in a_1, \dots, a_k, w_0 . Total number of constraints is at most $kN + N + 1$. Note that, for constraints of type (4), for each j , number of constraints equals $|Q_j \cup \overline{Q}_j| = N$. Hence total number of constraints of type (4) are kN . Similarly, for constraints of type (5), for each $S \subseteq [k]$, we have total of $|T(S)_1 \cup T(S)_2| = |T(S)|$ constraints. Hence total constraints of type (5) are $\sum_{\substack{S \subseteq [k], \\ S \neq \phi}} |T(S)| \leq N$ (This follows due to observation that $T(S), S \subseteq [k]$ is a partition of $[N]$). One more constraint is on w_0 . Hence total number of constraints is $(k+1)N + 1$. Since number of constraints and

variables are $\text{poly}(k, d, N)$ and objective is convex quadratic so we conclude that this program can be solved in $\text{poly}(N, k, d)$ time.

Finally, the total number of convex programs to be solved is $O(2^{k+1} \cdot N^{kd+(2^k-1)d})$.

A.2. Proof of Proposition 3.4

Before proving this proposition, we state a polynomial time algorithm (Theorem 1 of [7]) for training single hidden layer neural network.

Proposition A.1 *Let $f(x) = \sum_{j=1}^N w_j [a_j(x)]_+ + w_0$ be a single hidden layer neural network with N nodes. Let $x^i \in \mathbb{R}^d, i = 1, \dots, N$ be distinct data-points, and let $y^i \in \mathbb{R}, i = 1, \dots, N$ be arbitrary labels. Then there exists a $\text{poly}(N, d)$ -time algorithm to train this neural network such that $f(x^i) = y^i, i = 1, \dots, N$.*

Now we are ready to prove Proposition 3.4.

Note that a N-ReLU NN can be written as $c(x) = [\sum_{j=1}^N w_j [a_j(x)]_+ + w_0]_+$. Suppose $y = [y^1, \dots, y^N]^T \in \mathbb{R}^N$ be a vector of labels. We may assume that $y \geq \mathbf{0}$ since otherwise the answer to the training problem is a straight-forward “No”. Now, we prove that weights $w_i, i = 0, \dots, N$ and affine functions $a_j, j = 1, \dots, N$, satisfying $c(x^i) = y^i$ for all $i \in [N]$, can be computed in $\text{poly}(N, d)$ -time. Hence, providing a “Yes” answer to the training problem in $\text{poly}(N, d)$ -time.

Since $y^i \geq \mathbf{0}$ for all $i \in [N]$ and we want $y^i = c(x^i) = [f(x^i)]_+$, where $f(x) = \sum_{j=1}^N w_j [a_j(x)]_+ + w_0$, it suffices¹ to show that $f(x^i) = y^i$ for all $i \in [N]$. Using the fact that number of nodes in f matches the number of data points, N , then applying Proposition A.1, we obtain the result.

A.3. Proof of Lemma 4.6

Suppose (α_1, β_1) and (α_2, β_2) are solution satisfying condition for 2-hyperplane separability. Note that there is a data-point $\mathbf{0} \in S_1$ so we obtain $\beta_1, \beta_2 > 0$. Without loss of generality we can assume $\beta_1 = \beta_2 = 0.5$. This is due to the fact that scaling the original solution by any positive scalar yields a valid solution. Now we show that the solution of 2-hyperplane separability problem can be used to show hard-sorting of $S_0 \cup T_0 \cup S_1 \cup T_1$ with respect to $S_1 \cup T_1$. Hence in view of Lemma 4.5, we obtain the existence of a solution for problem **(P)**. Set $w_1 = w_2 = -1, c = 0$. Moreover, for $(x, y, z) \in \mathbb{R}^{d+2}$, consider the affine map $l_1(x, y, z) = -\alpha_1^T x - y - \beta_1$ and $l_2(x, y, z) = -\alpha_2^T x - z - \beta_2$. We claim that w_1, w_2, c, l_1, l_2 satisfy hard-sorting condition (3) for $S_0 \cup T_0 \cup S_1 \cup T_1$ with respect to $S_1 \cup T_1$. In particular, note that

¹Note that when $y^i = 0$, we can have $f(x^i) \leq 0$. However, assuming $f(x^i) = 0$ is sufficient since we can find a solution for arbitrary nonpositive assignments of y^i . This also keeps the proof clean as we can assume that $f(x^i) = y^i$ for all labels y^i uniformly (even for the 0-labels).

1. For $x \in S_1$, we have

$$-[-\alpha_1^T x - \beta_1]_+ - [-\alpha_2^T x - \beta_2]_+ = 0 = c.$$

2. For $x = (\mathbf{0}, l, m) \in T_1$, we have

$$-[-\beta_1 - l]_+ - [-\beta_2 - m]_+ = 0.$$

This follows since $\beta_1 = \beta_2 = 1/2$ and $l, m \in [0.75, 2.25]$ so the two ReLU terms inside are both zero for all $x \in T_1$.

3. For $x \in S_0$, we have

$$-[-\alpha_1^T x - \beta_1]_+ - [-\alpha_2^T x - \beta_2]_+ < 0.$$

This follows since at least one of $\alpha_1^T x + \beta_1$ and $\alpha_2^T x + \beta_2$ is strictly negative for $x \in S_0$ as (α_1, β_1) and (α_2, β_2) are solution for 2-hyperplane separability problem.

4. For $x = (\mathbf{0}, l, m) \in T_0$, we have

$$-[-\beta_1 - l]_+ - [-\beta_2 - m]_+ < 0.$$

This follows since $\beta_1 = \beta_2 = 1/2$ and either l or m equals -1 for $x \in T_0$.

This proves hard-sorting of $S_0 \cup T_0 \cup S_1 \cup T_1$ with respect to $S_1 \cup T_1$ and hence we have the existence of solution for training problem **(P)**.

A.4. Proof of Lemma 4.5

We first prove the forward direction. Suppose points are hard-sorted as required by the lemma.

Then define $\varepsilon := \min_{x \in S_0 \cup T_0} -w_1[l_1(x)]_+ - w_2[l_2(x)]_+ + c$. By definition, we have $\varepsilon > 0$. Then neural network $f(x) = \frac{2}{\varepsilon}[w_1[l_1(x)]_+ + w_2[l_2(x)]_+ - c + \varepsilon/2]_+$ solves training problem. This can be easily checked from the fact that

$$w_1[l_1(x)]_+ + w_2[l_2(x)]_+ - c \begin{cases} = 0 & \text{if } x \in S_1 \cup T_1; \\ < -\varepsilon & \text{if } x \in S_0 \cup T_0, \end{cases}$$

which holds under the assumption of hard-sorting.

Now we assume that points cannot be hard-sorted and conclude that there does not exist weight assignment solving training problem of 2-ReLU NN, hence proving the backward direction. Since the points cannot be hard-sorted so there does not exist any l_1, l_2, w_1, w_2, c satisfying condition (3). This fact along with Remark 4.2 implies that for all possible weights we either have

a) $w_1[l_1(x)]_+ + w_2[l_2(x)]_+$ is not constant for all $x \in S_1 \cup T_1$ or

b) If $w_1[l_1(x)]_+ + w_2[l_2(x)]_+ = c$ for all $x \in S_1 \cup T_1$ and some constant c , then same expression evaluated on $x \in S_0 \cup T_0$ is not strictly on same side of c .

If we choose l_1, l_2, w_1, w_2, c such that a) happens, then such weights will not solve training problem as their output of 2-ReLU NN for points $p \in S_1 \cup T_1$ will be at least two distinct numbers which is an undesirable outcome. Specifically, we want $[w_0 + w_1[l_1(x)]_+ + w_2[l_2(x)]_+]_+$ to evaluate to 1 for all $x \in S_1 \cup T_1$. Hence $w_1[l_1(x)]_+ + w_2[l_2(x)]_+$ must be a constant for all $x \in S_1 \cup T_1$. This requirement is violated in case a).

If we choose l_1, l_2, w_1, w_2, c such that b) happens, then we can set w_0, θ such that $F(x) = \theta[w_1[l_1(x)]_+ + w_2[l_2(x)]_+ + w_0]_+$, $w_0 + c > 0$ and $\theta = \frac{1}{w_0 + c}$. Here we introduced another parameter $\theta > 0$ in the definition of F for sake of convenience of argument but note that θ can be absorbed in the definition of l_1 and l_2 to obtain the original neural network function defined (2). Since not all $x \in S_0 \cup T_0$ are strictly on one side, we conclude there exist $x' \in S_0 \cup T_0$ such that $w_1[l_1(x')]_+ + w_2[l_2(x')]_+ = c' \geq c$ hence $F(x') := \theta[w_1[l_1(x')]_+ + w_2[l_2(x')]_+ + w_0]_+ \geq 1$ which is an undesirable outcome for a point with label 0.

Since all choices of l_1, l_2, w_1, w_2, c satisfy either a) or b), we conclude that there does not exist weights solving training problem of 2-ReLU NN.

A.5. Proof of Proposition 4.8

In order to prove Proposition 4.8, we need to prove one more technical result stated below. Proof of this new proposition is deferred to Appendix A.6 but here we state it and proceed with the proof of Proposition 4.8.

Proposition A.2 *Suppose affine functions $a_1, a_2 : \mathbb{R}^2 \rightarrow \mathbb{R}$ and $w_1, w_2 \in \{-1, 1\}$ be such that (i) $a_1(x)$ is a constant for all $x \in \mathbb{R}^2$ or (ii) $a_2(x)$ is a constant for all $x \in \mathbb{R}^2$ or (iii) $w_1 a_1(x) + w_2 a_2(x)$ is a constant for all $x \in \mathbb{R}^2$, then a_1, a_2, w_1, w_2 cannot satisfy hard-sorting of the data points $T_1 \cup T_0 \cup \{\mathbf{0}\}$ with respect to $T_1 \cup \{\mathbf{0}\}$.*

Now we are ready to prove Proposition 4.8.

Let a_1, a_2 satisfy hard-sorting of $T_1 \cup T_0 \cup \{\mathbf{0}\}$ with respect to $T_1 \cup \{\mathbf{0}\}$. Then due to Remark 4.3, we have that a_1, a_2 satisfy hard-sorting of $T_1 \cup T_0$ with respect to T_1 . We will show that any a_1, a_2 satisfying the above condition must satisfy the requirement of Proposition 4.8.

Let us partition the set of points \mathbb{R}^2 into four partitions $S_{0,0}, S_{+,0}, S_{0,+}$ and $S_{+,+}$ based on sign of $[a_1(x)]_+$ and $[a_2(x)]_+$. Then, we have to show that at least one element in T_1 lies in the partition $S_{0,0}$.

For sake of contradiction, assume that $T_1 \cap S_{0,0} = \emptyset$. Then, using pigeonhole principle, we have that at least one of $S_{+,0}, S_{0,+}$ and $S_{+,+}$ must contain three points from the set T_1 . Note that any three points in the set T_1 are not collinear. Moreover, the function $w_1[a_1]_+ + w_2[a_2]_+$ is affine in all three regions, $S_{+,0}, S_{0,+}$ and $S_{+,+}$ of \mathbb{R}^2 and is non-constant in view of Proposition

A.2. Hence, we cannot satisfy hard-sorting since those three points in T_1 will break the requirement in condition (3) for hard-sorting. Hence, we obtain a contradiction.

A.6. Proof of Proposition A.2

First observe that if $a_1, a_2 : \mathbb{R}^2 \rightarrow \mathbb{R}$ satisfy hard-sorting of $T_1 \cup T_0 \cup \{\mathbf{0}\}$ with respect to $T_1 \cup \{\mathbf{0}\}$, then neither of them can be a constant function. In particular, it is straightforward to see that both of them cannot be constant. If only one of them is constant, then data needs to be linearly separable which is not the case for gadget data-points $T_1 \cup T_0 \cup \{\mathbf{0}\}$. Therefore, we will assume that both of them are affine functions with non-zero normal vectors.

Note that in view of Remark 4.4 and the fact that $w_1 a_1(x) + w_2 a_2(x) = c$ for all $x \in \mathbb{R}^2$, we may assume that magnitude of the normal to these lines is equal i.e. $\|\nabla a_1\| = \|\nabla a_2\| \neq 0$. For the sake of this proof, we extend the definition of hard-sorting to include the condition

$$w_1 [a_1(x)]_+ + w_2 [a_2(x)]_+ \begin{cases} = c & \text{for all } x \in T_1 \cup \{\mathbf{0}\}; \\ > c & \text{for all } x \in T_0, \end{cases}$$

along with condition (3). Due to this extended definition and in view of Remark 4.2, we just need to check for case $(w_1, w_2) = (1, 1)$ and $(w_1, w_2) = (1, -1)$. More specifically, $(w_1, w_2) = (-1, -1)$ yields a hard-sorting solution iff there exists a hard-sorting solution for $(w_1, w_2) = (1, 1)$. Equivalent argument can be made about the case $(w_1, w_2) = (-1, 1)$ and $(w_1, w_2) = (1, -1)$.

Then, we have two possible situations here: a_1, a_2 satisfy 1) $a_1(x) + a_2(x) = c, \forall x \in \mathbb{R}^2$ when normals point in opposite directions and 2) $a_1(x) - a_2(x) = c, \forall x \in \mathbb{R}^2$ when normals point in same direction. We will consider both these cases separately and show that expression $w_1 [a_1]_+ + w_2 [a_2]_+$, for the choices of w_1, w_2 mentioned above, cannot hard-sort the data as required.

Case 1: Normals point in the opposite directions. Here $w_1 = w_2 = 1$ and we assume $a_1 + a_2 = c$. Suppose $c \geq 0$. Then it can be verified that

$$[a_1(x)]_+ + [a_2(x)]_+ = \begin{cases} c & \text{if } c \geq a_1(x) \geq 0 \\ a_1(x) & \text{if } a_1(x) \geq c \\ c - a_1(x) & \text{if } a_1(x) \leq 0. \end{cases}$$

By the extended hard-sorting requirement, we need all points in $T_1 \cup \{\mathbf{0}\}$ to be contained in the set $\{x : a_1(x) \in [0, c]\}$ and all points in T_0 should not be in this set. Now observe that if $c = 0$, then the set $\{x : a_1(x) = 0\}$ is one dimensional, and therefore cannot contain all the points of $T_1 \cup \{\mathbf{0}\}$. Hence we must have $c > 0$ and all points in $T_1 \cup \{\mathbf{0}\}$ lie inside the region of two parallel lines $a_1(x) = 0$ and $a_1(x) = c$ as $[a_1(x)]_+ + [a_2(x)]_+$ evaluates to the constant c

in this region. It can be seen that this separation of $T_1 \cup \{\mathbf{0}\}$ from T_0 is impossible to achieve by two parallel lines.

Similarly when $c < 0$, then it can be verified that

$$[a_1(x)]_+ + [a_2(x)]_+ = \begin{cases} 0 & \text{if } c \leq a_1(x) \leq 0 \\ a_1(x) & \text{if } a_1(x) \geq 0 \\ c - a_1(x) & \text{if } a_1(x) \leq c \end{cases}$$

Again, for the extended hard-sorting, as in the previous case, we need all points in $T_1 \cup \{\mathbf{0}\}$ should be in set $\{x : a_1(x) \in [c, 0]\}$ and all points in T_0 should not be in this set which cannot be achieved.

Case 2: Normals point in the same direction. Then $a_1(x) - a_2(x) = c$. Suppose $c \geq 0$. Then it can be verified that

$$[a_1(x)]_+ - [a_2(x)]_+ = \begin{cases} a_1(x) & \text{if } c \geq a_1(x) \geq 0 \\ c & \text{if } a_1(x) \geq c \\ 0 & \text{if } a_1(x) \leq 0 \end{cases}$$

If $c = 0$ then $[a_1(x)]_+ - [a_2(x)]_+ = 0$ for all $x \in \mathbb{R}^2$. So this cannot hard-sort data. Hence for hard-sorting we definitely need $c > 0$. Moreover, we need either 1) $T_1 \cup \{\mathbf{0}\} \subset \{x : a_1(x) \leq 0\}$ and $T_0 \subset \{x : a_1(x) > 0\}$ or 2) $T_1 \cup \{\mathbf{0}\} \subset \{x : a_1(x) \geq c\}$ and $T_0 \subset \{x : a_1(x) < c\}$. In both cases, we need the points in $T_1 \cup \{\mathbf{0}\}$ must be separable from the points in T_0 by a line. This is not possible.

Note that when $c < 0$, one can write $a_2 - a_1 = -c$ and write similar functional form for $[a_2]_+ - [a_1]_+$.

Since in both cases, we were unable to achieve hard-sorting $T_1 \cup T_0 \cup \{\mathbf{0}\}$ w.r.t. $T_1 \cup \{\mathbf{0}\}$, so we conclude the proof.

A.7. Proof of Lemma 4.9

Proposition 4.8 yields that any hard-sorting a_1, a_2 must satisfy $a_1(x) \leq 0, a_2(x) \leq 0$ for at least one $x \in T_1$.

Now, suppose sign of w_1, w_2 is different. Suppose $w_1 = 1, w_2 = -1$. Since a_1 and a_2 satisfy hard-sorting of gadget so we have $[a_1(x)]_+ - [a_2(x)]_+ = c, \forall x \in T_1$. Due to Proposition 4.8, we obtain $c = 0$. Then to fulfill hard-sorting condition, we need $[a_1(x)]_+ - [a_2(x)]_+ < 0 \forall x \in T_0$. (The case for $w_1 = -1, w_2 = 1$ will have same proof with all a_2 exchanged by a_1 in next 3 lines.) This implies $a_2(x) > 0$ for all $x \in T_0$. However note that $T_1 \subset \text{conv}(T_0)$. So we get a contradiction to the assumption that sign of weights w_1, w_2 is different.

Now note that if sign of w_1, w_2 is same then we cannot set $w_1 = w_2 = 1$ due to requirement (3) of hard-sorting. Hence we have that $w_1 = w_2 = -1$.

A.8. Proof of Corollary 4.11

The reduction is again from 2-affine separability problem. It involves the same gadget of 18 points in Figure 3 and similar labels except that labels need to be extended from \mathbb{R} to \mathbb{R}^j . Simply add $j - 1$ zeros to original output labels in (\mathbf{P}) . We call this instance (\mathbf{P}') .

First, we show that if there is a solution to 2-affine separability problem then there is a solution for (\mathbf{P}') . For the output of the first node, we can use the construction in Lemma 4.6 to obtain an exact fit for the first node. For rest $j - 1$ nodes, the output is 0 for all data-points so we can easily extend the solution to obtain an exact fit for all nodes. In particular, for $k \in [j]$, every k -th node in the second layer is connected to two nodes in the first layer by distinct edges whose weights are parameterized by $w_{k,1}, w_{k,2}$ and bias weight $w_{k,0}$. We can set $w_{k,1} = w_{k,2} = -1$ and $w_{k,0} = 0$ for all $k \in [j] \setminus \{1\}$. The output at k -th node can be written as $[w_{k,1}[a_1(x)]_+ + w_{k,2}[a_2(x)]_+ + w_{k,0}]_+$. In view of the above values of $w_{k,0}, w_{k,1}$ and $w_{k,2}$ for $k \in [j] \setminus \{1\}$, we note that $w_{k,1}[a_1(x)]_+ + w_{k,2}[a_2(x)]_+ + w_{k,0} \leq 0$ for all $k \in [j] \setminus \{1\}$. This yields the output 0 at all nodes $k \in [j] \setminus \{1\}$, irrespective of the affine functions a_1, a_2 in the first layer. Hence all nodes are satisfied to global optimality which completes the forward direction of the reduction.

Reverse direction follows immediately from the construction of the gadget and Lemma 4.10. To see this, note that solution for (\mathbf{P}') implies that first node is satisfied to global optimality. Ignoring the rest $j - 1$ nodes, we see that this solution is also a solution for (\mathbf{P}) . At this point, invoking Lemma 4.10 yields the solution corresponding to 2-hyperplane separability problem. Hence, we conclude that training problem of $(2, j)$ -ReLU NN is NP-hard.

References

- [1] A. Blum, R. L. Rivest, Training a 3-node neural network is np-complete, in: Proceedings of the First Annual Workshop on Computational Learning Theory, COLT '88, 1988, pp. 9–18.
- [2] N. Megiddo, On the complexity of polyhedral separability., *Discrete and Computational Geometry* 3 (4) (1988) 325–338.
- [3] G. E. Hinton, S. Osindero, Y.-W. Teh, A fast learning algorithm for deep belief nets, *Neural Computation* 18 (7) (2006) 1527–1554.
- [4] R. Collobert, J. Weston, A unified architecture for natural language processing: Deep neural networks with multitask learning, in: Proceedings of the 25th International Conference on Machine Learning, ICML '08, 2008, pp. 160–167.

- [5] A. Mohamed, G. E. Dahl, G. Hinton, Acoustic modeling using deep belief networks, *Trans. Audio, Speech and Lang. Proc.* (2012) 14–22.
- [6] A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with deep convolutional neural networks, in: *Proceedings of the 25th International Conference on Neural Information Processing Systems - Volume 1, NIPS'12, 2012*, pp. 1097–1105.
- [7] C. Zhang, S. Bengio, M. Hardt, B. Recht, O. Vinyals, Understanding deep learning requires rethinking generalization, *CoRR* (2016).
- [8] S. Hochreiter, Y. Bengio, P. Frasconi, J. Schmidhuber, Gradient flow in recurrent nets: the difficulty of learning long-term dependencies (2001).
- [9] A. R. Klivans, A. A. Sherstov, Cryptographic hardness for learning intersections of half-spaces, *J. Comput. Syst. Sci.* 75 (1) (2009) 2–12.
- [10] S. Shalev-Shwartz, S. Ben-David, *Understanding Machine Learning: From Theory to Algorithms*, Cambridge University Press, 2014.
- [11] R. Livni, S. Shalev-Shwartz, O. Shamir, On the computational efficiency of training neural networks, in: *Advances in Neural Information Processing Systems 27, 2014*, pp. 855–863.
- [12] B. DasGupta, H. T. Siegelmann, E. Sontag, On a learnability question associated to neural networks with continuous activations, in: *Proceedings of the Seventh Annual Conference on Computational Learning Theory, COLT '94, 1994*, pp. 47–56.
- [13] O. Shamir, Distribution-specific hardness of learning neural networks, *CoRR* (2016).
- [14] L. Song, S. Vempala, J. Wilmes, B. Xie, On the complexity of learning neural networks, *CoRR* (2017).
- [15] R. Arora, A. Basu, P. Mianjy, A. Mukherjee, Understanding deep neural networks with rectified linear units, 2016.
- [16] P. Manurangsi, D. Reichman, The computational complexity of training relu(s), *CoRR* abs/1810.04207 (2018). [arXiv:1810.04207](https://arxiv.org/abs/1810.04207).
URL <http://arxiv.org/abs/1810.04207>
- [17] S. S. Dey, G. Wang, Y. Xie, An approximation algorithm for training one-node relu neural network (2018). [arXiv:1810.03592](https://arxiv.org/abs/1810.03592).
- [18] H. Edelsbrunner, J. O'Rourke, R. Seidel, Constructing arrangements of lines and hyperplanes with applications, *SIAM J. Comput.* 15 (2) (1986) 341–363.