
Almost Global Convergence to Global Minima for Gradient Descent in Deep Linear Networks

Zhenyu Liao, Yacine Chitour, Romain Couillet

Laboratoire des Signaux et Systèmes
CentraleSupélec, Université Paris-Saclay
Gif-sur-Yvette, France

Abstract

In this article we prove the global convergence (as opposed to the less challenging local behavior), for almost all training data-target pairs and almost all initializations, of a linear deep network to a global minimum when using the classical gradient descent method with small step size. This global result is obtained through an original geometric framework relying on a key invariance property induced by the network structure and providing, as a fundamental side result, a clearer picture of the loss landscape. We further argue that the presented framework is sufficiently powerful to envision extensions to nonlinear deep networks.

1 Introduction

Despite the rapid growing list of successful applications of deep neural networks trained with back-propagation in various fields from computer vision [15] to speech recognition [19] and natural language processing [7], our theoretical understanding on these elaborate systems, however, is developing at a more modest pace.

One of the major difficulties in the design of deep neural networks today is that, to obtain networks with greater expressive power, we cascade more and more layers to make them “deeper” and hope to extract more “abstract” features from the (numerous) training data so as to improve the networks in terms of generalization performance. Nonetheless, from an optimization viewpoint, this “deeper” structure poses problems because it gives rise to non-convex loss functions and makes the optimization seemingly intractable. In general finding a global minimum of a *generic* non-convex function is an NP-complete problem [20] and it is unfortunately the case for neural networks as it was shown in [3] that even training a very simple network is indeed NP-complete.

Yet, many non-convex problems such as phase retrieval, independent component analysis and orthogonal tensor decomposition are known to obey the important properties [23] that 1) all local minima are also global; and 2) around any saddle point the objective function has a negative directional curvature (i.e., the possibility to continue to descend) and thus allow for the possibility to find some way to fall into a “basin” with a (comparably) low loss “with high probability”. In this regard, the loss surfaces of deep neural networks are receiving an unprecedented research interest: in the pioneering work of Baldi & Hornik [2] the landscape of mean square losses was studied in the case of linear auto-encoders (i.e., the same dimension for input data and output targets) of depth one; more recently in the work of Saxe et al. [22] the dynamics of the corresponding gradient descent system was first studied, by assuming the input data \mathbf{X} empirical correlation matrix $\mathbf{X}\mathbf{X}^\top$ to be identity, in a linear deep neural networks, so as to propose a novel initialization method. Then in [13] the author proved that under some appropriate rank condition on the (cascading) matrix product, all critical points of a deep linear neural networks are either global minima or saddle points with Hessian admitting eigenvalues with different signs, meaning that linear deep networks are somehow “close” to those examples mentioned at the beginning of this paragraph. Nonetheless, the results in [22, 13] are

incomplete in the sense that they do not provide enough (global) information regarding when and how can gradient descent trajectories result in these global minima (recall that, to escape from saddle points within a reasonable time one may alternatively use second-order methods with information from the Hessian, artificially perturb the gradient with noise as in [12], etc.). Concretely speaking, previous analyses in [13, 16] only focus on the *local* behavior of each critical point and a “global picture” on the *whole* space occupied by the network weights is still in demand.

In this paper, we elaborate on the model from [22, 13] and evaluate the dynamics of the associated gradient system in a “continuous” manner. We prove that, for almost every choice of training data-target pair (\mathbf{X}, \mathbf{Y}) and almost every initialization for the weight matrices \mathbf{W}_i , the corresponding trajectory of the gradient system converges to a global minimum of the loss function. Based on a cornerstone “invariant” in the parameter space induced by the network cascading structure, we establish a generic framework for the geometric understanding of deep neural networks and provide the aforementioned sought-for global picture of the gradient descent dynamics in the specific case of linear networks. Due to space limitation, only proof sketches are provided in the $H = 1$ layer scenario, along with generalization intuitions to the $H > 1$ case. The complete derivations for the case $H \geq 1$ are available in an extended version.

2 System Model and Main Result

2.1 Problem setup

We start with a deep linear neural network with H hidden layers as illustrated in Figure 1. To begin with, the network structure as well as associated notations are presented as follows.

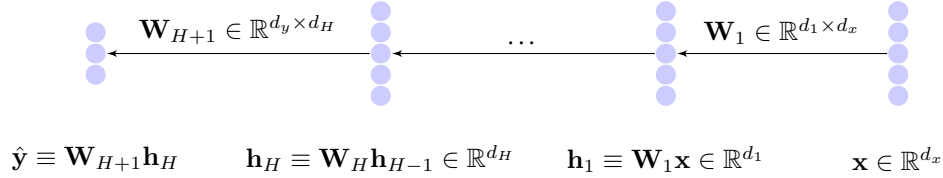


Figure 1: Illustration of the H -hidden-layer linear neural network

Let the pair (\mathbf{X}, \mathbf{Y}) denote the training data and associated targets, with $\mathbf{X} \equiv [\mathbf{x}_1, \dots, \mathbf{x}_m] \in \mathbb{R}^{d_x \times m}$ and $\mathbf{Y} \equiv [\mathbf{y}_1, \dots, \mathbf{y}_m] \in \mathbb{R}^{d_y \times m}$, where m denotes the number of instances in the training set and d_x, d_y the dimensions of data and targets, respectively. We denote the weight matrix $\mathbf{W}_i \in \mathbb{R}^{d_i \times d_{i-1}}$ that connects \mathbf{h}_{i-1} to \mathbf{h}_i for $i = 1, \dots, H+1$ and set $\mathbf{h}_0 = \mathbf{x}$, $\mathbf{h}_{H+1} = \hat{\mathbf{y}}$ as in Figure 1. The network output is thus given by $\hat{\mathbf{Y}} = \mathbf{W}_{H+1} \dots \mathbf{W}_1 \mathbf{X}$. We denote \mathbf{W} the $(H+1)$ -tuple of $(\mathbf{W}_1, \dots, \mathbf{W}_{H+1})$ for simplicity and work on the mean squared error $\mathcal{L}(\mathbf{W})$ given by the Frobenius norm below,

$$\mathcal{L}(\mathbf{W}) = \frac{1}{2} \|\mathbf{Y} - \hat{\mathbf{Y}}\|_F^2 = \frac{1}{2} \|\mathbf{Y} - \mathbf{W}_{H+1} \dots \mathbf{W}_1 \mathbf{X}\|_F^2 \quad (1)$$

under the following assumptions:

Assumption 1 (Dimension Condition). $m \geq d_x \geq \max(d_1, \dots, d_H) \geq \min(d_1, \dots, d_H) \geq d_y$. In particular in the case $H = 1$ this condition yields $m \geq d_x \geq d_1 \geq d_y$.

Assumption 2 (Full Rank Data and Targets). The matrices \mathbf{X} and \mathbf{Y} are of full (row) rank, i.e., of rank d_x and d_y , respectively, accordingly with Assumption 1.

Assumption 1 and 2 on the dimension and rank of the training data are realistic and practically easy to satisfy, as discussed in previous works [2, 13].¹

Under Assumptions 1 and 2, with the singular value decomposition on $\mathbf{X} = \mathbf{U}_\mathbf{X} \Sigma_\mathbf{X} \mathbf{V}_\mathbf{X}^\top$ with $\mathbf{V}_\mathbf{X} = [\mathbf{V}_\mathbf{X}^1 \mid \mathbf{V}_\mathbf{X}^2]$, $\mathbf{V}_\mathbf{X}^1 \in \mathbb{R}^{m \times d_x}$ and then on $\mathbf{Y} \mathbf{V}_\mathbf{X}^1 \equiv \bar{\mathbf{Y}} = \mathbf{U}_\mathbf{Y} \Sigma_\mathbf{Y} \mathbf{V}_\mathbf{Y}^\top$, together with an

¹Assumption 1 is demanded here for convenience and our results can be extended to handle more elaborate dimension settings. Similarly, when the training data is rank deficient, the learning problem can be reduced to a lower-dimensional one by removing these non-informative data in such a way that Assumption 2 holds.

immediate change of variable, we get $\mathcal{L}(\mathbf{W}) = L(\mathbf{W}) + \frac{1}{2}\|\mathbf{Y}\mathbf{V}_{\mathbf{X}}^2\|_F^2$ with

$$L(\mathbf{W}) \equiv \frac{1}{2}\|\Sigma_{\mathbf{Y}} - \overline{\mathbf{W}}_{H+1}\overline{\mathbf{W}}_H \dots \overline{\mathbf{W}}_2\overline{\mathbf{W}}_1\|_F^2 \quad (2)$$

where $\Sigma_{\mathbf{X}} \equiv [\mathbf{S}_{\mathbf{X}} \mid \mathbf{0}] \in \mathbb{R}^{d_x \times m}$, $\Sigma_{\mathbf{Y}} \in \mathbb{R}^{d_y \times d_x}$ and we denote $\overline{\mathbf{W}}_{H+1} \equiv \mathbf{U}_{\mathbf{Y}}^T \mathbf{W}_{H+1} \in \mathbb{R}^{d_y \times d_H}$, $\overline{\mathbf{W}}_1 \equiv \mathbf{W}_1 \mathbf{U}_{\mathbf{X}} \mathbf{S}_{\mathbf{X}} \mathbf{V}_{\mathbf{Y}} \in \mathbb{R}^{d_1 \times d_x}$ and $\overline{\mathbf{W}}_i = \mathbf{W}_i$ for $i = 2, \dots, H$. Therefore the state space² of $\Xi \equiv (\overline{\mathbf{W}}_{H+1}, \dots, \overline{\mathbf{W}}_1)$ is equal to $\mathcal{X} = \mathbb{R}^{d_y \times d_H} \times \dots \times \mathbb{R}^{d_1 \times d_x}$. In particular, for $H = 1$ we have $d_H = d_1$ and \mathcal{X} has dimension $d_1(d_x + d_y)$.

With the above notations, we demand in addition the following assumption on the target $\overline{\mathbf{Y}}$.

Assumption 3 (Distinct Singular Values). *The target $\overline{\mathbf{Y}}$ has d_y distinct singular values.*

Although seemingly restrictive, Assumption 3 actually holds for an open and dense subset of $\mathbb{R}^{d_y \times d_x}$.

The objective of this article is to study the gradient descent [5] dynamics (GDD) defined as

Definition 1 (GDD). *The Gradient Descent Dynamic of L is the dynamical system defined on \mathcal{X} by*

$$\frac{d\Xi}{dt} = -\nabla_{\Xi} L(\Xi) \quad (3)$$

where $\nabla_{\Xi} L(\Xi)$ denotes the gradient of the loss function L with respect to Ξ . A point $\Xi \in \mathcal{X}$ is a critical point of L if and only if $\nabla_{\Xi} L(\Xi) = \mathbf{0}$ and we denote $\text{Crit}(L)$ the set of critical points.

In the following, we work directly on the equivalent equation (2) and start by evaluating the gradient of L . With the previous notations, for $\xi \equiv (\overline{\mathbf{W}}_{H+1}, \dots, \overline{\mathbf{W}}_1)$, we expand the variation of $L(\Xi + \xi)$ as

$$L(\Xi + \xi) = L(\Xi) + \Delta_{\Xi}(\xi) + O(\|\xi\|^2)$$

with $\overline{\mathbf{M}} \equiv \Sigma_{\mathbf{Y}} - \overline{\mathbf{W}}_{H+1} \dots \overline{\mathbf{W}}_1$, $L(\Xi) = \frac{1}{2}\|\overline{\mathbf{M}}\|_F^2$ and the differential $\Delta_{\Xi}(\xi)$ given by $\Delta_{\Xi}(\xi) \equiv -\sum_{j=1}^{H+1} \text{tr}(\overline{\mathbf{M}}^T \overline{\mathbf{W}}_{H+1} \dots \overline{\mathbf{W}}_{j+1} \overline{\mathbf{w}}_j \overline{\mathbf{W}}_{j-1} \dots \overline{\mathbf{W}}_1)$. We thus derive from Definition 1 the dynamics of L , for $j = 1, \dots, H+1$, as

$$\frac{d\overline{\mathbf{W}}_j}{dt} \equiv -\nabla_{\overline{\mathbf{W}}_j} L(\Xi) = (\overline{\mathbf{W}}_{H+1} \dots \overline{\mathbf{W}}_{j+1})^T \overline{\mathbf{M}} (\overline{\mathbf{W}}_{j-1} \dots \overline{\mathbf{W}}_1)^T. \quad (4)$$

We first remark the following interesting (and crucial to what follows) property of the gradient system (4), inspired by [22] which essentially considered the case where all dimensions are equal to one.

Lemma 1 (Invariant in GDD). *Consider any trajectory of the gradient system given by (4). Then, for $j = 1, \dots, H$, the value of $\overline{\mathbf{W}}_{j+1}^T \overline{\mathbf{W}}_{j+1} - \overline{\mathbf{W}}_j \overline{\mathbf{W}}_j^T$ remains constant, i.e.,*

$$\overline{\mathbf{W}}_{j+1}^T(t) \overline{\mathbf{W}}_{j+1}(t) - \overline{\mathbf{W}}_j(t) \overline{\mathbf{W}}_j^T(t) = \mathbf{C}_j^0 \equiv (\overline{\mathbf{W}}_{j+1}^T \overline{\mathbf{W}}_{j+1} - \overline{\mathbf{W}}_j \overline{\mathbf{W}}_j^T) \Big|_{t=0}, \quad \forall t \geq 0.$$

In particular, in the case of $H = 1$ we get $\overline{\mathbf{W}}_2^T \overline{\mathbf{W}}_2 - \overline{\mathbf{W}}_1 \overline{\mathbf{W}}_1^T = \mathbf{C}^0 \equiv (\overline{\mathbf{W}}_2^T \overline{\mathbf{W}}_2 - \overline{\mathbf{W}}_1 \overline{\mathbf{W}}_1^T) \Big|_{t=0}$.

Proof. Simply check that $\frac{d}{dt} (\overline{\mathbf{W}}_{j+1}^T \overline{\mathbf{W}}_{j+1} - \overline{\mathbf{W}}_j \overline{\mathbf{W}}_j^T) = 0$. \square

Lemma 1 provides a key structural property of the GDD that is instrumental to ensure the boundedness of the gradient descent trajectories and thus in turn to prove the convergence to global minima.

2.2 Main Results

Our main result is Theorem 1 which provides information on the convergence of the GDD trajectories to global minima under reasonable conditions. This result is obtained as follows: we first prove that every trajectory of the GDD converges to a critical point of the loss function L . We then go on with a precise characterization and classification of the critical points, followed by a *local* analysis of

²The network (weight) parameters Ξ evolve through time and are considered to be *state variables* of the dynamical system, while the pair (\mathbf{X}, \mathbf{Y}) is fixed and thus referred as the “parameters” of the given system.

(4) around *each* critical point. More concretely, we first focus on the state space \mathcal{X} and figure out “how much” is occupied by the saddle points (it turns out that there is no local maximum [13]): we stratify the set of critical points $\text{Crit}(L)$ in d_y subsets, one of them corresponding to the set of global minima and the $d_y - 1$ others, $\text{Crit}_r(L)$, with $r = 1, \dots, d_y - 1$, corresponding to the set of saddle points. Then, we perform the aforementioned fine study on the *local* behavior of gradient descent trajectories “around” each saddle point, so as to measure “how much” of the space is converging towards a given saddle point. Precisely speaking, we determine, for $1 \leq r \leq d_y - 1$, an upper bound D_S^r on the “dimension” of the basin of attraction of each single saddle point in $\text{Crit}_r(L)$, which is the set of initializations for which the GDD trajectories converge to that saddle point, so as to evaluate the “chance” of being stuck when the corresponding trajectory gets close to such a given saddle point. To ensure *global* convergence from the *local* analysis, we next derive an upper bound D_C^r for the “dimension” of $\text{Crit}_r(L)$ and then show that the sum $D_S^r + D_C^r$ is smaller than the dimension of the state space \mathcal{X} minus two. This allows us to conclude by means of a transversality argument [10].

We start with the global convergence to critical points of all gradient descent trajectories. While one expects the gradient descent algorithm to converge to critical points, this may not always be the case. Two possible (undesirable) situations are 1) a trajectory is unbounded or 2) it oscillates “around” several critical points without convergence, i.e., along an ω -limit set made of a continuum of critical points (see [24] for notions on ω -limit sets). The property of an iterative algorithm (like gradient descent) to converge to a critical point for any initialization is referred to as “global convergence” [25]. However, it is very important to stress the fact that it does not imply (contrary to what the name might suggest) convergence to a global (or good) minimum for all initializations.

To answer the convergence question, we resort to Łojasiewicz’s theorem³ for the convergence of a gradient descent flow of the type of (4) with real analytic right-hand side [17]. Since the loss function $L(\Xi)$ is a polynomial of degree $(H + 1)^2$ in the components of Ξ , Łojasiewicz’s theorem ensures that if a given trajectory of the gradient descent flow is bounded (i.e., it remains in a compact set for every $t \geq 0$) it must converge to a critical point with a guaranteed rate of convergence. In particular, the previous phenomenon of “oscillation” cannot occur and we are left to ensure the absence of unbounded trajectories. Lemma 1 is the core argument to show that all trajectories of the GDD are indeed bounded, leading to the first result of this article as follows.

Proposition 1 (Global Convergence of GDD to Critical Points). *Let (\mathbf{X}, \mathbf{Y}) be a data-target pair satisfying Assumptions 1 and 2. Then, every trajectory of the corresponding gradient flow (4) converges to a critical point as $t \rightarrow \infty$, at rate at least $t^{-\alpha}$, for some fixed $\alpha > 0$ only depending on the dimensions of the problem.*

A first consequence of Proposition 1 is that it provides a rigorous justification for the appropriate discretization of the GDD given in (3). Indeed the step size can be chosen in terms of an a priori bound for the whole trajectory, which is explicitly determined only with the initial condition (see [4]). This is in contrast with [16] whose step size of the discretization of the GDD is determined with a bound on the Hessian norm of a critical point, however, the latter assumption supposes that the trajectory converges to a “well-known” critical point (with a prior information on the Hessian for example) while no such information is available at the initial stage of training.

To provide guarantees of global convergence to a “good” critical point, we then carry out the aforementioned analysis of the dimension of the sets $\text{Crit}_r(L)$ to obtain our main result as follows.

Theorem 1 (GDD Converges to a Global Minimum for Almost All Initializations). *Let Assumptions 1-3 hold. Then there exists an open and dense subset $\overline{\mathcal{P}}$ of the parameter (data) space \mathcal{P} so that, for every pair (\mathbf{X}, \mathbf{Y}) in $\overline{\mathcal{P}}$, there exists an open and dense subset $\overline{\mathcal{X}}$ of the state space \mathcal{X} such that every trajectory of the GDD in (4) corresponding to (\mathbf{X}, \mathbf{Y}) and starting in $\overline{\mathcal{X}}$ converges to a global minimum with $L(\Xi) = 0$.*

Previous works [13, 16] only studied *local* properties of critical points by establishing that the basin of attraction of each saddle point, i.e., the set of initializations of the GDD trajectories converging to that saddle point, is of measure zero. However, to obtain a *global* picture, one must estimate how “big” is the *union* of all these basins of attraction. For that purpose, first note that the set of saddle points, being an algebraic variety of positive dimension (see Item *iii*) of Proposition 2 below), is therefore uncountable. This is why the previous works of local nature left open the possibility that a *global* convergence result may not hold since the uncountable union of measure zero sets may sum

³We defer the readers to Section A in Supplementary Material for a detailed description of the theorem.

up to a set of positive measure. We solve this issue here by proving that the union of all the basins of attraction associated with the saddle points is in fact contained in a codimension two subset of the state space \mathcal{X} . In the next section, a more advanced sketch of proof of Theorem 1 is provided. For the sake of readability and to avoid cumbersome technical details, only the case of $H = 1$ is elaborated. This proof provides the main arguments for the more technical analysis of the $H \geq 1$ scenario, available in an extended version of this article.

3 Detailed Analysis of the case $H = 1$

In this section, we provide a detailed proof of Theorem 1 in the case of a single-hidden-layer linear network (i.e., $H = 1$). To this end, we start with Proposition 1 which states the global convergence of gradient flows in (4) to critical points, with a polynomial convergence rate in the worst case. In the following, when state variables are concerned, we frequently drop the argument t for simplicity.

3.1 Global Convergence to Critical Points in GDD

Proof of Proposition 1 for $H = 1$. First note that for $H = 1$ the loss function L is a polynomial of degree four in the elements of Ξ . According to Lojasiewicz’s theorem, it is enough to show that every trajectory is bounded. To this end, we note from (4) that

$$\begin{cases} \frac{d\bar{\mathbf{W}}_1}{dt} = \bar{\mathbf{W}}_2^T \bar{\mathbf{M}} \\ \frac{d\bar{\mathbf{W}}_2}{dt} = \bar{\mathbf{M}} \bar{\mathbf{W}}_1^T \end{cases} \Rightarrow \frac{d \operatorname{tr} (\bar{\mathbf{W}}_1^T \bar{\mathbf{W}}_1 + \bar{\mathbf{W}}_2^T \bar{\mathbf{W}}_2)}{dt} = 4 \operatorname{tr} (\bar{\mathbf{W}}_1^T \bar{\mathbf{W}}_2^T \bar{\mathbf{M}}),$$

where we recall $\bar{\mathbf{M}} \equiv \Sigma_Y - \bar{\mathbf{W}}_2 \bar{\mathbf{W}}_1$ and therefore (since $\operatorname{tr} \mathbf{A} \mathbf{A}^T = \|\mathbf{A}\|_F^2$)

$$\begin{aligned} \frac{d (\|\bar{\mathbf{W}}_1\|_F^2 + \|\bar{\mathbf{W}}_2\|_F^2)}{dt} &= 4 \operatorname{tr} \bar{\mathbf{W}}_1^T \bar{\mathbf{W}}_2^T \Sigma_Y - 2 \operatorname{tr} \bar{\mathbf{W}}_1^T \bar{\mathbf{W}}_2^T \bar{\mathbf{W}}_2 \bar{\mathbf{W}}_1 - 2 \operatorname{tr} \bar{\mathbf{W}}_1^T \bar{\mathbf{W}}_2^T \bar{\mathbf{W}}_2 \bar{\mathbf{W}}_1 \\ &= 4 \operatorname{tr} \bar{\mathbf{W}}_1^T \bar{\mathbf{W}}_2^T \Sigma_Y - \operatorname{tr} \bar{\mathbf{W}}_2^T \bar{\mathbf{W}}_2 (\bar{\mathbf{W}}_2^T \bar{\mathbf{W}}_2 - \mathbf{C}^0) - \operatorname{tr} \bar{\mathbf{W}}_1 \bar{\mathbf{W}}_1^T (\bar{\mathbf{W}}_1 \bar{\mathbf{W}}_1^T + \mathbf{C}^0) - 2 \|\bar{\mathbf{W}}_2 \bar{\mathbf{W}}_1\|_F^2 \\ &\leq 4 \operatorname{tr} \bar{\mathbf{W}}_1^T \bar{\mathbf{W}}_2^T \Sigma_Y - \frac{1}{d_1} (\|\bar{\mathbf{W}}_1\|_F^4 + \|\bar{\mathbf{W}}_2\|_F^4) - 2 \|\bar{\mathbf{W}}_2 \bar{\mathbf{W}}_1\|_F^2 + \lambda_{\max}(\mathbf{C}^0) (\|\bar{\mathbf{W}}_1\|_F^2 + \|\bar{\mathbf{W}}_2\|_F^2) \\ &\leq c_1 (\|\bar{\mathbf{W}}_1\|_F^2 + \|\bar{\mathbf{W}}_2\|_F^2) - c_2 (\|\bar{\mathbf{W}}_1\|_F^2 + \|\bar{\mathbf{W}}_2\|_F^2)^2 \end{aligned}$$

for some $c_1, c_2 > 0$, where we used Cauchy–Schwarz inequality $(\operatorname{tr} \mathbf{A} \mathbf{A}^T)^2 \leq \operatorname{tr}(\mathbf{A} \mathbf{A}^T)^2 \cdot \operatorname{tr} \mathbf{I}$ along with $|\operatorname{tr} \mathbf{A} \mathbf{A}^T \mathbf{B}| \leq \lambda_{\max}(\mathbf{B}) \operatorname{tr} \mathbf{A} \mathbf{A}^T$. Setting $F \equiv \|\bar{\mathbf{W}}_1\|_F^2 + \|\bar{\mathbf{W}}_2\|_F^2$ the above inequality reads $\frac{dF}{dt} \leq c_1 F - c_2 F^2$ with $F \geq 0$ and hence the sum $\|\bar{\mathbf{W}}_1\|_F^2 + \|\bar{\mathbf{W}}_2\|_F^2$ is uniformly bounded for all $t \geq 0$. With Lemma 1 we know that the difference $\|\bar{\mathbf{W}}_2\|_F^2 - \|\bar{\mathbf{W}}_1\|_F^2$ is also uniformly bounded, which further leads to the boundedness of all trajectories of both $\bar{\mathbf{W}}_1$ and $\bar{\mathbf{W}}_2$. Since the trajectories of $\bar{\mathbf{W}}_1, \bar{\mathbf{W}}_2$ (and thus $\bar{\mathbf{M}}$) are uniformly bounded for all $t \geq 0$, the norm of the gradient $\|\Delta_\Xi\|_F$ as well as all trajectories in the GDD are bounded. The guaranteed rate of convergence can be obtained from estimates associated with polynomial gradient systems [8]. \square

3.2 Characterization of Critical Points

Proposition 1 ensures, for all initializations, the convergence of the gradient descent to a critical point, i.e., a point Ξ in the state space \mathcal{X} verifying $\Delta_\Xi(\xi) = 0$. Nonetheless, the information on the “quality” of the solution achieved by the algorithm is still missing. To obtain a clearer picture, we now focus on the set of *all* critical points by further decomposing the loss L with $\Sigma_Y \equiv [\mathbf{S}_Y \mid \mathbf{0}]$ for diagonal $\mathbf{S}_Y \in \mathbb{R}^{d_y \times d_y}$ with $[\mathbf{S}_Y]_{ii} > 0$ as

$$L(\bar{\mathbf{W}}_1, \bar{\mathbf{W}}_2) = \frac{1}{2} \|\Sigma_Y - \bar{\mathbf{W}}_2 \bar{\mathbf{W}}_1\|_F^2 = \frac{1}{2} \|\mathbf{S}_Y - \mathbf{C} \mathbf{A}\|_F^2 + \frac{1}{2} \|\mathbf{C} \mathbf{B}\|_F^2 \quad (5)$$

with $\mathbf{C} \equiv \bar{\mathbf{W}}_2 \in \mathbb{R}^{d_y \times d_H}$, $\mathbf{A} \in \mathbb{R}^{d_1 \times d_y}$ and $\mathbf{B} \in \mathbb{R}^{d_1 \times (d_x - d_y)}$ such that $[\mathbf{A} \mid \mathbf{B}] \equiv \bar{\mathbf{W}}_1$.

Under the notations above, we further expand $L(\Xi + \xi)$ to obtain its higher order variation as

$$L(\mathbf{A} + \mathbf{a}, \mathbf{B} + \mathbf{b}, \mathbf{C} + \mathbf{c}) \equiv L(\Xi + \xi) = L(\Xi) + \Delta_\Xi(\xi) + H_\Xi(\xi) + O(\|\xi\|^3)$$

with $\mathbf{M} \equiv \mathbf{S}_Y - \mathbf{C}\mathbf{A}$, $L(\Xi) = \frac{1}{2}\|\mathbf{M}\|_F^2 + \frac{1}{2}\|\mathbf{C}\mathbf{B}\|_F^2$ and

$$\begin{aligned}\Delta_{\Xi}(\xi) &\equiv -\text{tr}(\mathbf{M}^T(\mathbf{C}\mathbf{a} + \mathbf{c}\mathbf{A})) + \text{tr}(\mathbf{B}^T\mathbf{C}^T(\mathbf{C}\mathbf{b} + \mathbf{c}\mathbf{B})) \\ H_{\Xi}(\xi) &\equiv -\text{tr}(\mathbf{M}^T\mathbf{c}\mathbf{a}) + \frac{1}{2}\|\mathbf{C}\mathbf{a} + \mathbf{c}\mathbf{A}\|_F^2 + \text{tr}(\mathbf{B}^T\mathbf{C}^T\mathbf{c}\mathbf{b}) + \frac{1}{2}\|\mathbf{C}\mathbf{b} + \mathbf{c}\mathbf{B}\|_F^2 = O(\|\xi\|^2)\end{aligned}$$

that give the differential and the Hessian of L , respectively.

Recall that $\text{Crit}(L) \equiv \{\Xi \mid \Delta_{\Xi}(\xi) = 0\}$ and denote $\mathbf{M} \equiv \mathbf{S}_Y - \mathbf{C}\mathbf{A}$, so that, by Definition 1,

$$\begin{cases} \frac{d\mathbf{A}}{dt} &\equiv -\nabla_{\mathbf{A}}L(\Xi) = \mathbf{C}^T\mathbf{M} = \mathbf{0} \\ \frac{d\mathbf{B}}{dt} &\equiv -\nabla_{\mathbf{B}}L(\Xi) = -\mathbf{C}^T\mathbf{C}\mathbf{B} = \mathbf{0} \\ \frac{d\mathbf{C}}{dt} &\equiv -\nabla_{\mathbf{C}}L(\Xi) = \mathbf{M}\mathbf{A}^T - \mathbf{C}\mathbf{B}\mathbf{B}^T = \mathbf{0} \end{cases} \Leftrightarrow \begin{cases} \mathbf{C}^T\mathbf{S}_Y = \mathbf{C}^T\mathbf{C}\mathbf{A} \\ \mathbf{C}\mathbf{B} = \mathbf{0} \\ \mathbf{A}\mathbf{S}_Y = \mathbf{A}\mathbf{A}^T\mathbf{C}^T. \end{cases} \quad (6)$$

Observing the symmetric structure of \mathbf{A}, \mathbf{C} in (6) we have the following lemma.

Lemma 2 (Same Kernel for \mathbf{A} and \mathbf{C}^T). *Let Assumption 1 and 2 hold. Then for all $\Xi \in \text{Crit}(L)$,*

$$\text{Ker } \mathbf{A} = \text{Ker } \mathbf{C}^T, \text{ with } \text{Ker } \mathbf{A} \equiv \{\mathbf{x}, \mathbf{A}\mathbf{x} = \mathbf{0}\}.$$

Moreover, denote r the common rank of \mathbf{A} and \mathbf{C} with $0 < r \leq d_y$. Then there exists some orthogonal matrix $\mathbf{U} \in \mathbb{R}^{d_y \times d_y}$ such that

$$\begin{cases} \mathbf{A}\mathbf{U} = [\bar{\mathbf{A}} \mid \mathbf{0}_{d_1 \times (d_y-r)}] \\ \mathbf{C}^T\mathbf{U} = [\bar{\mathbf{C}}^T \mid \mathbf{0}_{d_1 \times (d_y-r)}] \\ \mathbf{U}^{-1}\mathbf{S}_Y\mathbf{U} = \mathbf{S}_Y \end{cases} \quad (7)$$

with $\bar{\mathbf{A}}, \bar{\mathbf{C}}^T \in \mathbb{R}^{d_1 \times r}$. Moreover, if \mathbf{S}_Y has distinct eigenvalues (i.e., $\bar{\mathbf{Y}}$ has d_y distinct singular values, as demanded in Assumption 3), then \mathbf{U} is a permutation matrix.

Sketch of proof. It can be shown with basic algebraic manipulations that the eigenvectors of \mathbf{S}_Y^2 (thus of \mathbf{S}_Y) form a basis of both $\text{Ker } \mathbf{A}$ and $\text{Ker } \mathbf{C}^T$. Therefore $\text{Ker } \mathbf{A} = \text{Ker } \mathbf{C}^T$ and in particular $\dim \text{Ker } \mathbf{A} = \dim \text{Ker } \mathbf{C}^T$. We denote this dimension $d_y - r$ and \mathbf{A}, \mathbf{C} are thus both of rank r . Choose \mathbf{U}_2 from $\text{Ker } \mathbf{A}$ and $\mathbf{U}_1 \perp \text{Ker } \mathbf{A}$; we deduce $\mathbf{U} = [\mathbf{U}_1 \mid \mathbf{U}_2]$ so that (7) holds. \square

Remark from (7) in Lemma 2 that, for arbitrary \mathbf{S}_Y , there are infinitely many possibilities on the choice of \mathbf{U} with the risk of occupying too much of the state space \mathcal{X} , since, with the change of variable in Lemma 2 the state variable now becomes the tuple $(\bar{\mathbf{A}}, \mathbf{B}, \bar{\mathbf{C}}, \mathbf{U})$. Using Assumption 3, \mathbf{U} only takes a finite number of values (the 2^{d_y} permutation matrices) for a given $\Xi \in \text{Crit}_r(L)$, hence the state variable essentially becomes the tuple $(\bar{\mathbf{A}}, \mathbf{B}, \bar{\mathbf{C}})$.

For $\Xi \in \text{Crit}(L)$ with \mathbf{A}, \mathbf{C} of rank r with $0 < r \leq d_y$, rewriting the diagonal \mathbf{S}_Y in two blocks

$\mathbf{S}_Y = \begin{bmatrix} \mathbf{D}_Y & \mathbf{0} \\ \mathbf{0} & \mathbf{E}_Y \end{bmatrix}$, with $\mathbf{D}_Y \in \mathbb{R}^{r \times r}$ and $\mathbf{E}_Y \in \mathbb{R}^{(d_y-r) \times (d_y-r)}$, together with Lemma 2, we simplify (6) as

$$\begin{cases} \bar{\mathbf{C}}\bar{\mathbf{A}} = \mathbf{D}_Y \\ \bar{\mathbf{C}}\mathbf{B} = \mathbf{0} \end{cases}, \quad \mathbf{M} = \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{E}_Y \end{bmatrix} \quad (8)$$

with the fact that $\bar{\mathbf{C}}^T, \bar{\mathbf{A}}$ are both of full rank (equal to r). The loss $L(\Xi)$ (at critical points) can thus be simplified as $L(\Xi) = \frac{1}{2}\|\mathbf{E}_Y\|_F^2$ where \mathbf{E}_Y measures the “quality” of each critical points.

For any $\Xi \in \text{Crit}(L)$, with Lemma 2 we are allowed to “extract” the full rank (sub-)structures of \mathbf{A}, \mathbf{C} with \mathbf{S}_Y unchanged, via a simple change of basis. For $1 \leq r \leq d_y$, let $\text{Crit}_r(L)$ be the subset of $\text{Crit}(L)$ such that the rank of \mathbf{A} and of \mathbf{C} is equal to r . Then, one has the following disjoint union

$$\text{Crit}(L) = \cup_{r=1}^{d_y} \text{Crit}_r(L).$$

This precise characterization of critical points naturally leads to the following proposition on the loss function $L(\cdot)$, that can be further “visualized” as in Figure 2.

Proposition 2 (Landscape of Single-hidden-layer Linear Network). Under Assumptions 1-3, the loss function $L(\Xi)$ has the following properties:

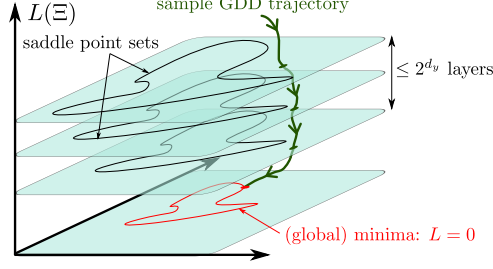


Figure 2: A geometric “vision” of the loss landscape.

- i) The set of possible limits of L along the GDD given by (3) is equal to the finite set made of the sum of the squares of any subset of the singular values of $\bar{\mathbf{Y}}$.
- ii) The set $\text{Crit}_{d_y}(L)$ is in fact the set of local (and global) minima, with $L = 0$ and $\mathbf{M} = \mathbf{0}$.
- iii) Every critical point in $\text{Crit}_r(L)$ with $1 \leq r \leq d_y - 1$ is a saddle point such that the Hessian has at least one negative eigenvalue. In particular, the set of saddle points is an algebraic variety of positive dimension, i.e., (up to a permutation matrix) the zero set of the polynomial functions given in (8), with $\mathbf{E}_{\mathbf{Y}} \neq \mathbf{0}$.

Proof. Item i) follows directly from the discussion preceding the proposition. As regard Item ii), we write the Hessian for a given $\Xi \in \text{Crit}(L)$ as $H_{\Xi}(\xi) = -\text{tr}(\mathbf{M}^T \mathbf{c} \mathbf{a}) + \frac{1}{2} \|\mathbf{C} \mathbf{a} + \mathbf{c} \mathbf{A}\|_F^2 + \frac{1}{2} \|\mathbf{C} \mathbf{b} + \mathbf{c} \mathbf{B}\|_F^2$, with $\Xi \equiv (\mathbf{A}, \mathbf{B}, \mathbf{C})$ such that (6) is satisfied. The fact that no critical point is a *local maximum* is easily checked by taking $\xi \equiv (\mathbf{a}, \mathbf{b}, \mathbf{c})$ such that $\mathbf{C} \mathbf{b} + \mathbf{c} \mathbf{B} \neq \mathbf{0}$ and $\mathbf{a} \mathbf{M}^T = \mathbf{0}$. Therefore, the Hessian has (strictly) positive eigenvalues and Ξ is not local maximum. Such ξ always exists for at least one of \mathbf{B}, \mathbf{C} away from $\mathbf{0}$, while the case $\mathbf{B} = \mathbf{C} = \mathbf{0}$ leads to $L(\Xi) = \frac{1}{2} \|\mathbf{M}\|_F^2 = \frac{1}{2} \|\mathbf{S}_{\mathbf{Y}}\|^2$ that is not of practical interest.

Moreover, note from Lemma 2 that, for Ξ with \mathbf{A}, \mathbf{C} both of full rank ($r = d_y$) we have $\mathbf{M} = \mathbf{0}$ resulting in equivalent global minima. Showing Item iii) amounts to prove that a critical point Ξ with $\mathbf{M} \neq \mathbf{0}$ is in fact a saddle point. For that purpose, we first take \mathbf{b} such that $\mathbf{C} \mathbf{b} + \mathbf{c} \mathbf{B} = \mathbf{0}$ and it then remains to show that there always exists a pair (\mathbf{a}, \mathbf{c}) such that the Hessian has *at least* one negative eigenvalue (i.e., $H_{\Xi}(\xi) < 0$, for $\xi \equiv (\mathbf{a}, \mathbf{b}, \mathbf{c})$). To this end, we hope to find (\mathbf{a}, \mathbf{c}) so that $\|\mathbf{M}^T \mathbf{c}\|_F > \|\mathbf{C} \mathbf{c}^T\|_F$ (for example $\mathbf{C} \mathbf{c}^T = \mathbf{0}$ while $\mathbf{M}^T \mathbf{c} \neq \mathbf{0}$). Taking $\mathbf{a} = c_0 \mathbf{c}^T \mathbf{M}$ results in $H_{\Xi}(\xi) = -c_0 \text{tr}(\mathbf{M}^T \mathbf{c} \mathbf{c}^T \mathbf{M}) + \frac{1}{2} \|c_0 \mathbf{C} \mathbf{c}^T \mathbf{M} + \mathbf{c} \mathbf{A}\|_F^2 = -c_0 \|\mathbf{c}^T \mathbf{M}\|_F^2 + \frac{1}{2} \|\mathbf{c} \mathbf{A}\|_F^2$, where we recall from (6) that $\mathbf{C}^T \mathbf{M} = \mathbf{0}$. If there exists such a \mathbf{c} we can *always* find positive c_0 such that $H_{\Xi}(\xi) < 0$. The existence of \mathbf{c} is guaranteed by the fact that $\mathbf{M} \neq \mathbf{0}$.⁴ This further indicates that all $\Xi \in \text{Crit}(L)$ with rank deficient \mathbf{A}, \mathbf{C} ($r \leq d_y - 1$) are in fact saddle points. \square

The fact that all local minima are equivalently global minima and all critical points that are not global minima are saddle points is in fact already known for single-hidden-layer linear networks [2] as well as for deep linear networks [13]. Here we provided an alternative and shorter proof.

3.3 Convergence to Global Minima for Almost All Initializations

Having characterized the critical points, we now show that the GDD almost always converges to a local (and thus global) minimum, thereby completing the proof of Theorem 1.

End of proof of Theorem 1. To complete the proof of Theorem 1 (for $H = 1$) it suffices to evaluate, for $r = 1, \dots, d_y - 1$: 1) D_C^r , an upper bound of the dimension of $\text{Crit}_r(L)$ and 2) U_S^r , a lower bound for a given $\Xi \in \text{Crit}_r(L)$ of the dimension of $\mathcal{C}(\Xi)$, the linear span of “variations” at Ξ , whose corresponding trajectories *do not converge* to Ξ . This measures the “possibility” for the GDD to escape from the saddle point Ξ .

⁴With the same change of basis as in Lemma 2, see for details in Proof 1 of Supplementary Material.

Then, by standard transversality arguments [10], one has 1) the basin of attraction of any $\Xi \in \text{Crit}_r(L)$ is contained in a set of dimension $D_S^r \leq \dim \mathcal{X} - U_S^r$, and 2) the set of initializations with corresponding trajectories of the GDD converge to an element of $\text{Crit}_r(L)$ is contained in a subset of \mathcal{X} of dimension $D_C^r + D_S^r$. We show $D_C^r + D_S^r \leq \dim \mathcal{X} - 2$ for $r = 1, \dots, d_y - 1$ and hence the conclusion. We refer the readers to Proof 2 in Supplementary Material for a detailed exposition. \square

As stated in Proposition 1 and Theorem 1, GDD achieves at least a polynomial convergence rate [8] to a global minimum (for almost all initializations). As a side and immediate aftermath, it can be shown that, upon proper initialization, exponential convergence can be achieved (here for $H = 1$).

Remark 1 (Exponential Convergence of GDD). *Let Assumptions 1 and 2 hold. Then, every trajectory of the GDD such that $\mathbf{C}^0 \equiv (\overline{\mathbf{W}}_2^\top \overline{\mathbf{W}}_2 - \overline{\mathbf{W}}_1 \overline{\mathbf{W}}_1^\top)|_{t=0}$ has at least d_y strictly positive eigenvalues, converges to a global minimum at the rate of $e^{-2\alpha t}$ with α the d_y -th smallest eigenvalue of \mathbf{C}^0 .*

Proof. Recalling that for $H = 1$ we have $L = \frac{1}{2} \|\Sigma_Y - \overline{\mathbf{W}}_2 \overline{\mathbf{W}}_1\|_F^2$, with (4) we deduce

$$\frac{d\|\overline{\mathbf{M}}\|_F^2}{dt} = \frac{d\text{tr}(\overline{\mathbf{M}}^\top \overline{\mathbf{M}})}{dt} = -2\text{tr}(\overline{\mathbf{M}}^\top \overline{\mathbf{M}} \overline{\mathbf{W}}_1^\top \overline{\mathbf{W}}_1 + \overline{\mathbf{W}}_2 \overline{\mathbf{W}}_2^\top \overline{\mathbf{M}} \overline{\mathbf{M}}^\top) \leq -2c_0 \|\overline{\mathbf{M}}\|_F^2 \quad (9)$$

with $\overline{\mathbf{M}} \equiv \Sigma_Y - \overline{\mathbf{W}}_2 \overline{\mathbf{W}}_1$ and $c_0 = \lambda_{\min}(\overline{\mathbf{W}}_1^\top \overline{\mathbf{W}}_1) + \lambda_{\min}(\overline{\mathbf{W}}_2 \overline{\mathbf{W}}_2^\top)$. Since $\overline{\mathbf{W}}_1^\top \overline{\mathbf{W}}_1 \in \mathbb{R}^{d_x \times d_x}$ is of maximum rank d_1 (with $d_1 \leq d_x$ from Assumption 1), we have $\lambda_{\min}(\overline{\mathbf{W}}_1^\top \overline{\mathbf{W}}_1) = 0$. Nonetheless, $\overline{\mathbf{W}}_2 \overline{\mathbf{W}}_2^\top \in \mathbb{R}^{d_y \times d_y}$ may be of full rank so that $\lambda_{\min}(\overline{\mathbf{W}}_2 \overline{\mathbf{W}}_2^\top) > 0$. To this end, we decompose $\overline{\mathbf{W}}_2 = [\overline{\mathbf{W}}_{21} \mid \overline{\mathbf{W}}_{22}]$, with $\overline{\mathbf{W}}_{21} \in \mathbb{R}^{d_y \times d_y}$. Then with the inclusion principle of Hermitian matrices (e.g., Theorem 4.3.28 in [11]) we deduce $\lambda_{\min}(\overline{\mathbf{W}}_{21}^\top \overline{\mathbf{W}}_{21}) \geq \lambda_{d_y}(\overline{\mathbf{W}}_2 \overline{\mathbf{W}}_2^\top)$. Moreover, since $\lambda_{\min}(\overline{\mathbf{W}}_{21}^\top \overline{\mathbf{W}}_{21}) = \lambda_{\min}(\overline{\mathbf{W}}_{21} \overline{\mathbf{W}}_{21}^\top)$, by Lemma 1 and Weyl's inequality (e.g., Corollary 4.3.12 in [11]), we have $\lambda_{\min}(\overline{\mathbf{W}}_2 \overline{\mathbf{W}}_2^\top) \geq \lambda_{\min}(\overline{\mathbf{W}}_{21} \overline{\mathbf{W}}_{21}^\top) \geq \lambda_{d_y}(\overline{\mathbf{W}}_2 \overline{\mathbf{W}}_2^\top) \geq \lambda_{d_y}(\mathbf{C}^0)$. As such, (9) yields $\frac{d\text{tr}(\overline{\mathbf{M}}^\top \overline{\mathbf{M}})}{dt} \leq -2\lambda_{d_y}(\mathbf{C}^0) \text{tr}(\overline{\mathbf{M}}^\top \overline{\mathbf{M}})$ which concludes the proof. \square

4 Concluding Remarks

To the best of the authors' knowledge, it is the first time that the *global* behavior of the gradient descent dynamics in linear neural networks is fully characterized, in the sense that we show a *global* convergence to critical points of all trajectories of the gradient flow via Łojasiewicz's theorem, which helps eliminate the possibility of divergence. Then with a fine local study of critical points we exclude the (possible) worries concerning the "accumulation" of saddle points together with associated basin of attractions so that they form "disjoint layers" that are of total measure zero in the total weight space. Interestingly, Łojasiewicz's theorem is more powerful than needed here and may enable extensions of the present results to more advanced dynamics than the simple GDD (see Remark 2 in Supplementary Material for more details).

It is interesting to note that the authors in [9, 6], made a strong case to warn against saddle points in deep learning, which is in sharp contrast with our conclusions. Yet, the analysis in [9, 6] is asymptotic in the network dimensions, where the present one is set for fixed network sizes. It would be of interest to conciliate both results to gain a even clearer picture of deep linear learning in practical scenarios.

When nonlinear networks are considered, obtaining an equivalent version of Lemma 1 would be a key enabler to achieve the global convergence to critical points as per Łojasiewicz's theorem and therefore would allow for a better understanding of the *nonlinear* deep networks performance. Exploring a random model setting for \mathbf{X}, \mathbf{Y} , the authors in [6] argue that the loss surfaces of these networks loosely recall (yet is formally quite different from) a spin-glass model, familiar to statistical physicists. In this case, as the network gets large, local minima gather in a thin "band" of similar losses isolated from the global minimum. Stating that the number of local minima outside that band diminishes exponentially with the size of the network, the authors argue that the gradient descent dynamics (in their case the stochastic gradient descent dynamics) converges to this band and therefore leads to deep *nonlinear* networks with good generalization performance. Taking advantage of a random nature for \mathbf{X}, \mathbf{Y} in our present setting would allow for a refinement of our proposed geometric vision, likely by means of a "statistical extension" of the key Lemma 1.

References

- [1] Pierre-Antoine Absil, Robert Mahony, and Benjamin Andrews. Convergence of the iterates of descent methods for analytic cost functions. *SIAM Journal on Optimization*, 16(2):531–547, 2005.
- [2] Pierre Baldi and Kurt Hornik. Neural networks and principal component analysis: Learning from examples without local minima. *Neural networks*, 2(1):53–58, 1989.
- [3] Avrim Blum and Ronald L Rivest. Training a 3-node neural network is np-complete. In *Advances in neural information processing systems*, pages 494–501, 1989.
- [4] Joseph-Frédéric Bonnans, Jean Charles Gilbert, Claude Lemaréchal, and Claudia A Sagastizábal. *Numerical optimization: theoretical and practical aspects*. Springer Science & Business Media, 2006.
- [5] Stephen Boyd and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- [6] Anna Choromanska, Mikael Henaff, Michael Mathieu, Gérard Ben Arous, and Yann LeCun. The loss surfaces of multilayer networks. In *Artificial Intelligence and Statistics*, pages 192–204, 2015.
- [7] Ronan Collobert and Jason Weston. A unified architecture for natural language processing: Deep neural networks with multitask learning. In *Proceedings of the 25th international conference on Machine learning*, pages 160–167. ACM, 2008.
- [8] Didier D’Acunto and Krzysztof Kurdyka. Explicit bounds for the Łojasiewicz exponent in the gradient inequality for polynomials. In *Annales Polonici Mathematici*, volume 1, pages 51–61, 2005.
- [9] Yann N Dauphin, Razvan Pascanu, Caglar Gulcehre, Kyunghyun Cho, Surya Ganguli, and Yoshua Bengio. Identifying and attacking the saddle point problem in high-dimensional non-convex optimization. In *Advances in neural information processing systems*, pages 2933–2941, 2014.
- [10] Mark Goresky and Robert MacPherson. Stratified morse theory. In *Stratified Morse Theory*, pages 3–22. Springer, 1988.
- [11] Roger A Horn and Charles R Johnson. *Matrix analysis*. Cambridge university press, 1990.
- [12] Chi Jin, Rong Ge, Praneeth Netrapalli, Sham M Kakade, and Michael I Jordan. How to escape saddle points efficiently. *arXiv preprint arXiv:1703.00887*, 2017.
- [13] Kenji Kawaguchi. Deep learning without poor local minima. In *Advances in Neural Information Processing Systems*, pages 586–594, 2016.
- [14] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [15] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [16] Jason D Lee, Ioannis Panageas, Georgios Piliouras, Max Simchowitz, Michael I Jordan, and Benjamin Recht. First-order methods almost always avoid saddle points. *arXiv preprint arXiv:1710.07406*, 2017.
- [17] S Łojasiewicz. Sur les trajectoires du gradient d’une fonction analytique. *Seminari di geometria*, 1983:115–117, 1982.
- [18] Stanisław Łojasiewicz. Ensembles semi-analytiques. *Lectures Notes IHES (Bures-sur-Yvette)*, 1965.

- [19] Abdel-rahman Mohamed, George E Dahl, and Geoffrey Hinton. Acoustic modeling using deep belief networks. *IEEE Transactions on Audio, Speech, and Language Processing*, 20(1):14–22, 2012.
- [20] Katta G Murty and Santosh N Kabadi. Some np-complete problems in quadratic and nonlinear programming. *Mathematical programming*, 39(2):117–129, 1987.
- [21] Ning Qian. On the momentum term in gradient descent learning algorithms. *Neural networks*, 12(1):145–151, 1999.
- [22] Andrew M Saxe, James L McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *arXiv preprint arXiv:1312.6120*, 2013.
- [23] Ju Sun, Qing Qu, and John Wright. When are nonconvex problems not scary? *arXiv preprint arXiv:1510.06096*, 2015.
- [24] Gerald Teschl. *Ordinary differential equations and dynamical systems*, volume 140. American Mathematical Society Providence, 2012.
- [25] Willard I Zangwill. Convergence conditions for nonlinear programming algorithms. *Management Science*, 16(1):1–13, 1969.

Supplementary Material

Almost Global Convergence to Global Minima for Gradient Descent in Deep Linear Networks

A Łojasiewicz’s theorem

We first recall Łojasiewicz’s theorem for the convergence of real analytic gradient flows, which is essentially the key enabler to prove the global convergence of the GDD trajectories.

Theorem 2 (Łojasiewicz’s theorem, [17]). *Let L be a real analytic function and let $\Xi(\cdot)$ be a solution trajectory of the gradient system given by Definition 1. Further assume that $\sup_{t \geq 0} \|\Xi(t)\| < \infty$. Then $\Xi(\cdot)$ converges to a critical point of L , as $t \rightarrow \infty$.⁵*

Remark 2. Since the fundamental (strict) gradient descent direction (as in Definition 1) in Łojasiewicz’s theorem can in fact be relaxed to a (more general) angle condition (see for example Theorem 2.2 in [1]), the line of argument developed in the core of the article may be similarly followed to prove the global convergence of more advanced optimizers (e.g., SGD, SGD-Momentum [21], ADAM [14], etc.), for which the direction of descent is not strictly the opposite of the gradient direction. This constitutes an important direction of future exploration.

B Proofs

In this section we provide detailed proofs of Proposition 2 and Theorem 1, in the following Proofs 1 and 2 respectively.

Proof 1. (*Complementary Proof of Proposition 2*). The existence of a matrix \mathbf{c} satisfying $\|\mathbf{M}^T \mathbf{c}\|_F > \|\mathbf{C} \mathbf{c}^T\|_F$ (for example $\mathbf{C} \mathbf{c}^T = \mathbf{0}$ while $\mathbf{M}^T \mathbf{c} \neq \mathbf{0}$, so that $H_{\Xi}(\xi) < 0$) is guaranteed by the fact that $\mathbf{M} \neq \mathbf{0}$, where we recall that $\mathbf{M} = \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{E}_Y \end{bmatrix}$. More concretely we would like to find \mathbf{c} such that

$$\begin{cases} \mathbf{C}^T \mathbf{M} = \mathbf{0} \\ \mathbf{C} \mathbf{c}^T = \mathbf{0} \\ \mathbf{c}^T \mathbf{M} \neq \mathbf{0} \end{cases} \quad (10)$$

are satisfied. The existence of \mathbf{c} is guaranteed by Lemma 2 as

$$\begin{cases} \mathbf{C}^T \mathbf{M} = \mathbf{C}^T \mathbf{U} \mathbf{U}^T \mathbf{M} \mathbf{U} \mathbf{U}^T = \begin{bmatrix} \overline{\mathbf{C}}^T & \mathbf{0}_{d_1 \times (d_y - r)} \end{bmatrix} \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{E}_Y \end{bmatrix} \mathbf{U}^T = \mathbf{0} \\ \mathbf{C} \mathbf{c}^T = \mathbf{U} \mathbf{U}^T \mathbf{C} \mathbf{c}^T \mathbf{U} \mathbf{U}^T = \mathbf{U} \begin{bmatrix} \overline{\mathbf{C}} \\ \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{c}_1^T & \mathbf{c}_2^T \end{bmatrix} \mathbf{U}^T = \mathbf{U} \begin{bmatrix} \overline{\mathbf{C}} \mathbf{c}_1^T & \overline{\mathbf{C}} \mathbf{c}_2^T \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{U}^T = \mathbf{0} \\ \mathbf{c}^T \mathbf{M} = \mathbf{c}^T \mathbf{U} \mathbf{U}^T \mathbf{M} \mathbf{U} \mathbf{U}^T = \begin{bmatrix} \mathbf{c}_1^T & \mathbf{c}_2^T \end{bmatrix} \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{E}_Y \end{bmatrix} \mathbf{U}^T = \begin{bmatrix} \mathbf{0} & \mathbf{c}_2^T \mathbf{E}_Y \end{bmatrix} \mathbf{U}^T \neq \mathbf{0} \end{cases}$$

with $\mathbf{c}^T \mathbf{U} = \begin{bmatrix} \mathbf{c}_1^T & \mathbf{c}_2^T \end{bmatrix}$. To fulfill (10) it suffices to take $\mathbf{c}_1 = \mathbf{0}$ and $\mathbf{c}_2^T \in \text{Ker } \overline{\mathbf{C}}$ with $\mathbf{c}_2^T \mathbf{E}_Y \neq \mathbf{0}$, which is possible since by definition $\overline{\mathbf{C}}^T \in \mathbb{R}^{d_1 \times r}$ is of full rank $r \leq d_1$, together with $\mathbf{E}_Y \neq \mathbf{0}$. \square

Proof 2. (*Proof of Theorem 1*). Recall that to complete the proof of Theorem 1, it remains to determine, for $1 \leq r \leq d_y - 1$ the following two quantities:

- 1) D_C^r , an upper bound of the dimension of $\text{Crit}_r(L)$,
- 2) U_S^r , a lower bound for any given $\Xi \in \text{Crit}_r(L)$ of the dimension of $\mathcal{C}(\Xi)$, the linear span of the non-converging “variations” around Ξ .

We start by counting the dimension of $\text{Crit}_r(L)$ for $1 \leq r \leq d_y - 1$ under the above setting, as described in the following lemma.

⁵This theorem is based on the fundamental Łojasiewicz’s inequality of analytic functions [18].

Lemma 3 (Upper Bound for the Dimension of $\text{Crit}_r(L)$). *Under Assumption 1-3, an upper bound for the dimension of the subset of critical points $\text{Crit}_r(L)$ for $1 \leq r \leq d_y - 1$ is $2d_1r - r^2 + (d_1 - r)(d_x - d_y)$.*

Proof. As discussed in Section 3.2, $\text{Crit}_r(L)$ is fully characterized by (8). The space described by (8) is in fact an algebraic variety (i.e., the set of zeros of a polynomial of degree three) made of a finite number of two by two disjoint smooth strata [10]. One can therefore attach a dimension to this algebraic variety as the largest dimension of each smooth stratum (the latter integer defined as the standard dimension of a differentiable manifold [10]). These strata are characterized according to the dimension of the span of the columns of $\bar{\mathbf{A}}, \bar{\mathbf{B}}, \bar{\mathbf{C}}^\top$. A lengthy but immediate computation shows that the largest dimension among these strata is obtained in the case where the columns of $\bar{\mathbf{A}}, \bar{\mathbf{B}}, \bar{\mathbf{C}}^\top$ are linearly independent. Hence, one can compute the upper bound D_C^r of $\dim \text{Crit}_r(L)$ and obtain that

$$D_C^r = \dim \bar{\mathbf{A}} + \dim \bar{\mathbf{C}} + \dim \bar{\mathbf{B}} - \dim \{\text{constraints in (8)}\} = 2d_1r - r^2 + (d_1 - r)(d_x - d_y). \quad (11)$$

□

Given a saddle point $\Xi \equiv (\mathbf{A}, \mathbf{B}, \mathbf{C}) \in \text{Crit}_r(L)$ with $1 \leq r \leq d_y - 1$, we provide in the next lemma a lower bound U_S^r for the dimension of $\mathcal{C}(\Xi)$, the linear span of “variations” at Ξ , for which the corresponding trajectories *do not converge* to Ξ .

Lemma 4 (Lower Bound for the Dimension of $\mathcal{C}(\Xi)$). *Under Assumption 1-3, for $\Xi \in \text{Crit}_r(L)$ with $1 \leq r \leq d_y - 1$, the dimension of $\mathcal{C}(\Xi)$ is larger than or equal to*

$$U_S^r = (d_1 - r)(d_x - d_y) + 2rd_1 - r^2 + (d_y - r)(d_1 - r) + 1. \quad (12)$$

Proof. We first note from (6) that, for any $\Xi \equiv (\mathbf{A}, \mathbf{B}, \mathbf{C}) \in \text{Crit}_r(L)$, all variation $\xi \equiv (\mathbf{a}, \mathbf{b}, \mathbf{c})$ that satisfies

$$\begin{aligned} (\mathbf{C}^\top + \mathbf{c}^\top)\mathbf{S}_\mathbf{Y} &= (\mathbf{C}^\top + \mathbf{c}^\top)(\mathbf{C} + \mathbf{c})(\mathbf{A} + \mathbf{a}) + O(\|\xi\|^2) \\ (\mathbf{C} + \mathbf{c})(\mathbf{B} + \mathbf{b}) &= O(\|\xi\|^2) \\ (\mathbf{A} + \mathbf{a})\mathbf{S}_\mathbf{Y} &= (\mathbf{A} + \mathbf{a})(\mathbf{A}^\top + \mathbf{a}^\top)(\mathbf{C}^\top + \mathbf{c}^\top) + O(\|\xi\|^2) \end{aligned}$$

leads to other critical points and thus shall be included in counting $\dim \mathcal{C}(\Xi)$. With (6), the above relations further simplify as

$$\begin{cases} \mathbf{c}^\top \mathbf{M} = \mathbf{C}^\top (\mathbf{c}\mathbf{A} + \mathbf{C}\mathbf{a}) \\ \mathbf{C}\mathbf{b} + \mathbf{c}\mathbf{B} = \mathbf{0} \\ \mathbf{a}\mathbf{M}^\top = \mathbf{A}(\mathbf{c}\mathbf{A} + \mathbf{C}\mathbf{a})^\top \end{cases} \quad (13)$$

where the higher order terms are removed. These equations actually represent the tangent space of the set of variations ξ yielding other critical points in a neighborhood of Ξ . The subspace of solutions of the second equation of (13) with fixed \mathbf{c} is of dimension $(d_1 - r)(d_x - d_y)$ since it is the kernel of the linear map $\mathbf{b} \mapsto \mathbf{C}\mathbf{b}$ for any given \mathbf{C} of rank r .

For the first and third equations of (13), with the change of basis from Lemma 2 we write

$$\mathbf{a}\mathbf{U} = [\mathbf{a}_1 \mid \mathbf{a}_2], \quad \mathbf{c}^\top \mathbf{U} = [\mathbf{c}_1^\top \mid \mathbf{c}_2^\top]$$

with $\mathbf{a}_1 \in \mathbb{R}^{d_1 \times r}$, $\mathbf{a}_2 \in \mathbb{R}^{d_1 \times (d_y - r)}$, $\mathbf{c}_1 \in \mathbb{R}^{r \times d_1}$ and $\mathbf{c}_2 \in \mathbb{R}^{(d_y - r) \times d_1}$, which further leads to the following equivalent linear system,

$$\begin{cases} \bar{\mathbf{C}}\mathbf{a}_1 + \mathbf{c}_1\bar{\mathbf{A}} = \mathbf{0} \\ \mathbf{c}_2^\top \mathbf{E}_\mathbf{Y}^2 = \bar{\mathbf{C}}^\top \bar{\mathbf{C}} \bar{\mathbf{A}} \bar{\mathbf{A}}^\top \mathbf{c}_2^\top \\ \mathbf{a}_2 = \bar{\mathbf{A}} \bar{\mathbf{A}}^\top \mathbf{c}_2^\top \mathbf{E}_\mathbf{Y}^{-1} \end{cases} \quad (14)$$

where we use the fact that $\mathbf{E}_\mathbf{Y}$ is invertible.

As a consequence we resort to counting the dimension of the kernel of $(\mathbf{a}_1, \mathbf{a}_2, \mathbf{c}_1, \mathbf{c}_2) \mapsto$ the linear system in (14). More concretely, the dimension of the kernel of $(\mathbf{a}_1, \mathbf{c}_1) \mapsto \bar{\mathbf{C}}\mathbf{a}_1 + \mathbf{c}_1\bar{\mathbf{A}}$ is equal to $2rd_1 - r^2$, while the dimension of the kernel of $\mathbf{c}_2 \mapsto \mathbf{c}_2^\top \mathbf{E}_\mathbf{Y}^2 - \bar{\mathbf{C}}^\top \bar{\mathbf{C}} \bar{\mathbf{A}} \bar{\mathbf{A}}^\top \mathbf{c}_2^\top$ gives $(d_y - r)(d_1 - r)$, since the matrix $\bar{\mathbf{C}}^\top \bar{\mathbf{C}} \bar{\mathbf{A}} \bar{\mathbf{A}}^\top \in \mathbb{R}^{d_1 \times d_1}$ is of rank r .

Finally, note from Item *iii*) of Proposition 2 that there always exists *at least* one negative eigenvalue for all $\Xi \in \text{Crit}_r(L)$ with $1 \leq r \leq d_y - 1$, the dimension of $\mathcal{C}(\Xi)$ is therefore at least larger than or equal to $(d_1 - r)(d_x - d_y) + 2rd_1 - r^2 + (d_y - r)(d_1 - r) + 1$, which completes the proof. \square

By adding D_C^r and $\dim \mathcal{X} - U_S^r$, with D_C^r and U_S^r given in (11) and (12) respectively, we obtain an upper bound for the dimension of a submanifold S_r of the state space \mathcal{X} containing all initializations, the GDD trajectories of which converge to an element of $\text{Crit}_r(L)$ for $1 \leq r \leq d_y - 1$. This dimension is equal to

$$\begin{aligned} & d_1(d_x + d_y) - (d_1 - r)(d_x - d_y) - 2rd_1 + r^2 - (d_y - r)(d_1 - r) - 1 + 2d_1r - r^2 \\ & + (d_1 - r)(d_x - d_y) = d_1(d_x + d_y) - (d_y - r)(d_1 - r) - 1 \leq d_1(d_x + d_y) - 2 \end{aligned}$$

since by Assumption 1 we have $1 \leq (d_y - r) \leq (d_1 - r)$. Recalling for $H = 1$ we have $\dim \mathcal{X} = d_1(d_x + d_y)$ and we thus deduce that the aforementioned submanifold S_r is of codimension greater than or equal to two. Since there is a finite number of submanifolds S_r , their union for $r = 1, \dots, d_y - 1$ is also of codimension two. This concludes the proof of Theorem 1. \square