

A Unifying view on implicit bias  
in training linear neural networks.  
- Chulhee Yun et al.

이종진

Seoul National University

*ga0408@snu.ac.kr*

May 20, 2021

# Table of Contents

1. Contributions

2. Preliminary

3. main theorems

## Contributions

- ▶ Implicit bias gradient flow of the linear tensor networks.
- ▶ Consider two cases(separable classification / undetermined regression)
- ▶ Subsume existing results without removing standard convergence assumptions.

## Contributions

- ▶ Linear tensor networks / classification  
→ Singular vectors of a tensor defined by the network.
- ▶ Orthogonally decomposable linear network / classification  
→ A solution of minimizing  $\ell_{2/L}$  max-margin problem in a "transformed" input space defined by the network.
- ▶ Orthogonally decomposable linear network / regression  
→ A solution of minimizing norm-like functions that interpolates between weighted  $\ell_1$  and  $\ell_2$  in a "transformed" input space.

## Contributions; in a separable classification

- ▶ (Thm 1), A linear tensor networks.
  - (Cor 1) A  $L$ -layer linear fully-connected network
- ▶ (Thm 2), A orthogonally decomposable linear network
  - (Cor 2), A  $L$ -layer linear diagonal network.
  - (Cor 3), A  $L$ -layer linear full-length convolution network.
- ▶ (Thm 3), A 2-layer linear network with a single data point  $(x,y)$ 
  - (Cor 4), A 2 - layer linear convolutional network with a single data point  $(x,y)$

## Contributions; in a undertermined regression

- ▶ (Thm 5), A orthogonally decomposable linear network
  - (Cor 5), A  $L$ -layer linear diagonal network.
  - (Cor 6), A  $L$ -layer linear full-length convolution network.
- ▶ (Thm 6), A 2-layer linear network with a single data point  $(x,y)$ .

## Problem settings

- ▶  $\{(x_i, y_i)\}_{i=1}^n$ , where  $x_i \in \mathbb{R}^d$  and  $y_i \in \mathbb{R}$
- ▶  $\mathbf{X} \in \mathbb{R}^{n \times d}$ ,  $y \in \mathbb{R}^n$
- ▶ For binary classification,
  - $y_i \in \{\pm 1\}$
  - Data is separable
  - Exponential loss,  $\ell(\hat{y}, y) = \exp(-\hat{y}y)$
- ▶ For regression
  - Undetermined case ( $n \geq d$ )
  - Squared error loss,  $\ell(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$

## Tensor networks

- ▶ A linear map  $M$  that maps  $x$  to an order-  $L$  tensor  $M(x) \in \mathbb{R}^{k_1 \times \dots \times k_L}$ , where  $L \geq 2$ .

- ▶ A tensor network with parameters  $v_l \in \mathbb{R}^{k_l}$  and activation  $\phi$ ,

$$H_1(\mathbf{x}) = \phi(M(\mathbf{x}) \circ (\mathbf{v}_1, \mathbf{l}_{k_2}, \dots, \mathbf{l}_{k_L})) \in \mathbb{R}^{k_2 \times \dots \times k_L}$$

$$H_l(\mathbf{x}) = \phi(H_{l-1}(\mathbf{x}) \circ (\mathbf{v}_l, \mathbf{l}_{k_{l+1}}, \dots, \mathbf{l}_{k_L})) \in \mathbb{R}^{k_{l+1} \times \dots \times k_L}, \text{ for } l = 2, \dots, L-1$$

$$f(\mathbf{x}; \Theta) = H_{L-1}(\mathbf{x}) \circ v_L \in \mathbb{R}$$

where  $\circ$  is a multilinear multiplication.

- ▶ Use  $\Theta$  to denote the collection of all parameters  $(v_1, \dots, v_L)$  and name  $M(x)$  as a data tensor.



## A multilinear multiplication

- Given a tensor  $A \in \mathbb{R}^{k_1 \times \cdots \times k_L}$  and linear maps  $B_l \in \mathbb{R}^{p_l \times k_l}$  for  $l \in [L]$ , the multilinear multiplication  $\circ$  between them is defined as

$$\begin{aligned} A \circ (\mathbf{B}_1^T, \mathbf{B}_2^T, \dots, \mathbf{B}_L^T) &= \sum_{j_1, \dots, j_L} [A]_{j_1, \dots, j_L} \left( e_{j_1}^{k_1} \otimes \cdots \otimes e_{j_L}^{k_L} \right) \circ (\mathbf{B}_1^T, \dots, \mathbf{B}_L^T) \\ &:= \sum_{j_1, \dots, j_L} [A]_{j_1, \dots, j_L} \left( \mathbf{B}_1 e_{j_1}^{k_1} \otimes \cdots \otimes \mathbf{B}_L e_{j_L}^{k_L} \right) \in \mathbb{R}^{p_1 \times \cdots \times p_L} \end{aligned}$$

## Linear tensor networks

- ▶ The tensor formulation includes
  1. Diagonal networks
  2. Convolution networks
  3. Fully-connected networks.
- ▶ Consider linear tensor networks, which means  $\phi(t) = t$ .

$$f(x; \Theta) = M(x) \circ (v_1, v_2, \dots, v_L)$$

- ▶ The output of the network can also be written as  $f(x; \Theta) = x^T \beta(\Theta)$ , where  $\beta(\Theta) \in \mathbb{R}^d$

## Diagonal networks

- ▶ An  $L$ -layer diagonal network can be written as

$$f_{\text{diag}}(x; \Theta_{\text{diag}}) = \phi(\cdots \phi(\phi(x \odot w_1) \odot w_2) \cdots \odot w_{L-1})^T w_L$$

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

- ▶  $M_{\text{diag}}(x) \in \mathbb{R}^{d \times \cdots \times d}$  and  $[M_{\text{diag}}(x)]_{j,j,\dots,j} = [x]_j$ , while other components are 0.
- ▶  $v_l = w_l$  for all  $l$

## Convolutional networks

- ▶ The convolutional networks can be written as

$$f_{\text{conv}}(x; \Theta_{\text{conv}}) = \phi(\cdots \phi(\phi(x \star w_1) \star w_2) \cdots \star w_{L-1})^T w_L,$$

where  $w_l \in \mathbb{R}^{k_l}$  with  $k_l \leq d$  and  $k_L = d$ , and  $\star$  defines the circular convolution.

- ▶  $a \star b \in \mathbb{R}^d$  defined as  $[a \star b]_i = \sum_{j=1}^k [a]_{(i+j-1) \bmod d} [b]_j$ , for  $i \in [d]$ . for any  $a \in \mathbb{R}^d$  and  $b \in \mathbb{R}^k (k \leq d)$
- ▶  $M_{\text{conv}}(x) \in \mathbb{R}^{k_1 \times \cdots \times k_L}$  as  $[M_{\text{conv}}(x)]_{j_1, j_2, \dots, j_L} = [x]_{(\sum_{l=1}^L j_l - L + 1) \bmod d}$  for  $j_l \in [k_l], l \in [L]$ .
- ▶  $v_l = w_l$  and  $M = M_{\text{conv}}$ .

## Fully-connected networks

- ▶ An  $L$ -layer fully-connected network is defined as

$$f_{\text{fc}}(x; \Theta_{\text{fc}}) = \phi \left( \cdots \phi \left( \phi \left( x^T W_1 \right) W_2 \right) \cdots W_{L-1} \right) w_L$$

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

- ▶ One can represent  $f_{\text{fc}}$  as the tensor form by
  - Defining parameters  $v_l = \text{vec}(W_l)$  for  $l \in [L-1]$  and  $v_L = w_L$
  - Constructing the tensor  $M_{\text{fc}}(x)$  by a recursive "block diagonal" manner.

## Singular value decomposition of tensor

- ▶ Given an order-  $L$  tensor  $A \in \mathbb{R}^{k_1 \times \dots \times k_L}$ , we define the singular vectors  $u_1, u_2, \dots, u_L$  and singular value  $s$  to be the solution of the following system of equations:

$$su_l = A \circ (u_1, \dots, u_{l-1}, I_{k_l}, u_{l+1}, \dots, u_L), \text{ for } l \in [L]$$

- ▶ We can characterize the limit direction of parameters after reaching 100% training accuracy.

## Theorem 1

### Theorem 1

Assume that the gradient flow satisfies  $\mathcal{L}(\Theta(t_0)) < 1$  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  $\mathbf{M}(-\mathbf{u}^\infty)$ . where  $\mathbf{r}(t) \in \mathbb{R}^n$  is defined as

$$[\mathbf{r}(t)]_i = \ell'(f(x_i; \Theta(t)), y_i) = \begin{cases} -y_i \exp(-y_i f(x_i; \Theta(t))) & \text{for classification,} \\ f(x_i; \Theta(t)) - y_i & \text{for regression.} \end{cases}$$

$$\blacktriangleright \dot{\mathbf{v}}_l = -\nabla_{\mathbf{v}_l} \mathcal{L}(\Theta) = \mathbf{M}(-\mathbf{X}^T \mathbf{r}) \circ (\mathbf{v}_1, \dots, \mathbf{v}_{l-1}, \mathbf{I}_{k_l}, \mathbf{v}_{l+1}, \dots, \mathbf{v}_L), \quad \forall l \in [L]$$

## Corollary 1

Corollary 1. (cf. Ji & Telgarsky, 2020)

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

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



## Theorem 2

- ▶ Theorem 1 is not a full characterization of the limit directions, because there are usually multiple solutions that satisfy a condition of singular value and vectors.
- ▶ Singular vectors of high order tensors are much less understood than the matrix counterparts, let alone orthogonal decompositions.
- ▶ The following assumptions define a class of orthogonally decomposable data tensors  $\mathbf{M}(x)$

## Condition for orthogonally decomposable data tensor.

### Assumption 1

For the data tensor  $M(x) \in \mathbb{R}^{k_1 \times \dots \times k_L}$  of a linear tensor network (6), there exist a full column rank matrix  $S \in \mathbb{C}^{m \times d}$  ( $d \leq m \leq \min_l k_l$ ) and matrices  $U_1 \in \mathbb{C}^{k_1 \times m}, \dots, U_L \in \mathbb{C}^{k_L \times m}$  such that  $U_l^H U_l = I_m$  for all  $l \in [L]$ , and the data tensor  $M(x)$  can be written as

$$M(x) = \sum_{j=1}^m [\mathbf{S}x]_j \left( [U_1]_{\cdot,j} \otimes [U_2]_{\cdot,j} \otimes \dots \otimes [U_L]_{\cdot,j} \right)$$

## Theorem 2

### Theorem 2

Suppose a linear tensor network satisfies Assumption 1. If there exists  $\lambda > 0$  such that the initial directions  $\bar{v}_1, \dots, \bar{v}_L$  of the network parameters satisfy  $\left| [U_\ell^T \bar{v}_\ell]_j \right|^2 - \left| [U_L^T \bar{v}_L]_j \right|^2 \geq \lambda$  for all  $\ell \in [L-1]$  and  $j \in [m]$ , then  $\beta(\Theta(t))$  converges in a direction that aligns with  $S^T \rho^\infty$ , where  $\rho^\infty \in \mathbb{C}^m$  denotes a stationary point of the following optimization problem

$$\underset{\rho \in \mathbb{C}^m}{\text{minimize}} \quad \|\rho\|_{2/L} \quad \text{subject to} \quad y_i x_i^T S^T \rho \geq 1, \quad \forall i \in [n]$$

## Theorem 2

- ▶ The gradient flow finds sparse  $\rho^\infty$  that minimizes the  $\ell_{2/L}$  norm in the "singular value space," where the data points  $x_i$  are transformed into vectors  $Sx_i$  consisting of singular values of  $M(x_i)$ .
- ▶ Also, the proof of Theorem 2 reveals that in case of  $L = 2$ , the parameters  $v_l(t)$  in fact converge to the top singular vectors of the data tensor  $M(-X^T r)$ ;
- ▶ Compared to Theorem 1. we have a more complete characterization of "which" singular vectors to converge to.

## Corollary 2

### Corollary 2

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

$$\underset{z \in \mathbb{R}^d}{\text{minimize}} \quad \|z\|_{2/L} \quad \text{subject to} \quad y_i x_i^T z \geq 1, \forall i \in [n]$$

## Corollary 3 (cf. Gunasekar et al., 2018b)

### Corollary 3

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

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

where  $\mathbf{F} \in \mathbb{C}^{d \times d}$  to be the matrix of discrete Fourier transform basis

$$[\mathbf{F}]_{j,k} = \frac{1}{\sqrt{d}} \exp\left(-\frac{\sqrt{-1} \cdot 2\pi(j-1)(k-1)}{d}\right).$$

## Corollary 3

- ▶ For full-length convolution networks ( $k_1 = \dots = k_L = d$ ) satisfy Assumption 1.
- ▶  $S = d^{\frac{L-1}{2}}$  and  $U_1 = \dots = U_L = F^*$

## Theorem 3

### Theorem 3

Suppose we have a 2-layer linear tensor network and a single data point  $(x, y)$ .

Consider the compact SVD  $M(x) = U_1 \text{diag}(s) U_2^T$ , where

$U_1 \in \mathbb{R}^{k_1 \times m}$ ,  $U_2 \in \mathbb{R}^{k_2 \times m}$ , and  $s \in \mathbb{R}^m$  for  $m \leq \min\{k_1, k_2\}$ .

Let  $\rho^\infty \in \mathbb{R}^m$  be a solution of the following optimization problem

$$\text{minimize}_{\rho \in \mathbb{R}^m} \quad \|\rho\|_1 \quad \text{subject to} \quad y s^T \rho \geq 1$$

Assume that there exists  $\lambda > 0$  such that the initial directions  $\bar{v}_1, \bar{v}_2$  of the network parameters satisfy  $[\mathbf{U}_1^T \bar{v}_1]_j^2 - [\mathbf{U}_2^T \bar{v}_2]_j^2 \geq \lambda$  for all  $j \in [m]$ . Then,  $v_1$  and  $v_2$  converge in direction to  $U_1 \eta_1^\infty$  and  $U_2 \eta_2^\infty$ , where

$|\eta_1^\infty| = |\eta_2^\infty| = |\rho^\infty|^{\odot 1/2}$ , and  $\text{sign}(\eta_1^\infty) = \text{sign}(y) \odot \text{sign}(\eta_2^\infty)$ .



## Underdetermined regression

- ▶ 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.
- ▶  $w_\ell(0) = \alpha \bar{w}_\ell$  for  $\ell \in [L - 1]$  and  $w_L(0) = 0$ .

## Lemma 4

- ▶ To define norm like function, 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  $c \in (0, \infty]$ .

Define  $h_L(t) = p_L(t)^{L-1}q_L(t)$ ; then,  $h_L(t)$  is odd and strictly increasing, satisfying  $\lim_{t \uparrow c} h_L(t) = \infty$  and  $\lim_{t \downarrow -c} h_L(t) = -\infty$ .

## Theorem 5

### Theorem 5

Suppose a linear tensor network satisfies Assumption 1. Assume further that the matrices  $U_1, \dots, U_L$  and  $S$  from Assumption 1 are all real matrices. For some  $\lambda > 0$ , choose any vector  $\bar{\eta} \in \mathbb{R}^m$  satisfying  $[\bar{\eta}]_j^2 \geq \lambda$  for all  $j \in [m]$ , and choose initial directions  $\bar{v}_\ell = U_\ell \bar{\eta}$  for  $\ell \in [L-1]$  and  $\bar{v}_L = 0$ . Then, the linear coefficients  $\beta(\Theta(t))$  converge to  $S^T \rho^\infty$ , where  $\rho^\infty$  is the solution of

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

where  $Q_{L,\alpha,\bar{\eta}} : \mathbb{R}^m \rightarrow \mathbb{R}$  is a norm-like function defined using

$$H_L(t) := \int_0^t h_L^{-1}(\tau) d\tau.$$

## Theorem 5

- ▶  $Q_{L,\alpha,\bar{\eta}}(\rho) := \alpha^2 \sum_{j=1}^m [\bar{\eta}]_j^2 H_L \left( \frac{[\rho]_j}{\alpha L |\bar{\eta}_j|^L} \right)$
- ▶  $Q_{L,\alpha,\bar{\eta}}(\rho)$  interpolates between the weighted  $\ell_1$  and weighted  $\ell_2$  norm of  $\rho$
- ▶  $H_L(t)$ :
  - grows like the absolute value function if  $t$  is large.
  - grows like a quadratic function if  $t$  is close to zero.

## Corollary 5 (cf. Woodworth et al., 2018b)

### Corollary 5

Consider an  $L$ -layer linear diagonal network. For some  $\lambda > 0$ , choose any vector  $\bar{\mathbf{w}} \in \mathbb{R}^d$  satisfying  $[\bar{\mathbf{w}}]_j^2 \geq \lambda$  for all  $j \in [d]$ , and choose initial directions  $\bar{\mathbf{w}}_l = \bar{\mathbf{w}}$  for  $l \in [L-1]$  and  $\bar{\mathbf{w}}_L = 0$ . Then, the linear coefficients  $\beta_{diag}(\Theta_{diag}(t))$  converge to the solution  $\mathbf{z}^\infty$  of

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

## Corollary 6

### Corollary 6

Consider an  $L$ -layer linear full-length convolutional network. Assume that the data points  $\{x_i\}_{i=1}^n$  are all even. For some  $\lambda > 0$ , choose any even vector  $\bar{w}$  satisfying  $[F\bar{w}]_j^2 \geq \lambda$  for all  $j \in [d]$ , and choose initial directions  $\bar{w}_l = \bar{w}$  for  $l \in [L-1]$  and  $\bar{w}_L = 0$ . Then, the linear coefficients  $\beta_{\text{conv}}(\Theta_{\text{conv}}(t))$  converge to the solution  $z^\infty$  of

$$\underset{z \in \mathbb{R}^d, \text{ even}}{\text{minimize}} Q_{L, \alpha, F\bar{w}}(\mathbf{F}z) := \alpha^2 \sum_{j=1}^d [F\bar{w}]_j^2 H_L \left( \frac{[Fz]_j}{\alpha^L |[F\bar{w}]_j|^L} \right) \quad \text{subject to} \quad \mathbf{X}z = \mathbf{y}$$

## Theorem 6

### Theorem 6

Suppose we have a 2-layer linear tensor network and a single data point  $(x, y)$ .

Consider the compact SVD  $M(x) = U_1 \text{diag}(s) U_2^T$ , where

$U_1 \in \mathbb{R}^{k_1 \times m}$ ,  $U_2 \in \mathbb{R}^{k_2 \times m}$ , and  $s \in \mathbb{R}^m$  for  $m \leq \min\{k_1, k_2\}$ . Assume that

there exists  $\lambda > 0$  such that the initial directions  $\bar{v}_1, \bar{v}_2$  of the network

parameters satisfy  $[U_1^T \bar{v}_1]_j^2 - [U_2^T \bar{v}_2]_j^2 \geq \lambda$  for all  $j \in [m]$ . Then, gradient flow

converges to a global minimizer of the loss  $\mathcal{L}$ , and  $v_1(t)$  and  $v_2(t)$  converge to

the limit points:

$$v_1^\infty = \alpha U_1 \left( U_1^T \bar{v}_1 \odot \cosh \left( g^{-1} \left( \frac{y}{\alpha^2} \right) s \right) + U_2^T \bar{v}_2 \odot \sinh \left( g^{-1} \left( \frac{y}{\alpha^2} \right) s \right) \right) + \alpha \left( I_{k_1} - U_1 U_1^T \right) \bar{v}_1$$

$$v_2^\infty = \alpha U_2 \left( U_1^T \bar{v}_1 \odot \cosh \left( g^{-1} \left( \frac{y}{\alpha^2} \right) s \right) + U_2^T \bar{v}_2 \odot \sinh \left( g^{-1} \left( \frac{y}{\alpha^2} \right) s \right) \right) + \alpha \left( I_{k_2} - U_2 U_2^T \right) \bar{v}_2$$

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

$$g(\nu) = \sum_{j=1}^m [s]_j \left( \frac{[U_1^T \bar{v}_1]_j^2 + [U_2^T \bar{v}_2]_j^2}{2} \sinh(2[s]_j \nu) + [U_1^T \bar{v}_1]_j [U_2^T \bar{v}_2]_j \cosh(2[s]_j \nu) \right)$$

The end

The End