

**GLOBAL CONVERGENCE AND GEOMETRIC
CHARACTERIZATION OF SLOW TO FAST WEIGHT
EVOLUTION IN NEURAL NETWORK TRAINING FOR
CLASSIFYING LINEARLY NON-SEPARABLE DATA**

ZIANG LONG*

Department of Mathematics
University of California, Irvine, CA 92697, USA

PENGHANG YIN

Department of Mathematics and Statistics
State University of New York at Albany, Albany, NY 12222, USA

JACK XIN

Department of Mathematics
University of California, Irvine, Irvine, CA 92697, USA

ABSTRACT. In this paper, we study the dynamics of gradient descent in learning neural networks for classification problems. Unlike in existing works, we consider the linearly non-separable case where the training data of different classes lie in orthogonal subspaces. We show that when the network has sufficient (but not exceedingly large) number of neurons, (1) the corresponding minimization problem has a desirable landscape where all critical points are global minima with perfect classification; (2) gradient descent is guaranteed to converge to the global minima. Moreover, we discovered a geometric condition on the network weights so that when it is satisfied, the weight evolution transitions from a slow phase of weight direction spreading to a fast phase of weight convergence. The geometric condition says that the convex hull of the weights projected on the unit sphere contains the origin.

1. Introduction. Deep neural networks (DNN) have achieved remarkable performances in image and speech classification tasks among other AI applications in recent years; for examples, see [8, 9, 15, 17]. Although there have been numerous theoretical contributions to understand their success, the learning process in the actual network training remains largely empirical. One interesting phenomenon is that over-parametrized DNN's trained by stochastic gradient descent generalize [13, 14] instead of overfitting the training data contrary to conventional statistical learning. Though several convergence results are proved in the over-parameterized regime for deep networks [6, 11, 1], the network weights move only in a small neighborhood of the random initialization and so their dynamics are very localized. Partly, this may be attributed to the exceedingly large number of neurons in convergence theory, far surpassing what is used in practice where the weights evolve significantly

2020 *Mathematics Subject Classification.* 90C26, 68W40.

Key words and phrases. Neural networks, learning dynamics, geometric condition, slow-to-fast convergence, classification, gradient descent.

* Corresponding author: Ziang Long.

from random start through hundreds of epochs in training to reach best prediction accuracy.

Our work here addresses how the weights evolve towards a global minimum of loss function as the number of neurons increases from the feature dimension (the least necessary) to the over-parametrized regime. To facilitate analysis, our model network structure is motivated by [4] on classifying linearly separable data. We instead study a linearly non-separable multi-category classification problem with an emphasis on the dynamics of weights in terms of the two time scales of evolution and a geometric characterization of the transition time. Our training data of the two classes will lie in orthogonal sub-spaces, which extends the data configuration in [3] where the subspace of each class is one dimensional for an XOR detection problem. Orthogonality of input data from the two classes implies that the training process in each class can be analyzed independently of the other. In the one-dimensional case [3], each weight update does not increase the loss on any sample point. In the multi-dimensional case here, we find that during gradient descent weight update, it is not possible that the loss is non-increasing in the point-wise sense (on each input data). Instead, the population loss is decreasing (i.e. in the sense of expectation). The population loss here is based on the hinge loss function and the network activation function is ReLU. Under a mild non-degenerate data condition, we prove that all critical points of our non-convex and non-smooth population loss function are global minima. Similar landscapes (a local minimum is a global minimum) are known for deep networks with activation functions that are either strictly convex [12], or real analytic and strictly increasing [14].

1.1. Prior works and our contributions. In DNN training, one observes that the network learning consists of alternating phases: plateaus where the validation error remains fairly constant and periods of rapid improvement where a lot of progress is made over a few epochs. Prior to our work, [5] studied slow and fast weight dynamics in a solvable model while minimizing a binary cross entropy or hinge loss function on linearly separable data. In the regression context, [2] came across such two time-scale phenomenon in training a two-linear-layer convolutional network with prescribed ground truth and unit Gaussian input data. This particular data assumption makes it possible to readily derive the closed-form expressions of the population loss and gradient, and then analyze the energy landscape and convergence of the gradient descent algorithm.

In this work, we study network weight dynamics in training a one-hidden-layer ReLU network via hinge loss minimization on multi-category classification of linearly non-separable data lying in n orthogonal sub-spaces. Our main contributions are:

- We discovered a geometric condition (GC) to characterize the transition time T from the first (slow) phase of weight evolution to the second fast weight convergence. The condition says that the convex hull of the weights on the unit sphere contains the origin, see Fig. 1 for an illustration. Equivalent geometric conditions are also derived (Lemma 1). In the first (slow) phase, the weight directions spread out over the unit sphere to satisfy GC.
- We obtain upper bound on T in terms of data distribution function provided that the network weights are uniformly bounded during training which we observed numerically.

- We give probabilistic bounds on the validity of geometric condition for random initialization, which suggests that the larger the number of neurons, the more likely GC holds and the earlier the fast phase of evolution begins.
- We prove the global convergence of gradient descent training algorithm under the uniformly bounded weight assumption. In case of positive network bias, we prove a global Lipschitz gradient property of the loss function and sub-sequential convergence of weights to a global minimum. In case of zero network bias, we prove that the loss function has Lipschitz gradient away from the origin and is piece-wise C^1 .
- We prove that all critical points of the population loss function are global minima under a non-degenerate data condition.
- We provide numerical examples to substantiate our theory, extend the data assumption, and illustrate the weight dynamics as the network size increases towards the over-parametrized regime. We visualize the feature and weight vectors in DNNs on MNIST data in connection with our model findings.

Organization. In section 2, we introduce the settings of the classification problem, including the assumptions on the data and network architecture. In section 3, we state the main results regarding the convergence guarantee of the gradient descent algorithm for training the neural net in the cases of with and without a bias term in the linear layer. In section 4, we present preliminaries about the landscape of the training loss function. The convergence analysis of main results will be sketched in section 5. In section 6, we substantiate our theoretical findings with numerical simulations. All the technical proofs are detailed in the appendix.

Notations. We denote by \mathcal{S}^{d-1} the unit sphere in \mathbb{R}^d , and $|\mathcal{S}^{d-1}|$ the area of the unit sphere in the corresponding dimension. For any finite dimensional linear space $V \subseteq \mathbb{R}^d$, we define V^k to be the collection of matrices of form $[\mathbf{x}_1, \dots, \mathbf{x}_k] \in \mathbb{R}^{d \times k}$, where $\mathbf{x}_j \in V$ is the j -th column vector. For any set \mathcal{X} , $\mathbf{1}_{\mathcal{X}}(x) = 1$ if $x \in \mathcal{X}$ else 0, is the indicator function of \mathcal{X} . For any vector $\mathbf{x} \in \mathbb{R}^d$, we denote $|\mathbf{x}|$ be the ℓ_2 norm of \mathbf{x} . For a matrix $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_k] \in \mathbb{R}^{d \times k}$, $|\mathbf{W}| := \sum_{j=1}^k |\mathbf{w}_j|$ is the column-wise ℓ_2 -norm sum.

2. Problem setup. In this section, we consider the multi-category classification problem in the d -dimensional space $\mathcal{X} = \mathbb{R}^d$. Let $\mathcal{Y} = [n] := \{1, 2, \dots, n\}$ be the set of labels, and $\{\mathcal{D}_i\}_{i=1}^n$ be n probabilistic distributions over $\mathcal{X} \times \mathcal{Y}$. Throughout the theoretical analysis of this paper, we make the following assumptions on the data:

1. **(Separability)** There are n orthogonal subspaces $V_i \subseteq \mathcal{X}$ for $i \in [n]$ with $\bigoplus_{i=1}^n V_i = \mathcal{X}$, such that

$$\mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}_i} [\mathbf{x} \in V_i \text{ and } y = i] = 1.$$

2. **(Boundedness of data)** For $i \in [n]$, There exist positive constants m_i and M_i , such that

$$\mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}_i} [m_i \leq |\mathbf{x}| \leq M_i] = 1.$$

3. **(Boundedness of p.d.f.)** For $i \in [n]$, let p_i be the probability density function of distribution \mathcal{D}_i restricted on V_i . For any $\mathbf{x} \in V_i$ with $m_i < |\mathbf{x}| < M_i$, it holds that

$$0 < p_{\min} \leq p_i(\mathbf{x}) \leq p_{\max} < \infty.$$

Later on, we denote \mathcal{D} to be the evenly mixed distribution of \mathcal{D}_i 's. For notation simplicity, we let $m = \min_{i \in [n]} m_i$, $M = \max_{i \in [n]} M_i$ and $d_i = \dim V_i$.

We consider a two-layer neural network with k hidden neurons. Denote by $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_k] \in \mathbb{R}^{d \times k}$ the weight matrix in the hidden layer. For any input data $\mathbf{x} \in \mathcal{X} = \mathbb{R}^d$, we have

$$h_j = \langle \mathbf{w}_j, \mathbf{x} \rangle - b_j \text{ and } f_i = \sum_{j=1}^n v_{i,j} \sigma(h_j)$$

and the neural net outputs

$$(1) \quad f(\mathbf{W}; \mathbf{x}) = \mathbf{V} \sigma(\mathbf{W}^\top \mathbf{x}) = [f_1, \dots, f_n]^\top,$$

where $\sigma := \max(\cdot, 0)$ is the ReLU function acting element-wise, and the bias $b_j \geq 0$ and $\mathbf{V} = (v_{i,j})$ are constants. Throughout this paper, we assume the following:

Assumption 1. $\mathbf{V} = (v_{i,j}) \in \mathbb{R}^{k \times n}$ satisfies

1. For any $i \in [n]$, there exists some $j \in [n]$ such that $v_{i,j} > 0$.
2. If $v_{i,j} > 0$ then $v_{r,j} < 0$ for all $r \neq i$ and $r \in [n]$.
3. There exists some constant $v > 0$ such that $|v_{i,j}| = v$.

One can show that as long as $k \geq n$, such \mathbf{V} can be constructed easily.

The prediction is given by the maximum coordinate index of the network output

$$\hat{y}(\mathbf{W}; \mathbf{x}) = \operatorname{argmax}_{i \in [n]} f_i,$$

ideally $\hat{y}(x) = i$ if $x \in V_i$. The classification accuracy in percentage is the frequency that this occurs (when network output label \hat{y} matches the true label) on a validation data set. Given the data sample $\{\mathbf{x}, y\}$, the associated hinge loss function reads

$$(2) \quad l(\mathbf{W}; \{\mathbf{x}, y\}) := \sum_{i \neq y} \max\{0, 1 - f_y + f_i\}.$$

For network training, we consider the gradient descent algorithm with step size $\eta > 0$

$$(3) \quad \mathbf{W}^t = \mathbf{W}^{t-1} - \eta \nabla l(\mathbf{W}^{t-1})$$

to solve following population loss minimization problem

$$(4) \quad \min_{\mathbf{W} \in \mathbb{R}^{d \times k}} l(\mathbf{W}) = \mathbb{E}_{\{\mathbf{x}, y\} \sim \mathcal{D}} [l(\mathbf{W}; \{\mathbf{x}, y\})],$$

where the sample loss function $l(\mathbf{W}; \{\mathbf{x}, y\})$ is given by (2). Let l_i be the population loss function of data type i . More precisely,

$$l_i(\mathbf{W}) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_i} [l(\mathbf{W}; \{\mathbf{x}, y\})] = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_i} \left[\sum_{r \neq i} \sigma(1 - f_i + f_r) \right].$$

Thus, we can rewrite the loss function as

$$l(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n l_i(\mathbf{W}).$$

Note that the population loss function

$$l_i(\mathbf{W}) = \sum_{r \neq i} \int_{\{f_i < f_r + 1\}} (1 - f_i + f_r) p_i(\mathbf{x}) d\mathbf{x}$$

has no closed-form solution even if p is a constant function on its support. We cannot use closed-form formula to analyze the learning process, which makes our work different from many other works.

3. Main results. Although (4) is a non-convex optimization problem, we show that under mild conditions, the gradient descent algorithm (3) converges to a global minimum with zero classification error. Specifically, we consider two different networks with a positive bias $b_j > 0$ (Theorem 3.1) and without a bias (Theorem 3.2), respectively. For both cases, we have the fact that any critical point of problem (4) is a global minimum (Proposition 1). The key difference between these two cases is that the population loss function has Lipschitz continuous gradient (Lemma 4.2) when $b_j > 0$, whereas this desirable property does not hold otherwise. For the latter case $b_j = 0$, we present a totally different analysis based on a geometric condition proved to emerge during the training process (Proposition 3). Under this geometric condition, the objective value converges zero (Proposition 4).

Theorem 3.1. *Assume $0 < \sum_{j=1}^k b_j < 1$ and assumption 1 holds in (1), and $\{\mathbf{W}^t\}$ generated by the algorithm (3) are bounded uniformly in t . If there exists $(\mathbf{x}, i) \sim \mathcal{D}_i$ and some indices $j \in [k]$, such that $v_{i,j} > 0$ and $|\langle \mathbf{w}_j^0, \mathbf{x} \rangle| > b_j$, then there exists some $\eta_0(v, k, b, p_{max}, n, M) > 0$ such that for the learning rate $\eta < \eta_0$, $\lim_{t \rightarrow \infty} l_i(\mathbf{W}^t) = 0$ and*

$$\lim_{t \rightarrow \infty} \mathbb{P}_{\{\mathbf{x}, i\} \sim \mathcal{D}_i} [\hat{y}(\mathbf{W}^t; \mathbf{x}) \neq i] = 0, \quad i \in [n], \forall n \geq 2.$$

Proof of Theorem 3.1. By Lemma 4.1, we only need to prove the convergence of the simplified network (6). From Lemma 4.2, we know $l_i(\mathbf{W})$ has Lipschitz gradient. We can assume for any $\mathbf{W}_1, \mathbf{W}_2 \in V_i^k$, we have

$$|\nabla l_i(\mathbf{W}_1) - \nabla l_i(\mathbf{W}_2)| \leq L |\mathbf{W}_1 - \mathbf{W}_2|.$$

As long as we take $\eta < \frac{L}{2}$ in algorithm (3), we know

$$(5) \quad l_i(\mathbf{W}^{t+1}) \leq l_i(\mathbf{W}^t) - \left(\eta - \frac{\eta^2 L}{2} \right) |\nabla l_i(\mathbf{W}^t)|^2 \leq l_i(\mathbf{W}^t).$$

Hence, $l_i(\mathbf{W}^t)$ is monotonically decreasing. Therefore, for any convergent subsequence $\{\mathbf{W}^{t_k}\}$ with the limit \mathbf{W}_0 , there exists $l_0 \geq 0$ such that

$$\lim_{t \rightarrow \infty} l(\mathbf{W}^t) = \lim_{k \rightarrow \infty} l(\mathbf{W}^{t_k}) = l_0.$$

Now, we can take subsequence and limit on both side of equation (5), we get

$$l_0 \leq l_0 - \left(\eta - \frac{\eta^2 L}{2} \right) |\nabla l_i(\mathbf{W}_0)|^2.$$

Now, we see that $\nabla l_i(\mathbf{W}_0) = \mathbf{0}$. □

Theorem 3.2. *Assume $b_j = 0$ and assumption 1 holds in (1), and $\{\mathbf{W}^t\}$ generated by algorithm (3) is bounded by R uniformly in t . If there exist some $(\mathbf{x}, i) \sim \mathcal{D}_i$ and some indices $j \in [k]$, such that $v_{i,j} > 0$ and $\langle \mathbf{w}_j^0, \mathbf{x} \rangle \neq 0$, then $\lim_{t \rightarrow \infty} l_i(\mathbf{W}^t) = 0$ and*

$$\lim_{t \rightarrow \infty} \mathbb{P}_{\{\mathbf{x}, i\} \sim \mathcal{D}_i} [\hat{y}(\mathbf{W}^t; \mathbf{x}) \neq i] = 0, \quad i \in [n], \forall n \geq 2.$$

Proof of Theorem 3.2. By Proposition 3, we know the iterates $\{\mathbf{W}^t\}$ stay in the first phase is bounded by $|T_1|$. Combining Proposition 4 which shows the summation of squared loss values in phase two is bounded, the summation of all loss values in the learning process is bounded

$$\sum_{t \in T_2} l_i(\mathbf{W}^t)^2 < \infty.$$

The desired result follows. \square

Remark 1. The assumptions on the initialization $|\langle \mathbf{w}_j^0, \mathbf{x} \rangle| > b_j$ in both theorems are natural. This assumption guarantees that the neuron \mathbf{w}_j is activated by some input data. Without this assumption, the algorithm suffers zero gradient and fails to update.

Remark 2. Theorem 3.2 does not explicitly require the learning rate η to be small. However, a larger learning rate will implicitly result in a larger bound in Propositions 3 and 4 which we used to prove Theorem 3.2.

4. Preliminaries.

4.1. Decomposition.

Lemma 4.1. 1. For any $i \in [n]$, if $\mathbf{W}_i^* \in \mathbb{R}^{d \times k}$ solves the optimization problem

$$\min_{\mathbf{W} \in V_i^k} l_i(\mathbf{W}),$$

then $\mathbf{W}^* = \sum_{i=1}^n \mathbf{W}_i^*$ solves the original problem

$$\min_{\mathbf{W} \in \mathbb{R}^{d \times k}} l(\mathbf{W}).$$

2. If $\mathbf{W}' = \mathbf{W} - \eta \nabla_{\mathbf{W}} l_i(\mathbf{W})$, then for any $r \neq i$, we have

$$l(\mathbf{W}'; \{\mathbf{x}, r\}) = l(\mathbf{W}; \{\mathbf{x}, r\})$$

for almost all $(\mathbf{x}, r) \sim \mathcal{D}_r$.

Proof of Lemma 4.1. Note that $\mathcal{X} = \mathbb{R}^d = \bigoplus_{i=1}^n V_i$, we decompose $\mathbf{w}_j = \sum_{i=1}^n \mathbf{w}_{j,i}$ where $\mathbf{w}_{j,i} \in V_i$.

Since V_i 's are orthogonal spaces, we have

$$\langle \mathbf{w}_j, \mathbf{x} \rangle = \langle \mathbf{w}_{j,i}, \mathbf{x} \rangle$$

if $\mathbf{x} \in V_i$. Now, assume $\mathbf{x} \in V_i$, let

$$\mathbf{W}_i = [\mathbf{w}_{1,i}, \dots, \mathbf{w}_{k,i}],$$

we have $f(\mathbf{W}; \mathbf{x}) = f(\mathbf{W}_i, \mathbf{x})$ and hence we get our first claim.

On the other hand, we have

$$\nabla_{\mathbf{w}_j} l(\mathbf{W}; \{\mathbf{x}, y\}) = - \sum_{i \neq y} (v_{y,j} - v_{i,j}) \mathbb{1}_{\Omega_{y,i}}(\mathbf{x}) \mathbb{1}_{\Omega_{\mathbf{w}_j}}(\mathbf{x}) \mathbf{x},$$

where

$$\Omega_{y,i} = \{\mathbf{x} : f_y < f_i + 1\} \text{ and } \Omega_{\mathbf{w}_j} = \{\mathbf{x} : \langle \mathbf{w}_j, \mathbf{x} \rangle > b_j\}.$$

Now, we see that $\nabla_{\mathbf{w}_j} l(\mathbf{W}; \{\mathbf{x}, y\}) \in V_i$ for almost all $(\mathbf{x}, i) \sim \mathcal{D}_i$ so that $\mathbf{W}'_r = \mathbf{W}_r$ for all $r \neq i$ and hence our desired result follows. \square

The optimization problem can be decomposed to n independent problems of the same form i.e. the optimization of l_i 's. Therefore, it suffices to consider only one subproblem. Let $\mathbf{W}_i = [\mathbf{w}_1, \dots, \mathbf{w}_k] \in V_i^k$, where $\mathbf{w}_j \in V_i = \mathbb{R}^{d_i}$, the network output for the input data $\mathbf{x} \in V_i = \mathbb{R}^{d_i}$ is given by

$$(6) \quad \tilde{f}^i(\mathbf{W}_i; \mathbf{x}) = [f_1, \dots, f_n]^\top = \mathbf{V} \sigma(\mathbf{W}_i^\top \mathbf{x}).$$

Networks (1) and (6) are different, since the parameters in (1) are $\mathbf{W} \in \mathbb{R}^{d \times k}$, whereas in (6), we have $\mathbf{W}_i \in V_i^k \cong \mathbb{R}^{d_i \times k}$. The corresponding input data are also in different intrinsic dimensions. From now on, we just focus on the loss function associated with data of Class i :

$$l_i(\mathbf{W}) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_i} \left[\sum_{r \neq i} \max\{0, 1 - f_i + f_r\} \right].$$

4.2. Landscape. The following Proposition 1 shows that while the loss function is non-convex, any critical point is in fact a global minimum, except for some degenerate cases.

Proposition 1. *Consider the neural network in (6). Assume $d > 1$, if \mathbf{W} is a critical point of $l_i(\mathbf{W})$ and there exists some $\mathbf{x} \in V_i$ such that $f(\mathbf{W}; \mathbf{x}) \neq \mathbf{0}$ then we have $l_i(\mathbf{W}) = 0$.*

Proof of Proposition 1. For any $j \in [k]$, we have

$$\nabla_{\mathbf{w}_j} l_i(\mathbf{W}) = - \sum_{r \neq i} (v_{i,j} - v_{r,j}) \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_i} \left[\mathbf{1}_{\Omega_{i,r}}(\mathbf{x}) \mathbf{1}_{\Omega_{\mathbf{w}_j}}(\mathbf{x}) \mathbf{x} \right] = \mathbf{0}.$$

Recall the definition of $\Omega_{\mathbf{w}_j}$, we know that each summand is a zero vector. By assumption 1, we know that either $v_{i,j} = v_{r,j}$ or

$$\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_i} \left[\mathbf{1}_{\Omega_{i,r}}(\mathbf{x}) \mathbf{1}_{\Omega_{\mathbf{w}_j}}(\mathbf{x}) \right] = \mathbf{0},$$

where the latter implies $\Omega_{i,r} \cap \Omega_{\mathbf{w}_j} = \emptyset$. Observe that

$$f_i - f_r = \sum_{j=1}^k (v_{i,j} - v_{r,j}) \sigma(h_j),$$

we see that if there exists some $(\mathbf{x}, y) \sim \mathcal{D}_i$ such that $f_i - f_r < 1$ then there must exist some $\mathbf{x} \in \Omega_{i,r}$ but this implies $\mathbf{x} \notin \Omega_{\mathbf{w}_j}$ for all $j \in [k]$ such that $v_{i,j} \neq v_{r,j}$ which gives $f_i = f_r = 0$. This contradicts with our assumption, so we get $f_i - f_r \geq 1$ for all $\mathbf{x} \sim \mathcal{D}_i$ and this implies $l_i(\mathbf{W}) = 0$. \square

The above result holds only when the global minimum of training loss function exists. The following proposition shows that the loss function has plenty of global minima.

Proposition 2. *Consider the network in (6). If the convex hull spanned by vertices $\{\mathbf{w}_j : v_{i,j} > 0\}$ contains a ball centered at the origin with radius $\max_{j \in [k]} \frac{1+b_j}{m_i}$, and $\{\mathbf{w}_j : v_{i,j} < 0\}$ lies in a ball with radius $\min_{j \in [k]} \frac{b_j}{M_i}$, then $l_i(\mathbf{W}) = 0$.*

The above proposition shows that if number of neurons is greater than the dimension of input data, then global minimum exists. Next, we study the smoothness of the loss function. The following proposition shows that as long as weights are bounded away from 0, then the loss function has Lipschitz gradient.

Lemma 4.2. *Consider the network in (1) with positive bias $0 < \sum_{j=1}^k b_j < 1$. The loss function $l(\mathbf{W})$ is Lipschitz differentiable, i.e, there exists some constant $L > 0$ depending on $k, \mathbf{b}, p_{max}, M, \mathbf{V}$, such that*

$$|\nabla l(\mathbf{W}_1) - \nabla l(\mathbf{W}_2)| \leq L |\mathbf{W}_1 - \mathbf{W}_2|.$$

Proof of Lemma 4.2. Note that $l(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n l_i(\mathbf{W})$, it suffice to show that each l_i has Lipschitz gradient. Note that

$$l_i(\mathbf{W}) = \sum_{r \neq i} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_i} [\sigma(1 - f_i + f_r)],$$

it suffice to show each summand has Lipschitz gradient. Now, we write the gradient of the summands as

$$\begin{aligned} & \nabla_{\mathbf{w}_j} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_i} [\sigma(1 - f_i + f_r)] \\ &= \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_i} [\nabla_{\mathbf{w}_j} \sigma(1 - f_i + f_r)] \\ &= \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_i} [\mathbb{1}_{\Omega_{i,r}}(\mathbf{x}) \mathbb{1}_{\Omega_{\mathbf{w}_j}}(\mathbf{x}) (v_{r,j} - v_{i,j}) \mathbf{x}]. \end{aligned}$$

On one hand, if $v_{r,j} = v_{i,j}$, then the formula above is obviously zero and Lipschitz gradient follows. On the other hand, we can without loss of generality assume $|v_{i,j} - v_{r,j}| = 1$. For notation simplicity, we fix i and r and denote

$$\varphi(\mathbf{W}) := \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_i} [\mathbb{1}_{\Omega_{i,r}}(\mathbf{x}) \mathbb{1}_{\Omega_{\mathbf{w}_j}}(\mathbf{x}) \mathbf{x}] \text{ and } \phi(\mathbf{W}; \mathbf{x}) = f_i(\mathbf{W}; \mathbf{x}) - f_r(\mathbf{W}; \mathbf{x}).$$

Now, our desired result becomes $\varphi(\mathbf{W})$ is Lipschitz in \mathbf{W} . Denote $\Omega_1 = \Omega_{i,r}(\mathbf{W}_1)$ and $\Omega_2 = \Omega_{i,r}(\mathbf{W}_2)$ where $\mathbf{W}_i = [\mathbf{w}_1^i, \dots, \mathbf{w}_k^i]$, we only need to show there exists some constant L such that

$$|\varphi(\mathbf{W}_1) - \varphi(\mathbf{W}_2)| \leq L |\mathbf{W}_1 - \mathbf{W}_2|.$$

Note that

$$\begin{aligned} & |\varphi(\mathbf{W}_1) - \varphi(\mathbf{W}_2)| \\ &= \left| \mathbb{E} [\mathbb{1}_{\Omega_1}(\mathbf{x}) \mathbb{1}_{\Omega_{\mathbf{w}_j^1}}(\mathbf{x}) \mathbf{x}] - \mathbb{E} [\mathbb{1}_{\Omega_2}(\mathbf{x}) \mathbb{1}_{\Omega_{\mathbf{w}_j^2}}(\mathbf{x}) \mathbf{x}] \right| \\ &\leq \underbrace{\mathbb{E} [\mathbb{1}_{\Omega_1 \Delta \Omega_2}(\mathbf{x}) |\mathbf{x}|]}_{\textcircled{1}} + \underbrace{\mathbb{E} [\mathbb{1}_{\Omega_{\mathbf{w}_j^1} \Delta \Omega_{\mathbf{w}_j^2}}(\mathbf{x}) |\mathbf{x}|]}_{\textcircled{2}}, \end{aligned}$$

we can deal with $\textcircled{1}$ and $\textcircled{2}$ respectively.

W.l.o.g, we assume $\epsilon = |\mathbf{W}_1 - \mathbf{W}_2| \leq \frac{1}{2}$. On one hand, if $\mathbf{x} \in \Omega_1 \Delta \Omega_2$, then we claim

$$1 - \epsilon |\mathbf{x}| \leq \phi(\mathbf{W}_1; \mathbf{x}) \leq 1 + \epsilon |\mathbf{x}|$$

because

$$\begin{aligned} |\phi(\mathbf{W}_1; \mathbf{x}) - \phi(\mathbf{W}_2; \mathbf{x})| &= \left| \sum_{j=1}^k (v_{i,j} - v_{r,j}) [\sigma(h_j^1) - \sigma(h_j^2)] \right| \\ &\leq \left| \sum_{j=1}^k \langle \mathbf{w}_j^1 - \mathbf{w}_j^2, \mathbf{x} \rangle \right| \leq |\mathbf{W}_1^\top \mathbf{x} - \mathbf{W}_2^\top \mathbf{x}| \leq \epsilon |\mathbf{x}|. \end{aligned}$$

Furthermore, we claim for such \mathbf{x} 's the gradient of ϕ is bounded away from zero. More precisely, with $\mathbf{x} = r\omega$ where $r = |\mathbf{x}|$, we have

$$\begin{aligned} \frac{d}{dr} \phi(\mathbf{W}; \mathbf{x}) &= \left[\sum_{j=1}^k (v_{i,j} - v_{r,j}) \mathbb{1}_{\Omega_{\mathbf{w}_j}}(\mathbf{x}) \mathbf{w}_j \right] \omega \\ &\geq \left(\phi(\mathbf{W}; \mathbf{x}) - \sum_{j=1}^k b_j \right) / r \geq \frac{1}{M} \left(\frac{1}{2} - \sum_{j=1}^k b_j \right) =: C > 0. \end{aligned}$$

Now, we have

$$\begin{aligned} \textcircled{1} &\leq \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_i} [\mathbb{1}_{1-\epsilon|\mathbf{x}| \leq \phi(\mathbf{W}_1, \mathbf{x}) \leq 1+\epsilon|\mathbf{x}|}] (\mathbf{x}) |\mathbf{x}| \\ &\leq \int_{1-\epsilon|\mathbf{x}| \leq \phi(\mathbf{W}_1, \mathbf{x}) \leq 1+\epsilon|\mathbf{x}|} |\mathbf{x}| p_{max} d\mathbf{x} \leq 2 \left(|\mathcal{S}^{d_i-1}| \frac{M^{d_i} p_{max}}{C} \right) \epsilon. \end{aligned}$$

As for $\textcircled{2}$, w.l.o.g. we can assume $|\mathbf{w}_j^1| \geq |\mathbf{w}_j^2|$. Note that when $|\mathbf{w}_j^1| \leq \frac{b_j}{M}$ then $h_j \leq 0$ so that $\nabla_{\mathbf{w}_j} l_i(\mathbf{W}) = \mathbf{0}$ and this $\textcircled{2} = 0$. Hence, we only need to take care of the case when $|\mathbf{w}_j^1| \geq \frac{b_j}{M}$. Note that $|\mathbf{w}_j^1 - \mathbf{w}_j^2| \leq \epsilon$ we know

$$\sin \theta \leq \frac{\epsilon M}{b_j},$$

where θ denotes the acute angle between \mathbf{w}_j^1 and \mathbf{w}_j^2 . We have the following estimate

$$\textcircled{2} \leq p_{max} \frac{\epsilon M}{b_j} |\mathcal{S}^{d_i-1}|.$$

Combine with $\textcircled{1}$, we get our desired result. \square

Note that Lipschitz differentiability in Lemma 4.2 does not hold for the case $b_j = 0$, as the gradient might be volatile near the origin.

5. Convergence analysis for non-bias case. With the Lipschitz differentiability shown in Lemma 4.2 in the case $b_j > 0$, it is not hard to prove the convergence result in Theorem 3.1. In this section, we focus on the non-bias case ($b_j = 0$) where the Lipschitz differentiability fails and sketch the convergence analysis.

Lemma 5.1. *Consider the network (6), $|\mathbf{w}_j^t|$ is non-decreasing in t if $v_{i,j} > 0$. For any $r > 0$, choosing learning rate $\eta < \min \left\{ \frac{r}{C_p M_i^2}, \frac{r}{2\eta n M_i} \right\}$, then $|\mathbf{w}_j^t|$ is non-increasing if $v_{i,j} < 0$ and $|\mathbf{w}_j^t| > r$, where C_p is a constant satisfying*

$$C_p = \max_{\mathbf{v} \in V_i, a \in \mathbb{R}} \int_{\{(\mathbf{v}, \mathbf{x})=a\}} p_i(\mathbf{x}) d\mathbf{x} = O(p_{max} M^{d_i}).$$

Proof of Lemma 5.1. We first define

$$C_p = \max_{\mathbf{v} \in V_1, a \in \mathbb{R}} \int_{\langle \mathbf{v}, \mathbf{x} \rangle = a} p_1(\mathbf{x}) dS \leq M^{d-1} p_{\max}.$$

Recall the definition of $\Omega_{\mathbf{w}_j}$, for all $\mathbf{x} \in \Omega_{\mathbf{w}_j}$, we have $\langle \tilde{\mathbf{w}}_j, \mathbf{x} \rangle > 0$. For $v_{i,j} > 0$,

$$\langle \tilde{\mathbf{w}}_j, -\nabla_{\mathbf{w}_j} l_i(\mathbf{W}) \rangle = 2v \sum_{r \neq i} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_i} \left[\mathbb{1}_{\Omega_{i,r}}(\mathbf{x}) \mathbb{1}_{\Omega_{\mathbf{w}_j}}(\mathbf{x}) \langle \tilde{\mathbf{w}}_j, \mathbf{x} \rangle \right] \geq 0.$$

Note we have from (3) that

$$\mathbf{w}_j^{t+1} = \mathbf{w}_j^t - \eta \nabla_{\mathbf{w}_j} l_i(\mathbf{W}^t).$$

So,

$$|\mathbf{w}_j^{t+1}| = \langle \mathbf{w}_j^{t+1}, \tilde{\mathbf{w}}_j^{t+1} \rangle \geq \langle \mathbf{w}_j^{t+1}, \tilde{\mathbf{w}}_j^t \rangle \geq \langle \mathbf{w}_j^t, \tilde{\mathbf{w}}_j^t \rangle = |\mathbf{w}_j^t|.$$

Let $\Omega_{i,r}^j = \Omega_{i,r} \cap \Omega_{\mathbf{w}_j}$. For $v_{i,j} < 0$, we know

$$|\nabla_{\mathbf{w}_j} l_i(\mathbf{W})| = 2v \left| \sum_{r \neq i} \mathbb{E} \left[\mathbb{1}_{\Omega_{i,r}^j}(\mathbf{x}) \mathbf{x} \right] \right| \leq 2vM \sum_{r \neq i} \mathbb{P} \left[\Omega_{i,r}^j \right]^2,$$

where we omit the distribution $\mathbf{x} \sim \mathcal{D}_i$.

On the other hand, by definition of C_p , we know

$$|\langle \nabla_{\mathbf{w}_j} l_i(\mathbf{W}), \tilde{\mathbf{w}}_j \rangle| = \left| 2v \sum_{r \neq i} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_i} \left[\mathbb{1}_{\Omega_{i,r}^j}(\mathbf{x}) \langle \tilde{\mathbf{w}}_j, \mathbf{x} \rangle \right] \right| \geq \frac{v}{C_p} \sum_{r \neq i} \mathbb{P} \left[\Omega_{i,r}^j \right]^2.$$

When $0 < \eta < \frac{r}{2vnM}$, we have

$$\langle \mathbf{w}_j - \eta \nabla_{\mathbf{w}_j} l_i(\mathbf{W}), \tilde{\mathbf{w}}_j \rangle = \langle \mathbf{w}_j, \tilde{\mathbf{w}}_j \rangle - \eta \langle \nabla_{\mathbf{w}_j} l_i(\mathbf{W}), \tilde{\mathbf{w}}_j \rangle > r - 2v\eta M \sum_{r \neq i} \mathbb{P} \left[\Omega_{i,r}^j \right] > 0.$$

Now, we decompose $\nabla_{\mathbf{w}_j} l_i(\mathbf{W}^t)$ into two parts,

$$\nabla_{\mathbf{w}_j} l_i(\mathbf{W}^t) = \underbrace{\langle \tilde{\mathbf{w}}_j^t, \nabla_{\mathbf{w}_j} l_i(\mathbf{W}^t) \rangle \tilde{\mathbf{w}}_j^t}_{\mathbf{n}} + \underbrace{\left(\nabla_{\mathbf{w}_j} l_i(\mathbf{W}^t) - \langle \tilde{\mathbf{w}}_j^t, \nabla_{\mathbf{w}_j} l_i(\mathbf{W}^t) \rangle \tilde{\mathbf{w}}_j^t \right)}_{\boldsymbol{\nu}}.$$

So that when $\eta < \frac{r}{M^2 C_p}$, we have

$$\begin{aligned} & |\mathbf{w}_j^{t+1}|^2 \\ &= |\mathbf{w}_j^t - \eta \nabla_{\mathbf{w}_j} l_i(\mathbf{W}^t)|^2 \\ &= |\mathbf{w}_j^t - \eta \mathbf{n}|^2 + |\eta \boldsymbol{\nu}|^2 \\ &= |\mathbf{w}_j^t|^2 - \eta \left(2 |\mathbf{w}_j^t| |\mathbf{n}| - \eta (|\mathbf{n}|^2 + |\boldsymbol{\nu}|^2) \right) \\ &\leq |\mathbf{w}_j^t|^2 - 2v\eta \left(\frac{|\mathbf{w}_j^t|}{C_p} - \eta M^2 \right) \sum_{r \neq i} \mathbb{P} \left[\Omega_{i,r}^j \right]^2 \\ &\leq |\mathbf{w}_j^t|^2, \end{aligned}$$

as long as $|\mathbf{w}_j| \geq r$. Now, we get the desired result. \square

The above theorem provides the dynamic information of the weights. When we input data with distribution \mathcal{D}_i , the weights $\{\mathbf{w}_j : v_{i,j} > 0\}$ become more useful for classification as the norm of those \mathbf{w}_j 's grows larger on every iteration. On the other hand, $\{\mathbf{w}_j : v_{i,j} < 0\}$ serve only as noise when $v_{i,j} < 0$. On one hand, when $\mathbf{w} \in \{\mathbf{w}_j : v_{i,j} < 0\}$ has large magnitude, the learning process guarantees the decreasing of $|\mathbf{w}|$. On the other hand, when $\mathbf{w} \in \{\mathbf{w}_j : v_{i,j} < 0\}$ has a small magnitude, it shall be trapped into a small region near the origin and contribute little to classification.

The learning process consists of two phases. In the beginning, since the weights are randomly distributed, there may well be some data that does not activate any neurons. The weights then automatically spread out. We call this process the first (slow learning) phase, during which the learning process is rather slow. The next Lemma lists four equivalent statements of a geometric condition. When the geometrical condition holds, we say that the learning process enters the second (fast learning) phase.

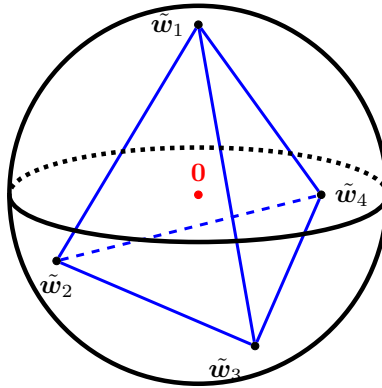


FIGURE 1. Geometric Condition in Lemma 5.2 ($d = 3$)

Lemma 5.2. Let $\tilde{\mathbf{w}}_j = \frac{\mathbf{w}_j}{|\mathbf{w}_j|} \in \mathcal{S}^{d-1}$ be k points on the unit sphere, where $1 \leq j \leq k$. Let Λ_W be the convex hull of $\{\tilde{\mathbf{w}}_j\}$ and $\{H_i\}_{i=1}^l$ be the facets of convex hull. The following statements are equivalent:

1. For any unit vector $\mathbf{n} \in \mathcal{S}^{d-1}$, there exist some j_1 and j_2 such that $\langle \mathbf{n}, \tilde{\mathbf{w}}_{j_1} \rangle > 0$ and $\langle \mathbf{n}, \tilde{\mathbf{w}}_{j_2} \rangle < 0$.
2. There exists no closed hemisphere that contains all $\tilde{\mathbf{w}}_j$.
3. $\mathbf{0}$ lies in the interior of Λ_W .
4. For each $0 \leq i \leq l$, we have $\mathbf{0}$ and all $\tilde{\mathbf{w}}_j$ lie on the same side of H_i .

Proof of Lemma 5.2. (1 \Leftrightarrow 2) is trivial.

(1 \Rightarrow 3) Proof by contradiction. Assume $\mathbf{0} \notin \Lambda_W^\circ$. Since Λ_W° is an open convex set, and $\{\mathbf{0}\}$ is a convex set, we know from geometric form of Hahn-Banach Theorem, that there exist a closed hyper-plane that separates $\{\mathbf{0}\}$ and Λ_W° . Hence, there exist a unit vector $\mathbf{n} \in \mathcal{S}^{d-1}$, such that $\langle \mathbf{n}, \tilde{\mathbf{w}}_j \rangle \geq \langle \mathbf{n}, \mathbf{0} \rangle = 0$ for all j , but this contradicts with our assumptions in 1.

(3 \Rightarrow 4) Assume that $H_i = \{ \langle \mathbf{v}_i, \mathbf{w} \rangle = \alpha_i \}$, where $\alpha_i > 0$. Note that the convex hull Λ_W is a polytope with faces H_i . We know, Λ_W° all lies on one side of H_i . Since

$\mathbf{0} \in \Lambda_W^\circ$, we know for all $\mathbf{w} \in \Lambda_W^\circ$ we have $\langle \mathbf{v}_i, \mathbf{w} \rangle < \alpha_i$, and hence $\langle \mathbf{v}_i, \tilde{\mathbf{w}}_j \rangle \leq \alpha_i$ for all j .

(4 \Rightarrow 3) is trivial.

(3 \Rightarrow 1) $\mathbf{0} \in \Lambda_W^\circ$ implies there exist positive numbers λ_j , such that

$$\sum_{j=1}^k \lambda_j \tilde{\mathbf{w}}_j = \mathbf{0}.$$

Hence for any unit vector \mathbf{n} , we have

$$\sum_{j=1}^k \lambda_j \langle \mathbf{n}, \tilde{\mathbf{w}}_j \rangle = 0.$$

Since $\tilde{\mathbf{w}}_j$ are in general position, so $\langle \mathbf{n}, \tilde{\mathbf{w}}_j \rangle$ cannot all be 0. Hence, there must be both positive and negative terms. Now, we get the desired result. \square

Remark 3. If \mathbf{w}_j is initialized such that $\tilde{\mathbf{w}}_j$ is uniformly distributed on \mathcal{S}^{d-1} , then the probability that the geometric condition (GC) in Lemma 5.2 holds is

$$P_{\text{gc}} := \text{Prob}(\text{GC holds}) = 2^{1-k} \sum_{j=d}^{k-1} \binom{k-1}{j}.$$

In particular, at any fixed feature dimension d ,

$$\lim_{k \rightarrow \infty} P_{\text{gc}} = 1.$$

Proof of Remark 3. For any $J \in \{\pm 1\}^k$, we let $S_J(\{\tilde{\mathbf{w}}_j\}) = \{J_j \tilde{\mathbf{w}}_j : 1 \leq j \leq k\}$. From the choice of $\tilde{\mathbf{w}}_j$, we know all S_J have the same distribution, so that

$$\text{Prob}_{\{\tilde{\mathbf{w}}_j\}}(\text{GC holds}) = \text{Prob}_{S_J}(\text{GC holds}).$$

Hence, we can simplify the probability as follows

$$\text{Prob}(\text{GC holds}) = \frac{1}{2^n} \mathbb{E} \left[\sum_J \mathbf{1}_{\text{gc}}(S_J(\{\tilde{\mathbf{w}}_j\})) \right].$$

We claim that $\sum_J \mathbf{1}_{\text{gc}}(S_J(\{\tilde{\mathbf{w}}_j\}))$ is a constant independent of choice of $\tilde{\mathbf{w}}_j$ as long as they are in general position.

Let $\tilde{\mathbf{w}}_j \in \mathcal{S}^{d-1}$ where $1 \leq j \leq k$. Each $\tilde{\mathbf{w}}_j$ corresponds to a subspace H_j with codimension 1 in \mathbb{R}^d , that is

$$H_j = \{\mathbf{x} \in \mathbb{R}^d : \langle \tilde{\mathbf{w}}_j, \mathbf{x} \rangle = 0\}.$$

Note that any connected region of $(\cup_J H_J)^c$ corresponds to a choice of J such that the geometric condition fails. More precisely, for any given connected region D of $(\cup_J H_J)^c$, we know $\langle \tilde{\mathbf{w}}_j, \mathbf{x} \rangle$ is one sign for all $\mathbf{x} \in D$, so with

$$J_j = \text{sign}(\langle \tilde{\mathbf{w}}_j, \mathbf{x} \rangle),$$

we know \mathbf{x} correspond to $S_J(\{\tilde{\mathbf{w}}_j\})$ where the geometric condition fails. Hence, the number of connected regions of $(\cup_J H_J)^c$ is same as $2^n - \sum_J \mathbf{1}_{\text{gc}}(S_J(\{\tilde{\mathbf{w}}_j\}))$. By [7], we have

$$\sum_J \mathbf{1}_{\text{gc}}(S_J(\{\tilde{\mathbf{w}}_j\})) = 2^n - 2 \sum_{j=0}^{d-1} \binom{k-1}{j} = 2 \sum_{j=d}^k \binom{k-1}{j}.$$

The desired result follows. \square

Per Remark 3, the more neurons the network has, higher possibility the geometric condition in Lemma 5.2 holds upon initialization. As a consequence, the learning process skips the first phase and goes straight to the fast learning phase. This explains why gradient descent for learning an over-parameterized network converges rapidly to a nearby critical point from random initialization.

The following proposition gives an upper bound on the maximum number of iterations for the learning process to enter the second phase.

Proposition 3. *Let $b_j = 0$ in (6), and assume that $|\mathbf{W}^t| \leq R$ for all t . Let T_1 be the set of t such that $\{\mathbf{w}_j^t : v_{i,j} > 0\}$ does not satisfy the geometric condition in Lemma 5.2, then $|T_1| \leq \frac{C_p R}{v\eta p_R^2}$, where p_R is a positive constant. Also, the following estimate holds for p_R :*

$$p_R = \Omega \left(\frac{p_{\min}}{\sqrt{d_i} (M_i R)^{d_i}} \right)$$

where C_p is the constant in Lemma 5.1. More precisely,

$$|T_1| = O \left(\frac{C_p d_i R^{2d_i+1} M_i^{2d_i}}{v\eta p_{\min}^2} \right) = O \left(\frac{C_p d R^{2d+1} M^{2d}}{v\eta p_{\min}^2} \right).$$

Proof of Proposition 3. W.l.o.g, assume $i = 1$ and $v_{1,j} > 0$ if and only if $j \in [k_1]$. Let $t \in T_1$, by Lemma 5.2, we know, for any fixed $t \in T_1$, there exists a $\alpha \in [0, \frac{\pi}{2}]$ and a unit vector $\mathbf{v} \in \mathbb{R}^{d_1}$, such that $\langle \mathbf{v}, \tilde{\mathbf{w}}_j^t \rangle \geq \sin \alpha$ for all $j \in [k_1]$ and $\langle \mathbf{v}, \tilde{\mathbf{w}}_1^t \rangle = \sin \alpha$. W.l.o.g, we assume $\mathbf{v} = (1, 0, \dots, 0)$ and $\langle \mathbf{v}, \tilde{\mathbf{w}}_1^t \rangle = \sin \alpha$. For any non-zero $\mathbf{x} \in V_1$, we write $\tilde{\mathbf{x}} = \frac{\mathbf{x}}{|\mathbf{x}|} \in \mathcal{S}^{d_1-1}$.

Note that $|\mathbf{W}|$ is bounded by R , we know there exist $\beta \in (0, \frac{\pi}{2} - \alpha)$ such that $\sum_{j=1}^k |\mathbf{w}_j^t| \leq R = \frac{1}{2vM_1 \sin \beta}$. Now, w.l.o.g, we can assume $\tilde{\mathbf{w}}_1^t = (\cos \alpha, \sin \alpha, 0, \dots, 0)$. Take $\mathbf{n} = (-\cos \alpha, \sin \alpha, 0, \dots, 0)$, we know $\langle \mathbf{n}, \tilde{\mathbf{w}}_1^t \rangle = 0$. For all $\mathbf{x} \in V_1$ such that $\langle \mathbf{n}, \tilde{\mathbf{x}} \rangle > \cos \beta$, we have

$$\cos \beta < \langle \mathbf{n}, \tilde{\mathbf{x}} \rangle = -\tilde{x}_1 \cos \alpha + \tilde{x}_2 \sin \alpha \leq -\tilde{x}_1 \cos \alpha + \sqrt{1 - \tilde{x}_1^2} \sin \alpha.$$

So, we have

$$\tilde{x}_1 < -\cos(\alpha + \beta).$$

Now, for all $\mathbf{x} \in V_1$ such that $\tilde{x}_1 < -\cos(\alpha + \beta)$, we have

$$f_1(\mathbf{W}^t; \mathbf{x}) - f_r(\mathbf{W}^t; \mathbf{x}) \leq 2vM_1 \sum_{j=1}^k \sigma(\langle \mathbf{w}_j^t, \tilde{\mathbf{x}} \rangle) \leq 2vM_1 \sum_{j=1}^k |\mathbf{w}_j^t| \sin \beta < 1.$$

So, with

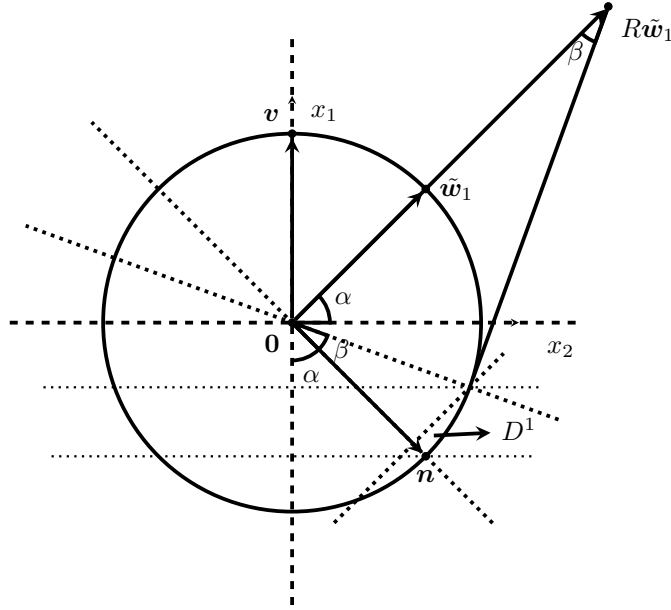
$$D^1 := \{\mathbf{x} \in V_1 : \langle \mathbf{n}, \tilde{\mathbf{x}} \rangle > \cos \beta \text{ and } \langle \tilde{\mathbf{w}}_1, \tilde{\mathbf{x}} \rangle > 0\},$$

we have for any $r > 1$

$$D^1 \subset \Omega_{1,r}.$$

See Fig 2 for a intuition of D^1 . Note that \mathbf{n} and $\tilde{\mathbf{w}}_1$ perpendicular to each other, and the probability distribution p_1 is bounded from below, so there exists a constant p_R , such that

$$\mathbb{P} \left[\Omega_{1,r} \cap \Omega_{\mathbf{w}_j^t} \right] \geq \mathbb{P} \left[D^1 \right] = p_R > 0,$$

FIGURE 2. 2-dim section of \mathbb{R}^d spanned by \tilde{w}_1 and \mathbf{n}

and we have the following estimate for p_R

$$\begin{aligned}
 p_R &= \int_{D^1} \langle \tilde{w}_1^t, \nabla_{\mathbf{w}_1} l(\mathbf{W}) \rangle p_1(\mathbf{x}) \, d\mathbf{x} \\
 &\geq \frac{p_{\min} |\mathcal{S}^{d_1-2}|}{|\mathcal{S}^{d_1-1}|} \int_{\cos \beta}^1 (1-y^2)^{\frac{d_1-3}{2}} \, dy \\
 &= \frac{p_{\min} |\mathcal{S}^{d_1-2}|}{|\mathcal{S}^{d_1-1}|} \int_0^\beta (\sin \theta)^{d_1-2} \, d\theta \\
 &= \Omega \left(\frac{p_{\min} (\sin \beta)^{d_1-1} |\mathcal{S}^{d_1-2}|}{(d_1-1) |\mathcal{S}^{d_1-1}|} \right) \\
 &= \Omega \left(\frac{p_{\min}}{\sqrt{d_1} (M_1 R)^{d_1}} \right).
 \end{aligned}$$

Now, we know the gradient of \mathbf{w}_1 on \tilde{w}_1 direction is bounded below, with same arguments in proof of Lemma 5.1, we have

$$|\langle \tilde{w}_1^t, \nabla_{\mathbf{w}_1} l_1(\mathbf{W}^t) \rangle| \geq \frac{v p_R^2}{C_p}.$$

Note that

$$\sum_{j=1}^{k_1} |\mathbf{w}_j^{t+1}| - |\mathbf{w}_j^t| \geq |\mathbf{w}_1^{t+1}| - |\mathbf{w}_1^t| \geq \eta |\langle \tilde{w}_1^t, \nabla_{\mathbf{w}_1} l(\mathbf{W}^t) \rangle| \geq \frac{v \eta p_R^2}{C_p}$$

and since $|\mathbf{w}_j^t|$ is non-decreasing for $j \in [k_1]$, we have

$$\sum_{t \in T_1} \sum_{j=1}^{k_1} |\mathbf{w}_j^{t+1}| - |\mathbf{w}_j^t| \leq R.$$

Combining the above two equations, it follows $|T_1| \leq \frac{C_p R}{v \eta \rho_R^2}$. \square

At the beginning of the second phase, $\tilde{\mathbf{w}}_j$ are already evenly distributed. That is, the convex hull of $\tilde{\mathbf{w}}_j$ contains the origin, and any input data must at least activate some of the neurons. As long as an input data contributes to the loss, it also contributes to the gradient. So learning process becomes faster during the second phase. The following proposition shows the loss decays faster than before.

Proposition 4. *Let $b_j = 0$ in (6). Let T_2 be the set of t such that $\{\mathbf{w}_j^t : v_{i,j} > 0\}$ satisfies the geometric condition in Lemma 5.2, and that $|\mathbf{W}^t|$ is upper-bounded by R at all t . Then:*

$$\sum_{t \in T_2} l_i(\mathbf{W}^t)^2 \leq 4\eta^{-1} v n^2 C_p R^2 M_i^2 R.$$

Proof of Proposition 4. First, we estimate $l_i(\mathbf{W}^t)$ as follows

$$\begin{aligned} l_i(\mathbf{W}^t) &= \sum_{r \neq i} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_i} [\sigma(1 - f_i + f_r)] \\ &\leq \sum_{r \neq i} (2vRM_i) \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}_i} [f_i < 1 + f_r] \\ &= (2vRM_i) \sum_{r \neq i} \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}_i} [\Omega_{i,r}]. \end{aligned}$$

Now, we have

$$\sum_{r \neq i} \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}_i} [\Omega_{i,r}] \geq \frac{l_i(\mathbf{W}^t)}{2vRM_i}.$$

Second, we estimate the gradient. Let $\Omega_{i,r}^j = \Omega_{i,r} \cap \Omega_{\mathbf{w}_j}$ and assume $v_{i,j} > 0$, we have

$$\nabla_{\mathbf{w}_j} l_i(\mathbf{W}^t) = \sum_{r \neq i} 2v \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_i} [\mathbf{1}_{\Omega_{i,r}^j}(\mathbf{x}) \mathbf{x}].$$

By the same arguments in proof of Lemma 5.1, we know

$$\langle \nabla_{\mathbf{w}_j} l_i(\mathbf{W}^t), \tilde{\mathbf{w}}_j^t \rangle \geq 2v \frac{\mathbb{P}[\Omega_{i,r}^j]^2}{2C_p} = \frac{v \mathbb{P}[\Omega_{i,r}^j]^2}{C_p}.$$

Next, we utilize the geometric condition which implies

$$\sum_{v_{i,j} > 0} \mathbb{P}[\Omega_{i,r}^j] \geq \mathbb{P}[\Omega_{i,r}]$$

and get

$$\begin{aligned}
\sum_{v_{i,j}>0} |\mathbf{w}_j^{t+1}| - |\mathbf{w}_j^t| &\geq \sum_{\substack{v_{i,j}>0 \\ r \neq i}} \frac{v\eta \mathbb{P}[\Omega_{i,r}^j]^2}{C_p} = \frac{v\eta}{C_p} \sum_{\substack{v_{i,j}>0 \\ r \neq i}} \mathbb{P}[\Omega_{i,r}^j]^2 \\
&\geq \frac{v\eta}{C_p} \frac{1}{k_i} \sum_{r \neq i} \mathbb{P}[\Omega_{i,r}]^2 \geq \frac{v\eta}{C_p} \frac{1}{k_i(n-1)} \left(\sum_{r \neq i} \mathbb{P}[\Omega_{i,r}] \right)^2 \\
&\geq \frac{v\eta}{n^2 C_p} \frac{l_i(\mathbf{W}^t)^2}{4v^2 R^2 M_i^2} = \frac{\eta}{4vn^2 C_p R^2 M_i^2} l_i(\mathbf{W}^t)^2.
\end{aligned}$$

where k_i is the number of j 's such that $V_{i,j} > 0$.

Finally, we combine the inequalities above and the assumption $|\mathbf{W}^t| < R$ and get

$$\sum_{t \in T_2} l_i(\mathbf{W}^t)^2 \leq \frac{4vn^2 C_p R^2 M_i^2}{\eta} \sum_{\substack{v_{i,j}>0 \\ t \in T_2}} |\mathbf{w}_j^{t+1}| - |\mathbf{w}_j^t| \leq 4\eta^{-1} vn^2 C_p R^2 M_i^2 R.$$

□

From the above proposition, we see that in order to get $l_i(\mathbf{W}^t) < \epsilon$ for some $t \in T_2$, we only need

$$|T_2| \geq \frac{4vn^2 C_p R^2 M_i^2 R}{\eta \epsilon}.$$

Compared with Proposition 3, we see that $|T_2|$ does not rely on the dimension d_i , whereas the upper bound of $|T_1|$ is exponential in d_i .

6. Experiments. In this section, we report the results of our experiments on both synthetic and MNIST data. The experiments on synthetic data aim to show that *convergence to global minimum continues to hold if data subspaces form acute angles, going beyond the theoretical orthogonality assumption* under which convergence is observed to be the fastest. Our theoretical results and the geometric conditions are supported by simulations. Experiments on MNIST dataset exhibit subspace structures in data flow and slow-to-fast training dynamics on LeNet-5. These phenomena from our model are worth further study in deep networks.

6.1. Synthetic data. Let $\{\mathbf{e}_j\}_{j \in [4]}$ be orthonormal basis of \mathbb{R}^4 , θ be an acute angle and $\mathbf{v}_1 = \mathbf{e}_1$, $\mathbf{v}_2 = \sin \theta \mathbf{e}_2 + \cos \theta \mathbf{e}_3$, $\mathbf{v}_3 = \mathbf{e}_3$, $\mathbf{v}_4 = \mathbf{e}_4$. Now, we have two linearly independent subspaces of \mathbb{R}^4 namely $V_1 = \text{Span}(\{\mathbf{v}_1, \mathbf{v}_2\})$ and $V_2 = \text{Span}(\{\mathbf{v}_3, \mathbf{v}_4\})$. We can easily calculate that the angle between V_1 and V_2 is θ . Next, we define

$$\hat{\mathcal{X}}_1 = \{r(\cos \varphi \mathbf{v}_1 + \sin \varphi \mathbf{v}_2) : r \in S_r, \varphi \in S_\varphi\},$$

$$\hat{\mathcal{X}}_2 = \{r(\cos \varphi \mathbf{v}_3 + \sin \varphi \mathbf{v}_4) : r \in S_r, \varphi \in S_\varphi\},$$

where

$$S_r = \left\{ \frac{20}{j} : j \in [20] - [9] \right\}, \quad S_\varphi = \left\{ \frac{j\pi}{40} : j \in [80] \right\}.$$

Let \mathcal{X}_1 corresponds to label $y = 1$ and \mathcal{X}_2 corresponds to label $y = -1$. Since we are considering a binary classification problem, the neural network structure can be

simplified as (7):

$$(7) \quad \tilde{f}(\mathbf{W}; \mathbf{x}) = \sum_{j=1}^k \sigma(h_j) - \sum_{j=k+1}^{2k} \sigma(h_j),$$

and the prediction is given by $\hat{y}(\mathbf{x}) = \text{sign}f(\tilde{\mathbf{W}}; \mathbf{x})$. Now, the population loss becomes

$$l_i(\mathbf{W}) := \frac{1}{|\mathcal{X}_i|} \sum_{\mathbf{x} \in \mathcal{X}_i} \max \left\{ 0, 1 + (-1)^i \tilde{f}(\mathbf{W}; \mathbf{x}) \right\}.$$

In our first simulation, we set $k = 4$ in (7) and run gradient descent (3) on $l_1 + l_2$ with learning rate $\eta = 0.1$. Fig. 3 shows the iterations it takes to converge to global minima given θ and a Gaussian noise added to \mathcal{X}_i 's. From this simulation, we see that the orthogonal data assumptions are only technically needed and our convergence result holds in more general settings.

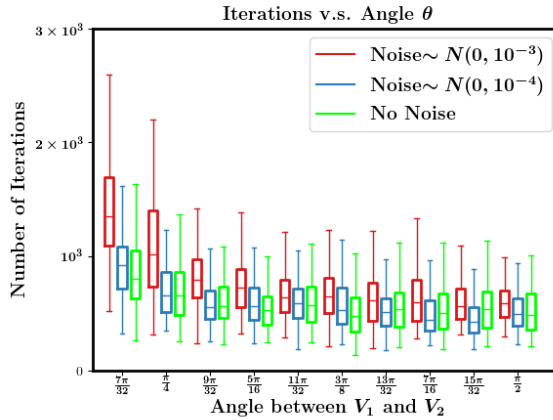


FIGURE 3. Number of iterations to convergence v.s. θ , the angle between subspaces V_1 and V_2 .

In our second and remaining simulations of this subsection, we take $\theta = \frac{\pi}{2}$ so that V_1 and V_2 are orthogonal. Lemma 4.1 suggests the learning process of l_1 and l_2 are independent, so we only simulate the training process of l_1 and assume $\mathbf{w}_j \in V_1 \cong \mathbb{R}^2$. We take entries of \mathbf{W}^0 to be i.i.d. standard normal i.e. $\mathbf{w}_j^0 \sim N(\mathbf{0}, I_2)$. We train the network (7) with gradient descent (3) in all our simulations, where learning rate $\eta = 0.1$. The left plot of Fig. 4 shows how many iterations algorithm (3) takes in searching for a global minima from the random initialization mentioned above. For each box, the red mark indicates the median, and the bottom and top of the box indicate the 25th and 75th percentiles, respectively. As we can see from the graph, as number of hidden neurons ($2k$) becomes larger, the algorithm (3) tends to need less iterations in searching for a global minima.

In the third simulation, we compare the convergence speed with and without the geometric condition being satisfied. We introduce two initialization method: random initialization and half space initialization i.e. with $\hat{w}_{j,i} \sim N(0, 1)$, random initialization takes $w_{j,i}^0 = \hat{w}_{j,i}$ whereas half space initialization takes $w_{j,1}^0 = |\hat{w}_{j,1}|$ and $w_{j,2}^0 = \hat{w}_{j,2}$. We run the algorithm for 100 times with different numbers of

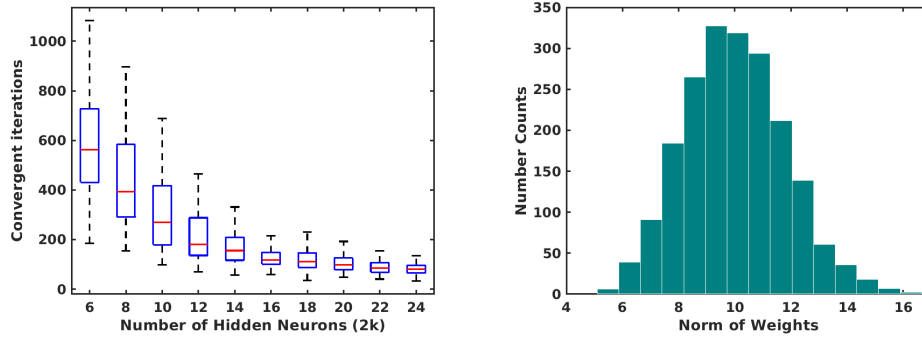


FIGURE 4. Left: convergent iterations vs. number of neurons ($d = 2$). Right: histogram of norm of weights: $\max_t |\mathbf{W}^t|$ ($d = 2$ and $k = 4$).

hidden neurons using initialization methods, and report the means and standard variances of the number of iterations in Table 1. We see from Remark 3 how the P_{gc} increases when the number of hidden neurons grows. However, the half space initialization never satisfies the geometric condition, as all the weights lie in the same half space. A widely believed explanation on why a neural network can fit all training labels is that the neural network is over-parameterized. Our work explained one of the reasons why over-parameterization helps convergence: it helps the weights to spread more 'evenly' and quickly after initialization. Table 1 shows that when we randomly initialize, the iterations for convergence in gradient descent (3) come down a lot as the number of hidden neurons increases; much less so in half space initialization.

TABLE 1. Iterations taken (mean \pm std) to convergence with random and half space initializations.

# of Neurons ($2k$)	Random Init.	Half Space Init.
6	578.90 \pm 205.43	672.41 \pm 226.53
8	423.96 \pm 190.91	582.16 \pm 200.81
10	313.29 \pm 178.67	550.19 \pm 180.59
12	242.72 \pm 178.94	517.26 \pm 172.46
14	183.53 \pm 108.60	500.42 \pm 215.87
16	141.00 \pm 80.66	487.42 \pm 220.48
18	126.52 \pm 62.07	478.25 \pm 202.71
20	102.09 \pm 32.32	412.46 \pm 195.92
22	90.65 \pm 28.01	454.08 \pm 203.00
24	82.93 \pm 26.76	416.82 \pm 216.58

Our fourth simulation take specifically $2k = 8$. With 2000 runs we did a histogram of the maximum norm of \mathbf{W} during the training process shown in the right plot of Fig. 4. In fact, our third simulation suggests our boundedness assumption on \mathbf{W} in Theorem 3.1 and Theorem 3.2 are reasonable.

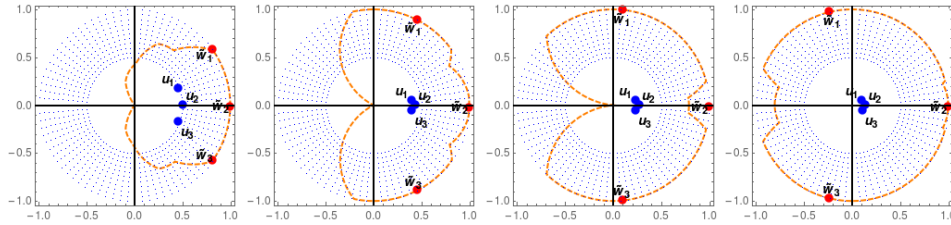


FIGURE 5. Dynamics of weights: \tilde{w}_j and u_j

In our last simulation, we take $k = 3$ so that there are in total 6 hidden neurons. For notation simplicity, we denote $u_j = w_{j+3}$. For $j \in [3]$, we plot \tilde{w}_j 's and u_j 's in Fig. 5, where we plot \tilde{w}_j 's instead of w_j 's since some of $|w_j|$'s are greater than one. Before algorithm (3) starts, the parameters in neural network (6) are initialized to be

$$w_j^0 = u_j^0 = \frac{3}{4} \left(\cos \frac{(2-j)\pi}{6}, \sin \frac{(2-j)\pi}{6} \right)$$

for $j \in \{1, 2, 3\}$. In Fig. 5, the tiny blue points are input data under Kelvin transformation: $\mathbf{x} \rightarrow \mathbf{x}^* = \frac{\mathbf{x}}{|\mathbf{x}|^2}$. Take $\mathbf{x} = \frac{1}{r}(\cos \theta, \sin \theta)$ so that under Kelvin transformation $\mathbf{x}^* = r(\cos \theta, \sin \theta)$. For convenience, we let $\tilde{\mathbf{x}} = r\mathbf{x} = (\cos \theta, \sin \theta)$. The orange dashed curve has expression in polar coordinates:

$$\rho(\theta) = \min \left\{ 1, \sigma \left(\tilde{f}(\mathbf{W}; \tilde{\mathbf{x}}) \right) \right\}.$$

Note that we are taking Hinge loss $l(\mathbf{W}; \{\mathbf{x}, 1\}) = 0$ if and only if $\tilde{f}(\mathbf{W}; \mathbf{x}) \geq 1$, i.e.

$$\tilde{f}(\mathbf{W}; \tilde{\mathbf{x}}) = r\tilde{f}(\mathbf{W}; \mathbf{x}) \geq r = |\mathbf{x}^*|.$$

Here, in our data set $\hat{\mathcal{X}}$, all data point have norm less than one under Kelvin transformation, so $l(\mathbf{W}; \{\mathbf{x}, 1\}) = 0$ if and only if $\rho(\theta) \geq |\mathbf{x}^*|$. This means, the blue points when surrounded by the orange dashed curve provide zero loss. In particular, when $\rho(\theta) = 1$, the population loss is 0.

6.2. MNIST experiments. The two-phase dynamics we proved in our model does appear in deep network training on real (non-synthetic) data sets. In experiments on MNIST, we used a simplified version of LeNet-5 [16] with 2 convolutional layers and two fully-connected (fc) layers; see [10] for the full version where F6 and output layers correspond to our fc layers. The simplified Lenet-5 is trained via stochastic gradient descent (SGD) at constant learning rate 0.01, batch size 1000 and without momentum or regularization. We show in Fig. 6 (left) the loss value vs. iterations during training. At the early stage, the loss decays slowly, then the fast phase sets in after 400 iterations. Fig. 6 (left) clearly supports our theory on the slow and fast dynamics of gradient descent.

Network visualization helps understand its geometric properties. In Fig. 6 (right), we plot 2D projections of feature vectors at input to fc layer extracted by a neural network [18] consisting of 6 convolutional layers followed by an fc layer. Projected features from different classes cluster around linearly independent subspaces. The plot suggests that in trained deep networks, the linearly independent subspace assumption approximately holds for the input to the fully connected layer before classification output. Similar subspace structure of high level feature vectors

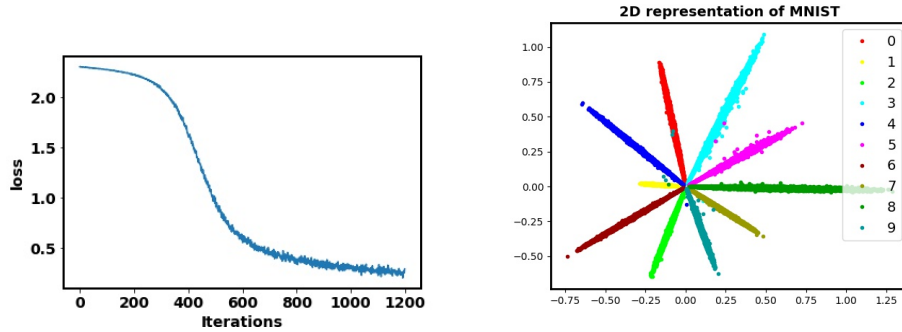


FIGURE 6. **Left:** Slow-to-Fast transition during LeNet [16] training on MNIST dataset. **Right:** 2D projections of MNIST features from a trained convolutional neural network [18]. The 10 classes are color coded, the feature points cluster near linearly independent subspaces.

on CIFAR-10 and enlargement of subspace angles to improve classification accuracy have been studied in [19].

In Fig. 7, we show four projections onto unit sphere \mathcal{S}^2 (inside randomly selected 3D subspaces) of the weight vectors of the first and second fc layers of LeNet [16]. Visual inspection on these and others (not shown) suggests that our geometric condition holds with high probability.

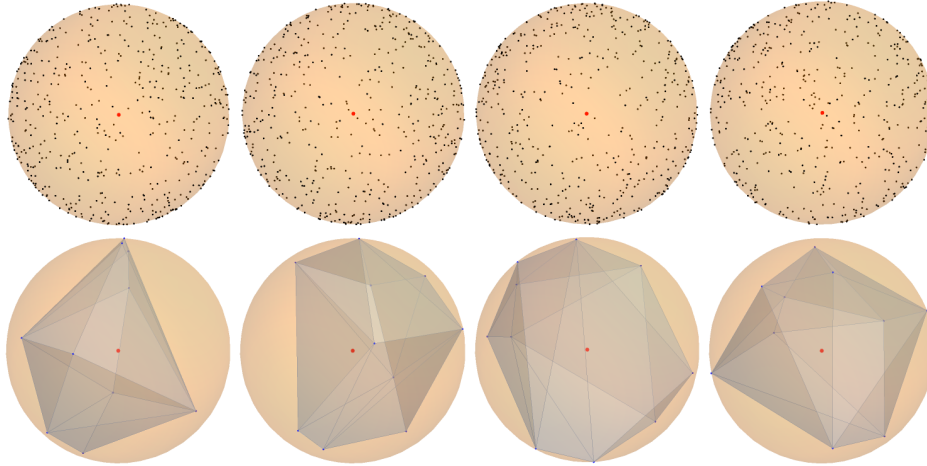


FIGURE 7. **Top row:** Projections onto \mathcal{S}^2 (inside randomly selected 3D subspaces) of weight vectors in the first fully connected layer of a trained LeNet. **Bottom row:** Projections onto \mathcal{S}^2 (inside randomly selected 3D subspaces) of weight vectors and their convex hull in the second fully connected layer of a trained LeNet.

7. Conclusions and future work. The slow and fast dynamics of neural network weights under gradient descent is critical for understanding the learning process. We performed the first theoretical study on training neural networks to classify linearly inseparable data sets away from the over-parameterized regime. We discovered a two time-scale phenomenon of network weights during gradient descent training: a slow phase where the weights spread out to satisfy a geometric condition, and a subsequent fast phase where the weights converge to a global minimum. Both the two scale dynamics and geometric conditions are supported by LeNet-5 training on MNIST data.

A future direction is to provide a concrete relation between the number of weights and the rate of convergence, and quantify the effect of over-parameterization on the rate of convergence. Another direction is to devise efficient method to verify the geometric condition computationally for weight vectors in training deep neural networks on general linearly non-separable data sets.

Acknowledgments. This work was partially supported by NSF grants IIS-1632935, DMS-1854434, DMS-1924548, and DMS-1924935.

REFERENCES

- [1] Z. Allen-Zhu, Y. Li and Z. Song, A convergence theory for deep learning via over-parameterization, preprint, [arXiv:1811.03962](https://arxiv.org/abs/1811.03962).
- [2] A. Brutzkus and A. Globerson, Globally optimal gradient descent for a ConvNet with Gaussian inputs, preprint, [arXiv:1702.07966](https://arxiv.org/abs/1702.07966).
- [3] A. Brutzkus and A. Globerson, Over-parameterization improves generalization in the XOR detection problem, preprint.
- [4] A. Brutzkus, A. Globerson, E. Malach and S. Shalev-Shwartz, SGD learns over-parameterized networks that provably generalize on linearly separable data, 6th International Conference on Learning Representations, Vancouver, BC, Canada, 2018, preprint, [arXiv:1710.10174](https://arxiv.org/abs/1710.10174).
- [5] R. T. des Combes, M. Pezeshki, S. Shabanian, A. Courville and Y. Bengio, Convergence properties of deep neural networks on separable data, 2019. Available from: <https://openreview.net/forum?id=HJfQrsOqt7>.
- [6] S. S. Du, X. Zhai, B. Póczós and A. Singh, Gradient descent provably optimizes over-parameterized neural networks, preprint, [arXiv:1810.02054](https://arxiv.org/abs/1810.02054).
- [7] C. Ho and S. Zimmerman, On the number of regions in an m -dimensional space cut by n hyperplanes, *Austral. Math. Soc. Gaz.*, **33** (2006), 260–264.
- [8] S. Hochreiter and J. Schmidhuber, [Long short-term memory](#), *Neural Comput.*, **9** (1997), 1735–1780.
- [9] A. Krizhevsky, I. Sutskever and G. E. Hinton, [ImageNet classification with deep convolutional neural networks](#), in *Advances in Neural Information Processing Systems*, 2012, 1097–1105.
- [10] *LeNet-5 – A Classic CNN Architecture*. Available from: <https://engmrk.com/lenet-5-a-classic-cnn-architecture/>.
- [11] Y. Li and Y. Liang, Learning overparameterized neural networks via stochastic gradient descent on structured data, preprint, [arXiv:1808.01204](https://arxiv.org/abs/1808.01204).
- [12] S. Liang, R. Sun, Y. Li and R. Srikant, Understanding the loss surface of neural networks for binary classification, preprint, [arXiv:1803.00909](https://arxiv.org/abs/1803.00909).
- [13] B. Neyshabur, Z. Li, S. Bhojanapalli, Y. LeCun and N. Srebro, Towards understanding the role of over-parameterization in generalization of neural networks, preprint, [arXiv:1805.12076](https://arxiv.org/abs/1805.12076).
- [14] Q. Nguyen, M. C. Mukkamala and M. Hein, On the loss landscape of a class of deep neural networks with no bad local valleys, preprint, [arXiv:1809.10749](https://arxiv.org/abs/1809.10749).
- [15] S. Ren, K. He, R. Girshick and J. Sun, [Faster R-CNN: Towards real-time object detection with region proposal networks](#), *IEEE Trans. Pattern Anal. Machine Intell.*, **39** (2017), 1137–1149.
- [16] A. Rosebrock, [LeNet – Convolutional neural network in Python](#), 2016. Available from: <https://www.pyimagesearch.com/2016/08/01/lenet-convolutional-neural-network-in-python/>.

- [17] D. Silver, A. Huang, C. J. Maddison, A. Guez and L. Sifre, et al., [Mastering the game of Go with deep neural networks and tree search](#), *Nature*, **529** (2016), 484–489.
- [18] H. Wang, Y. Wang, Z. Zhou, X. Ji and D. Gong, et al., CosFace: Large margin cosine loss for deep face recognition, preprint, [arXiv:1801.09414](#).
- [19] P. Yin, J. Xin and Y. Qi, [Linear feature transform and enhancement of classification on deep neural network](#), *J. Sci. Comput.*, **76** (2018), 1396–1406.

Received November 2019; 1st revision July 2020; final revision October 2020.

E-mail address: zlong6@uci.edu

E-mail address: pyin@albany.edu

E-mail address: jxin@math.uci.edu