

# A UNIFYING VIEW ON IMPLICIT BIAS IN TRAINING LINEAR NEURAL NETWORKS

**Chulhee Yun\***

MIT

chulheey@mit.edu

**Shankar Krishnan**

Google Research

skrishnan@google.com

**Hossein Mobahi**

Google Research

hmobahi@google.com

## ABSTRACT

We study the implicit bias of gradient flow (i.e., gradient descent with infinitesimal step size) on linear neural network training. We propose a tensor formulation of neural networks that includes fully-connected, diagonal, and convolutional networks as special cases, and investigate the linear version of the formulation called linear tensor networks. With this formulation, we can characterize the convergence direction of the network parameters as singular vectors of a tensor defined by the network. For  $L$ -layer linear tensor networks that are orthogonally decomposable, we show that gradient flow on separable classification finds a stationary point of the  $\ell_{2/L}$  max-margin problem in a “transformed” input space defined by the network. For underdetermined regression, we prove that gradient flow finds a global minimum which minimizes a norm-like function that interpolates between weighted  $\ell_1$  and  $\ell_2$  norms in the transformed input space. Our theorems subsume existing results in the literature while removing standard convergence assumptions. We also provide experiments that corroborate our analysis.

## 1 INTRODUCTION

Overparametrized neural networks have infinitely many solutions that achieve zero training error, and such global minima have different generalization performance. Moreover, training a neural network is a high-dimensional nonconvex problem, which is typically intractable to solve. However, the success of deep learning indicates that first-order methods such as gradient descent or stochastic gradient descent (GD/SGD) not only (a) succeed in finding global minima, but also (b) are biased towards solutions that generalize well, which largely has remained a mystery in the literature.

To explain part (a) of the phenomenon, there is a growing literature studying the convergence of GD/SGD on overparametrized neural networks (e.g., Du et al. (2018a;b); Allen-Zhu et al. (2018); Zou et al. (2018); Jacot et al. (2018); Oymak & Soltanolkotabi (2020), and many more). There are also convergence results that focus on linear networks, without nonlinear activations (Bartlett et al., 2018; Arora et al., 2019a; Wu et al., 2019; Du & Hu, 2019; Hu et al., 2020). These results typically focus on the convergence of loss, hence do not address *which* of the many global minima is reached.

Another line of results tackles part (b), by studying the implicit bias or regularization of gradient-based methods on neural networks or related problems (Gunasekar et al., 2017; 2018a;b; Arora et al., 2018; Soudry et al., 2018; Ji & Telgarsky, 2019a; Arora et al., 2019b; Woodworth et al., 2020; Chizat & Bach, 2020; Gissin et al., 2020). These results have shown interesting progress that even without explicit regularization terms in the training objective, algorithms such as GD applied on neural networks have an *implicit bias* towards certain solutions among the many global minima. However, in proving such results, many results rely on *convergence assumptions* such as global convergence of loss to zero and/or directional convergence of parameters and gradients. Ideally, such convergence assumptions should be removed because they cannot be tested *a priori* and there are known examples where GD does not converge to global minima under certain initializations (Bartlett et al., 2018; Arora et al., 2019a).

---

\*Based on work performed during internship at Google Research

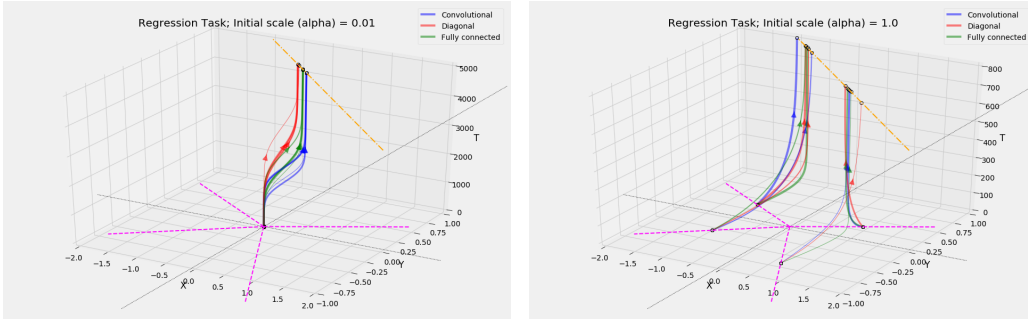


Figure 1: Gradient descent trajectories of linear coefficients of linear fully-connected, diagonal, and convolutional networks on a regression task, initialized with different initial scales  $\alpha = 0.01, 1$ . Networks are initialized at the same coefficients (circles on purple lines), but follow different trajectories due to implicit biases of networks induced from their architecture. The figures show that our theoretical predictions on limit points (circles on yellow line, the set of global minima) agree with the solution found by GD. For details of the experimental setup, see Section 6.

### 1.1 SUMMARY OF OUR CONTRIBUTIONS

We study the implicit bias of gradient flow (GD with infinitesimal step size) on linear neural networks. Following recent progress on this topic, we consider classification and regression problems that have multiple solutions with zero training error. Our analyses apply to a **general class of networks**, and prove **both convergence and implicit bias**, providing a more *complete characterization* of the algorithm trajectory without relying on convergence assumptions.

- We propose a general *tensor formulation* of nonlinear neural networks which includes many network architectures considered in the literature. In this paper, we focus on the linear version of this formulation (i.e., no nonlinear activations), called *linear tensor networks*.
- For linearly separable classification, we prove that linear tensor network parameters converge in direction to *singular vectors* of a tensor defined by the network. As a *corollary*, we show that linear fully-connected networks converge to the  $\ell_2$  max-margin solution (Ji & Telgarsky, 2020).
- For separable classification, we further show that if the linear tensor network is orthogonally decomposable (Assumption 1), the gradient flow finds the  $\ell_{2/\text{depth}}$  max-margin solution in the singular value space, leading the parameters to converge to the *top singular vectors* of the tensor when depth = 2. This theorem subsumes known results on linear convolutional networks and diagonal networks proved in Gunasekar et al. (2018b), *without* using convergence assumptions.
- For underdetermined linear regression, we study the limit points of gradient flow on orthogonally decomposable networks (Assumption 1), and provide a full characterization of the limit points. This theorem covers results on deep matrix sensing (Arora et al., 2019b) as a special case, and extends a similar recent result (Woodworth et al., 2020) to a broader class of networks.
- For underdetermined linear regression with deep linear fully-connected networks, we prove that the network converges to the minimum  $\ell_2$  norm solutions as we scale the initialization to zero.
- Lastly, we present simple experiments that corroborate our theoretical analysis. Figure 1 shows that our predictions of limit points match with solutions found by GD.

## 2 PROBLEM SETTINGS AND RELATED WORKS

We first define notation used in the paper. Given a positive integer  $a$ , let  $[a] := \{1, \dots, a\}$ . We use  $I_d$  to denote the  $d \times d$  identity matrix. Given a matrix  $\mathbf{A}$ , we use  $\text{vec}(\mathbf{A})$  to denote its vectorization, i.e., the concatenation of all columns of  $\mathbf{A}$ . For two vectors  $\mathbf{a}$  and  $\mathbf{b}$ , let  $\mathbf{a} \otimes \mathbf{b}$  be their tensor product,  $\mathbf{a} \cdot \mathbf{b}$  be their element-wise product, and  $\mathbf{a}^{\otimes k}$  be the element-wise  $k$ -th power of  $\mathbf{a}$ . Given an order- $L$  tensor  $\mathbf{A} \in \mathbb{R}^{k_1 \times \dots \times k_L}$ , we use  $[\mathbf{A}]_{j_1, \dots, j_L}$  to denote the  $(j_1, j_2, \dots, j_L)$ -th element of  $\mathbf{A}$ , where  $j_l \in [k_l]$  for all  $l \in [L]$ . In element indexing, we use  $\mathbf{A}_{i_1, \dots, i_L}$  to denote all indices in the corresponding dimension, and  $\mathbf{A}_{a:b}$  to denote all indices from  $a$  to  $b$ . For example, for a matrix  $\mathbf{A}$ ,  $[\mathbf{A}]_{:,4:6}$  denotes a submatrix that consists of 4th–6th columns of  $\mathbf{A}$ . The square bracket notation for indexing overloads

with  $[a]$  when  $a \geq \mathbb{N}$ , but they will be distinguishable from the context. Since element indices start from 1, we re-define the modulo operation  $a \bmod d := a - \lfloor \frac{a-1}{d} \rfloor d \in [d]$  for  $a > 0$ . We use  $\mathbf{e}_j^k$  to denote the  $j$ -th standard basis vector of the vector space  $\mathbb{R}^k$ . Lastly, we define the multilinear multiplication between a tensor and linear maps, which can be viewed as a generalization of left- and right-multiplication on a matrix. Given a tensor  $\mathbf{A} \in \mathbb{R}^{k_1 \times \dots \times k_L}$  and linear maps  $\mathbf{B}_l \in \mathbb{R}^{p_l \times k_l}$  for  $l \in [L]$ , we define the multilinear multiplication between them as

$$\begin{aligned} \mathbf{A} (\mathbf{B}_1^T, \mathbf{B}_2^T, \dots, \mathbf{B}_L^T) &= \prod_{j_1, \dots, j_L} [\mathbf{A}]_{j_1, \dots, j_L} (\mathbf{e}_{j_1}^{k_1} \quad \dots \quad \mathbf{e}_{j_L}^{k_L}) (\mathbf{B}_1^T, \dots, \mathbf{B}_L^T) \\ &:= \prod_{j_1, \dots, j_L} [\mathbf{A}]_{j_1, \dots, j_L} (\mathbf{B}_1 \mathbf{e}_{j_1}^{k_1} \quad \dots \quad \mathbf{B}_L \mathbf{e}_{j_L}^{k_L}) \in \mathbb{R}^{p_1 \times \dots \times p_L}. \end{aligned}$$

## 2.1 PROBLEM SETTINGS

We are given a dataset  $f(\mathbf{x}_i, y_i)_{i=1}^n$ , where  $\mathbf{x}_i \in \mathbb{R}^d$  and  $y_i \in \mathbb{R}$ . We let  $\mathbf{X} \in \mathbb{R}^{n \times d}$  and  $\mathbf{y} \in \mathbb{R}^n$  be the data matrix and the label vector, respectively. We study binary classification and linear regression in this paper, focusing on the settings where there exist *many* global solutions. For binary classification, we assume  $y_i \in \{-1, 1\}$  and that the data is separable: there exists a unit vector  $\mathbf{z}$  and a constant  $\gamma > 0$  such that  $y_i \mathbf{x}_i^T \mathbf{z} \geq \gamma$  for all  $i \in [n]$ . For regression, we consider the underdetermined case ( $n > d$ ) where there are many parameters  $\mathbf{z} \in \mathbb{R}^d$  such that  $\mathbf{X}\mathbf{z} = \mathbf{y}$ . Throughout the paper, we assume that  $\mathbf{X}$  has full row rank.

We use  $f(\cdot; \theta) : \mathbb{R}^d \rightarrow \mathbb{R}$  to denote a neural network parametrized by  $\theta$ . Given the network and the dataset, we consider minimizing the training loss  $L(\theta) := \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i; \theta), y_i)$  over  $\theta$ . Following previous results (e.g., [Lyu & Li \(2020\)](#); [Ji & Telgarsky \(2020\)](#)), we use the exponential loss  $\ell(\hat{y}, y) = \exp(-\hat{y}y)$  for classification problems. For regression, we use the squared error loss  $\ell(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$ . On the algorithm side, we minimize  $L$  using gradient flow, which can be viewed as GD with infinitesimal step size. The gradient flow dynamics is defined as  $\frac{d\theta}{dt} = -\nabla_{\theta} L(\theta)$ .

## 2.2 RELATED WORKS

**Gradient flow/descent in separable classification.** For linear models  $f(\mathbf{x}; \mathbf{z}) = \mathbf{x}^T \mathbf{z}$  with separable data, [Soudry et al. \(2018\)](#) show that the GD run on  $L$  drives  $\|\mathbf{z}\|$  to  $\infty$ , but  $\mathbf{z}$  converges in direction to the  $\ell_2$  max-margin classifier. The limit direction of  $\mathbf{z}$  is aligned with the solution of

$$\text{minimize}_{\mathbf{z} \in \mathbb{R}^d} \|\mathbf{z}\| \quad \text{subject to} \quad y_i \mathbf{x}_i^T \mathbf{z} \geq 1 \text{ for } i \in [n], \quad (1)$$

where the norm in the cost is the  $\ell_2$  norm. [Nacson et al. \(2019b;c\)](#); [Gunasekar et al. \(2018a\)](#); [Ji & Telgarsky \(2019b;c\)](#) extend these results to other (stochastic) algorithms and non-separable settings.

[Gunasekar et al. \(2018b\)](#) study the same problem on linear neural networks and show that GD exhibits different implicit bias depending on the architecture. The authors show that the linear coefficients of the network converges in direction to the solution of (1) with different norms:  $\ell_2$  norm for linear fully-connected networks,  $\ell_{2/L}$  (quasi-)norm for diagonal networks, and DFT-domain  $\ell_{2/L}$  (quasi-)norm for convolutional networks with full-length filters. Here,  $L$  denotes the depth. We note that [Gunasekar et al. \(2018b\)](#) assume that GD globally minimizes the loss, and the network parameters and the gradient with respect to the linear coefficients converge in direction. Subsequent results ([Ji & Telgarsky, 2019a; 2020](#)) remove such assumptions for linear fully-connected networks.

A recent line of results ([Nacson et al., 2019a; Lyu & Li, 2020; Ji & Telgarsky, 2020](#)) studies general homogeneous models and show divergence of parameters to infinity, monotone increase of smoothed margin, directional convergence and alignment of parameters (see Section 4 for details). [Lyu & Li \(2020\)](#) also characterize the limit direction of parameters as the KKT point of a nonconvex max-margin problem similar to (1), but this characterization does not provide useful insights for the functions  $f(\cdot; \theta)$  represented by specific architectures, because the formulation is in the parameter space  $\theta$ . Also, these results require that gradient flow/descent has already reached 100% training accuracy. Although we study a more restrictive set of networks (i.e., deep linear), we provide a more complete characterization of the implicit bias for the functions  $f(\cdot; \theta)$ , without assuming 100% training accuracy.

**Gradient flow/descent in linear regression.** It is known that for linear models  $f(\mathbf{x}; \mathbf{z}) = \mathbf{x}^T \mathbf{z}$ , GD converges to the global minimum that is closest in  $\ell_2$  distance to the initialization (see e.g.,

Gunasekar et al. (2018a)). However, relatively less is known for deep networks, even for linear networks. This is partly because the parameters do not diverge to infinity, hence making limit points highly dependent on the initialization, this dependency renders analysis difficult. A related problem of matrix sensing aims to minimize  $\sum_{i=1}^n (y_i - \mathbf{h} \mathbf{A}_i \mathbf{W}_1 \cdots \mathbf{W}_L)^2$  over  $\mathbf{W}_1, \dots, \mathbf{W}_L \in \mathbb{R}^{d \times d}$ . It is shown in Gunasekar et al. (2017); Arora et al. (2019b) that if the sensor matrices  $\mathbf{A}_i$  commute and we initialize all  $\mathbf{W}_i$ 's to  $\alpha \mathbf{I}$ , GD finds the minimum nuclear norm solution as  $\alpha \neq 0$ .

Chizat et al. (2019) show that if a network is zero at initialization, and we scale the network output by a factor of  $\alpha \neq 1$ , then the GD dynamics enters a ‘‘lazy regime’’ where the network behaves like a first-order approximation at its initialization, as also seen in results studying kernel approximations of neural networks and convergence of GD in the corresponding RKHS (e.g., Jacot et al. (2018)).

Woodworth et al. (2020) study linear regression with a diagonal network of the form  $f(\mathbf{x}; \mathbf{w}_+, \mathbf{w}_-) = \mathbf{x}^T (\mathbf{w}_+^{\odot L} \mathbf{w}_-^{\odot L})$ , where  $\mathbf{w}_+$  and  $\mathbf{w}_-$  are identically initialized  $\mathbf{w}_+(0) = \mathbf{w}_-(0) = \alpha \mathbf{w}$ . The authors show that the global minimum reached by GD minimizes a norm-like function which interpolates between (weighted)  $\ell_1$  norm ( $\alpha \neq 0$ ) and  $\ell_2$  norm ( $\alpha \neq 1$ ). In our paper, we consider a more general class of orthogonally decomposable networks, and obtain similar results interpolating between weighted  $\ell_1$  and  $\ell_2$  norms. We also remark that our results include the results in Arora et al. (2019b) as a special case, and we do not assume convergence to global minima, as done in Gunasekar et al. (2017); Arora et al. (2019b); Woodworth et al. (2020).

### 3 TENSOR FORMULATION OF NEURAL NETWORKS

In this section, we present a general tensor formulation of neural networks. Given an input  $\mathbf{x} \in \mathbb{R}^d$ , the network uses a linear map  $\mathbf{M}$  that maps  $\mathbf{x}$  to an order- $L$  tensor  $\mathbf{M}(\mathbf{x}) \in \mathbb{R}^{k_1 \times \dots \times k_L}$ , where  $L \geq 2$ . Using parameters  $\mathbf{v}_l \in \mathbb{R}^{k_l}$  and activation  $\phi$ , the network computes its layers as the following:

$$\begin{aligned} \mathbf{H}_1(\mathbf{x}) &= \phi(\mathbf{M}(\mathbf{x}) (\mathbf{v}_1, \mathbf{I}_{k_2}, \dots, \mathbf{I}_{k_L})) \in \mathbb{R}^{k_2 \times \dots \times k_L}, \\ \mathbf{H}_l(\mathbf{x}) &= \phi(\mathbf{H}_{l-1}(\mathbf{x}) (\mathbf{v}_l, \mathbf{I}_{k_{l+1}}, \dots, \mathbf{I}_{k_L})) \in \mathbb{R}^{k_{l+1} \times \dots \times k_L}, \text{ for } l = 2, \dots, L-1, \\ f(\mathbf{x}; \mathbf{v}) &= \mathbf{H}_{L-1}(\mathbf{x}) \mathbf{v}_L \in \mathbb{R}. \end{aligned} \quad (2)$$

We use  $\mathbf{v}$  to denote the collection of all parameters  $(\mathbf{v}_1, \dots, \mathbf{v}_L)$ . We call  $\mathbf{M}(\mathbf{x})$  the *data tensor*. Although this new formulation may look a bit odd in the first glance, it is general enough to capture many network architectures considered in the literature, including fully-connected networks, diagonal networks, and circular convolutional networks. We formally define these architectures below.

**Diagonal networks.** An  $L$ -layer diagonal network is written as

$$f_{\text{diag}}(\mathbf{x}; \mathbf{v}_{\text{diag}}) = \phi(\phi(\phi(\mathbf{x} \mathbf{w}_1) \mathbf{w}_2) \cdots \mathbf{w}_{L-1})^T \mathbf{w}_L, \quad (3)$$

where  $\mathbf{w}_l \in \mathbb{R}^d$  for  $l \in [L]$ . The representation of  $f_{\text{diag}}$  as the tensor form (2) is straightforward. Let  $\mathbf{M}_{\text{diag}}(\mathbf{x}) \in \mathbb{R}^{d \times \dots \times d}$  have  $[\mathbf{M}_{\text{diag}}(\mathbf{x})]_{j_1, j_2, \dots, j_L} = [\mathbf{x}]_{j_l}$ , while all the remaining entries of  $\mathbf{M}_{\text{diag}}(\mathbf{x})$  are set to zero. We can set  $\mathbf{v}_l = \mathbf{w}_l$  for all  $l$ , and  $\mathbf{M} = \mathbf{M}_{\text{diag}}$  to verify that (2) and (3) are equivalent.

**Circular convolutional networks.** The tensor formulation (2) includes convolutional networks

$$f_{\text{conv}}(\mathbf{x}; \mathbf{v}_{\text{conv}}) = \phi(\phi(\phi(\mathbf{x} \star \mathbf{w}_1) \star \mathbf{w}_2) \cdots \star \mathbf{w}_{L-1})^T \mathbf{w}_L, \quad (4)$$

where  $\mathbf{w}_l \in \mathbb{R}^{k_l}$  with  $k_l = d$  and  $k_L = d$ , and  $\star$  defines the circular convolution: for any  $\mathbf{a} \in \mathbb{R}^d$  and  $\mathbf{b} \in \mathbb{R}^k$  ( $k = d$ ), we have  $\mathbf{a} \star \mathbf{b} \in \mathbb{R}^d$  defined as  $[\mathbf{a} \star \mathbf{b}]_i = \sum_{j=1}^k [\mathbf{a}]_{(i+j-1) \bmod d} [\mathbf{b}]_j$ , for  $i \in [d]$ . Define  $\mathbf{M}_{\text{conv}}(\mathbf{x}) \in \mathbb{R}^{k_1 \times \dots \times k_L}$  as  $[\mathbf{M}_{\text{conv}}(\mathbf{x})]_{j_1, j_2, \dots, j_L} = [\mathbf{x}]_{\sum_{l=1}^L j_l - L + 1 \bmod d}$  for  $j_l \in [k_l]$ ,  $l \in [L]$ . Setting  $\mathbf{v}_l = \mathbf{w}_l$  and  $\mathbf{M} = \mathbf{M}_{\text{conv}}$ , we can verify that (2) and (4) are identical.

**Fully-connected networks.** An  $L$ -layer fully-connected network is defined as

$$f_{\text{fc}}(\mathbf{x}; \mathbf{v}_{\text{fc}}) = \phi(\phi(\phi(\mathbf{x}^T \mathbf{W}_1 \mathbf{W}_2) \cdots \mathbf{W}_{L-1}) \mathbf{w}_L), \quad (5)$$

where  $\mathbf{W}_l \in \mathbb{R}^{d_l \times d_{l+1}}$  for  $l \in [L-1]$  (we use  $d_1 = d$ ) and  $\mathbf{w}_L \in \mathbb{R}^{d_L}$ . One can represent  $f_{\text{fc}}$  as the tensor form (2) by defining parameters  $\mathbf{v}_l = \text{vec}(\mathbf{W}_l)$  for  $l \in [L-1]$  and  $\mathbf{v}_L = \mathbf{w}_L$ , and constructing the tensor  $\mathbf{M}_{\text{fc}}(\mathbf{x})$  by a recursive ‘‘block diagonal’’ manner. For example, if  $L = 2$ , we can define  $\mathbf{M}_{\text{fc}}(\mathbf{x}) \in \mathbb{R}^{d_1 d_2 \times d_2}$  to be the Kronecker product of  $\mathbf{I}_{d_2}$  and  $\mathbf{x}$ . For deeper networks, we defer the full description of  $\mathbf{M}_{\text{fc}}(\mathbf{x})$  to Appendix B.

**Our focus: linear tensor networks.** Throughout this section, we have used the activation  $\phi$  to motivate our tensor formulation (2) for neural networks with nonlinear activations. For the remaining of the paper, we study the case whose activation is *linear*, i.e.,  $\phi(t) = t$ . In this case,

$$f(\mathbf{x}; \mathbf{v}) = \mathbf{M}(\mathbf{x}) (\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_L). \quad (6)$$

We will refer to (6) as *linear tensor networks*, where ‘‘linear’’ is to indicate that the activation is linear. Note that as a function of parameters  $\mathbf{v}_1, \dots, \mathbf{v}_L$ ,  $f(\mathbf{x}; \mathbf{v})$  is in fact multilinear. We also remark that when depth  $L = 2$ , the data tensor  $\mathbf{M}(\mathbf{x})$  is a  $k_1 \times k_2$  matrix and the network formulation boils down to  $f(\mathbf{x}; \mathbf{v}) = \mathbf{v}_1^T \mathbf{M}(\mathbf{x}) \mathbf{v}_2$ .

Since the data tensor  $\mathbf{M}(\mathbf{x})$  is a linear function of  $\mathbf{x}$ , the linear tensor network is also a linear function of  $\mathbf{x}$ . Thus, the output of the network can also be written as  $f(\mathbf{x}; \mathbf{v}) = \mathbf{x}^T \mathbf{c}(\mathbf{v})$ , where  $\mathbf{c}(\mathbf{v}) \in \mathbb{R}^d$  denotes the *linear coefficients* computed as a function of the network parameters  $\mathbf{v}$ . Since the linear tensor network  $f(\mathbf{x}; \mathbf{v})$  is linear in  $\mathbf{x}$ , the expressive power of  $f$  is at best a linear model  $\mathbf{x}^T \mathbf{z}$ . However, even though the models have the same expressive power, their architectural differences lead to different implicit biases in training, which is the focus of our investigation in this paper. Studying separable classification and underdetermined regression is useful for highlighting such biases because there are *infinitely many* coefficients that perfectly classify or fit the dataset.

For our linear tensor network, the evolution of the parameters  $\mathbf{v}_l$  via gradient flow reads

$$\begin{aligned} \dot{\mathbf{v}}_l &= -\mathbf{r}_{\mathbf{v}_l} \mathcal{L}(\mathbf{v}) = -\sum_{i=1}^n \ell'(f(\mathbf{x}_i; \mathbf{v}), y_i) \mathbf{M}(\mathbf{x}_i) (\mathbf{v}_1, \dots, \mathbf{v}_{l-1}, \mathbf{I}_{k_l}, \mathbf{v}_{l+1}, \dots, \mathbf{v}_L) \\ &= \mathbf{M}(\mathbf{X}^T \mathbf{r}) (\mathbf{v}_1, \dots, \mathbf{v}_{l-1}, \mathbf{I}_{k_l}, \mathbf{v}_{l+1}, \dots, \mathbf{v}_L), \quad \mathcal{L} \in [L], \end{aligned}$$

where we initialize  $\mathbf{v}_l(0) = \alpha \mathbf{v}_l$ , for  $l \in [L]$ . We refer to  $\alpha$  and  $\mathbf{v}_l$  as the *initial scale* and *initial direction*, respectively. We note that we do not restrict  $\mathbf{v}_l$ 's to be unit vectors, in order to allow different scaling (at initialization) over different layers. The vector  $\mathbf{r} \in \mathbb{R}^n$  is the *residual vector*, and each component of  $\mathbf{r}$  is defined as

$$[\mathbf{r}]_i = \ell'(f(\mathbf{x}_i; \mathbf{v}), y_i) = \begin{cases} y_i \exp(-y_i f(\mathbf{x}_i; \mathbf{v})) & \text{for classification,} \\ f(\mathbf{x}_i; \mathbf{v}) - y_i & \text{for regression.} \end{cases} \quad (7)$$

## 4 IMPLICIT BIAS OF GRADIENT FLOW IN SEPARABLE CLASSIFICATION

In this section, we present our results on the implicit bias of gradient flow in binary classification with linearly separable data. Recent papers (Lyu & Li, 2020; Ji & Telgarsky, 2020) on this separable classification setup prove that after 100% training accuracy has been achieved by gradient flow (along with other technical conditions), the parameters of  $L$ -homogeneous models diverge to infinity, while converging in direction that aligns with the direction of the negative gradient. Mathematically,

$$\lim_{t \rightarrow \infty} k(\mathbf{v}(t))k = 1, \quad \lim_{t \rightarrow \infty} \frac{\mathbf{v}(t)}{\|\mathbf{v}(t)\|} = \frac{-\nabla_{\mathbf{v}} \mathcal{L}(\mathbf{v}(t))}{\|\nabla_{\mathbf{v}} \mathcal{L}(\mathbf{v}(t))\|} = \mathbf{1}.$$

Since the linear tensor network satisfies the technical assumptions in the prior works, we apply these results to our setting and develop a new characterization of the limit directions of the parameters. Here, we present theorems on separable classification with general linear tensor networks. Corollaries for specific networks are deferred to Appendix A.

### 4.1 LIMIT DIRECTIONS OF PARAMETERS ARE SINGULAR VECTORS

Consider the singular value decomposition (SVD) of a matrix  $\mathbf{A} = \sum_{j=1}^m s_j (\mathbf{u}_j \mathbf{v}_j^T)$ , where  $m$  is the rank of  $\mathbf{A}$ . Note that the tuples  $(\mathbf{u}_j, \mathbf{v}_j, s_j)$  are solutions to the system of equations  $s\mathbf{u} = \mathbf{A}\mathbf{v}$  and  $s\mathbf{v} = \mathbf{A}^T \mathbf{u}$ . Lim (2005) generalizes this definition of singular vectors and singular values to higher-order tensors: given an order- $L$  tensor  $\mathbf{A} \in \mathbb{R}^{k_1 \times \dots \times k_L}$ , we define the singular vectors  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_L$  and singular value  $s$  to be the solution of the following system of equations:

$$s\mathbf{u}_l = \mathbf{A} (\mathbf{u}_1, \dots, \mathbf{u}_{l-1}, \mathbf{I}_{k_l}, \mathbf{u}_{l+1}, \dots, \mathbf{u}_L), \quad \text{for } l \in [L]. \quad (8)$$

Using the definition of the singular vectors of tensors, we can characterize the limit direction of parameters after reaching 100% training accuracy. In Appendix C, we prove the following:



**Theorem 1.** Assume that the gradient flow satisfies  $L(\mathbf{v}(t_0)) < 1$  for some  $t_0 \geq 0$  and  $\mathbf{X}^T \mathbf{r}(t)$  converges in direction, say  $\mathbf{u}^\infty := \lim_{t \rightarrow \infty} \frac{\mathbf{X}^T \mathbf{r}(t)}{\|\mathbf{X}^T \mathbf{r}(t)\|_2}$ . Then,  $\mathbf{v}_1, \dots, \mathbf{v}_L$  converge to the singular vectors of  $\mathbf{M}(\mathbf{u}^\infty)$ .

For this theorem, we make some convergence assumptions, because the network is fully general; this is the *only* result where we assume convergence. In fact, for the special case of linear fully-connected networks, the directional convergence assumption is *not* required, and the linear coefficients  $\mathbf{v}_i(t)$  converge in direction to the  $\ell_2$  max-margin classifier. We state this corollary in Appendix A.1; this result also appears in Ji & Telgarsky (2020), but we provide an alternative proof.

#### 4.2 LIMIT DIRECTIONS FOR ORTHOGONALLY DECOMPOSABLE NETWORKS

Admittedly, Theorem 1 is not a *full* characterization of the limit directions, because there are usually multiple solutions that satisfy (8). For example, in case of  $L = 2$ , the data tensor  $\mathbf{M}(\mathbf{u}^\infty)$  is a matrix and the number of possible limit directions (up to scaling) of  $(\mathbf{v}_1, \mathbf{v}_2)$  is at least the rank of  $\mathbf{M}(\mathbf{u}^\infty)$ . Singular vectors of high order tensors are much less understood than the matrix counterparts, and are much harder to deal with. Although their existence is implied from the variational formulation (Lim, 2005), they are intractable to compute. Testing if a given number is a singular value, approximating the corresponding singular vectors, and computing the best rank-1 approximation are all NP-hard (Hillar & Lim, 2013); let alone orthogonal decompositions.

Given this intractability, it might be reasonable to make some assumptions on the “structure” of the data tensor  $\mathbf{M}(\mathbf{x})$ , so that they are easier to handle. The following assumption defines a class of *orthogonally decomposable* data tensors, which includes **linear diagonal networks** and **linear full-length convolutional networks** as special cases (for the proof, see Appendix D.2 and D.3).

**Assumption 1.** For the data tensor  $\mathbf{M}(\mathbf{x}) \in \mathbb{R}^{k_1 \times \dots \times k_L}$  of a linear tensor network (6), there exist a full column rank matrix  $\mathbf{S} \in \mathbb{C}^{m \times d}$  ( $d = m \cdot \min_l k_l$ ) and matrices  $\mathbf{U}_1 \in \mathbb{C}^{k_1 \times m}, \dots, \mathbf{U}_L \in \mathbb{C}^{k_L \times m}$  such that  $\mathbf{U}_l^H \mathbf{U}_l = \mathbf{I}_m$  for all  $l \in [L]$ , and the data tensor  $\mathbf{M}(\mathbf{x})$  can be written as

$$\mathbf{M}(\mathbf{x}) = \prod_{j=1}^m [\mathbf{S}\mathbf{x}]_j ([\mathbf{U}_1]_{\cdot,j} \quad [\mathbf{U}_2]_{\cdot,j} \quad \dots \quad [\mathbf{U}_L]_{\cdot,j}). \quad (9)$$

In this assumption, we allow  $\mathbf{U}_1, \dots, \mathbf{U}_L$  and  $\mathbf{S}$  to be complex matrices, although  $\mathbf{M}(\mathbf{x})$  and parameters  $\mathbf{v}_i$  stay real, as defined earlier. For a complex matrix  $\mathbf{A}$ , we use  $\mathbf{A}^*$  to denote its entry-wise complex conjugate,  $\mathbf{A}^T$  to denote its transpose (without conjugating), and  $\mathbf{A}^H$  to denote its conjugate transpose. In case of  $L = 2$ , Assumption 1 requires that the data tensor  $\mathbf{M}(\mathbf{x})$  (now a matrix) has singular value decomposition  $\mathbf{M}(\mathbf{x}) = \mathbf{U}_1 \text{diag}(\mathbf{S}\mathbf{x}) \mathbf{U}_2^T$ ; i.e., the left and right singular vectors are independent of  $\mathbf{x}$ , and the singular values are linear in  $\mathbf{x}$ . Using Assumption 1, the following theorem characterizes the limit directions.

**Theorem 2.** Suppose a linear tensor network satisfies Assumption 1. If there exists  $\lambda > 0$  such that the initial directions  $\mathbf{v}_1, \dots, \mathbf{v}_L$  of the network parameters satisfy  $\|\mathbf{U}_l^T \mathbf{v}_l\|_j^2 \geq \lambda \|\mathbf{U}_l^T \mathbf{v}_l\|_j^2$  for all  $l \in [L-1]$  and  $j \in [m]$ , then  $(\mathbf{v}(t))$  converges in a direction that aligns with  $\mathbf{S}^T \mathbf{z}^\infty$ , where  $\mathbf{z}^\infty \in \mathbb{C}^m$  denotes a stationary point of the following optimization problem

$$\text{minimize}_{\mathbf{z} \in \mathbb{C}^m} \|\mathbf{z}\|_{2/L} \quad \text{subject to} \quad \|\mathbf{y}_i \mathbf{x}_i^T \mathbf{S}^T \mathbf{z}\|_1 \leq 1, \quad \forall i \in [n].$$

If  $\mathbf{S}$  is invertible, then  $(\mathbf{v}(t))$  converges in a direction that aligns with a stationary point  $\mathbf{z}^\infty$  of

$$\text{minimize}_{\mathbf{z} \in \mathbb{R}^d} \|\mathbf{z}\|_{2/L} \quad \text{subject to} \quad \|\mathbf{y}_i \mathbf{x}_i^T \mathbf{z}\|_1 \leq 1, \quad \forall i \in [n].$$

Theorem 2 shows that the gradient flow finds sparse  $\mathbf{z}^\infty$  that minimizes the  $\ell_{2/L}$  norm in the “singular value space,” where the data points  $\mathbf{x}_i$  are transformed into vectors  $\mathbf{S}\mathbf{x}_i$  consisting of singular values of  $\mathbf{M}(\mathbf{x}_i)$ . Also, the proof of Theorem 2 reveals that in case of  $L = 2$ , the parameters  $\mathbf{v}_i(t)$  in fact converge to the *top* singular vectors of the data tensor  $\mathbf{M}(\mathbf{X}^T \mathbf{r})$ ; thus, compared to Theorem 1, we have a more complete characterization of “which” singular vectors to converge to.

The proof of Theorem 2 is in Appendix D. Since the orthogonal decomposition (Assumption 1) of  $\mathbf{M}(\mathbf{x})$  tells us that the singular vectors  $\mathbf{M}(\mathbf{x})$  in  $\mathbf{U}_1, \dots, \mathbf{U}_L$  are independent of  $\mathbf{x}$ , we can transform the network parameters  $\mathbf{v}_i$  to  $\mathbf{U}_l^T \mathbf{v}_i$  and show that the network behaves similar to a linear diagonal network. This observation comes in handy in the characterization of limit directions.

**Remark 1** (Necessity of initialization assumptions). In order to remove the assumption that the loss converges to zero, at least some condition on initialization is *necessary*, because there are examples showing non-convergence of gradient flow for certain initializations (Bartlett et al., 2018; Arora et al., 2019a). In our theorems, we pose assumptions on initial directions  $\mathbf{v}_l$  that are sufficient conditions for the loss  $L(\cdot(t))$  to converge to zero. Although such sufficient conditions are “stronger” than assuming  $L(\cdot(t)) \neq 0$ , they are useful because they can be easily checked *a priori*, i.e., before running gradient flow. We note an important fact that in Theorems 2 and onwards, the conditions on initialization are used **solely to prove convergence of the loss to zero, and our statements on the implicit bias hold whenever the loss converges to zero**, even for initializations that do not satisfy our conditions. In addition, we argue that our assumptions are not too restrictive;  $\lambda$  can be arbitrarily small, so the conditions are satisfied *with probability 1* if we set  $\mathbf{v}_L = \mathbf{0}$  and randomly sample other  $\mathbf{v}_l$ ’s. Setting one layer to zero to prove convergence is also studied in Wu et al. (2019). Lastly, the condition that  $\mathbf{v}_L$  is “small” can be replaced with any layer; e.g., convergence still holds if  $\|\mathbf{U}_l^T \mathbf{v}_l\|_j^2 \leq \lambda \|\mathbf{U}_1^T \mathbf{v}_1\|_j^2$  for all  $l = 2, \dots, L$  and  $j \geq [m]$ .

**Remark 2** (Comparison to existing results). Theorem 2 leads to corollaries (stated in Appendix A.2) on linear diagonal and full-length convolutional networks, showing that diagonal (or convolutional) networks converge to the stationary point of the max-margin problem with respect to the  $\ell_{2/L}$  norm (or DFT-domain  $\ell_{2/L}$  norm). Theorem 2 recovers the results in Gunasekar et al. (2018b) without relying on assumptions such as directional convergence of parameters and gradients.

**Remark 3** (Implications to architecture design). Theorem 2 shows that the gradient flow finds a solution that is sparse in a “transformed” input space where all data points are transformed with  $\mathbf{S}$ . This implies something interesting about architecture design: if the sparsity of the solution under a certain linear transformation  $T$  is needed, one can design a network using Assumption 1 by setting  $\mathbf{S} = T$ . Training such a network will give us a solution that has the desired sparsity property.

Other than Assumption 1, there is another setting where we can prove a full characterization of limit directions: when there is one data point ( $n = 1$ ) and the network is 2-layer ( $L = 2$ ). This “extremely overparametrized” case is motivated by an experimental paper (Zhang et al., 2019) which studies generalization performance of different architectures when there is only one training data point.

**Theorem 3.** *Suppose we have a 2-layer linear tensor network (6) and a single data point  $(\mathbf{x}, y)$ . Consider the compact SVD  $\mathbf{M}(\mathbf{x}) = \mathbf{U}_1 \text{diag}(\mathbf{s}) \mathbf{U}_2^T$ , where  $\mathbf{U}_1 \in \mathbb{R}^{k_1 \times m}$ ,  $\mathbf{U}_2 \in \mathbb{R}^{k_2 \times m}$ , and  $\mathbf{s} \in \mathbb{R}^m$  for  $m = \min\{k_1, k_2\}$ . Let  $\mathbf{v}^\infty \in \mathbb{R}^m$  be a solution of the following optimization problem*

$$\text{minimize}_{\mathbf{v} \in \mathbb{R}^m} \|\mathbf{v}\|_{k_1} \quad \text{subject to} \quad \mathbf{y} \mathbf{s}^T \mathbf{v} = 1.$$

*Assume that there exists  $\lambda > 0$  such that the initial directions  $\mathbf{v}_1, \mathbf{v}_2$  of the network parameters satisfy  $\|\mathbf{U}_1^T \mathbf{v}_1\|_j^2 \leq \lambda \|\mathbf{U}_2^T \mathbf{v}_2\|_j^2$  for all  $j \geq [m]$ . Then,  $\mathbf{v}_1$  and  $\mathbf{v}_2$  converge in direction to  $\mathbf{U}_1 \mathbf{v}^\infty$  and  $\mathbf{U}_2 \mathbf{v}^\infty$ , where  $\|\mathbf{v}^\infty\|_1 = \|\mathbf{v}^\infty\|_2 = \|\mathbf{v}^\infty\|_{k_1} = \|\mathbf{v}^\infty\|_{k_2} = 1$ , and  $\text{sign}(\mathbf{v}^\infty) = \text{sign}(y) \text{sign}(\mathbf{s})$ .*

The proof of Theorem 3 can be found in Appendix E. Since  $\mathbf{v}^\infty$  is the minimum  $\ell_1$  norm solution in the singular value space, the parameters  $\mathbf{v}_1$  and  $\mathbf{v}_2$  converge in direction to the top singular vectors. We would like to emphasize that this theorem can be applied to *any* network architecture that can be represented as a linear tensor network. Recall that the previous result (Gunasekar et al., 2018b) only considers full-length filters ( $k_1 = d$ ), hence providing limited insights on networks with small filters, e.g.,  $k_1 = 2$ . In light of this, we present a corollary in Appendix A.3 showing that linear coefficients of convolutional networks converge in direction to a “filtered” version of  $\mathbf{x}$ .

## 5 IMPLICIT BIAS OF GRADIENT FLOW IN UNDERDETERMINED REGRESSION

In Section 4, the limit directions of parameters we characterized do not depend on initialization. This is due to the fact that the parameters diverge to infinity in separable classification problems, so that the initialization becomes unimportant in the limit. This is not the case in regression setting, because parameters do not diverge to infinity. As we show in this section, the limit points are closely tied to initialization, and our analyses characterize the dependency between them.

### 5.1 LIMIT POINT CHARACTERIZATION FOR ORTHOGONALLY DECOMPOSABLE NETWORKS

For the orthogonally decomposable networks satisfying Assumption 1 with real  $\mathbf{S}$  and  $\mathbf{U}_l$ ’s, we consider how limit points of gradient flow change according to initialization. We consider a specific

initialization scheme that, in the special case of diagonal networks, corresponds to setting  $w_l(0) = \alpha \mathbf{W}$  for  $l \geq [L - 1]$  and  $w_L(0) = \mathbf{0}$ . We use the following lemma on a relevant system of ODEs:

**Lemma 4.** Consider the system of ODEs, where  $p, q : \mathbb{R} \rightarrow \mathbb{R}$ :

$$\dot{p} = p^{L-2}q, \quad \dot{q} = p^{L-1}, \quad p(0) = 1, \quad q(0) = 0.$$

Then, the solutions  $p_L(t)$  and  $q_L(t)$  are continuous on their maximal interval of existence of the form  $(-c, c) \subset \mathbb{R}$  for some  $c \geq (0, 1]$ . Define  $h_L(t) = p_L(t)^{L-1}q_L(t)$ ; then,  $h_L(t)$  is odd and strictly increasing, satisfying  $\lim_{t \uparrow c} h_L(t) = 1$  and  $\lim_{t \downarrow -c} h_L(t) = -1$ .

Using the function  $h_L(t)$  from Lemma 4, we can obtain the following theorem that characterizes the limit points as the minimizer of a norm-like function  $Q_{L,\alpha,-}$  among the global minima.

**Theorem 5.** Suppose a linear tensor network satisfies Assumption 1. Assume further that the matrices  $\mathbf{U}_1, \dots, \mathbf{U}_L$  and  $\mathbf{S}$  from Assumption 1 are all real matrices. For some  $\lambda > 0$ , choose any vector  $\mathbf{z} \in \mathbb{R}^m$  satisfying  $[\mathbf{z}]_j^2 \geq \lambda$  for all  $j \geq [m]$ , and choose initial directions  $\mathbf{v}_l = \mathbf{U}_l$  for  $l \geq [L - 1]$  and  $\mathbf{v}_L = \mathbf{0}$ . Then, the linear coefficients  $(\cdot)(t)$  converge to  $\mathbf{S}^T \mathbf{z}^\infty$ , where  $\mathbf{z}^\infty$  is the solution of

$$\text{minimize}_{\mathbf{z} \in \mathbb{R}^m} Q_{L,\alpha,-}(\mathbf{z}) := \alpha^2 \prod_{j=1}^m [\mathbf{z}]_j^2 H_L \left( \frac{[\mathbf{z}]_j}{\alpha^{L-1} |[\mathbf{z}]_j|^{L-1}} \right) \quad \text{subject to} \quad \mathbf{X} \mathbf{S}^T \mathbf{z} = \mathbf{y},$$

where  $Q_{L,\alpha,-} : \mathbb{R}^m \rightarrow \mathbb{R}$  is a norm-like function defined using  $H_L(t) := \int_0^t h_L^{-1}(\tau) d\tau$ . If  $\mathbf{S}$  is invertible, then  $(\cdot)(t)$  converges to the solution  $\mathbf{z}^\infty$  of

$$\text{minimize}_{\mathbf{z} \in \mathbb{R}^d} Q_{L,\alpha,-}(\mathbf{S}^{-T} \mathbf{z}) \quad \text{subject to} \quad \mathbf{X} \mathbf{z} = \mathbf{y}.$$

The proofs of Lemma 4 and Theorem 5 are deferred to Appendix F.

**Remark 4** (Interpolation between  $\ell_1$  and  $\ell_2$ ). It can be checked that  $H_L(t)$  grows like the absolute value function if  $t$  is large, and grows like a quadratic function if  $t$  is close to zero. This means that

$$\lim_{\alpha \rightarrow 0} Q_{L,\alpha,-}(\mathbf{z}) \propto \prod_{j=1}^m \frac{|[\mathbf{z}]_j|}{|[\mathbf{z}]_j|^{L-2}}, \quad \lim_{\alpha \rightarrow \infty} Q_{L,\alpha,-}(\mathbf{z}) \propto \prod_{j=1}^m \frac{[\mathbf{z}]_j^2}{|[\mathbf{z}]_j|^{2L-2}},$$

so  $Q_{L,\alpha,-}$  interpolates between the weighted  $\ell_1$  and weighted  $\ell_2$  norms of  $\mathbf{z}$ . Also, the weights in the norm are *dependent* on the initialization direction unless  $L = 2$  and  $\alpha \neq 0$ . In general,  $Q_{L,\alpha,-}$  interpolates the standard  $\ell_1$  and  $\ell_2$  norms only if  $\prod_{j=1}^m [\mathbf{z}]_j$  is the same for all  $j \geq [m]$ . This result is similar to the observations made in [Woodworth et al. \(2020\)](#) which considers a diagonal network with a ‘‘differential’’ structure  $f(\mathbf{x}; \mathbf{w}_+, \mathbf{w}_-) = \mathbf{x}^T (\mathbf{w}_+^{\odot L} - \mathbf{w}_-^{\odot L})$ . In contrast, our results apply to a more general class of networks, without the need to have the differential structure. In Appendix A.4, we state corollaries of Theorem 5 for linear diagonal networks and linear full-length convolutional networks with even data points. There, we also show that deep matrix sensing with commutative sensor matrices ([Arora et al., 2019b](#)) is a special case of our setting.

Next, we present the regression counterpart of Theorem 3, for 2-layer linear tensor networks with a single data point. For this extremely overparametrized setup, we can fully characterize the limit points as functions of initialization  $\mathbf{v}_1(0) = \alpha \mathbf{v}_1$  and  $\mathbf{v}_2(0) = \alpha \mathbf{v}_2$ , for *any* linear tensor networks including linear convolutional networks with filter size smaller than input dimension.

**Theorem 6.** Suppose we have a 2-layer linear tensor network (6) and a single data point  $(\mathbf{x}, y)$ . Consider the compact SVD  $\mathbf{M}(\mathbf{x}) = \mathbf{U}_1 \text{diag}(\mathbf{s}) \mathbf{U}_2^T$ , where  $\mathbf{U}_1 \in \mathbb{R}^{k_1 \times m}$ ,  $\mathbf{U}_2 \in \mathbb{R}^{k_2 \times m}$ , and  $\mathbf{s} \in \mathbb{R}^m$  for  $m = \min\{k_1, k_2\}$ . Assume that there exists  $\lambda > 0$  such that the initial directions  $\mathbf{v}_1, \mathbf{v}_2$  of the network parameters satisfy  $[\mathbf{U}_1^T \mathbf{v}_1]_j^2 \geq [\mathbf{U}_2^T \mathbf{v}_2]_j^2 \geq \lambda$  for all  $j \geq [m]$ . Then, gradient flow converges to a global minimizer of the loss  $L$ , and  $\mathbf{v}_1(t)$  and  $\mathbf{v}_2(t)$  converge to the limit points:

$$\begin{aligned} \mathbf{v}_1^\infty &= \alpha \mathbf{U}_1 \left( \mathbf{U}_1^T \mathbf{v}_1 \cosh \left( g^{-1} \frac{y}{\alpha^2} \mathbf{s} + \mathbf{U}_2^T \mathbf{v}_2 \sinh \left( g^{-1} \frac{y}{\alpha^2} \mathbf{s} + \alpha (\mathbf{I}_{k_1} - \mathbf{U}_1 \mathbf{U}_1^T) \mathbf{v}_1 \right) \right) \right), \\ \mathbf{v}_2^\infty &= \alpha \mathbf{U}_2 \left( \mathbf{U}_1^T \mathbf{v}_1 \sinh \left( g^{-1} \frac{y}{\alpha^2} \mathbf{s} + \mathbf{U}_2^T \mathbf{v}_2 \cosh \left( g^{-1} \frac{y}{\alpha^2} \mathbf{s} + \alpha (\mathbf{I}_{k_2} - \mathbf{U}_2 \mathbf{U}_2^T) \mathbf{v}_2 \right) \right) \right), \end{aligned}$$

where  $g^{-1}$  is the inverse of the following strictly increasing function

$$g(\nu) = \prod_{j=1}^m [\mathbf{s}]_j \frac{[\mathbf{U}_1^T \mathbf{v}_1]_j^2 + [\mathbf{U}_2^T \mathbf{v}_2]_j^2}{2} \sinh(2[\mathbf{s}]_j \nu) + [\mathbf{U}_1^T \mathbf{v}_1]_j [\mathbf{U}_2^T \mathbf{v}_2]_j \cosh(2[\mathbf{s}]_j \nu).$$

The proof can be found in Appendix G. We can observe that as  $\alpha \rightarrow 0$ , we have  $g^{-1} \frac{y}{\alpha^2} \rightarrow 1$ , which results in exponentially faster growth of the  $\sinh(\cdot)$  and  $\cosh(\cdot)$  for the top singular values. As a result, the top singular vectors dominate the limit points  $\mathbf{v}_1^\infty$  and  $\mathbf{v}_2^\infty$  as  $\alpha \rightarrow 0$ , and they do not depend on the initial directions  $\mathbf{v}_1, \mathbf{v}_2$ . Experiment results in Section 6 support this observation.



## 5.2 IMPLICIT BIAS IN FULLY-CONNECTED NETWORKS: THE $\alpha \rightarrow 0$ LIMIT

We state our last theoretical element of this paper, which proves that the linear coefficients  $\mathbf{w}_{fc}$  of deep linear fully-connected networks converge to the minimum  $\ell_2$  norm solution as  $\alpha \rightarrow 0$ . We assume for simplicity that  $d_1 = d_2 = \dots = d_L = d$  in this section, but we can extend it for  $d_l \neq d$  without too much difficulty. Recall  $f_{fc}(\mathbf{x}; \mathbf{w}_{fc}) = \mathbf{x}^T \mathbf{W}_1 \dots \mathbf{W}_{L-1} \mathbf{w}_L$ . We minimize the training loss  $L$  with initialization  $\mathbf{W}_l(0) = \alpha \mathbf{W}_l$  for  $l \geq [L-1]$  and  $\mathbf{w}_L(0) = \alpha \mathbf{w}_L$ .

**Theorem 7.** *Assume that initial directions  $\mathbf{W}_1, \dots, \mathbf{W}_{L-1}, \mathbf{w}_L$  satisfy (1)  $\mathbf{W}_l^T \mathbf{W}_l = \mathbf{W}_{l+1} \mathbf{W}_{l+1}^T$  for  $l \geq [L-2]$ , and (2) there exists  $\lambda > 0$  such that  $\mathbf{W}_{L-1}^T \mathbf{W}_{L-1} = \mathbf{w}_L \mathbf{w}_L^T = \lambda \mathbf{I}_d$ . Then, the gradient flow converges to a global minimum, and  $\lim_{\alpha \rightarrow 0} \lim_{t \rightarrow \infty} \mathbf{w}_{fc}(\mathbf{w}_{fc}(t)) = \mathbf{X}^T (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{y}$ .*

The proof is presented in Appendix H. Theorem 7 shows that in the limit  $\alpha \rightarrow 0$ , linear fully-connected networks have bias towards the minimum  $\ell_2$  norm solution, regardless of the depth. This is consistent with the results shown for classification. We also note that the convergence to a global minimum holds for any  $\alpha > 0$ , and our sufficient conditions ( $\mathbf{W}_l^T \mathbf{W}_l = \mathbf{W}_{l+1} \mathbf{W}_{l+1}^T$  and  $\mathbf{W}_{L-1}^T \mathbf{W}_{L-1} = \mathbf{w}_L \mathbf{w}_L^T = \lambda \mathbf{I}_d$ ) for global convergence is a generalization of the zero-asymmetric initialization scheme ( $\mathbf{W}_1 = \dots = \mathbf{W}_{L-1} = \mathbf{I}_d$  and  $\mathbf{w}_L = \mathbf{0}$ ) proposed in Wu et al. (2019).

## 6 EXPERIMENTS

**Regression.** To fully visualize the trajectory of linear coefficients, we run simple experiments with 2-layer linear fully-connected/diagonal/convolutional networks with a single 2-dimensional data point  $(\mathbf{x}, y) = ([1 \ 2], 1)$ . For this dataset, the minimum  $\ell_2$  norm solution (corresponding to fully-connected networks) of the regression problem is  $[0.2 \ 0.4]$ , whereas the minimum  $\ell_1$  norm solution (corresponding to diagonal) is  $[0 \ 0.5]$  and the minimum DFT-domain  $\ell_1$  norm solution (corresponding to convolutional) is  $[0.33 \ 0.33]$ . We randomly pick four directions  $\mathbf{z}_1, \dots, \mathbf{z}_4 \in \mathbb{R}^2$ , and choose initial directions of the network parameters in a way that their linear coefficients at initialization are exactly  $\mathbf{w}(0) = \alpha^2 \mathbf{z}_j$ . With varying initial scales  $\alpha \in \{0.01, 0.5, 1\}$ , we run GD with small step size  $\eta = 10^{-3}$  for large enough number of iterations  $T = 5 \cdot 10^3$ . Figures 1 and 2 plot the trajectories of  $\mathbf{w}(t)$  (appropriately clipped for visual clarity) as well as the predicted limit points (Theorem 6). We observe that even though the networks start at the same linear coefficients  $\alpha^2 \mathbf{z}_j$ , they evolve differently due to different architectures. Note that the prediction of limit points is accurate, and the solution found by GD is less dependent on initial directions when  $\alpha$  is small.

**Classification.** It is shown in the existing works as well as in Section 4 that the limit directions of linear coefficients are independent of the initialization. Is this also true in practice? To see this, we run a set of toy experiments on classification with two data points  $(\mathbf{x}_1, y_1) = ([1 \ 2], +1)$  and  $(\mathbf{x}_2, y_2) = ([0 \ 3], -1)$ . One can check that the max-margin classifiers for this problem are in the same directions to the corresponding min-norm solutions in the regression problem above. We use the same networks as in regression, and the same set of initial directions satisfying  $\mathbf{w}(0) = \alpha^2 \mathbf{z}_j$ . With initial scales  $\alpha \in \{0.01, 0.5, 1\}$ , we run GD with step size  $\eta = 5 \cdot 10^{-4}$  for  $T = 2 \cdot 10^6$  iterations. All experiments reached  $L(\mathbf{w}) \lesssim 10^{-5}$  at the end. The trajectories are plotted in Figure 2 in the Appendix. We find that, in contrast to our theoretical characterization, the actual coefficients are quite dependent on initialization, because we do not train the network all the way to zero loss. This observation is also consistent with a recent analysis (Moroshko et al., 2020) for diagonal networks, and suggests that understanding the behavior of iterates after a finite number of steps is an important future work.

## 7 CONCLUSION

This paper studies the implicit bias of gradient flow on training linear tensor networks. Under a general tensor formulation of linear networks, we provide theorems characterizing how the network architectures and initializations affect the limit directions/points of gradient flow. Our work provides a unified framework that connects multiple existing results on implicit bias of gradient flow as special cases.

## REFERENCES

- Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. *arXiv preprint arXiv:1811.03962*, 2018.
- Sanjeev Arora, Nadav Cohen, and Elad Hazan. On the optimization of deep networks: Implicit acceleration by overparameterization. In *International Conference on Machine Learning*, pp. 244–253, 2018.
- Sanjeev Arora, Nadav Cohen, Noah Golowich, and Wei Hu. A convergence analysis of gradient descent for deep linear neural networks. In *International Conference on Learning Representations*, 2019a.
- Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix factorization. In *Advances in Neural Information Processing Systems*, pp. 7413–7424, 2019b.
- Peter Bartlett, Dave Helmbold, and Philip Long. Gradient descent with identity initialization efficiently learns positive definite linear transformations by deep residual networks. In *International Conference on Machine Learning*, pp. 521–530, 2018.
- Lenaic Chizat and Francis Bach. Implicit bias of gradient descent for wide two-layer neural networks trained with the logistic loss. *arXiv preprint arXiv:2002.04486*, 2020.
- Lenaic Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming. In *Advances in Neural Information Processing Systems*, pp. 2937–2947, 2019.
- Simon S Du and Wei Hu. Width provably matters in optimization for deep linear neural networks. *arXiv preprint arXiv:1901.08572*, 2019.
- Simon S Du, Jason D Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks. *arXiv preprint arXiv:1811.03804*, 2018a.
- Simon S Du, Xiyu Zhai, Barnabas Póczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks. *arXiv preprint arXiv:1810.02054*, 2018b.
- Daniel Gissin, Shai Shalev-Shwartz, and Amit Daniely. The implicit bias of depth: How incremental learning drives generalization. In *International Conference on Learning Representations*, 2020.
- Suriya Gunasekar, Blake E Woodworth, Srinadh Bhojanapalli, Behnam Neyshabur, and Nati Srebro. Implicit regularization in matrix factorization. In *Advances in Neural Information Processing Systems*, pp. 6151–6159, 2017.
- Suriya Gunasekar, Jason Lee, Daniel Soudry, and Nathan Srebro. Characterizing implicit bias in terms of optimization geometry. In *International Conference on Machine Learning*, pp. 1832–1841, 2018a.
- Suriya Gunasekar, Jason D Lee, Daniel Soudry, and Nati Srebro. Implicit bias of gradient descent on linear convolutional networks. In *Advances in Neural Information Processing Systems*, pp. 9461–9471, 2018b.
- Christopher J Hillar and Lek-Heng Lim. Most tensor problems are NP-hard. *Journal of the ACM (JACM)*, 60(6):1–39, 2013.
- Wei Hu, Lechao Xiao, and Jeffrey Pennington. Provable benefit of orthogonal initialization in optimizing deep linear networks. In *International Conference on Learning Representations*, 2020.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*, pp. 8571–8580, 2018.
- Ziwei Ji and Matus Telgarsky. Gradient descent aligns the layers of deep linear networks. In *International Conference on Learning Representations*, 2019a.
- Ziwei Ji and Matus Telgarsky. The implicit bias of gradient descent on nonseparable data. In *Conference on Learning Theory*, pp. 1772–1798, 2019b.

- Ziwei Ji and Matus Telgarsky. A refined primal-dual analysis of the implicit bias. *arXiv preprint arXiv:1906.04540*, 2019c.
- Ziwei Ji and Matus Telgarsky. Directional convergence and alignment in deep learning. *arXiv preprint arXiv:2006.06657*, 2020.
- Lek-Heng Lim. Singular values and eigenvalues of tensors: a variational approach. In *1st IEEE International Workshop on Computational Advances in Multi-Sensor Adaptive Processing, 2005.*, pp. 129–132. IEEE, 2005.
- Kaifeng Lyu and Jian Li. Gradient descent maximizes the margin of homogeneous neural networks. In *International Conference on Learning Representations*, 2020.
- Edward Moroshko, Suriya Gunasekar, Blake Woodworth, Jason D Lee, Nathan Srebro, and Daniel Soudry. Implicit bias in deep linear classification: Initialization scale vs training accuracy. *arXiv preprint arXiv:2007.06738*, 2020.
- Mor Shpigel Nacson, Suriya Gunasekar, Jason Lee, Nathan Srebro, and Daniel Soudry. Lexicographic and depth-sensitive margins in homogeneous and non-homogeneous deep models. In *International Conference on Machine Learning*, pp. 4683–4692, 2019a.
- Mor Shpigel Nacson, Jason Lee, Suriya Gunasekar, Pedro Henrique Pamplona Savarese, Nathan Srebro, and Daniel Soudry. Convergence of gradient descent on separable data. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 3420–3428. PMLR, 2019b.
- Mor Shpigel Nacson, Nathan Srebro, and Daniel Soudry. Stochastic gradient descent on separable data: Exact convergence with a fixed learning rate. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 3051–3059. PMLR, 2019c.
- Samet Oymak and Mahdi Soltanolkotabi. Towards moderate overparameterization: global convergence guarantees for training shallow neural networks. *IEEE Journal on Selected Areas in Information Theory*, 2020.
- Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit bias of gradient descent on separable data. *The Journal of Machine Learning Research*, 19(1):2822–2878, 2018.
- Blake Woodworth, Suriya Gunasekar, Jason D Lee, Edward Moroshko, Pedro Savarese, Itay Golan, Daniel Soudry, and Nathan Srebro. Kernel and rich regimes in overparametrized models. In *Conference On Learning Theory*, 2020.
- Lei Wu, Qingcan Wang, and Chao Ma. Global convergence of gradient descent for deep linear residual networks. In *Advances in Neural Information Processing Systems*, pp. 13389–13398, 2019.
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Michael C Mozer, and Yoram Singer. Identity crisis: Memorization and generalization under extreme overparameterization. *arXiv preprint arXiv:1902.04698*, 2019.
- Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Stochastic gradient descent optimizes over-parameterized deep ReLU networks. *arXiv preprint arXiv:1811.08888*, 2018.

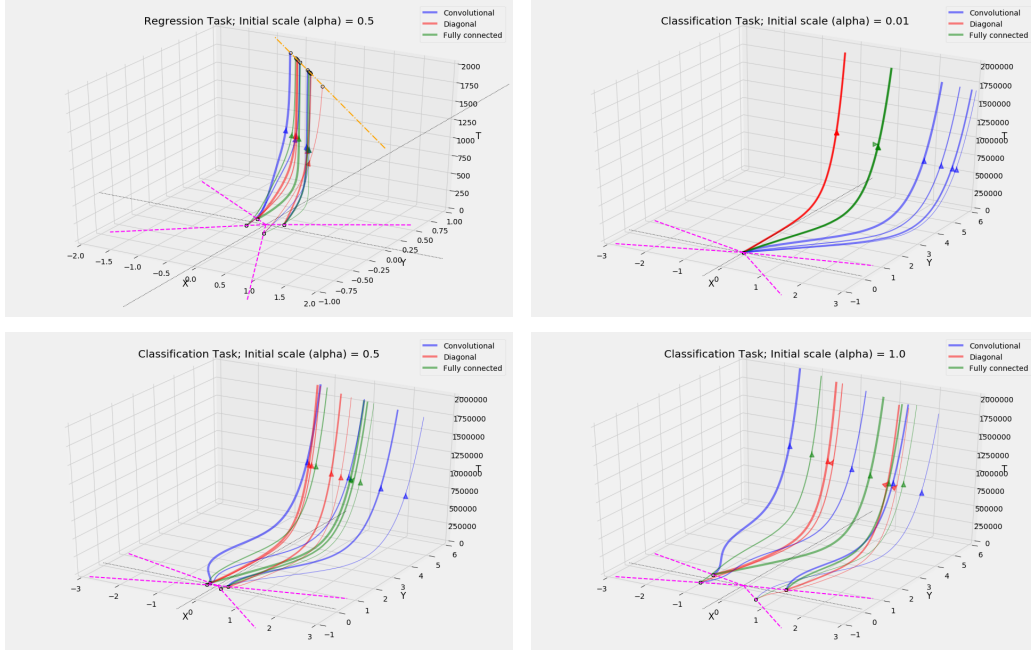


Figure 2: Gradient descent trajectories of linear coefficients of linear fully-connected, diagonal, and convolutional networks on a regression task with initial scale  $\alpha = 0.5$  (top left), and networks on a classification task with initial scales  $\alpha = 0.01, 0.5, 1$  (rest). Networks are initialized at the same coefficients (circles on purple lines), but follow different trajectories due to different implicit biases of networks induced from their architecture. The top left figure shows that our theoretical predictions on limit points (circles on yellow line, the set of global minima) agree with the solution found by GD. For details of the experimental setup, please refer to Section 6.

## A COROLLARIES ON SPECIFIC NETWORK ARCHITECTURES

We present corollaries obtained by specializing the theorems in the main text to specific network architectures. We briefly review the linear neural network architectures studied in this section.

**Linear fully-connected networks.** An  $L$ -layer linear fully-connected network is defined as

$$f_{\text{fc}}(\mathbf{x}; \mathbf{w}_{\text{fc}}) = \mathbf{x}^T \mathbf{W}_1 \cdots \mathbf{W}_{L-1} \mathbf{w}_L, \quad (10)$$

where  $\mathbf{W}_l \in \mathbb{R}^{d_l \times d_{l+1}}$  for  $l \in [L-1]$  (we use  $d_1 = d$ ) and  $\mathbf{w}_L \in \mathbb{R}^{d_L}$ .

**Linear diagonal networks.** An  $L$ -layer linear diagonal network is written as

$$f_{\text{diag}}(\mathbf{x}; \mathbf{w}_{\text{diag}}) = (\mathbf{x} \mathbf{W}_1 \cdots \mathbf{W}_{L-1})^T \mathbf{w}_L, \quad (11)$$

where  $\mathbf{w}_l \in \mathbb{R}^{d}$  for  $l \in [L]$ .

**Linear (circular) convolutional networks.** An  $L$ -layer linear convolutional network is written as

$$f_{\text{conv}}(\mathbf{x}; \mathbf{w}_{\text{conv}}) = ((\mathbf{x} \star \mathbf{w}_1) \star \mathbf{w}_2) \cdots \star \mathbf{w}_{L-1})^T \mathbf{w}_L, \quad (12)$$

where  $\mathbf{w}_l \in \mathbb{R}^{k_l}$  with  $k_l \leq d$  and  $k_L = d$ , and  $\star$  defines the circular convolution: for any  $\mathbf{a} \in \mathbb{R}^d$  and  $\mathbf{b} \in \mathbb{R}^k$  ( $k \leq d$ ), we have  $\mathbf{a} \star \mathbf{b} \in \mathbb{R}^d$  defined as  $[\mathbf{a} \star \mathbf{b}]_i = \sum_{j=1}^k [\mathbf{a}]_{(i+j-1) \bmod d} [\mathbf{b}]_j$ , for  $i \in [d]$ . In case of  $k_l = d$  for all  $l \in [L]$ , we refer to this network as *full-length* convolutional networks.

**Deep matrix sensing.** The deep matrix sensing problem considered in [Gunasekar et al. \(2017\)](#); [Arora et al. \(2019b\)](#) aims to minimize the following problem

$$\underset{\mathbf{W}_1, \dots, \mathbf{W}_L \in \mathbb{R}^{d \times d}}{\text{minimize}} \quad L_{\text{ms}}(\mathbf{W}_1, \dots, \mathbf{W}_L) := \sum_{i=1}^n (y_i - \mathbf{h}(\mathbf{A}_i, \mathbf{W}_1, \dots, \mathbf{W}_L))^2, \quad (13)$$

where the sensor matrices  $\mathbf{A}_1, \dots, \mathbf{A}_n \in \mathbb{R}^{d \times d}$  are symmetric. Following [Gunasekar et al. \(2017\)](#); [Arora et al. \(2019b\)](#), we consider sensor matrices  $\mathbf{A}_1, \dots, \mathbf{A}_n \in \mathbb{R}^{d \times d}$  that commute. To make the problem underdetermined, we assume that  $n > d$ , and  $\mathbf{A}_i$ 's are linearly independent.

## A.1 COROLLARY OF THEOREM 1

**Corollary 1.** Consider an  $L$ -layer linear fully-connected network (10). If the training loss satisfies  $L(\text{fc}(t_0)) < 1$  for some  $t_0 > 0$ , then  $\text{fc}(\text{fc}(t))$  converges in a direction that aligns with the solution of the following optimization problem

$$\text{minimize}_{\mathbf{z} \in \mathbb{R}^d} \|\mathbf{z}\|_2^2 \quad \text{subject to} \quad y_i \mathbf{x}_i^T \mathbf{z} = 1, \quad \forall i \in [n].$$

Corollary 1 shows that whenever the network separates the data correctly, the linear coefficients  $\text{fc}(\text{fc})$  convergence in direction to the  $\ell_2$  max-margin classifier. Note that this corollary does not require the directional convergence of  $\mathbf{X}^T \mathbf{r}$ , which is different from Theorem 1. In fact, this corollary also appears in Ji & Telgarsky (2020), but we provide an alternative proof based on our tensor formulation. The proof of Corollary 1 can be found in Appendix C.

## A.2 COROLLARIES OF THEOREM 2

**Corollary 2.** Consider an  $L$ -layer linear diagonal network (11). If there exists  $\lambda > 0$  such that the initial directions  $\mathbf{w}_1, \dots, \mathbf{w}_L$  of the network parameters satisfy  $[\mathbf{w}_l]_j^2 \geq \lambda [\mathbf{w}_L]_j^2$  for all  $l \in [L-1]$  and  $j \in [d]$ , then  $\text{diag}(\text{diag}(t))$  converges in a direction that aligns with a stationary point  $\mathbf{z}^\infty$  of

$$\text{minimize}_{\mathbf{z} \in \mathbb{R}^d} \|\mathbf{z}\|_{2/L} \quad \text{subject to} \quad y_i \mathbf{x}_i^T \mathbf{z} = 1, \quad \forall i \in [n].$$

For full-length convolutional networks, we define  $\mathbf{F} \in \mathbb{C}^{d \times d}$  to be the matrix of discrete Fourier transform basis  $[\mathbf{F}]_{j,k} = \frac{1}{\sqrt{d}} \exp(\frac{\sqrt{-1} \cdot 2\pi(j-1)(k-1)}{d})$ . Note that  $\mathbf{F}^* = \mathbf{F}^{-1}$ , and both  $\mathbf{F}$  and  $\mathbf{F}^*$  are symmetric, but not Hermitian.

**Corollary 3.** Consider an  $L$ -layer linear full-length convolutional network (12). If there exists  $\lambda > 0$  such that the initial directions  $\mathbf{w}_1, \dots, \mathbf{w}_L$  of the network parameters satisfy  $\sum_j [\mathbf{F} \mathbf{w}_l]_j^2 \geq \lambda \sum_j [\mathbf{F} \mathbf{w}_L]_j^2$  for all  $l \in [L-1]$  and  $j \in [d]$ , then  $\text{conv}(\text{conv}(t))$  converges in a direction that aligns with a stationary point  $\mathbf{z}^\infty$  of

$$\text{minimize}_{\mathbf{z} \in \mathbb{R}^d} \|\mathbf{F} \mathbf{z}\|_{2/L} \quad \text{subject to} \quad y_i \mathbf{x}_i^T \mathbf{z} = 1, \quad \forall i \in [n].$$

Corollary 2 shows that in the limit, linear diagonal network finds a sparse solution  $\mathbf{z}$  that is a stationary point of the  $\ell_{2/L}$  max-margin classification problem. Corollary 3 has a similar conclusion except that the standard  $\ell_{2/L}$  norm is replaced with DFT-domain  $\ell_{2/L}$  norm. By specifying mild conditions on initialization (see Remark 1), these corollaries remove the convergence assumptions required in Gunasekar et al. (2018b). The proofs of Corollaries 2 and 3 are in Appendix D.

## A.3 COROLLARY OF THEOREM 3

Recall that Theorem 3 can be applied to any 2-layer networks that can be represented as linear tensor networks. Examples include the convolutional networks that are not full-length (i.e., filter size  $k_1 < d$ ), which are not covered by the previous result (Gunasekar et al., 2018b). Here, we present the characterization of convergence directions of  $\text{conv}(\text{conv}(t))$  for small-filter cases:  $k_1 = 1$  and  $k_1 = 2$ .

**Corollary 4.** Consider a 2-layer linear convolutional network (12) with  $k_1 = 1$  and a single data point  $(\mathbf{x}, y)$ . If there exists  $\lambda > 0$  such that the initial directions  $\mathbf{w}_1$  and  $\mathbf{w}_2$  of the network parameters satisfy  $k \|\mathbf{x}\|_2^2 \|\mathbf{v}_1\|_2^2 \geq (\mathbf{x}^T \mathbf{v}_2)^2 \geq k \|\mathbf{x}\|_2^2 \lambda$ , then  $\text{conv}(\text{conv}(t))$  converges in direction that aligns with  $y \mathbf{x}$ .

Consider a 2-layer linear convolutional network (12) with  $k_1 = 2$  and a single data point  $(\mathbf{x}, y)$ . Let  $\mathbf{x} := [[\mathbf{x}]_2 \quad \dots \quad [\mathbf{x}]_d \quad [\mathbf{x}]_1]$ , and  $\mathbf{x}' := [[\mathbf{x}]_d \quad [\mathbf{x}]_1 \quad \dots \quad [\mathbf{x}]_{d-1}]$ . If there exists  $\lambda > 0$  such that the initial directions  $\mathbf{w}_1$  and  $\mathbf{w}_2$  of the network parameters satisfy

$$([\mathbf{v}_1]_1 + [\mathbf{v}_1]_2)^2 \frac{((\mathbf{x} + \mathbf{x}')^T \mathbf{v}_2)^2}{k \|\mathbf{x}\|_2^2 + \mathbf{x}^T \mathbf{x}} \geq \lambda, \quad \text{and} \quad ([\mathbf{v}_1]_1 - [\mathbf{v}_1]_2)^2 \frac{((\mathbf{x} - \mathbf{x}')^T \mathbf{v}_2)^2}{k \|\mathbf{x}\|_2^2 - \mathbf{x}^T \mathbf{x}} \geq \lambda,$$

then  $\text{conv}(\text{conv}(t))$  converges in a direction that aligns with a filtered version of  $\mathbf{x}$ :

$$\lim_{t \rightarrow \infty} \frac{\text{conv}(\text{conv}(t))}{k \|\text{conv}(\text{conv}(t))\|_2} \propto \begin{cases} 2y \mathbf{x} + y \mathbf{x}' + y \mathbf{x} & \text{if } \mathbf{x}^T \mathbf{x} > 0, \\ 2y \mathbf{x} - y \mathbf{x}' - y \mathbf{x} & \text{if } \mathbf{x}^T \mathbf{x} < 0. \end{cases}$$



Corollary 4 shows that if the filter size is  $k_1 = 1$ , then the limit direction of  $\text{conv}(\text{conv})$  is always the  $l_2$  max-margin classifier. Note that this is quite different from the case  $k_1 < d$  which converges to the DFT-domain  $l_1$  max-margin classifier. However, for  $k_1 < d$ , it is difficult to characterize the limit direction as the max-margin classifier of some common norms. Rather, the limit directions of  $\text{conv}(\text{conv})$  correspond to a “ltered” version of the data point, and the weights of the filter depend on the data point. For  $k_1 = 2$ , the filter is a low-pass filter if the autocorrelation of  $x$  is positive, and high-pass if the autocorrelation is negative. For  $k_1 > 2$ , the filter weights are more complicated to characterize in terms of  $x$  and the filter length increases as  $k_1$  increases. We prove Corollary 4 in Appendix E.

#### A.4 COROLLARIES OF THEOREM 5

In this subsection, we apply Theorem 5 to linear diagonal networks, linear full-length convolutional networks with even data, and deep matrix sensing. The proofs of the corollaries can be found in Appendix F.

Corollary 5. Consider an  $L$ -layer linear diagonal network (11). For some  $\epsilon > 0$ , choose any vector  $w \in \mathbb{R}^d$  satisfying  $|w_j| \geq \epsilon$  for all  $j \in [d]$ , and choose initial direction  $w_1 = w$  for  $l \in [L-1]$  and  $w_L = 0$ . Then, the linear coefficients  $\text{diag}(\text{diag}(t))$  converge to the solution  $z^*$  of

$$\text{minimize}_{z \in \mathbb{R}^d} Q_L; w(z) := \sum_{j=1}^d [w_j]^2 H_L \frac{|z_j|}{|w_j|} \quad \text{subject to} \quad Xz = y:$$

Recall that the original statement of Assumption 1 allows the matrices  $S_1, \dots, U_L$  to be complex, but Theorem 5 poses another assumption that these matrices are real. In applying Theorem 2 to convolutional networks to get Corollary 3, we used the fact that the data  $\text{Mso}(x)$  of a linear full-length convolutional network satisfies Assumption 1 with  $w = d^{\frac{L-1}{2}} F$  and  $U_1 = \dots = U_L = \frac{1}{d} F$ , where  $F \in \mathbb{C}^{d \times d}$  is the matrix of discrete Fourier transform basis  $[F]_{j,k} = \frac{1}{d} \exp(i \frac{2\pi}{d} (j-1)(k-1))$  and  $F^*$  is the complex conjugate of  $F$ . Note that these are complex matrices, so one cannot directly apply Theorem 5 to convolutional networks. However, it turns out that if the data and initialization are even, we can derive a corollary for convolutional networks.

We say that a vector  $x \in \mathbb{R}^d$  is even when it satisfies the even symmetry, as in even functions. More concretely, a vector  $x \in \mathbb{R}^d$  is even if  $[x]_{j+2} = [x]_{d-j}$  for  $j = 0, \dots, \lfloor \frac{d-3}{2} \rfloor$ ; i.e., the vector has the even symmetry around its “origin”  $[x]_1$ . From the definition of the matrix  $F \in \mathbb{C}^{d \times d}$ , it is straightforward to check that if  $x$  is real and even, then its DFT  $Fx$  is also real and even (see Appendix F.4 for details).

Corollary 6. Consider an  $L$ -layer linear full-length convolutional network (12). Assume that the data points  $x_i, g_{i=1}^n$  are all even. For some  $\epsilon > 0$ , choose any even vector  $w$  satisfying  $|Fw_j| \geq \epsilon$  for all  $j \in [d]$ , and choose initial direction  $w_1 = w$  for  $l \in [L-1]$  and  $w_L = 0$ . Then, the linear coefficients  $\text{conv}(\text{conv}(t))$  converge to the solution  $z^*$  of

$$\text{minimize}_{z \in \mathbb{R}^d; \text{even}} Q_L; Fw(Fz) := \sum_{j=1}^d [Fw_j]^2 H_L \frac{|Fz_j|}{|Fw_j|} \quad \text{subject to} \quad Xz = y:$$

Corollaries 5 and 6 show that the interpolation between minimum weighted  $l_1$  and weighted  $l_2$  solutions occurs for diagonal networks, and also for convolutional networks (in DFT domain, with the restriction of even symmetry). The conclusion of Corollary 5 is similar to the results in Woodworth et al. (2020), but the network architecture (11) considered in our corollary is a slightly different from the “differential” network  $f(x; w_+, w_-) = x^T (w_+^L - w_-^L)$  in Woodworth et al. (2020).

As mentioned in the main text, we can actually show that the matrix sensing result in Arora et al. (2019b) is a special case of our Theorem 5. Given any symmetric matrix  $M \in \mathbb{R}^{d \times d}$ , let  $\text{eig}(M) \in \mathbb{R}^d$  be the  $d$ -dimensional vector of eigenvalues of  $M$ .

Corollary 7. Consider the depth- $L$  deep matrix sensing problem (13). Let  $A_i$ 's be symmetric, and assume  $A_1, \dots, A_n$  commute. For  $\epsilon > 0$ , choose initialization  $W_1(0) = I_d$  for  $l \in [L-1]$  and  $W_L(0) = 0$ . Then, the product  $W_L(t) \dots W_1(t)$  converge to the solution  $M^*$  of

$$\text{minimize}_{M \in \mathbb{R}^{d \times d}; \text{symmetric}} Q_L; (\text{eig}(M)) := \sum_{j=1}^d H_L \frac{|\text{eig}(M)_j|}{L} \quad \text{subject to} \quad L_{\text{ms}}(M) = 0:$$

Under an additional assumption that  $A_i$ 's are positive semidefinite, Theorem 2 in Arora et al. (2019b) studies the initialization  $W_l(0) = I_{d_l}$  for all  $l \in [L]$ , and shows that the limit point of  $W_L \cdots W_1$  converges to the minimum nuclear norm solution as  $\delta \rightarrow 0$ . We remove the assumption of positive definiteness of  $A_i$ 's and let  $W_L(0) = 0$ , to show a complete characterization of the solution found by gradient flow, which interpolates between the minimum nuclear norm (i.e., Schatten 1-norm) solution (when  $\delta = 0$ ) and the minimum Frobenius norm (i.e., Schatten 2-norm) solution (when  $\delta \rightarrow \infty$ ).

## B TENSOR REPRESENTATION OF FULLY-CONNECTED NETWORKS

In Section 3, we only defined the data tensor  $M_{fc}(x)$  of fully-connected networks for  $d_r = 2$ . Here, we describe an iterative procedure constructing the data tensor for deep fully-connected networks.

We start with  $T_1(x) := x \otimes R^{d_1}$ . Next, define a block diagonal matrix  $T_2(x) \in \mathbb{R}^{d_1 d_2 \times d_1 d_2}$  where the ‘‘diagonals’’  $[T_2(x)]_{d_1(j-1)+1:d_1 j, d_1(j-1)+1:d_1 j} = T_1(x)$  for  $j \in [d_2]$ , while all the other entries are filled with 0. We continue this ‘‘block diagonal’’ procedure, as the following. Having defined  $T_{l-1}(x) \in \mathbb{R}^{d_1 d_2 \times \dots \times d_1 d_2}$ ,

1. Define  $T_l(x) \in \mathbb{R}^{d_1 d_2 \times \dots \times d_1 d_1 d_1}$ .
2. Set  $[T_l(x)]_{d_1(j-1)+1:d_1 j, d_1(j-1)+1:d_1 j} = T_{l-1}(x)$ ;  $j \in [d_1]$ .
3. Set all the remaining entries of  $T_l(x)$  to zero.

We repeat this process for  $l = 2, \dots, L$ , and set  $M_{fc}(x) := T_L(x)$ . By defining the parameters of the tensor formulation  $v_l = \text{vec}(W_l)$  for  $l \in [L-1]$  and  $v_L = w_L$ , and using the tensor  $M(x) = M_{fc}(x)$ , we can check the equivalence of (2) and (5).

## C PROOFS OF THEOREM 1 AND COROLLARY 1

### C.1 PROOF OF THEOREM 1

The proof of Theorem 1 is outlined as follows. First, using the directional convergence and alignment results in Ji & Telgarsky (2020), we prove that each of our network parameters converges in direction, and it aligns with its corresponding negative gradient  $-v_l$ . Then, we prove that the directions of  $v_l$ 's are actually singular vectors of  $M(u^1)$ , where  $u^1 := \lim_{t \rightarrow 1} \frac{X^T r(t)}{\|X^T r(t)\|_2}$ .

Since a linear tensor network is an homogeneous polynomial of  $v_1, \dots, v_L$ , it satisfies the assumptions required for Theorems 3.1 and 4.1 in Ji & Telgarsky (2020). These theorems imply that if the gradient flow satisfies  $L'(t_0) < 1$  for some  $t_0 > 0$ , then  $(t)$  converges in direction, and the direction aligns with  $-L'(t)$ ; that is,

$$\lim_{t \rightarrow 1} \|L'(t)\|_2 = 1; \quad \lim_{t \rightarrow 1} \frac{L'(t)}{\|L'(t)\|_2} = -1; \quad \lim_{t \rightarrow 1} \frac{(L'(t))^T r}{\|L'(t)\|_2 \|r\|_2} = -1. \quad (14)$$

For linear tensor networks (6), the parameter is the concatenation of all parameter vectors  $v_1, \dots, v_L$ , so (14) holds for  $v = v_1^T \cdots v_L^T$ .

Now, recall that by the definition of the linear tensor network, we have the following gradient flow

$$\dot{v}_l = M(X^T r) (v_1; \dots; v_{l-1}; I_{k_l}; v_{l+1}; \dots; v_L);$$

Note that we can apply this to calculate the rate of growth of  $\|v_l\|_2^2$ :

$$\begin{aligned} \frac{d}{dt} \|v_l\|_2^2 &= 2v_l^T \dot{v}_l = 2v_l^T M(X^T r) (v_1; \dots; v_{l-1}; I_{k_l}; v_{l+1}; \dots; v_L) \\ &= 2M(X^T r) (v_1; \dots; v_{l-1}; v_l; v_{l+1}; \dots; v_L) \\ &= \frac{d}{dt} \|v_l\|_2^2 \quad \text{for any } l \in [L], \end{aligned}$$

so the rate at which  $k_2^2$  grows over time is the same for all layers  $[L]$ . By the definition of (14), we have

$$k_2^2 = \prod_{l=1}^L kv_l k_2^2 \neq 1 ;$$

which then implies

$$\lim_{t \rightarrow \infty} kv_l(t)k_2 \neq 1 ; \quad \lim_{t \rightarrow \infty} \frac{k_2(t)k_2}{kv_l(t)k_2} = \frac{k_2(t)k_2^2}{kv_l(t)k_2^2} = \frac{1}{L} ;$$

for all  $l \in [L]$ . Now, let  $I_1$  be the set of indices that correspond to the components of  $r$ . It follows from (14) that

$$\lim_{t \rightarrow \infty} \frac{v_l(t)}{kv_l(t)k_2} = \lim_{t \rightarrow \infty} \frac{v_l(t) k_2(t)k_2}{k_2(t)k_2 kv_l(t)k_2} = \lim_{t \rightarrow \infty} \frac{[v_l(t)]_{I_1} k_2(t)k_2}{k_2(t)k_2 kv_l(t)k_2} = \frac{1}{L} \mathbb{1}_{I_1} ;$$

thus showing the directional convergence of  $r$ .

Next, it follows from directional convergence of  $r$  and its alignment with  $L(r)$  (14) that  $r - L(r)$  also converges in direction, in the opposite direction of  $r$ . By comparing the components in  $I_1$ 's, we get that  $v_l L(r)$  converges in the opposite direction of  $r$ .

For any  $l \in [L]$ , now let  $v_l^1 := \lim_{t \rightarrow \infty} \frac{v_l(t)}{kv_l(t)k_2}$ . Also recall the assumption that  $X^T r(t)$  converges in direction; let the unit vector  $u^1 := \lim_{t \rightarrow \infty} \frac{X^T r(t)}{kX^T r(t)k_2}$  be the limit direction. By the gradient flow dynamics of  $v_l$ , we have

$$v_l^1 / r - v_l L(r) = M(u^1) (v_1^1 ; \dots ; v_{l-1}^1 ; l_{k_l} ; v_{l+1}^1 ; \dots ; v_L^1) ;$$

for all  $l \in [L]$ . Note that this equation has the same form as (8), the definition of singular vectors in tensors. So this proves that  $(v_1^1 ; \dots ; v_L^1)$  are singular vectors of  $M(u^1)$ .

### C.2 PROOF OF COROLLARY 1

The proof proceeds as follows. First, we will show using the structure of the data matrix that the limit direction of linear coefficients  $f_c(\frac{1}{f_c})$  is proportional to  $cu^1$ , where  $c$  is a nonzero scalar and  $u^1$  is the limit direction of  $X^T r$ . Then, through a closer look at  $f_c$  and  $c$ , we will prove that  $f_c(\frac{1}{f_c})$  is in fact a conic combination of the support vectors (i.e., the data points with the minimum margins). Finally, we will compare  $f_c(\frac{1}{f_c})$  with the KKT conditions of the  $\ell_2$  max-margin classification problem and conclude that  $f_c(\frac{1}{f_c})$  must be in the same direction as the max-margin classifier.

Due to the way how the data tensor  $f_c$  is constructed for fully-connected networks (Appendix B), we always have

$$r - v_l L(f_c) = M_{f_c}(X^T r) (l_{k_1} ; v_2 ; \dots ; v_L) \in \text{span} \left\{ \begin{matrix} 0 & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{matrix} \right\} \cup \text{span} \left\{ \begin{matrix} 0 & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{matrix} \right\} ;$$

From Theorem 1, we have directional convergence of  $r$  and its alignment with  $v_l L(f_c)$ . This means that the limit direction  $u^1$ , which is a fixed vector, must be also in the span of vectors written above. This implies that  $X^T r$  must also converge to some direction, say  $u^1 := \lim_{t \rightarrow \infty} \frac{X^T r(t)}{kX^T r(t)k_2}$ .

Now recall the definition of  $v_l$  in case of the fully-connected network  $k_1 = \text{vec}(W_1)$ . So, by reshaping  $v_l^1$  into its original  $d_1 \times d_2$  matrix form  $W_1^1$ , we have

$$W_1^1 / u^1 q^T ;$$

for some  $q \in \mathbb{R}^{d_2}$ . This implies that the linear coefficients  $f_c(\frac{1}{f_c})$  of the network converge in direction to

$$f_c(\frac{1}{f_c}) = W_1^1 W_2^1 \dots W_{L-1}^1 W_L^1 / u^1 q^T W_2^1 \dots W_{L-1}^1 W_L^1 = cu^1 ; \quad (15)$$

where  $c$  is some nonzero real number.

Let us now take a closer look at the vector  $u^1$ , the limit direction of  $X^T r$ . Recall from Section 2.1 that for any  $i \in [n]$ ,

$$[r]_i = y_i \exp(-y_i f_{fc}(x_i; \frac{1}{fc})) = y_i \exp(-y_i x_i^T f_{fc}(\frac{1}{fc}));$$

in case of classification. Recall that  $\|f_{fc}(t)\|_2 \rightarrow 1$  while converging to a certain direction  $f_{fc}(\frac{1}{fc})$ . This means that if

$$y_j x_j^T f_{fc}(\frac{1}{fc}) > y_i x_i^T f_{fc}(\frac{1}{fc})$$

for any  $i, j \in [n]$ , then

$$\lim_{t \rightarrow \infty} \frac{\exp(-y_j x_j^T f_{fc}(\frac{1}{fc}(t)))}{\exp(-y_i x_i^T f_{fc}(\frac{1}{fc}(t)))} = 0: \tag{16}$$

Take  $i$  to be the index of any support vector, i.e., any that attains the minimum  $y_i x_i^T f_{fc}(\frac{1}{fc})$  among all data points. Using such an observation (16) implies that  $\lim_{t \rightarrow \infty} [r(t)]_j = 0$  for any  $x_j$  that is not a support vector. Thus, by the argument above can in fact be written as

$$u^1 = \lim_{t \rightarrow \infty} \frac{\sum_{i=1}^n x_i [r(t)]_i}{\sum_{i=1}^n x_i [r(t)]_i} = \sum_{i=1}^n y_i x_i; \tag{17}$$

where  $\alpha_i \geq 0$  for all  $i \in [n]$ , and  $\alpha_j = 0$  for  $x_j$ 's that are not support vectors. Combining (17) and (15),

$$f_{fc}(\frac{1}{fc}) / c = \sum_{i=1}^n y_i x_i; \tag{18}$$

Recall that we do not yet know whether  $c$  introduced in (15), is positive or negative; we will now show that  $c$  has to be negative. From Lyu & Li (2020), we know that  $\|f_{fc}(t)\|_2 \rightarrow 1$ , which implies that  $y_i x_i^T f_{fc}(\frac{1}{fc}) > 0$  for all  $i \in [n]$ . However, if  $c > 0$ , then (18) implies that  $f_{fc}(\frac{1}{fc})$  is inside a cone  $K$  defined as

$$K := \left\{ \sum_{i=1}^n y_i x_i \mid \alpha_i \geq 0; \forall i \in [n] \right\};$$

Note that the polar cone of  $K$ , denoted as  $K^\circ$ , is

$$K^\circ := \{ z \mid z^T z \geq 0; \forall z \in K \} = \{ z \mid \sum_{i=1}^n y_i x_i^T z \geq 0; \forall i \in [n] \};$$

It is known that  $K \cap K^\circ = \{0\}$  for any convex cone  $K$  and its polar cone  $K^\circ$ . Therefore, having  $c > 0$  implies that  $f_{fc}(\frac{1}{fc}) \in K \cap K^\circ$ , which means that there exists some  $i \in [n]$  such that  $y_i x_i^T f_{fc}(\frac{1}{fc}) < 0$ ; this contradicts the fact that the loss goes to zero as  $t \rightarrow \infty$ . Therefore,  $c$  in (15) and (18) must be negative:

$$f_{fc}(\frac{1}{fc}) / c = - \sum_{i=1}^n y_i x_i; \tag{19}$$

for  $\alpha_i \geq 0$  for all  $i \in [n]$  and  $\alpha_j = 0$  for all  $x_j$ 's that are not support vectors.

Finally, compare (19) with the KKT conditions of the following optimization problem:

$$\text{minimize}_z \|z\|_2^2 \quad \text{subject to} \quad y_i x_i^T z \leq 1; \quad \forall i \in [n];$$

The KKT conditions of this problem are

$$z = \sum_{i=1}^n \alpha_i y_i x_i; \quad \text{and} \quad \alpha_i \geq 0; \quad \alpha_i (1 - y_i x_i^T z) = 0 \quad \text{for all } i \in [n];$$

where  $\alpha_1, \dots, \alpha_n$  are the dual variables. Note that this is (up to scaling) satisfied by  $f_{fc}(\frac{1}{fc})$  (19), if we replace  $\alpha_i$ 's with  $\alpha_i$ 's. This finishes the proof that  $f_{fc}(\frac{1}{fc})$  is aligned with the  $\ell_2$  max-margin classifier.

## D PROOFS OF THEOREM 2 AND COROLLARIES 2 & 3

### D.1 PROOF OF THEOREM 2

#### D.1.1 CONVERGENCE OF LOSS TO ZERO

Since Theorem 2 does not assume the existence of  $\theta$  satisfying  $L(\theta) < 1$ , we need to first show that given the conditions on initialization, the training loss  $\mathcal{L}(t)$  converges to zero. Recall from Section 2.1 that

$$\mathbf{v}_l = \mathbf{r} \cdot \mathbf{v}_l L(\theta) = \mathbf{M}(\mathbf{X}^T \mathbf{r}) \cdot (\mathbf{v}_1; \dots; \mathbf{v}_{l-1}; I_{k_l}; \mathbf{v}_{l+1}; \dots; \mathbf{v}_L):$$

Applying the structure (9) in Assumption 1, we get

$$\begin{aligned} \mathbf{v}_l &= \mathbf{M}(\mathbf{X}^T \mathbf{r}) \cdot (\mathbf{v}_1; \dots; \mathbf{v}_{l-1}; I_{k_l}; \mathbf{v}_{l+1}; \dots; \mathbf{v}_L) \\ &= \prod_{j=1}^n [\mathbf{S} \mathbf{X}^T \mathbf{r}]_j \cdot (\mathbf{v}_1^T [\mathbf{U}_1]_{\cdot j} \quad \mathbf{v}_{l-1}^T [\mathbf{U}_{l-1}]_{\cdot j} \quad [\mathbf{U}_l]_{\cdot j} \quad \mathbf{v}_{l+1}^T [\mathbf{U}_{l+1}]_{\cdot j} \quad \mathbf{v}_L^T [\mathbf{U}_L]_{\cdot j}) \\ &= \prod_{j=1}^n [\mathbf{S} \mathbf{X}^T \mathbf{r}]_j \cdot \prod_{k \in I} [\mathbf{U}_k^T \mathbf{v}_k]_{\cdot j} \cdot [\mathbf{U}_l]_{\cdot j} \end{aligned}$$

Left-multiplying  $\mathbf{U}_l^H$  (the conjugate transpose of  $\mathbf{U}_l$ ) to both sides, we get

$$\mathbf{U}_l^H \mathbf{v}_l = \mathbf{S} \mathbf{X}^T \mathbf{r} \cdot \prod_{k \in I} \mathbf{U}_k^T \mathbf{v}_k; \quad (20)$$

where  $\cdot$  denotes the product using entry-wise multiplication

Now consider the rate of growth for the absolute value squared of the component of  $\mathbf{U}_l^T \mathbf{v}_l$ :

$$\begin{aligned} \frac{d}{dt} j[\mathbf{U}_l^T \mathbf{v}_l]_j^2 &= \frac{d}{dt} [\mathbf{U}_l^T \mathbf{v}_l]_j [\mathbf{U}_l^T \mathbf{v}_l]_j = \frac{d}{dt} [\mathbf{U}_l^T \mathbf{v}_l]_j [\mathbf{U}_l^H \mathbf{v}_l]_j \\ &= [\mathbf{U}_l^T \mathbf{v}_l]_j [\mathbf{U}_l^H \mathbf{v}_l]_j + [\mathbf{U}_l^H \mathbf{v}_l]_j [\mathbf{U}_l^T \mathbf{v}_l]_j \\ &= 2 \operatorname{Re} [\mathbf{U}_l^H \mathbf{v}_l]_j [\mathbf{U}_l^T \mathbf{v}_l]_j \\ &= 2 \operatorname{Re} [\mathbf{S} \mathbf{X}^T \mathbf{r}]_j \cdot \prod_{k=1}^L [\mathbf{U}_k^T \mathbf{v}_k]_j \\ &= \frac{d}{dt} j[\mathbf{U}_l^T \mathbf{v}_l]_j^2 \quad \text{for any } l \in [L], \end{aligned}$$

so for any  $j \in [m]$ , the squared absolute value of the components in  $\mathbf{U}_l^T \mathbf{v}_l$  grow at the same rate for each layer  $l \in [L]$ . This means that the gap between any two different layers stays constant for all  $t \geq 0$ . Combining this with our conditions on initial directions, we have

$$\begin{aligned} j[\mathbf{U}_l^T \mathbf{v}_l(t)]_j^2 - j[\mathbf{U}_L^T \mathbf{v}_L(t)]_j^2 &= j[\mathbf{U}_l^T \mathbf{v}_l(0)]_j^2 - j[\mathbf{U}_L^T \mathbf{v}_L(0)]_j^2 \\ &= \epsilon_j^2 j[\mathbf{U}_l^T \mathbf{v}_l]_j^2 - \epsilon_j^2 j[\mathbf{U}_L^T \mathbf{v}_L]_j^2 \leq \epsilon_j^2; \end{aligned} \quad (21)$$

for any  $j \in [m]$ ,  $l \in [L-1]$ , and  $t \geq 0$ . This inequality also implies

$$j[\mathbf{U}_l^T \mathbf{v}_l(t)]_j^2 - j[\mathbf{U}_L^T \mathbf{v}_L(t)]_j^2 + \epsilon_j^2 \leq \epsilon_j^2; \quad (22)$$

Let us now consider the time derivative of  $\mathcal{L}(t)$ . We have the following chain of upper bounds on the time derivative:

$$\begin{aligned} \frac{d}{dt} \mathcal{L}(t) &= \mathbf{r} \cdot \mathcal{L}(t)^T - \mathcal{L}(t) = k \mathbf{r} \cdot \mathcal{L}(t) k_2^2 \\ &= k \mathbf{v}_L \mathcal{L}(t) k_2^2 = k \mathbf{v}_L(t) k_2^2 \\ &\stackrel{(a)}{=} k \mathbf{U}_L^H \mathbf{v}_L(t) k_2^2 \stackrel{(b)}{=} \mathbf{S} \mathbf{X}^T \mathbf{r}(t) \cdot \prod_{k \in L} \mathbf{U}_k^T \mathbf{v}_k(t) k_2^2 \\ &= \prod_{j=1}^n j[\mathbf{S} \mathbf{X}^T \mathbf{r}(t)]_j^2 \cdot \prod_{k \in L} j[\mathbf{U}_k^T \mathbf{v}_k(t)]_j^2 \end{aligned}$$



$$\begin{aligned}
 & \text{(c)} \quad \sum_{j=1}^m \| [S X^T r(t)]_j \|^2 \\
 & = \| S X^T r(t) \|_2^2 \\
 & \text{(d)} \quad \sum_{i=1}^m s_{\min}(S)^2 \| X^T r(t) \|_2^2; \tag{23}
 \end{aligned}$$

where (a) used the fact that  $\| U_L U_L^H v_L(t) \|_2^2 = \| v_L(t) \|_2^2$  because it is a projection onto a subspace, and  $\| U_L U_L^H v_L(t) \|_2^2 = \| U_L^H v_L(t) \|_2^2$  because  $U_L^H U_L = I_{k_L}$ ; (b) is due to (20); (c) is due to (22); and (d) used the fact that  $S \in \mathbb{R}^{m \times d}$  is a matrix that has full column rank, so for any  $v \in \mathbb{R}^d$ , we can use  $\| S v \|_2 = s_{\min}(S) \| v \|_2$  where  $s_{\min}(S)$  is the minimum singular value of  $S$ .

We now prove a lower bound on the quantity  $\| X^T r(t) \|_2^2$ . Recall from Section 2.1 the definition of  $[r(t)]_i = y_i \exp(-y_i f(x_i; \theta(t)))$  for classification problems. Also, recall the assumption that the dataset is linearly separable, which means that there exists a unit vector  $z \in \mathbb{R}^d$  such that

$$y_i x_i^T z > 0$$

holds for all  $i \in [n]$ , for some  $\gamma > 0$ . Using these,

$$\begin{aligned}
 \| X^T r(t) \|_2^2 & = \left\| \sum_{i=1}^n y_i x_i \exp(-y_i f(x_i; \theta(t))) \right\|_2^2 \\
 & \geq \left[ z^T \sum_{i=1}^n y_i x_i \exp(-y_i f(x_i; \theta(t))) \right]^2 \\
 & \geq \left[ \sum_{i=1}^n \exp(-y_i f(x_i; \theta(t))) \right]^2 = 2L(\theta(t))^2.
 \end{aligned}$$

Combining this with (23), we get

$$\frac{d}{dt} L(\theta(t)) \geq \sum_{i=1}^m s_{\min}(S)^2 2L(\theta(t))^2;$$

which implies

$$L(\theta(t)) \geq \frac{L(\theta(0))}{1 - \sum_{i=1}^m s_{\min}(S)^2 2t}.$$

Therefore,  $L(\theta(t)) \rightarrow \infty$  as  $t \rightarrow \frac{1}{2 \sum_{i=1}^m s_{\min}(S)^2}$ .

### D.1.2 CHARACTERIZING THE LIMIT DIRECTION

Since we proved that  $L(\theta(t)) \rightarrow \infty$ , the argument in the proof of Theorem 1 applies to this case, and shows that the parameters converge in direction and align with  $u = r_{v_1} L(\theta)$ . Let  $v_1^1 := \lim_{t \rightarrow \infty} \frac{v_1(t)}{\|v_1(t)\|_2}$  be the limit direction of  $v_1$ .

The remaining steps of the proof are as follows. We first prove  $S X^T r(t)$  converges in direction  $u^1$ . Using this  $u^1$ , we derive a number of conditions that has to be satisfied by the limit directions of the parameters. Finally, we compare these conditions with the KKT conditions of the minimization problem, and finish the proof.

By Assumption 1, we have

$$\begin{aligned}
 f(x; \theta) & = M(x; (v_1; \dots; v_L)) = \sum_{j=1}^n [S x]_j \sum_{l=1}^L [U_l^T v_l]_j \\
 & = \sum_{j=1}^n \sum_{l=1}^L [U_l^T v_l]_j [S]_{j, x} = x^T S^T \sum_{l \in [L]} U_l^T v_l = x^T S^T \theta.
 \end{aligned}$$

Here, we defined  $\theta := \sum_{l \in [L]} U_l^T v_l \in \mathbb{R}^m$ . Since the linear coefficients must be real, we have  $S^T \theta \in \mathbb{R}^d$  for any real  $v_l$ 's. Since  $v_l$ 's converge in direction,  $\theta$  also converges in direction, to  $\theta^1 := \sum_{l \in [L]} U_l^T v_l^1$ . So we can express the limit direction of  $(\theta)$  as

$$(\theta^1) / \| \theta^1 \| = S^T \theta^1 / \| S^T \theta^1 \|; \tag{24}$$

From (20) and alignment of  $\mathbf{u}_l$  and  $\mathbf{v}_l$ , we have

$$\lim_{t \uparrow 1} \mathbf{U}_l^H \mathbf{v}_l(t) = \lim_{t \uparrow 1} (\mathbf{U}_l^T \mathbf{v}_l(t)) / \lim_{t \uparrow 1} \sum_{k \in I} \mathbf{S} \mathbf{X}^T \mathbf{r}(t) \mathbf{U}_k^T \mathbf{v}_k(t); \quad (25)$$

Since all vectors  $\mathbf{U}_l^T \mathbf{v}_l(t)$  converge in direction, the term  $\sum_{k \in I} \mathbf{S} \mathbf{X}^T \mathbf{r}(t) \mathbf{U}_k^T \mathbf{v}_k(t)$  should also converge in direction. Let  $\mathbf{u}^1 := \lim_{t \uparrow 1} \frac{\sum_{k \in I} \mathbf{S} \mathbf{X}^T \mathbf{r}(t) \mathbf{U}_k^T \mathbf{v}_k(t)}{\|\sum_{k \in I} \mathbf{S} \mathbf{X}^T \mathbf{r}(t) \mathbf{U}_k^T \mathbf{v}_k(t)\|_2}$ . One can use the same argument as in Appendix C.2, more specifically (16) and (17), to show that  $\mathbf{u}^1$  can be written as

$$\mathbf{u}^1 = \lim_{t \uparrow 1} \frac{\sum_{i=1}^n \mathbf{S} \mathbf{P}_{i=1}^n \mathbf{x}_i[r(t)]_i}{\|\sum_{i=1}^n \mathbf{S} \mathbf{P}_{i=1}^n \mathbf{x}_i[r(t)]_i\|_2} = \sum_{i=1}^n y_i \mathbf{x}_i; \quad (26)$$

where  $y_i \geq 0$  for all  $i \in [n]$ , and  $y_j = 0$  for  $\mathbf{x}_j$ 's that are not support vectors, i.e., those satisfying  $y_j \mathbf{x}_j^T \mathbf{S}^T \mathbf{1} > \min_{i \in [n]} y_i \mathbf{x}_i^T \mathbf{S}^T \mathbf{1}$ .

Using  $\mathbf{u}^1$ , we can rewrite (25) as

$$\mathbf{U}_l^H \mathbf{v}_l^1 / \|\mathbf{u}^1\| = \sum_{k \in I} \mathbf{U}_k^T \mathbf{v}_k^1;$$

for all  $l \in [L]$ . Element-wise multiplying  $\mathbf{U}_l^T \mathbf{v}_l^1$  to both sides gives

$$\mathbf{U}_l^T \mathbf{v}_l^1 \|\mathbf{u}^1\| = \sum_{k \in I} \mathbf{U}_l^T \mathbf{v}_l^1 \|\mathbf{u}^1\| = \|\mathbf{u}^1\| \|\mathbf{U}_l^T \mathbf{v}_l^1\|; \quad (27)$$

where  $\|\cdot\|$  denotes element-wise  $l^2$  norm of the vector  $\cdot$ . Since the LHS of (27) is a positive real number, we have

$$\arg(\|\mathbf{U}_l^T \mathbf{v}_l^1\|) = 0 = \arg(\|\mathbf{u}^1\|) + \arg(\|\mathbf{U}_l^T \mathbf{v}_l^1\|); \quad (28)$$

so using this, (27) becomes

$$\|\mathbf{U}_l^T \mathbf{v}_l^1\| \|\mathbf{u}^1\| = \|\mathbf{u}^1\| \|\mathbf{U}_l^T \mathbf{v}_l^1\|; \quad (29)$$

Now element-wise multiply (29) for all  $l \in [L]$ , then we get

$$\|\mathbf{u}^1\| \|\mathbf{U}_l^T \mathbf{v}_l^1\| = \|\mathbf{U}_l^T \mathbf{v}_l^1\| \|\mathbf{u}^1\|; \quad (30)$$

A close look at (30) reveals that if  $\|\mathbf{u}^1\| \neq 0$  and  $\|\mathbf{U}_l^T \mathbf{v}_l^1\| \neq 0$  must satisfy that

$$\|\mathbf{U}_l^T \mathbf{v}_l^1\| \|\mathbf{u}^1\| = \|\mathbf{u}^1\| \|\mathbf{U}_l^T \mathbf{v}_l^1\|; \quad (31)$$

for all  $l \in [L]$ . There is another condition that has to be satisfied when

$$\|\mathbf{U}_l^T \mathbf{v}_l^1\| = 0; \|\mathbf{u}^1\| \neq 0 \Rightarrow \|\mathbf{U}_l^T \mathbf{v}_l^1\| = 0; \quad (32)$$

for any  $l \in [L]$ ; let us prove why. First, consider the time derivative  $\frac{d}{dt} \|\mathbf{U}_l^T \mathbf{v}_l^1\| = \|\mathbf{U}_l^T \dot{\mathbf{v}}_l^1\| \|\mathbf{U}_l^T \mathbf{v}_l^1\|$ .

$$\begin{aligned} \frac{d}{dt} \|\mathbf{U}_l^T \mathbf{v}_l^1\| &= \|\mathbf{U}_l^T \dot{\mathbf{v}}_l^1\| \|\mathbf{U}_l^T \mathbf{v}_l^1\| + \|\mathbf{U}_l^T \mathbf{v}_l^1\| \frac{d}{dt} \|\mathbf{U}_l^T \mathbf{v}_l^1\| \\ &\stackrel{(a)}{=} \|\mathbf{S} \mathbf{X}^T \mathbf{r}(t)\|_j (\|\mathbf{U}_l^T \mathbf{v}_l^1(t)\|_j^2 + \|\mathbf{U}_l^T \mathbf{v}_l^2(t)\|_j^2); \end{aligned} \quad (33)$$

where (a) used (20). Now consider

$$\frac{\frac{d}{dt} \|\mathbf{U}_l^T \mathbf{v}_l^1\|_j}{\|\mathbf{S} \mathbf{X}^T \mathbf{r}(t)\|_j \|\mathbf{U}_l^T \mathbf{v}_l^1\|_j} = \frac{\|\mathbf{S} \mathbf{X}^T \mathbf{r}(t)\|_j (\|\mathbf{U}_l^T \mathbf{v}_l^1(t)\|_j^2 + \|\mathbf{U}_l^T \mathbf{v}_l^2(t)\|_j^2)}{\|\mathbf{S} \mathbf{X}^T \mathbf{r}(t)\|_j \|\mathbf{U}_l^T \mathbf{v}_l^1\|_j^2}; \quad (34)$$

We want to compare this quantity for different  $l \in [L]$ . Before we do that, we take a look at the last term in the RHS of (34). Recall from (21) that

$$\|\mathbf{U}_l^T \mathbf{v}_l^1(t)\|_j^2 = \|\mathbf{U}_l^T \mathbf{v}_l^2(t)\|_j^2 + \|\mathbf{U}_l^T \mathbf{v}_l^1(0)\|_j^2 - \|\mathbf{U}_l^T \mathbf{v}_l^2(0)\|_j^2; \quad (35)$$

For simplicity, let  $\beta_j := \|\mathbf{U}_l^T \mathbf{v}_l^1(0)\|_j^2 - \|\mathbf{U}_l^T \mathbf{v}_l^2(0)\|_j^2$ , which is a positive number due to our assumption on initialization. Then, we can use (35) in (34) to get  $\frac{d}{dt} \|\mathbf{U}_l^T \mathbf{v}_l^1\|_j = \|\mathbf{U}_l^T \mathbf{v}_l^1(t)\|_j \|\mathbf{U}_l^T \mathbf{v}_l^2(t)\|_j$  to show that

$$\frac{\|\mathbf{U}_l^T \mathbf{v}_l^1(t)\|_j^2 + \|\mathbf{U}_l^T \mathbf{v}_l^2(t)\|_j^2}{\|\mathbf{U}_l^T \mathbf{v}_l^1(t)\|_j} = \frac{2\|\mathbf{U}_l^T \mathbf{v}_l^2(t)\|_j^2 + \beta_j}{\|\mathbf{U}_l^T \mathbf{v}_l^2(t)\|_j \|\mathbf{U}_l^T \mathbf{v}_l^1(t)\|_j + \beta_j} \geq 2;$$

$$\lim_{t \rightarrow \infty} \frac{j[U_1^T v_1(t)]_j^2 + j[U_2^T v_2(t)]_j^2}{j[(t)]_j} = 2 \quad \text{if} \quad \lim_{t \rightarrow \infty} j[U_2^T v_2(t)]_j = 1 :$$

Recall that we want to prove that (32) should necessarily hold. For the sake of contradiction, suppose that there exist  $\epsilon > 0$  that satisfies  $j[u^1]_j > j[u^1]_{oj}$ , for some  $j \in [m]$  satisfying  $j[u^1]_{oj} \neq 0$ . Note that having  $j[u^1]_j > j[u^1]_{oj}$  and  $j[u^1]_{oj} \neq 0$  implies that  $j[(t)]_j \neq 1$  and  $\frac{j[(t)]_j}{j[(t)]_{oj}} \neq 0$ . We now want to compare the ratio of (34) for  $j$  and  $j^o$ . First, note that

$$\lim_{t \rightarrow \infty} \frac{j[SX^T r(t)]_j = kSX^T r(t)k_2}{j[SX^T r(t)]_{j^o} = kSX^T r(t)k_2} = \frac{j[u^1]_j}{j[u^1]_{oj}} > 1: \quad (36)$$

Next, using  $\frac{j[(t)]_j}{j[(t)]_{oj}} \neq 0$  and the fact that  $f(x) = \frac{2x^2 + \epsilon}{x^2 + \epsilon}$  is a decreasing function of  $x > 0$  for any  $\epsilon > 0$ , we have

$$\frac{(j[U_1^T v_1(t)]_j^2 + j[U_2^T v_2(t)]_j^2) / j[(t)]_j}{(j[U_1^T v_1(t)]_{j^o}^2 + j[U_2^T v_2(t)]_{j^o}^2) / j[(t)]_{oj}} > 1; \quad (37)$$

for any  $t_0$ , when  $t_0$  is large enough. Combining (36) and (37) to compare the ratio of (34) for  $j$  and  $j^o$ , we get that there exists some  $\epsilon_0 > 0$  such that for any  $t > t_0$ , we have

$$\frac{\frac{d}{dt} j[(t)]_j}{j[(t)]_j} > \frac{\frac{d}{dt} j[(t)]_{j^o}}{j[(t)]_{j^o}} > 1: \quad (38)$$

This implies that the ratio of the absolute value of time derivative of  $j[(t)]_j$  to the absolute value of current value of  $j[(t)]_j$  is strictly bigger than that of  $j[(t)]_{j^o}$ . Moreover, we saw in (33) that the phase of  $\frac{d}{dt} j[(t)]_j$  converges to that of  $[u^1]_j$ . Since this holds for all  $t > t_0$ , (38) results in a growth of  $j[(t)]_j$  that is exponentially faster than that of  $j[(t)]_{j^o}$ , so  $j[(t)]_j$  becomes a dominant component in  $(t)$  as  $t \rightarrow \infty$ . This contradicts that  $j[u^1]_j = 0$ , hence the condition (32) has to be satisfied.

So far, we have characterized a number of conditions (26), (28), (31), (32) that have to be satisfied by the limit directions  $u^1$  and  $u^2$  of  $X^T r$  and  $v$ . We now consider the following optimization problem and prove that these conditions are in fact the KKT conditions of the optimization problem. Consider

$$\underset{C^m}{\text{minimize}} \quad k_1 k_2 = L \quad \text{subject to} \quad y_i x_i^T S^T = 1; \quad \forall i \in [n]; \quad (39)$$

The KKT conditions of this problem are

$$\exists k_1, k_2 \geq 0, \lambda_i \geq 0; \quad \lambda_i (1 - y_i x_i^T S^T) = 0 \quad \text{for all } i \in [n];$$

where  $\lambda_1, \dots, \lambda_n$  are the dual variables. The symbol  $\partial_{k_2=L}$  denotes the (local) subdifferential of the  $\| \cdot \|_{2=L}$  norm<sup>1</sup>, which can be written as

$\partial_{k_1} k_1 = \{ u \in C^m \mid \|u\|_j = 1 \text{ for all } j \in [m]; \text{ and } [u]_j \neq 0 \Rightarrow [u]_j = \exp(\frac{p-1}{2} \arg([u]_j)) \}$ ;  
if  $L = 2$  (in this case  $\partial_{k_1} k_1$  is the global subdifferential), and

$$\partial_{k_2=L} k_2 = \{ u \in C^m \mid [u]_j \neq 0 \Rightarrow [u]_j = \frac{2}{L} [u]_j^{\frac{2}{L-1}} \exp(\frac{p-1}{2} \arg([u]_j)) \};$$

if  $L > 2$ . By replacing  $\lambda_i$ 's with  $\mu_i$ 's defined in (26), we can check from (26), (28), (31), (32) that the  $\mu^1$  and  $\mu^2$  satisfy the KKT conditions up to scaling. Therefore, by (24),  $(t)$  converges in direction aligned with  $S^T \mu^1$ , where  $\mu^1$  is again aligned with a stationary point (global minimum in case of  $L = 2$ ) of the optimization problem (39).

If  $S$  is invertible, we can get  $\mu^T = (\mu^1)^T / \mu^1$ . Plugging this into the optimization problem (39) gives the last statement of the theorem.

<sup>1</sup>the definition of subdifferentials used here is taken from Gunasekar et al. (2018b).

## D.2 PROOF OF COROLLARY 2

It suffices to prove that linear diagonal networks satisfy Assumption 1 with  $S = I_d$ . The proof is very straightforward, since  $M_{\text{diag}}(x) \in \mathbb{R}^{d \times d}$  has  $[M_{\text{diag}}(x)]_{j,j} = [x]_j$  while all the remaining entries are zero. It is straightforward to verify that  $M_{\text{diag}}(x)$  satisfies Assumption 1 with  $S = U_1 = \dots = U_L = I_d$ . A direct substitution into Theorem 2 gives the corollary.

## D.3 PROOF OF COROLLARY 3

For full-length convolutional networks ( $k_1 = \dots = k_L = d$ ), we will prove that they satisfy Assumption 1 with  $S = d^{\frac{L-1}{2}} F$  and  $U_1 = \dots = U_L = F$ , where  $F \in \mathbb{C}^{d \times d}$  is the matrix of discrete Fourier transform basis  $[F]_{j,k} = \frac{1}{d} \exp\left(\frac{-i 2\pi (j-1)(k-1)}{d}\right)$  and  $F^*$  is the complex conjugate of  $F$ .

For simplicity of notation, define  $\omega = \exp\left(\frac{-i 2\pi}{d}\right)$ . With such matrices  $S$  and  $U_1, \dots, U_L$ , we can write  $M(x)$  as

$$\begin{aligned} M(x) &= \prod_{j=1}^L [Sx]_j ([U_1]_{:,j} \quad [U_2]_{:,j} \quad \dots \quad [U_L]_{:,j}) \\ &= \prod_{j=1}^L d^{\frac{L-2}{2}} \prod_{k=1}^d [x]_k \begin{pmatrix} \omega^{(j-1)(k-1)} \\ \omega^{2(j-1)(k-1)} \\ \vdots \\ \omega^{(d-1)(j-1)(k-1)} \end{pmatrix} \end{aligned}$$

where  $\omega^L$  denotes the  $L$ -times tensor product of  $\omega$ . We will show that  $M(x) = M_{\text{conv}}(x)$ .

For any  $j_1, \dots, j_L \in [d]$ ,

$$\begin{aligned} [M(x)]_{j_1, \dots, j_L} &= \prod_{l=1}^L \prod_{k=1}^d [x]_k \omega^{(j_l-1)(k-1)} \omega^{(j_{l+1}-1)(k-1)} \dots \omega^{(j_q-1)(k-1)} \dots \omega^{(j_L-1)(k-1)} \\ &= \prod_{k=1}^d [x]_k \omega^{(j_1-1)(k-1) + \dots + (j_L-1)(k-1)} \end{aligned}$$

Recall that

$$\prod_{l=1}^L \omega^{(j_l-1)(k-1)} = \begin{cases} d & \text{if } k-1 + \sum_{q=1}^L j_q + L \text{ is a multiple of } d; \\ 0 & \text{otherwise} \end{cases}$$

Using this, we have

$$\begin{aligned} [M(x)]_{j_1, \dots, j_L} &= \prod_{k=1}^d [x]_k \omega^{(j_1-1)(k-1) + \dots + (j_L-1)(k-1)} \\ &= [x]_{k_1 + \dots + k_L + 1 \pmod d} = [M_{\text{conv}}(x)]_{j_1, \dots, j_L} \end{aligned}$$

Hence, linear full-length convolutional networks satisfy Assumption 1 with  $S = d^{\frac{L-1}{2}} F$ . A direct substitution into Theorem 2 and then using the fact  $[Fz]_j = [F^*z]_j$  for any real vector  $z \in \mathbb{R}^d$  gives the corollary.

## E PROOFS OF THEOREM 3 AND COROLLARY 4

### E.1 PROOF OF THEOREM 3

#### E.1.1 CONVERGENCE OF LOSS TO ZERO

Since Theorem 3 does not assume the existence of  $\epsilon_0$  satisfying  $L(\epsilon_0) < 1$ , we need to first show that given the conditions on initialization, the training loss  $\mathcal{L}(t)$  converges to zero. Since

$L = 2$  and  $M(x) = U_1 \text{diag}(s) U_2^T$ , we can write the gradient flow dynamics from Section 2.1 as

$$\begin{aligned} \dot{v}_1 &= -M(x^T r) (I_{k_1}; v_2) = -r U_1 \text{diag}(s) U_2^T v_2; \\ \dot{v}_2 &= -M(x^T r) (v_1; I_{k_2}) = -r U_2 \text{diag}(s) U_1^T v_1; \end{aligned} \quad (40)$$

where  $r(t) = -y \exp(-y f(x; (t)))$  is the residual of the data point  $(x; y)$ . From (40) we get

$$U_1^T \dot{v}_1 = -r s \quad U_2^T \dot{v}_2; \quad U_2^T \dot{v}_2 = -r s \quad U_1^T \dot{v}_1; \quad (41)$$

Now consider the rate of growth for the  $j$ th component of  $U_1^T v_1$  squared:

$$\frac{d}{dt} [U_1^T v_1]_j^2 = 2[U_1^T v_1]_j [U_1^T \dot{v}_1]_j = -2r [s]_j [U_1^T v_1]_j [U_2^T v_2]_j = \frac{d}{dt} [U_2^T v_2]_j^2; \quad (42)$$

So for any  $j \in [m]$ ,  $[U_1^T v_1]_j^2$  and  $[U_2^T v_2]_j^2$  grow at the same rate. This means that the gap between the two layers stays constant for all  $t \geq 0$ . Combining this with our conditions on initial directions,

$$\begin{aligned} [U_1^T v_1(t)]_j^2 - [U_2^T v_2(t)]_j^2 &= [U_1^T v_1(0)]_j^2 - [U_2^T v_2(0)]_j^2 \\ &= -2[U_1^T v_1]_j^2 + 2[U_2^T v_2]_j^2 \leq 0; \end{aligned} \quad (43)$$

for any  $j \in [m]$  and  $t \geq 0$ . This inequality implies

$$[U_1^T v_1(t)]_j^2 \leq [U_2^T v_2(t)]_j^2 + 2[U_2^T v_2(0)]_j^2 - 2[U_1^T v_1(0)]_j^2; \quad (44)$$

Let us now consider the time derivative of  $\|L(\cdot)\|$ . We have the following chain of upper bounds on the time derivative:

$$\begin{aligned} \frac{d}{dt} \|L(\cdot)\| &= r \|L(\cdot)\|^T - \|L(\cdot)\| = kr \|L(\cdot)\| k_2^2 \\ &\leq kr \|v_2\| k_2^2 = k \|v_2(t)\| k_2^2 \\ &\stackrel{(a)}{\leq} k \|U_2^T v_2(t)\| k_2^2 \stackrel{(b)}{=} r(t)^2 s \|U_1^T v_1(t)\|_2^2 \\ &= r(t)^2 \sum_{j=1}^m [s]_j^2 [U_1^T v_1(t)]_j^2 \\ &\stackrel{(c)}{\leq} 2r(t)^2 \sum_{j=1}^m [s]_j^2 \\ &= 2ksk_2^2 \|L(\cdot)\|^2; \end{aligned}$$

where (a) used the fact that  $\|v_2(t)\| k_2^2 \leq k \|U_2^T v_2(t)\| k_2^2$  because it is a projection onto a subspace, and  $k \|U_2^T v_2(t)\| k_2^2 = k \|U_2^T v_2(t)\| k_2^2$  because  $U_2^T U_2 = I_{k_2}$ ; (b) is due to (41); (c) is due to (44). From this, we get

$$\|L(\cdot)\| \leq \frac{\|L(\cdot)\|}{1 + 2ksk_2^2 t}.$$

Therefore  $\|L(\cdot)\| \rightarrow 0$  as  $t \rightarrow \infty$ .

### E.1.2 CHARACTERIZING THE LIMIT DIRECTION

Since we proved that  $\|L(\cdot)\| \rightarrow 0$ , the argument in the proof of Theorem 1 applies to this case, and shows that the parameters converge in direction and align with  $u = r \|v_1\| L(\cdot)$ . Let  $v_1^1 := \lim_{t \rightarrow \infty} \frac{v_1(t)}{\|v_1(t)\| k_2}$  be the limit direction of  $v_1$ . As done in the proof of Theorem 2, define  $u = U_1^T v_1(t) - U_2^T v_2(t)$  and  $u^1 = U_1^T v_1^1 - U_2^T v_2^1$ .

It follows from  $r(t) = -y \exp(-y f(x; (t)))$  that we have  $\text{sign}(r(t)) = \text{sign}(y)$ . Using this, (41), and alignment of  $v_1$  and  $v_2$ , we have

$$U_1^T v_1^1 / y s \quad U_2^T v_2^1; \quad U_2^T v_2^1 / y s \quad U_1^T v_1^1; \quad (45)$$

Element-wise multiplying  $U_1^T v_1^1$  to both sides gives

$$(U_1^T v_1^1)^2 / y s \quad 1; \quad (U_2^T v_2^1)^2 / y s \quad 1; \quad (46)$$



Since the LHSs are positive and  $s_j$  is positive, the following equations have to be satisfied for all  $j \in [m]$ :

$$\text{sign}(y) = \text{sign}([v_1]_j): \quad (47)$$

Now, multiplying both sides of the two equations (46), we get

$$([v_1]_j)^2 / s_j^2 = ([v_2]_j)^2: \quad (48)$$

From (48),  $[v_1]_j$  must satisfy that

$$[v_1]_j \in \{0; [v_1]_j \in \mathbb{R}\} \Rightarrow j [s]_j = j [s]_j; \quad (49)$$

for all  $j \in [m]$ . As in the proof of Theorem 2, there is another condition that has to be satisfied:

$$[v_1]_j = 0; [v_1]_j \in \mathbb{R} \Rightarrow j [s]_j \neq j [s]_j; \quad (50)$$

for any  $j \in [m]$ ; let us prove why. First, consider the time derivative  $\dot{v} = [U_1^T v_1] [U_2^T v_2]$ .

$$\begin{aligned} \frac{d}{dt} [v]_j &= [U_1^T v_1(t)]_j \frac{d}{dt} [U_2^T v_2(t)]_j + [U_2^T v_2(t)]_j \frac{d}{dt} [U_1^T v_1(t)]_j \\ &\stackrel{(a)}{=} r(t) [s]_j ([U_1^T v_1(t)]_j^2 + [U_2^T v_2(t)]_j^2); \end{aligned}$$

where (a) used (41). Now consider

$$\frac{\frac{d}{dt} [v]_j}{j r(t) j [v]_j} = j [s]_j \frac{[U_1^T v_1(t)]_j^2 + [U_2^T v_2(t)]_j^2}{j [v]_j}. \quad (51)$$

We want to compare this quantity for different  $j \in [m]$ . Before we do that, we take a look at the last term in the RHS of (51). Recall from (43) that

$$[U_1^T v_1(t)]_j^2 = [U_2^T v_2(t)]_j^2 + [U_1^T v_1(0)]_j^2 - [U_2^T v_2(0)]_j^2. \quad (52)$$

For simplicity, let  $j := [U_1^T v_1(0)]_j^2 - [U_2^T v_2(0)]_j^2$ , which is a positive number due to our assumption on initialization. Then, we can use (52) and  $\dot{v}_j = j [U_1^T v_1(t)]_j j [U_2^T v_2(t)]_j$  to show that

$$\begin{aligned} \frac{[U_1^T v_1(t)]_j^2 + [U_2^T v_2(t)]_j^2}{j [v]_j} &= \frac{2[U_2^T v_2(t)]_j^2 + j}{j[U_2^T v_2(t)]_j [U_2^T v_2(t)]_j + j} > 2; \\ \lim_{t \rightarrow \infty} \frac{[U_1^T v_1(t)]_j^2 + [U_2^T v_2(t)]_j^2}{j [v]_j} &= 2 \quad \text{if} \quad \lim_{t \rightarrow \infty} j [U_2^T v_2(t)]_j = 1: \end{aligned}$$

Recall that we want to prove that (50) should necessarily hold. For the sake of contradiction, suppose that there exists  $j \in [m]$  that satisfies  $[v_1]_j = 0$  but  $j [s]_j > j [s]_j$ , for some  $j \in [m]$  satisfying  $[v_1]_j \in \mathbb{R}$ . Note that having  $[v_1]_j = 0$  and  $[v_1]_j \in \mathbb{R}$  implies that  $j [v]_j \neq 1$  and  $\frac{j [v]_j}{j [v]_j} \neq 0$ . We now want to compare the ratio of (51) for  $j$  and  $j^0$ . Using  $\frac{j [v]_j}{j [v]_j} \neq 0$  and the fact that  $x \mapsto \frac{2x^2 + j}{x^2 + j}$  is a decreasing function of  $x$  for any  $j > 0$ , we have

$$\frac{([U_1^T v_1(t)]_j^2 + [U_2^T v_2(t)]_j^2) / j [v]_j}{([U_1^T v_1(t)]_{j^0}^2 + [U_2^T v_2(t)]_{j^0}^2) / j [v]_{j^0}} > 1; \quad (53)$$

for any  $t \geq t_0$ , when  $t_0$  is large enough. Combining  $\frac{j [s]_j}{j [s]_{j^0}} > 1$  and (53) to compare the ratio of (51) for  $j$  and  $j^0$ , there exists some  $t_0 > 0$  such that for any  $t \geq t_0$ , we have

$$\frac{\frac{d}{dt} [v]_j}{\frac{d}{dt} [v]_{j^0}} = \frac{j [v]_j}{j [v]_{j^0}} > 1; \quad (54)$$

This implies that the ratio of the absolute value of time derivative  $|\dot{v}_j|$  to the absolute value of current value of  $[v]_j$  is strictly bigger than that of  $[v]_{j^0}$ . Moreover, by the definition of  $[v]$ ,  $\frac{d}{dt} [v]_j$  does not change sign over time. Since this holds for all  $t_0$ , (54) results in a growth of  $j [v]_j$  that is exponentially faster than that of  $j [v]_{j^0}$ , so  $[v]_j$  becomes a dominant component in  $[v]$  as  $t \rightarrow \infty$ . This contradicts that  $[v_1]_j = 0$ , hence the condition (50) has to be satisfied.

So far, we have characterized some conditions (47), (49), (50) that have to be satisfied by the limit direction  $\hat{v}_1$  of  $\hat{v}$ . We now consider the following optimization problem and prove that these conditions are in fact the KKT conditions of the optimization problem. Consider

$$\underset{v \in \mathbb{R}^m}{\text{minimize}} \quad \|k - k_1\| \quad \text{subject to} \quad y^T v = 1; \quad (55)$$

The KKT condition of this problem is

$$\lambda (k - k_1) + \mu y = 0;$$

where the global subdifferential  $\partial \|k - k_1\|$  is defined as

$$\partial \|k - k_1\| = \{u \in \mathbb{R}^m \mid \|u\| = 1 \text{ for all } j \in [m]; \text{ and } [u]_j \leq 0 \Rightarrow [u]_j = \text{sign}([k - k_1]_j)\};$$

We can check from (47), (49), (50) that the  $\hat{v}_1$  satisfies the KKT condition up to scaling.

Now, how do we characterize  $\hat{v}_1$  and  $\hat{v}_2$  in terms of  $\hat{v}$ ? Let  $\hat{v}_1 := U_1^T v_1^1$  and  $\hat{v}_2 := U_2^T v_2^1$ . Then,  $v_1^1 = U_1 \hat{v}_1 = U_1 U_1^T v_1^1$  holds because any component orthogonal to the column space of  $U_1$  stays unchanged while the component in the column space diverges to infinity. By (42),  $\hat{v}_1 = \hat{v}$  and  $\hat{v}_2 = \hat{v} - \hat{v}_1$ . By (45), we have  $\text{sign}(\hat{v}_1) = \text{sign}(y) = \text{sign}(\hat{v}_2)$ .

## E.2 PROOF OF COROLLARY 4

The proof of Corollary 4 boils down to characterizing the SVD of  $M_{\text{conv}}(x)$ .

### E.2.1 THE $k_1 = 1$ CASE

First, it is straightforward to check that for  $r=2$  and  $k_1 = 1$ , we have

$$M_{\text{conv}}(x) = v_1 v_2^T;$$

For  $k_1 = 1$ , the data tensor is simply  $M_{\text{conv}}(x) = x x^T$ . Thus, we have  $U_1 = 1$ ,  $U_2 = \frac{x}{\|x\|}$ , and  $s = \|x\|^2$ . Substituting  $U_1$  and  $U_2$  to the theorem gives the condition on initial directions in Corollary 4. Also, the theorem implies us that the limit direction of  $v_2$  satisfies  $v_2^T / \|v_2\| = y / \|y\|$ . Using this, it is easy to check that

$$M_{\text{conv}}(\hat{v}_1) / \|\hat{v}_1\| = \hat{v}_1 \hat{v}_1^T / \|x\|;$$

### E.2.2 THE $k_1 = 2$ CASE

First, it is straightforward to check that for  $r=2$  and  $k_1 = 2$ , we have

$$M_{\text{conv}}(x) = \begin{bmatrix} [v_1]_1 & 0 & 0 & 0 & [v_1]_2 \\ [v_1]_2 & [v_1]_1 & 0 & 0 & 0 \\ 0 & [v_1]_2 & [v_1]_1 & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & [v_1]_1 & 0 \\ 0 & 0 & 0 & [v_1]_2 & [v_1]_1 \end{bmatrix} v_2; \quad (56)$$

For  $k_1 = 2$ , by definition, the data tensor is

$$M_{\text{conv}}(x) = \begin{bmatrix} x^T \\ x \end{bmatrix};$$

and it is straightforward to check that the SVD of this matrix is

$$M_{\text{conv}}(x) = \begin{bmatrix} x^T \\ x \end{bmatrix} = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} \frac{\|x\|^2 + x^T x}{2} & 0 \\ 0 & \frac{\|x\|^2 - x^T x}{2} \end{bmatrix} \begin{bmatrix} \frac{x^T + x}{\sqrt{\|x\|^2 + x^T x}} \\ \frac{x^T - x}{\sqrt{\|x\|^2 - x^T x}} \end{bmatrix};$$

so

$$U_1 = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}; U_2 = \begin{bmatrix} \frac{x^T + x}{\sqrt{\|x\|^2 + x^T x}} \\ \frac{x^T - x}{\sqrt{\|x\|^2 - x^T x}} \end{bmatrix}; s = \begin{bmatrix} \frac{\|x\|^2 + x^T x}{2} \\ \frac{\|x\|^2 - x^T x}{2} \end{bmatrix};$$

Substituting  $U_1$  and  $U_2$  to the theorem gives the conditions on initial directions. Also, note that the maximum singular value depends on the sign of  $x$ . Consider the optimization problem in the theorem statement:

$$\text{minimize}_{\mathbb{R}^m} \|k\|_1 \quad \text{subject to} \quad y^T v = 1:$$

If  $x^T x > 0$ , then the solution  $v^1$  to this problem is in the direction of  $[0 \ y]$ . Therefore, the limit directions  $v_1^1$  and  $v_2^1$  will be of the form

$$v_1^1 / c_1 \begin{bmatrix} 1 \\ 1 \end{bmatrix}; \quad v_2^1 / c_2(x + x);$$

where  $\text{sign}(c_1) \text{sign}(c_2) = \text{sign}(y)$ . Using (56), it is straightforward to check that

$$\text{conv} \left( \begin{matrix} 1 \\ \text{conv} \end{matrix} \right) / y = \begin{matrix} 2 & 1 & 0 & 0 & 0 & 1^3 \\ 6 & 1 & 1 & 0 & 0 & 0^2 \\ 4 & 0 & 1 & 1 & 0 & 0^2 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 4 & 0 & 0 & 0 & 1 & 0^5 \\ 0 & 0 & 0 & 0 & 1 & 1 \end{matrix} (x + x) = y(2x + x + \dots x):$$

Similarly, if  $x^T x < 0$ , then the solution  $v^1$  is in the direction of  $[0 \ -y]$ . Using (56), we have

$$\text{conv} \left( \begin{matrix} 1 \\ \text{conv} \end{matrix} \right) / y = \begin{matrix} 2 & 1 & 0 & 0 & 0 & 1^3 \\ 6 & 1 & 1 & 0 & 0 & 0^2 \\ 4 & 0 & 1 & 1 & 0 & 0^2 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 4 & 0 & 0 & 0 & 1 & 0^5 \\ 0 & 0 & 0 & 0 & 1 & 1 \end{matrix} (x - x) = y(2x - x + \dots x):$$

## F PROOFS OF THEOREM 5, COROLLARIES 5, 6 & 7, AND LEMMA 4

### F.1 PROOF OF LEMMA 4

In this subsection, we restate Lemma 4 and prove it.

Lemma 4. Consider the system of ODEs, where  $p, q: \mathbb{R} \rightarrow \mathbb{R}$ :

$$\dot{p} = p^L - 2q; \quad \dot{q} = p^L - 1; \quad p(0) = 1; \quad q(0) = 0:$$

Then, the solutions  $p_L(t)$  and  $q_L(t)$  are continuous on their maximal interval of existence of the form  $(-c; c) \subset \mathbb{R}$  for some  $c \in (0; 1]$ . Define  $h_L(t) = p_L(t)^L - 1 - q_L(t)$ ; then,  $h_L(t)$  is odd and strictly increasing, satisfying  $\lim_{t \rightarrow -c} h_L(t) = -1$  and  $\lim_{t \rightarrow c} h_L(t) = 1$ .

Proof First, continuity (and also continuous differentiability) of  $p(t)$  and  $q(t)$  is straightforward because the RHSs of the ODEs are differentiable in  $p$  and  $q$ . Next, define  $\tilde{p}(t) = p(-t)$  and  $\tilde{q}(t) = -q(-t)$ . Then, one can show that  $\tilde{p}$  and  $\tilde{q}$  are also the solution of the ODE because

$$\begin{aligned} \frac{d}{dt} \tilde{p}(t) &= \frac{d}{dt} p(-t) = -\dot{p}(-t) = -p(-t)^L + 2q(-t) = \tilde{p}(t)^L - 2\tilde{q}(t); \\ \frac{d}{dt} \tilde{q}(t) &= \frac{d}{dt} -q(-t) = \dot{q}(-t) = p(-t)^L - 1 = \tilde{p}(t)^L - 1; \end{aligned}$$

However, by the Picard-Lindelöf theorem, the solution has to be unique; this means  $\tilde{p}(t) = p(t)$  and  $\tilde{q}(t) = -q(t) = -q(-t)$ , which proves that  $p$  is even and  $q$  is odd and also implies that the domain of  $p$  and  $q$  has to be of the form  $(-c; c)$  (i.e. symmetric around the origin) and  $h = p^L - 1 - q$  is odd.

To show that  $h$  is strictly increasing, it suffices to show that  $p$  and  $q$  are both strictly increasing on  $[0; c)$ . To this end, we show that  $p(t) > 1$  for all  $t \in [0; c)$ . First, due to the initial condition  $p(0) = 1$  and continuity of  $p$ , there exists  $\delta_1 > 0$  such that  $p(t) > 1$  for all  $t \in [0; \delta_1) =: I_1$ . This implies that  $\dot{q}(t) = p(t)^L - 1 > 0$  for  $t \in I_1 \cap [0; c)$ , so  $q$  is strictly increasing on  $I_1$ . Since  $q(0) = 0$ , we have  $q(t) > 0$  for  $t \in I_1 \cap [0; c)$ , which then implies that  $\dot{p}(t) = p(t)^L - 2q(t) > 0$ . Therefore,  $p$  is

also strictly increasing on  $[0; c]$ ; this then means  $p(t) < 1$  for  $t \in [0; c]$  because  $p(0) = 1$ . Now, due to  $p'(t) < 0$  and continuity of  $p$ , there exists  $c_2 > c_1$  such that  $p(t) > 0$  for all  $t \in [c_1; c_2] =: I_2$ . Using the argument above for  $I_2$  results in  $p(t) < 1$  for  $t \in [0; c_2]$ . Repeating this until the end of the domain, we can show that  $p(t) < 1$  holds for all  $t \in [0; c]$ . By  $p < 1$ , we have  $q = p^{L-1} < 1$  on  $[0; c]$ , so  $q$  is strictly increasing on  $[0; c]$ . Also,  $q(t) > 0$  on  $(0; c)$ , so  $p = p^{L-2}q > 0$  on  $(0; c)$  and  $p$  is also strictly increasing on  $[0; c]$ . This proves that  $h$  is strictly increasing on  $[0; c]$ , and also on  $(-c; c)$  by oddity of  $h$ .

Finally, it is left to show  $\lim_{t \rightarrow c} h(t) = 1$  and  $\lim_{t \rightarrow -c} h(t) = -1$ . If  $c < 1$ , then this together with monotonicity implies that the limits hold. To see why, suppose  $\lim_{t \rightarrow c} h(t) < 1$ . Then,  $p$  and  $q$  can be extended beyond  $c$ , which contradicts the fact that  $(c; c)$  is the maximal interval of existence of the solution. Next, consider the case  $c = 1$ . From  $p(t) < 1$ , we have  $q(t) < 1$  for  $t > 0$ . This implies that  $p'(t) < -t$  for  $t > 0$ . Now,  $p'(t) = p(t)^{L-2}q'(t) < -t$ , which gives  $p(t) < \frac{t^2}{2} + 1$  for  $t > 0$ . Therefore, we have

$$\lim_{t \rightarrow 1} h(t) = \lim_{t \rightarrow 1} p(t)^{L-1}q(t) = \lim_{t \rightarrow 1} \left( \frac{t^2}{2} + 1 \right)^{L-1} t = 1;$$

hence finishing the proof.  $\square$

## F.2 PROOF OF THEOREM 5

### F.2.1 CONVERGENCE OF LOSS TO ZERO

We first show that given the conditions on initialization, the training losses  $\mathcal{L}_l(t)$  converges to zero. Recall from Section 2.1 that

$$\underline{v}_l = \text{r}_{v_l}(\underline{v}_l) = M(X^T r) (v_1; \dots; v_{l-1}; I_{k_l}; v_{l+1}; \dots; v_L);$$

Applying the structure (9) in Assumption 1, we get

$$\begin{aligned} \underline{v}_l &= M(X^T r) (v_1; \dots; v_{l-1}; I_{k_l}; v_{l+1}; \dots; v_L) \\ &= \prod_{j=1}^n [S X^T r]_j (v_1^T [U_1]_{\cdot j} \quad v_{l-1}^T [U_{l-1}]_{\cdot j} \quad [U_l]_{\cdot j} \quad v_{l+1}^T [U_{l+1}]_{\cdot j} \quad v_L^T [U_L]_{\cdot j}) \\ &= \prod_{j=1}^n [S X^T r]_j \prod_{k \in I} [U_k^T v_k]_{\cdot j} [U_l]_{\cdot j}; \end{aligned}$$

Left-multiplying  $U_l^T$  to both sides, we get

$$U_l^T \underline{v}_l = S X^T r \prod_{k \in I} U_k^T v_k; \tag{57}$$

where  $\prod$  denotes the product using entry-wise multiplication

Now consider the rate of growth for the second power of the component of  $U_l^T \underline{v}_l$ :

$$\frac{d}{dt} [U_l^T \underline{v}_l]_j^2 = 2 [U_l^T \underline{v}_l]_j [U_l^T \dot{\underline{v}}_l]_j = 2 [S X^T r]_j \prod_{k=1}^L [U_k^T v_k]_j = \frac{d}{dt} [U_l^T v_l]_j^2$$

for any  $l \in [L]$ . Thus, for any  $j \in [m]$ , the second power of the  $j$ th components in  $U_l^T \underline{v}_l$  grow at the same rate for each layer  $l \in [L]$ . This means that the gap between any two different layers stays constant for all  $t \geq 0$ . Combining this with our conditions on initial directions, we have

$$[U_l^T \underline{v}_l(t)]_j^2 - [U_l^T \underline{v}_l(0)]_j^2 = [U_l^T \underline{v}_l(0)]_j^2 - [U_l^T \underline{v}_l(0)]_j^2 = 2 [U_l^T \underline{v}_l(0)]_j^2 - 2 [U_l^T \underline{v}_l(0)]_j^2;$$

for any  $j \in [m]$ ,  $l \in [L-1]$ , and  $t \geq 0$ . This inequality also implies

$$[U_l^T \underline{v}_l(t)]_j^2 - [U_l^T \underline{v}_l(t)]_j^2 + 2 [U_l^T \underline{v}_l(0)]_j^2 - 2 [U_l^T \underline{v}_l(0)]_j^2; \tag{58}$$

Let us now consider the time derivative of  $L(\mathbf{t})$ . We have the following chain of upper bounds on the time derivative:

$$\begin{aligned}
 \frac{d}{dt}L(\mathbf{t}) &= \mathbf{r}^\top \mathbf{L}(\mathbf{t}) - \mathbf{r}^\top \mathbf{L}(\mathbf{t}) = -\|\mathbf{r} - \mathbf{L}(\mathbf{t})\|_2^2 \\
 &\leq -\|\mathbf{v}_L \mathbf{L}(\mathbf{t})\|_2^2 = -\|\mathbf{U}_L^\top \mathbf{L}(\mathbf{t})\|_2^2 \\
 &\stackrel{(a)}{=} -\sum_{k \in [L]} \|\mathbf{U}_k^\top \mathbf{v}_k(\mathbf{t})\|_2^2 \\
 &= -\sum_{j=1}^m [\mathbf{S}^\top \mathbf{r}(\mathbf{t})]_j^2 - \sum_{k \in [L]} [\mathbf{U}_k^\top \mathbf{v}_k(\mathbf{t})]_j^2 \\
 &\stackrel{(c)}{\leq} -\sum_{j=1}^m [\mathbf{S}^\top \mathbf{r}(\mathbf{t})]_j^2 \\
 &= -\sum_{j=1}^m \lambda_j(\mathbf{S})^2 \|\mathbf{X}^\top \mathbf{r}(\mathbf{t})\|_2^2 \\
 &\stackrel{(d)}{\leq} -\sum_{j=1}^m \lambda_j(\mathbf{S})^2 \lambda_j(\mathbf{X})^2 \|\mathbf{r}(\mathbf{t})\|_2^2 \\
 &= -\sum_{j=1}^m \lambda_j(\mathbf{S})^2 \lambda_j(\mathbf{X})^2 L(\mathbf{t}); \tag{59}
 \end{aligned}$$

where (a) used the fact that  $\|\mathbf{v}_L \mathbf{L}(\mathbf{t})\|_2^2 = \|\mathbf{U}_L^\top \mathbf{L}(\mathbf{t})\|_2^2$  because it is a projection onto a subspace, and  $\|\mathbf{U}_L^\top \mathbf{L}(\mathbf{t})\|_2^2 = \sum_{k \in [L]} \|\mathbf{U}_k^\top \mathbf{v}_k(\mathbf{t})\|_2^2$  because  $\mathbf{U}_L^\top \mathbf{U}_L = \mathbf{I}_{k_L}$ ; (b) is due to (57); (c) is due to (58); and (d) used the fact that  $\mathbf{S} \in \mathbb{R}^{m \times d}$  and  $\mathbf{X} \in \mathbb{R}^{d \times n}$  are matrices that have full column rank, so for any  $\mathbf{z} \in \mathbb{R}^n$ , we can use  $\|\mathbf{S}^\top \mathbf{z}\|_2 \geq \lambda_{\min}(\mathbf{S}) \|\mathbf{z}\|_2$  and  $\|\mathbf{X}^\top \mathbf{z}\|_2 \geq \lambda_{\min}(\mathbf{X}) \|\mathbf{z}\|_2$  where  $\lambda_{\min}(\cdot)$  denotes the minimum singular value of a matrix.

From (59), we get

$$L(\mathbf{t}) \leq L(\mathbf{0}) \exp(-\sum_{j=1}^m \lambda_j(\mathbf{S})^2 \lambda_j(\mathbf{X})^2 t); \tag{60}$$

so that  $L(\mathbf{t}) \rightarrow 0$  as  $t \rightarrow \infty$ .

### F.2.2 CHARACTERIZING THE LIMIT POINT

Now, we move on to characterize the limit points of the gradient flow. First, by defining a ‘‘transformed’’ version of the parameters  $\mathbf{s}(\mathbf{t}) := \mathbf{U}_L^\top \mathbf{v}_L(\mathbf{t})$  and using (57), one can define an equivalent system of ODEs:

$$\begin{aligned}
 \dot{\mathbf{s}}_l &= -\sum_{k \in [L]} \mathbf{S}^\top \mathbf{r}(\mathbf{t}) \mathbf{U}_k^\top \mathbf{v}_k(\mathbf{t}) \quad \text{for } l \in [L]; \\
 \mathbf{s}_l(0) &= \mathbf{0} \quad \text{for } l \in [L-1]; \quad \mathbf{s}_L(0) = \mathbf{0}. \tag{61}
 \end{aligned}$$

Using Lemma 4, it is straightforward to verify that the solution to (61) has the following form. For odd  $L$ , we have

$$\begin{aligned}
 \mathbf{s}_l(t) &= -\sum_{j=1}^{L-2} \lambda_j(\mathbf{S})^2 \lambda_j(\mathbf{X})^2 \int_0^t \mathbf{S}^\top \mathbf{r}(\tau) d\tau \quad \text{for } l \in [L-1]; \\
 \mathbf{s}_L(t) &= -\sum_{j=1}^{L-2} \lambda_j(\mathbf{S})^2 \lambda_j(\mathbf{X})^2 \int_0^t \mathbf{S}^\top \mathbf{r}(\tau) d\tau; \tag{62}
 \end{aligned}$$

Similarly, for even  $L$ , the solution for (61) satisfies

$$\begin{aligned}
 \mathbf{s}_l(t) &= -\sum_{j=1}^{L-2} \lambda_j(\mathbf{S})^2 \lambda_j(\mathbf{X})^2 \int_0^t \mathbf{S}^\top \mathbf{r}(\tau) d\tau \quad \text{for } l \in [L-1]; \\
 \mathbf{s}_L(t) &= -\sum_{j=1}^{L-2} \lambda_j(\mathbf{S})^2 \lambda_j(\mathbf{X})^2 \int_0^t \mathbf{S}^\top \mathbf{r}(\tau) d\tau; \tag{63}
 \end{aligned}$$

Now that we know how the solutions look like, let us see how these relate to the linear coefficients of the network. By Assumption 1, we have

$$f(\mathbf{x}; \mathbf{s}) = M(\mathbf{x}) (\mathbf{v}_1; \dots; \mathbf{v}_L) = \sum_{j=1}^m [\mathbf{S}^\top \mathbf{x}]_j \sum_{l=1}^L [\mathbf{U}_l^\top \mathbf{v}_l]_j$$



$$= \sum_{j=1}^n \sum_{l=1}^L [W_j]_l [S]_j; \quad x = x^T S^T \quad Y \quad \text{I2[L]} \quad = x^T S^T :$$

Here, we define  $Q := \sum_{l=1}^L \sum_{j=1}^n [W_j]_l^2 [S]_j^2 \in \mathbb{R}^m$ . Therefore, the linear coefficients of the network can be written as  $(t) = S^T (t)$ . From the solutions (62) and (63), we can write

$$(t) = \sum_{i=1}^n (t) = \sum_{j=1}^n \sum_{l=1}^L h_L^{(j,l)} \sum_{j=1}^n \sum_{l=1}^L S X^T \int_0^{Z_t} r(\cdot) d \cdot ;$$

where  $h_L := p_L^{-1} q_L$ , defined in Lemma 4. By the convergence of the loss to zero (60), we have  $\lim_{t \rightarrow \infty} X^{-1} (t) = y$ . Therefore,

$$X S^T \sum_{j=1}^n \sum_{l=1}^L h_L^{(j,l)} \sum_{j=1}^n \sum_{l=1}^L S X^T \int_0^{Z_1} r(\cdot) d \cdot = y; \quad (64)$$

Next, we will show that  $z^1$  is in fact the solution of the following optimization problem

$$\text{minimize}_{z \in \mathbb{R}^m} Q_L; (z) \quad \text{subject to} \quad X S^T z = y; \quad (65)$$

where  $Q_L; : \mathbb{R}^m \rightarrow \mathbb{R}$  is a norm-like function defined using  $H_L(t) := \int_0^t h_L^{-1}(\cdot) d \cdot$ :

$$Q_L; (z) = \sum_{j=1}^n \sum_{l=1}^L [W_j]_l^2 H_L \left( \frac{[W_j]_l}{[S]_j} z \right) :$$

Note that the KKT conditions for (65) are

$$X S^T z = y; \quad r Q_L; (z) = S X^T \lambda ;$$

for some  $\lambda \in \mathbb{R}^n$ . It is clear from (64) that  $z^1$  satisfies the first condition (primal feasibility), so let us check the other one. Through a straightforward calculation, we get

$$r Q_L; (z) = \sum_{j=1}^n \sum_{l=1}^L [W_j]_l^2 h_L^{-1} \left( \frac{[W_j]_l}{[S]_j} z \right) :$$

Equating this with  $S X^T \lambda$  gives

$$\begin{aligned} \sum_{j=1}^n \sum_{l=1}^L [W_j]_l^2 h_L^{-1} \left( \frac{[W_j]_l}{[S]_j} z \right) &= S X^T \lambda \\ \sum_{j=1}^n \sum_{l=1}^L h_L^{-1} \left( \frac{[W_j]_l}{[S]_j} z \right) &= \sum_{j=1}^n \sum_{l=1}^L [W_j]_l^{-2} [S]_j^2 S X^T \lambda \\ \sum_{j=1}^n \sum_{l=1}^L h_L^{-1} \left( \frac{[W_j]_l}{[S]_j} z \right) &= \sum_{j=1}^n \sum_{l=1}^L [W_j]_l^{-2} [S]_j^2 S X^T \lambda : \end{aligned}$$

Hence, by setting  $\lambda = \int_0^{R_1} r(\cdot) d \cdot$ ,  $z^1$  satisfies this condition as well. Also,  $\mathcal{B}$  is invertible, we can substitute  $z = S^{-T} z$  to (65) to get the last statement of the theorem. This finishes the proof.

### F.3 PROOF OF COROLLARY 5

The proof is a direct consequence of the fact that Assumption 1 holds with  $U_1 = I_d = U_L$  for linear diagonal networks. Hence, the proof is the same as Corollary 2, proved in Appendix D.2.

### F.4 PROOF OF COROLLARY 6

We start by showing the DFT of a real and even vector is also real and even. Suppose  $x$  is real and even. First,

$$[F x]_j = \frac{1}{d} \sum_{k=1}^d x_k \exp \left( \frac{p-1}{d} 2(j-1)(k-1) \right)$$

$$\begin{aligned}
 &= \frac{1}{d} \sum_{k=1}^{X^d} [x]_k \cos \frac{2(j-1)(k-1)}{d} + \frac{1}{d} \sum_{k=1}^{X^d} [x]_k \sin \frac{2(j-1)(k-1)}{d} \\
 &= \frac{1}{d} \sum_{k=1}^{X^d} [x]_k \cos \frac{2(j-1)(k-1)}{d} \quad 2R;
 \end{aligned}$$

for all  $j \in [d]$ . To prove that  $Fx$  is even, for  $j = 0; \dots; \lfloor \frac{d}{2} \rfloor$ , we have

$$\begin{aligned}
 [Fx]_{j+2} &= \frac{1}{d} \sum_{k=1}^{X^d} [x]_k \cos \frac{2(j+1)(k-1)}{d} \\
 &= \frac{1}{d} \sum_{k=1}^{X^d} [x]_k \cos \frac{2(k-1)}{d} \frac{2(j+1)(k-1)}{d} \\
 &= \frac{1}{d} \sum_{k=1}^{X^d} [x]_k \cos \frac{2(d-j-1)(k-1)}{d} \\
 &= \frac{1}{d} \sum_{k=1}^{X^d} [x]_k \cos \frac{2(d-j-1)(k-1)}{d} \\
 &= [Fx]_{d-j};
 \end{aligned}$$

It is proved in Appendix D.3 that linear full-length convolutional networks  $\mathcal{F} = \{F_1, \dots, F_L\}$  ( $L = k_L = d$ ) satisfy Assumption 1 with  $S = d^{\frac{L-1}{2}} F$  and  $U_1 = \dots = U_L = F$ , where  $F \in \mathbb{C}^{d \times d}$  is the matrix of discrete Fourier transform basis  $[F]_{j;k} = \frac{1}{d} \exp(\frac{j-1}{d} \frac{2\pi i}{d} (k-1))$  and  $F$  is the complex conjugate of  $F^T$ .

The proof of convergence of loss to zero in Appendix F.2.1 is written for real matrices  $S; U_1; \dots; U_L$ , but we can actually apply the same argument as in Appendix D.1.1 and prove that the loss converges to zero, even in the case where  $S; U_1; \dots; U_L$  are complex.

Next, since  $U_l$ 's are complex, we can write the system of ODE as (see (20) for its derivation)

$$\dot{F} w_l = d^{\frac{L-1}{2}} F X^T r + \sum_{k \in I} F w_k; \tag{66}$$

Since all data points  $x_i$  and initialization  $w_l(0)$  are real and even, we have that  $F X^T r$  is real and even, and  $F w_l(0) = F w_k(0)$ 's are real and even. By (66), we see that the time derivatives of  $F w_l(t)$  are also real and even. Thus, the parameters  $w_l(t)$  are all real and even for all  $t \geq 0$ . From this observation, we can define  $ne_l(t) := F w_l(t)$ ,  $ne := F w$ , and  $S := d^{\frac{L-1}{2}} \text{Re}(F)$ , which are all real by the even symmetry. Then, starting from (61), the proof goes through.

### F.5 PROOF OF COROLLARY 7

Since the sensor matrices  $A_1; \dots; A_n$  commute, they are simultaneously diagonalizable with a real unitary matrix  $U \in \mathbb{R}^{d \times d}$ , i.e.,  $U^T A_i U$ 's are diagonal matrices. From the deep matrix sensing problem (13), we can compute  $w_l L_{ms}$ , which gives the gradient flow dynamics  $\dot{w}_l$ .

$$\dot{w}_l = r - w_l L_{ms} = W_l^T \left( r - \sum_{i=1}^n r_i A_i \right) W_l^T - W_{l+1}^T;$$

where  $r_i = h A_i$ ;  $W_l^T w_l - y_i$  is the residual for the  $i$ -th sensor matrix. If we left-multiply  $U^T$  and right-multiply  $U$  to both sides, we get

$$U^T \dot{w}_l U = U^T W_l^T \left( r - \sum_{i=1}^n r_i A_i \right) U^T W_l^T U - U^T W_{l+1}^T U; \tag{67}$$

If  $U^T W_k^T U$  is a diagonal matrix for all  $k \in I$ , then  $U^T W_l U$  is also a diagonal matrix. Note also that, since  $W_l(0) = I_d = U U^T$  for  $l \in [L-1]$ , the product  $U^T W_l U$  is a diagonal matrix at initialization. These observations imply that  $w_l(t)$ 's are all diagonalizable with  $U$  for all  $t \geq 0$ .

Now, denote  $v_l(t) = \text{eig}(W_l(t))$ , i.e.,  $U^T W_l U = \text{diag}(v_l)$ . Also, let  $x_i = \text{eig}(A_i)$ . Then, (67) can be written as

$$v_l = \left( \prod_{i=1}^n r_i x_i \right)_{k \in l} v_k.$$

Therefore, this is equivalent to the regression problem with linear diagonal networks, initialized at  $v_l(0) = 1$  for  $l \geq 2$  and  $v_l(0) = 0$ . Given this equivalence, Corollary 7 can be implied from Corollary 5.

## G PROOF OF THEOREM 6

### G.1 CONVERGENCE OF LOSS TO ZERO

We first show that given the conditions on initialization, the training loss  $\mathcal{L}(t)$  converges to zero. Since  $L = 2$  and  $M(x) = U_1 \text{diag}(s) U_2^T$ , we can write the gradient flow dynamics from Section 2.1 as

$$\begin{aligned} \dot{v}_1 &= -M(x^T r) (l_{k_1}; v_2) = -r U_1 \text{diag}(s) U_2^T v_2; \\ \dot{v}_2 &= -M(x^T r) (v_1; l_{k_2}) = -r U_2 \text{diag}(s) U_1^T v_1; \end{aligned} \quad (68)$$

where  $r(t) = f(x; y)$  is the residual of the data point  $(x; y)$ . From (68) we get

$$U_1^T \dot{v}_1 = -r s \quad U_2^T \dot{v}_2 = -r s \quad U_1^T v_1; \quad (69)$$

Now consider the rate of growth for the  $j$ th component of  $U_1^T v_1$  squared:

$$\frac{d}{dt} [U_1^T v_1]_j^2 = 2[U_1^T v_1]_j [U_1^T \dot{v}_1]_j = -2r [s]_j [U_1^T v_1]_j [U_2^T v_2]_j = \frac{d}{dt} [U_2^T v_2]_j^2.$$

So for any  $j \geq 2$ ,  $[U_1^T v_1]_j^2$  and  $[U_2^T v_2]_j^2$  grow at the same rate. This means that the gap between the two layers stays constant for all  $t \geq 0$ . Combining this with our conditions on initial directions,

$$\begin{aligned} [U_1^T v_1(t)]_j^2 - [U_2^T v_2(t)]_j^2 &= [U_1^T v_1(0)]_j^2 - [U_2^T v_2(0)]_j^2 \\ &= -2[U_1^T v_1]_j^2 + 2[U_2^T v_2]_j^2; \end{aligned}$$

for any  $j \geq 2$  and  $t \geq 0$ . This inequality implies

$$[U_1^T v_1(t)]_j^2 \leq [U_2^T v_2(t)]_j^2 + 2[U_1^T v_1]_j^2 - 2[U_2^T v_2]_j^2. \quad (70)$$

Let us now consider the time derivative of  $\mathcal{L}(t)$ . We have the following chain of upper bounds on the time derivative:

$$\begin{aligned} \frac{d}{dt} \mathcal{L}(t) &= -r \mathcal{L}(t) - \sum_{k=2}^m k v_k(t) k_2^2 \\ &\leq -r \mathcal{L}(t) - k v_2(t) k_2^2 \\ &\stackrel{(a)}{\leq} -r \mathcal{L}(t) - \sum_{j=2}^m [s]_j^2 [U_1^T v_1(t)]_j^2 \\ &\stackrel{(b)}{=} -r \mathcal{L}(t) - \sum_{j=2}^m [s]_j^2 [U_2^T v_2(t)]_j^2 \\ &\stackrel{(c)}{\leq} -2r \mathcal{L}(t) - \sum_{j=2}^m [s]_j^2 \\ &= -2r \mathcal{L}(t) - \sum_{j=2}^m [s]_j^2 k_2^2; \end{aligned}$$

where (a) used the fact that  $k v_k(t) k_2^2 \leq k U_2 U_2^T v_k(t) k_2^2$  because it is a projection onto a subspace, and  $k U_2 U_2^T v_k(t) k_2^2 = k U_2^T v_k(t) k_2^2$  because  $U_2^T U_2 = I_{k_2}$ ; (b) is due to (69); (c) is due to (70). From this, we get

$$\mathcal{L}(t) \leq \mathcal{L}(0) \exp(-2r t) - \sum_{j=2}^m [s]_j^2 k_2^2 t. \quad (71)$$

Therefore,  $\mathcal{L}(t) \rightarrow 0$  as  $t \rightarrow \infty$ .

## G.2 CHARACTERIZING THE LIMIT POINT

Now, we move on to characterize the limit points of the gradient flow. First, note that any changes made in  $v_1$  over time are in the subspace spanned by the columns of  $U_1$ . Therefore, any component in the initialization  $v_1(0) = v_1$  that is orthogonal to the column space of  $U_1$  stays constant.

So, we can focus on the evolution of  $v_1$  in the column space of  $U_1$ ; this can be done by defining a “transformed” version of the parameter  $z(t) := U_1^T v_1(t)$  and using (69), one can define an equivalent system of ODEs:

$$\begin{aligned} \dot{z}_1 &= -r s_1 z_1; & \dot{z}_2 &= -r s_2 z_2; \\ z_1(0) &= z_1; & z_2(0) &= z_2; \end{aligned} \quad (72)$$

where  $z_1 := U_1^T v_1$ ,  $z_2 := U_2^T v_2$ . It is straightforward to verify that the solution to (72) has the following form.

$$\begin{aligned} z_1(t) &= z_1 \cosh(-s_1 r t) + z_2 \sinh(-s_1 r t); \\ z_2(t) &= z_1 \sinh(-s_2 r t) + z_2 \cosh(-s_2 r t); \end{aligned} \quad (73)$$

By the convergence of the loss to zero (71), we have  $f(x; z(t)) = y$ . Note that  $f(x; z(t))$  can be written as

$$\begin{aligned} f(x; z(t)) &= M(x) (v_1(t); v_2(t)) = v_1(t)^T M(x) v_2(t) \\ &= v_1(t)^T U_1 \text{diag}(s) U_2^T v_2(t) = s^T (z_1(t) \quad z_2(t)); \end{aligned}$$

Therefore,

$$\begin{aligned} \lim_{t \rightarrow \infty} f(x; z(t)) &= \lim_{t \rightarrow \infty} s^T (z_1(t) \quad z_2(t)) \\ &= s^T (z_1^2 + z_2^2) \cosh(-s_1 r t) + s_1 z_1 z_2 \sinh(-s_1 r t) \\ &\quad + (z_1 z_2) \cosh(-s_2 r t) + s_2 z_1 z_2 \sinh(-s_2 r t) \\ &= s^T \frac{z_1^2 + z_2^2}{2} \sinh(2s_1 r t) + (z_1 z_2) \cosh(2s_2 r t) \\ &= \sum_{j=1}^n s_j \frac{[z_1]_j^2 + [z_2]_j^2}{2} \sinh(2[s]_j) + [z_1]_j [z_2]_j \cosh(2[s]_j) \\ &= y; \end{aligned} \quad (74)$$

where we defined  $s := \begin{bmatrix} s_1 \\ \vdots \\ s_n \end{bmatrix}$ . Consider the function  $g(s) = a \sinh(s) + b \cosh(s)$ . This is a strictly increasing function if  $a > |b|$ . Note also that

$$\frac{[z_1]_j^2 + [z_2]_j^2}{2} > |[z_1]_j [z_2]_j|; \quad (75)$$

which holds with equality if and only if  $[z_1]_j = j [z_2]_j$ . However, recall from our assumptions on initialization that  $[z_1]_j^2 + [z_2]_j^2 > 0$ , so (75) can only hold with strict inequality. Therefore,

$$g(s) := \sum_{j=1}^n s_j \frac{[z_1]_j^2 + [z_2]_j^2}{2} \sinh(2[s]_j) + [z_1]_j [z_2]_j \cosh(2[s]_j)$$

is a strictly increasing (hence invertible) function because it is a sum of strictly increasing functions. Using this, (74) can be written as  $g(s) = y$ , and by using the inverse of  $g$  we have

$$s = g^{-1} \left( \frac{y}{2} \right); \quad (76)$$

Plugging (76) into (73), we get

$$\begin{aligned}
& \lim_{t \uparrow 1} v_1(t) \\
&= U_1 \lim_{t \uparrow 1} v_1(t) + (I_{k_1} - U_1 U_1^T) v_1 \\
&= U_1^{-1} \cosh g^{-1} \frac{y}{2} s + \cosh g^{-1} \frac{y}{2} s + (I_{k_1} - U_1 U_1^T) v_1; \\
& \lim_{t \uparrow 1} v_2(t) \\
&= U_2 \lim_{t \uparrow 1} v_2(t) + (I_{k_2} - U_2 U_2^T) v_2 \\
&= U_2^{-1} \sinh g^{-1} \frac{y}{2} s + \cosh g^{-1} \frac{y}{2} s + (I_{k_2} - U_2 U_2^T) v_2;
\end{aligned}$$

This finishes the proof.

## H PROOF OF THEOREM 7

### H.1 CONVERGENCE OF LOSS TO ZERO

We first show that given the conditions on initialization, the training loss  $\mathcal{L}(t)$  converges to zero. Recall from (10) that the linear fully-connected network can be written as

$$f_{fc}(x; \mathbf{w}_{fc}) = x^T W_1 W_2 \dots W_{L-1} W_L;$$

From the definition of the training loss, it is straightforward to check that the gradient flow dynamics read

$$\begin{aligned}
\dot{W}_l &= -r_{w_l} L(f_{fc}) = -W_l^T X^T r w_L^T W_{l-1}^T - W_{l+1}^T, \text{ for } l \in [L-1]; \\
\dot{W}_L &= -r_{w_L} L(f_{fc}) = -W_L^T X^T r; \\
W_l(0) &= W_l \text{ for } l \in [L-1]; \\
w_L(0) &= w_L;
\end{aligned} \tag{77}$$

where  $r \in \mathbb{R}^n$  is the residual vector satisfying  $r_i = f_{fc}(x_i; \mathbf{w}_{fc}) - y_i$ , as defined in Section 2.1. From (77), we have

$$\begin{aligned}
W_l^T \dot{W}_l &= W_{l+1}^T W_{l+1}^T = -W_l^T X^T r w_L^T W_{l-1}^T - W_{l+1}^T; \\
W_l^T \dot{W}_l &= W_{l+1}^T W_{l+1}^T = -W_{l+1}^T W_{L-1} w_L r^T X W_l - W_l;
\end{aligned}$$

for any  $l \in [L-2]$ . From this, we have

$$\frac{d}{dt} W_l^T W_l = \frac{d}{dt} W_{l+1}^T W_{l+1};$$

and thus

$$\begin{aligned}
W_l(t)^T W_l(t) - W_{l+1}(t) W_{l+1}(t)^T &= W_l(0)^T W_l(0) - W_{l+1}(0) W_{l+1}(0)^T \\
&= -2W_l^T W_l - 2W_{l+1}^T W_{l+1};
\end{aligned} \tag{78}$$

for any  $l \in [L-2]$ . Similarly, we have

$$\begin{aligned}
W_{L-1}(t)^T W_{L-1}(t) - w_L(t) w_L(t)^T &= W_{L-1}(0)^T W_{L-1}(0) - w_L(0) w_L(0)^T \\
&= -2W_{L-1}^T W_{L-1} - 2w_L w_L^T;
\end{aligned} \tag{79}$$

Let us now consider the time derivative of  $\mathcal{L}(f_{fc}(t))$ . We have the following chain of upper bounds on the time derivative:

$$\begin{aligned}
\frac{d}{dt} \mathcal{L}(f_{fc}(t)) &= r_{fc}^T L(f_{fc}(t)) - \dot{f}_{fc}(t) = k r_{fc}^T L(f_{fc}(t)) k_2^2 \\
&\leq k r_{w_L}^T L(f_{fc}(t)) k_2^2 = k w_L(t) k_2^2 \\
&= k W_{L-1}^T X^T r k_2^2;
\end{aligned} \tag{80}$$

Note from (80) that if  $W_{L-1}^T W_1^T$  is full-rank, its minimum singular value is positive, and one can bound

$$k W_{L-1}^T W_1^T X^T r k_2 \geq \min(W_{L-1}^T W_1^T) k X^T r k_2 \quad (81)$$

We now prove that the matrix  $W_{L-1}^T W_1^T$  is full-rank, and its minimum singular value is bounded from below by  $\epsilon^{L-1} (L-1)^{-2}$  for any  $\epsilon > 0$ . To show this, it suffices to show that

$$W_{L-1}^T W_1^T W_1 W_{L-1} \geq \epsilon^{2L-2} (L-1) I_d \quad (82)$$

Now,

$$\begin{aligned} & W_{L-1}^T W_1^T W_1 W_{L-1} \\ \stackrel{(a)}{=} & W_{L-1}^T W_1^T (W_2 W_2^T + \epsilon^2 W_1^T W_1 + \epsilon^2 W_2 W_2^T) W_1 W_{L-1} \\ \stackrel{(b)}{=} & W_{L-1}^T W_1^T W_2^T W_2 W_2^T W_2 W_3 W_{L-1} \\ \stackrel{(a)}{=} & W_{L-1}^T W_1^T W_3^T (W_3 W_3^T + \epsilon^2 W_2^T W_2 + \epsilon^2 W_3 W_3^T) W_3 W_{L-1} \\ \stackrel{(b)}{=} & W_{L-1}^T W_1^T W_3^T (W_3 W_3^T)^2 W_3 W_{L-1} \\ = & (W_{L-1}^T W_1 W_{L-1})^{L-1}; \end{aligned}$$

where equalities marked in (a) used (78), and inequalities marked in (b) used the initialization conditions  $W_i^T W_i = W_{i+1}^T W_{i+1}$ . Next, it follows from (79) that

$$\begin{aligned} (W_{L-1}^T W_1 W_{L-1})^{L-1} &= (\epsilon^{2L-2} (W_{L-1}^T W_1 W_{L-1} + \epsilon^2 W_{L-1}^T W_1 W_{L-1} + \epsilon^2 W_L W_L^T)^{L-1} \\ &\geq \epsilon^{2L-2} (W_{L-1}^T W_1 W_{L-1} + \epsilon^2 W_L W_L^T)^{L-1} \\ \stackrel{(c)}{\geq} & \epsilon^{2L-2} (L-1) I_d; \end{aligned}$$

where (c) used the assumption that  $W_L^T W_L = W_L W_L^T = I_d$ . This proves (82). Applying (82) to (80) then gives

$$\begin{aligned} \frac{d}{dt} L(\mathbf{f}_c(t)) &\leq k W_{L-1}^T W_1^T X^T r k_2^2 \\ &\leq \min(W_{L-1}^T W_1^T)^2 k X^T r k_2^2 \\ &\leq \epsilon^{2L-2} (L-1) k X^T r k_2^2 \\ \stackrel{(d)}{\leq} & \epsilon^{2L-2} (L-1) \min(X)^2 k r k_2^2 \\ = & \epsilon^{2L-2} (L-1) \min(X)^2 L(\mathbf{f}_c(t)); \end{aligned}$$

where (d) used the fact that  $X^T$  is a full column rank matrix to apply a bound similar to (81). From this, we get

$$L(\mathbf{f}_c(t)) \leq L(\mathbf{f}_c(0)) \exp(-\epsilon^{2L-2} (L-1) \min(X)^2 t);$$

hence proving  $L(\mathbf{f}_c(t)) \rightarrow 0$  as  $t \rightarrow \infty$ .

## H.2 CHARACTERIZING THE LIMIT POINT: $\epsilon \rightarrow 0$ CASE

Now, we move on to characterize the limit points of the gradient flow, for the “active regime” case  $\epsilon \rightarrow 0$ . This part of the proof is motivated from the analysis in Ji & Telgarsky (2019a).

Let  $u_l$  and  $v_l$  be the top left and right singular vectors of  $W_l$ , for  $l \in [L-1]$ . Note that since  $W_l$  varies over time, the singular vectors and singular value also vary over time. Similarly, let the largest singular value of  $W_l$ . We will show that the linear coefficients  $\mathbf{f}_c(\mathbf{f}_c) = W_1 W_{L-1} W_L$  align with  $u_1$  as  $\epsilon \rightarrow 0$ , and  $u_1$  is in the subspace  $\text{span}(X)$  in the limit  $\epsilon \rightarrow 0$ , hence proving that  $\mathbf{f}_c(\mathbf{f}_c)$  is the minimum  $\ell_2$  norm solution in the limit  $\epsilon \rightarrow 0$ .

First, note from (78) and (79) that if we take trace of both sides, we get

$$k W_l k_F^2 \leq k W_{l+1} k_F^2 = \epsilon^2 (W_l^2 + W_{l+1}^2) \text{ for } l \in [L-2];$$

$$k\mathbf{W}_{L-1}k_{\text{F}}^2 \quad k\mathbf{W}_Lk_2^2 = \alpha^2(\mathbf{W}_{L-1}^2 \quad k\mathbf{W}_Lk_2^2).$$

Summing the equations above for  $l, l+1, \dots, L-1$ , we get

$$k\mathbf{W}_l k_{\text{F}}^2 \quad k\mathbf{W}_L k_2^2 = \alpha^2(\mathbf{W}_l^2 \quad k\mathbf{W}_L k_2^2). \quad (83)$$

Next, consider the operator norms (i.e., the maximum singular values), denoted as  $k k_2$ , of the matrices.

$$\begin{aligned} k\mathbf{W}_l k_2^2 &= \mathbf{u}_{l+1}^T \mathbf{W}_l^T \mathbf{W}_l \mathbf{u}_{l+1} \\ &\stackrel{(e)}{=} \mathbf{u}_{l+1}^T \mathbf{W}_{l+1} \mathbf{W}_{l+1}^T \mathbf{u}_{l+1} + \alpha^2 \mathbf{u}_{l+1}^T (\mathbf{W}_l^T \mathbf{W}_l \quad \mathbf{W}_{l+1} \mathbf{W}_{l+1}^T) \mathbf{u}_{l+1} \\ &= k\mathbf{W}_{l+1}k_2^2 + \alpha^2 \mathbf{u}_{l+1}^T (\mathbf{W}_l^T \mathbf{W}_l \quad \mathbf{W}_{l+1} \mathbf{W}_{l+1}^T) \mathbf{u}_{l+1} \\ &\quad k\mathbf{W}_{l+1}k_2^2 \quad \alpha^2 k\mathbf{W}_l^T \mathbf{W}_l \quad \mathbf{W}_{l+1} \mathbf{W}_{l+1}^T k_2 \quad \text{for } l \geq [L-2], \\ k\mathbf{W}_{L-1}k_2^2 &= \frac{\mathbf{W}_L}{k\mathbf{W}_L k_2} \mathbf{W}_{L-1}^T \mathbf{W}_{L-1} \frac{\mathbf{W}_L}{k\mathbf{W}_L k_2} \\ &\stackrel{(f)}{=} \frac{\mathbf{W}_L}{k\mathbf{W}_L k_2} \mathbf{W}_L \mathbf{W}_L^T \frac{\mathbf{W}_L}{k\mathbf{W}_L k_2} + \alpha^2 \frac{\mathbf{W}_L}{k\mathbf{W}_L k_2} (\mathbf{W}_{L-1}^T \mathbf{W}_{L-1} \quad \mathbf{W}_L \mathbf{W}_L^T) \frac{\mathbf{W}_L}{k\mathbf{W}_L k_2} \\ &\quad k\mathbf{W}_L k_2^2 \quad \alpha^2 k\mathbf{W}_{L-1}^T \mathbf{W}_{L-1} \quad \mathbf{W}_L \mathbf{W}_L^T k_2. \end{aligned}$$

where (e) used (78) and (f) used (79). Summing the inequalities gives

$$k\mathbf{W}_l k_2^2 \quad k\mathbf{W}_L k_2^2 \quad \alpha^2 \sum_{k=l}^{L-1} k\mathbf{W}_k^T \mathbf{W}_k \quad \mathbf{W}_{k+1} \mathbf{W}_{k+1}^T k_2. \quad (84)$$

From (83) and (84), we get a bound on the gap between the second powers of the Frobenius norm (or the  $\ell_2$  norm of singular values) and operator norm (or the maximum singular value  $s_l$ ) of  $\mathbf{W}_l$ :

$$k\mathbf{W}_l(t)k_{\text{F}}^2 \quad k\mathbf{W}_l(t)k_2^2 \quad \alpha^2(\mathbf{W}_l^2 \quad k\mathbf{W}_L k_2^2) + \alpha^2 \sum_{k=l}^{L-1} k\mathbf{W}_k^T \mathbf{W}_k \quad \mathbf{W}_{k+1} \mathbf{W}_{k+1}^T k_2, \quad (85)$$

which holds for any  $t \geq 0$ . The gap (85) implies that each  $\mathbf{W}_l$ , for  $l \geq [L-1]$ , can be written as

$$\mathbf{W}_l(t) = s_l(t) \mathbf{u}_l(t) \mathbf{v}_l(t)^T + O(\alpha^2). \quad (86)$$

Next, we show that the ‘‘adjacent’’ singular vectors  $\mathbf{v}_l$  and  $\mathbf{u}_{l+1}$  align with each other as  $\alpha \rightarrow 0$ . To this end, we will get lower and upper bounds for a quantity  $\mathbf{v}_l^T \mathbf{W}_{l+1} \mathbf{W}_{l+1}^T \mathbf{v}_l$ .

$$\begin{aligned} \mathbf{v}_l^T \mathbf{W}_{l+1} \mathbf{W}_{l+1}^T \mathbf{v}_l &= \mathbf{v}_l^T \mathbf{W}_l^T \mathbf{W}_l \mathbf{v}_l \quad \alpha^2 \mathbf{v}_l^T \mathbf{W}_l^T \mathbf{W}_l \mathbf{v}_l + \alpha^2 \mathbf{v}_l^T \mathbf{W}_{l+1} \mathbf{W}_{l+1}^T \mathbf{v}_l \\ &\quad k\mathbf{W}_l k_2^2 \quad \alpha^2 \mathbf{W}_l^T \mathbf{W}_l \quad \mathbf{W}_{l+1} \mathbf{W}_{l+1}^T \quad 2 \\ &= s_l^2 \quad \alpha^2 \mathbf{W}_l^T \mathbf{W}_l \quad \mathbf{W}_{l+1} \mathbf{W}_{l+1}^T \quad 2, \end{aligned} \quad (87)$$

$$\begin{aligned} \mathbf{v}_l^T \mathbf{W}_{l+1} \mathbf{W}_{l+1}^T \mathbf{v}_l &= \mathbf{v}_l^T (s_{l+1}^2 \mathbf{u}_{l+1} \mathbf{u}_{l+1}^T + \mathbf{W}_{l+1} \mathbf{W}_{l+1}^T \quad s_{l+1}^2 \mathbf{u}_{l+1} \mathbf{u}_{l+1}^T) \mathbf{v}_l \\ &= s_{l+1}^2 (\mathbf{v}_l^T \mathbf{u}_{l+1})^2 + \mathbf{v}_l^T (\mathbf{W}_{l+1} \mathbf{W}_{l+1}^T \quad s_{l+1}^2 \mathbf{u}_{l+1} \mathbf{u}_{l+1}^T) \mathbf{v}_l \\ &\quad s_{l+1}^2 (\mathbf{v}_l^T \mathbf{u}_{l+1})^2 + k\mathbf{W}_{l+1}k_{\text{F}}^2 \quad k\mathbf{W}_{l+1}k_2^2. \end{aligned} \quad (88)$$

Combining (87), (88), and (85), we get

$$\begin{aligned} s_l^2 \quad s_{l+1}^2 (\mathbf{v}_l^T \mathbf{u}_{l+1})^2 + \alpha^2 \mathbf{W}_l^T \mathbf{W}_l \quad \mathbf{W}_{l+1} \mathbf{W}_{l+1}^T \quad 2 + k\mathbf{W}_{l+1}k_{\text{F}}^2 \quad k\mathbf{W}_{l+1}k_2^2 \\ s_{l+1}^2 (\mathbf{v}_l^T \mathbf{u}_{l+1})^2 + \alpha^2(\mathbf{W}_{l+1}^2 \quad k\mathbf{W}_L k_2^2) + \alpha^2 \sum_{k=l}^{L-1} k\mathbf{W}_k^T \mathbf{W}_k \quad \mathbf{W}_{k+1} \mathbf{W}_{k+1}^T k_2. \end{aligned} \quad (89)$$

Next, by a similar reasoning as (87), we have

$$s_l^2 \quad \mathbf{u}_{l+1}^T \mathbf{W}_l^T \mathbf{W}_l \mathbf{u}_{l+1} \quad s_{l+1}^2 \quad \alpha^2 \mathbf{W}_l^T \mathbf{W}_l \quad \mathbf{W}_{l+1} \mathbf{W}_{l+1}^T \quad 2. \quad (90)$$

Combining (89) and (90) and dividing both sides by  $s_{l+1}^2$ , we get

$$(\mathbf{v}_l(t)^T \mathbf{u}_{l+1}(t))^2 \geq 1 \quad \alpha^2 \frac{G_l}{s_{l+1}(t)^2} \quad (91)$$



for  $t \geq 0$ , where

$$G_l := \mathbf{W}_l^T \mathbf{W}_l - \mathbf{W}_{l+1} \mathbf{W}_{l+1}^T + (\mathbf{W}_{l+1}^T \mathbf{F} - \mathbf{F}^T \mathbf{W}_{l+1})^2 + \sum_{k=l}^{L-1} k \mathbf{W}_k \mathbf{W}_k^T - \mathbf{W}_{k+1} \mathbf{W}_{k+1}^T k_2.$$

By a similar argument, we can also get

$$\frac{(\mathbf{v}_{L-1}(t)^T \mathbf{w}_L(t))^2}{k \mathbf{W}_L(t) k_2^2} \geq 1 - \alpha^2 \frac{G_{L-1}}{k \mathbf{W}_L(t) k_2^2}, \quad (92)$$

where

$$G_{L-1} := 2 \mathbf{W}_{L-1}^T \mathbf{W}_{L-1} - \mathbf{w}_L \mathbf{w}_L^T k_2.$$

From (91) and (92), we can note that as  $\alpha \rightarrow 0$ , the inner product between the adjacent singular vectors converges to 1, unless  $s_2, \dots, s_{L-1}, k \mathbf{W}_L k_2$  also diminish to zero. So it is left to show that the singular values do not diminish to zero as  $\alpha \rightarrow 0$ . To this end, recall that we proved in the previous subsection that

$$\lim_{t \rightarrow \infty} \mathbf{X} \mathbf{W}_1(t) - \mathbf{W}_{L-1}(t) \mathbf{w}_L(t) = \mathbf{y}.$$

A necessary condition for this to hold is that

$$\frac{k \mathbf{y} k_2}{k \mathbf{X} k_2} \leq \lim_{t \rightarrow \infty} k \mathbf{W}_1(t) - \mathbf{W}_{L-1}(t) \mathbf{w}_L(t) k_2 \leq \lim_{t \rightarrow \infty} \sum_{l=1}^{L-1} s_l(t) k \mathbf{W}_L(t) k_2.$$

This means that after converging to the global minimum solution of the problem (i.e.,  $t \rightarrow \infty$ ), the product of the singular values must be at least greater than some constant independent of  $\alpha$ . Moreover, we can see from (87) and (90) that the gap between singular values squared of adjacent layers is bounded by  $O(\alpha^2)$ , for all  $t \geq 0$ ; so the maximum singular values become closer and closer to each other as  $\alpha$  diminishes. This implies that

$$\lim_{\alpha \rightarrow 0} \lim_{t \rightarrow \infty} s_l(t) \geq \frac{k \mathbf{y} k_2^{1/L}}{k \mathbf{X} k_2^{1/L}} \text{ for } l \geq [L-1], \quad \lim_{\alpha \rightarrow 0} \lim_{t \rightarrow \infty} k \mathbf{W}_L(t) k_2 \geq \frac{k \mathbf{y} k_2^{1/L}}{k \mathbf{X} k_2^{1/L}}.$$

Therefore, we have the alignment of singular vectors at convergence as  $\alpha \rightarrow 0$ :

$$\lim_{\alpha \rightarrow 0} \lim_{t \rightarrow \infty} (\mathbf{v}_l(t)^T \mathbf{u}_{l+1}(t))^2 = 1, \text{ for } l \geq [L-2], \quad \lim_{\alpha \rightarrow 0} \lim_{t \rightarrow \infty} \frac{(\mathbf{v}_{L-1}(t)^T \mathbf{w}_L(t))^2}{k \mathbf{W}_L(t) k_2^2} = 1. \quad (93)$$

So far, we saw from (86) that  $\mathbf{W}_l(t)$ 's become rank-1 matrices as  $\alpha \rightarrow 0$ , and from (93) that the top singular vectors align with each other as  $t \rightarrow \infty$  and  $\alpha \rightarrow 0$ . These imply that, as  $t \rightarrow \infty$  and  $\alpha \rightarrow 0$ ,  $\mathbf{f}_c(\mathbf{f}_c)$  is a scalar multiple of the  $\mathbf{u}_1$ , the top left singular vector of  $\mathbf{W}_1$ :

$$\lim_{\alpha \rightarrow 0} \lim_{t \rightarrow \infty} \mathbf{f}_c(\mathbf{f}_c(t)) = c \lim_{\alpha \rightarrow 0} \lim_{t \rightarrow \infty} \mathbf{u}_1(t), \quad (94)$$

for some  $c \in \mathbb{R}$ .

In light of (94), it remains to take a close look at  $\mathbf{u}_1(t)$ . Note from the gradient flow dynamics of  $\mathbf{W}_1$  that  $\mathbf{W}_1$  is always a rank-1 matrix whose columns are in the row space of  $\mathbf{X}$ , since  $\mathbf{X}^T \mathbf{r} \in \text{row}(\mathbf{X})$ . This implies that, if we decompose  $\mathbf{W}_1$  into two orthogonal components  $\mathbf{W}_1^\perp$  and  $\mathbf{W}_1^\parallel$  so that the columns in  $\mathbf{W}_1^\parallel$  are in  $\text{row}(\mathbf{X})$  and the columns in  $\mathbf{W}_1^\perp$  are in the orthogonal subspace  $\text{row}(\mathbf{X})^\perp$ , we have

$$\mathbf{W}_1^\perp = \mathbf{0}, \quad \mathbf{W}_1^\parallel = \mathbf{W}_1.$$

That is, any component  $\mathbf{W}_1^\perp(0)$  orthogonal to  $\text{row}(\mathbf{X})$  remains unchanged for all  $t \geq 0$ , while the component  $\mathbf{W}_1^\parallel$  changes by the gradient flow. Since we have

$$\mathbf{W}_1^\perp(t) = \mathbf{W}_1^\perp(0) - \alpha \mathbf{W}_t^\perp,$$

the component in  $\mathbf{W}_1$  that is orthogonal to  $\text{row}(\mathbf{X})$  diminishes to zero as  $\alpha \rightarrow 0$ . This means that at the limit  $\alpha \rightarrow 0$ , the columns of  $\mathbf{W}_1$  are entirely from  $\text{row}(\mathbf{X})$ , which also means that

$$\lim_{\alpha \rightarrow 0} \lim_{t \rightarrow \infty} \mathbf{f}_c(\mathbf{f}_c(t)) \in \text{row}(\mathbf{X}).$$

However, recall that there is only one unique global minimum of  $\mathbf{X} \mathbf{z} = \mathbf{y}$  in  $\text{row}(\mathbf{X})$ : namely,  $\mathbf{z} = \mathbf{X}^T (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{y}$ , the minimum  $\ell_2$  norm solution. This finishes the proof.