Foundations of Computational Mathematics (FoCM) 2023

The Law of Parsimony in Gradient Descent for Learning Deep Linear Networks

Qing Qu

EECS, University of Michigan

June 30, 2023



・ 何 ト ・ ヨ ト ・ ヨ ト

Law of Parsimony in GD

1/34

Main Message



Throughout training of deep linear networks, the gradient descent (GD) dynamics possesses certain parsimonious structures.

✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 ✓ □→
 <l

Main Message



The parsimonious structures in GD dynamics leads to

- Efficient training via network compression
- Better understandings of hierarchical representations

→ ∃ →

Outline

1 Law of Parsimony in Gradient Dynamics

2 Efficient Deep Matrix Completion via Network Compression

3 Understanding Hierarchical Representations in Deep Neural Networks

(日) (四) (문) (문) (문)

4 Conclusion

Setup on Deep Linear Networks

• Training data $\{({m x}_i, {m y}_i)\}_{i=1}^N \subset \mathbb{R}^{d_x} imes \mathbb{R}^{d_y}$ with

 $oldsymbol{X} = [oldsymbol{x}_1 \ oldsymbol{x}_2 \ \dots \ oldsymbol{x}_N] \in \mathbb{R}^{d_x imes N}, \quad oldsymbol{Y} = [oldsymbol{y}_1 \ oldsymbol{y}_2 \ \dots \ oldsymbol{y}_N] \in \mathbb{R}^{d_y imes N}$

4/34

イロト 不得下 イヨト イヨト 二日

Setup on Deep Linear Networks

• Training data
$$\{({m x}_i,{m y}_i)\}_{i=1}^N \subset \mathbb{R}^{d_x} imes \mathbb{R}^{d_y}$$
 with

$$oldsymbol{X} = [oldsymbol{x}_1 \; oldsymbol{x}_2 \; \ldots \; oldsymbol{x}_N] \in \mathbb{R}^{d_x imes N}, \quad oldsymbol{Y} = [oldsymbol{y}_1 \; oldsymbol{y}_2 \; \ldots \; oldsymbol{y}_N] \in \mathbb{R}^{d_y imes N}$$

• Deep linear network (DLN):

$$f_{\boldsymbol{\Theta}}(\boldsymbol{x}) := \boldsymbol{W}_{L} \cdots \boldsymbol{W}_{1} \boldsymbol{x} = \boldsymbol{W}_{L:1} \boldsymbol{x},$$

where $W_l \in \mathbb{R}^{d_l \times d_{l-1}}$ and $\Theta = \{W_l\}_{l=1}^L$.

э

4/34

Setup on Deep Linear Networks

• Training data $\{({m x}_i, {m y}_i)\}_{i=1}^N \subset \mathbb{R}^{d_x} imes \mathbb{R}^{d_y}$ with

$$oldsymbol{X} = [oldsymbol{x}_1 \; oldsymbol{x}_2 \; \ldots \; oldsymbol{x}_N] \in \mathbb{R}^{d_x imes N}, \quad oldsymbol{Y} = [oldsymbol{y}_1 \; oldsymbol{y}_2 \; \ldots \; oldsymbol{y}_N] \in \mathbb{R}^{d_y imes N}$$

• Deep linear network (DLN):

$$f_{\Theta}(\boldsymbol{x}) := \boldsymbol{W}_{L} \cdots \boldsymbol{W}_{1} \boldsymbol{x} = \boldsymbol{W}_{L:1} \boldsymbol{x},$$

where $W_l \in \mathbb{R}^{d_l \times d_{l-1}}$ and $\Theta = \{W_l\}_{l=1}^L$.

Loss function:

$$\min_{\boldsymbol{\Theta}} \ \ell(\boldsymbol{\Theta}) = \frac{1}{2} \sum_{i=1}^{N} \|f_{\boldsymbol{\Theta}}(\boldsymbol{x}_{i}) - \boldsymbol{y}_{i}\|_{F}^{2} = \frac{1}{2} \|\boldsymbol{W}_{L:1}\boldsymbol{X} - \boldsymbol{Y}\|_{F}^{2}$$

4/34

.

3

 Orthogonal initialization. We use ε-scale orthogonal matrices for some ε > 0, with

$$\boldsymbol{W}_l^{\top}(0)\boldsymbol{W}_l(0) = \varepsilon^2 \boldsymbol{I} \quad \text{or} \quad \boldsymbol{W}_l(0)\boldsymbol{W}_l^{\top}(0) = \varepsilon^2 \boldsymbol{I}, \quad \forall l \in [L],$$

depending on the size of W_l .

5/34

<日

<</p>

 Orthogonal initialization. We use ε-scale orthogonal matrices for some ε > 0, with

$$\boldsymbol{W}_l^\top(0)\boldsymbol{W}_l(0) = \varepsilon^2\boldsymbol{I} \quad \text{or} \quad \boldsymbol{W}_l(0)\boldsymbol{W}_l^\top(0) = \varepsilon^2\boldsymbol{I}, \quad \forall l \in [L],$$

depending on the size of W_l .

• Learning dynamics of GD. We update all weights via GD for $t = 1, 2, \ldots$ as

$$\boldsymbol{W}_{l}(t) = (1 - \eta \lambda) \boldsymbol{W}_{l}(t - 1) - \eta \nabla_{\boldsymbol{W}_{l}} \ell(\boldsymbol{\Theta}(t - 1)), \ \forall \ l \in [L],$$

where $\eta > 0$ is the learning rate and $\lambda \ge 0$ controls weight decay.

イロト イヨト イヨト ・

We study the GD iterates for training DLNs under the following assumptions:

- The weight matrices are square except the last layer, i.e., $d_x = d_1 = d_2 = \cdots = d_{L-1} = d$ for some $d \in \mathbb{N}_+$.
- The input data is *whitened* in the sense that $oldsymbol{X}oldsymbol{X}^ op=oldsymbol{I}_{d_x}.^1$
- The cross correlation matrix YX[⊤] has certain *low-dimensional* structures (e.g., low-rank or wide matrix).

¹For any full rank $X \in \mathbb{R}^{d_x \times N}$ with $N \ge d_x$, whitened data can always be obtained with a data pre-processing step such as preconditioning.

Qing Qu (EECS, University of Michigan)

We study the GD iterates for training DLNs under the following assumptions:

- The weight matrices are square except the last layer, i.e., $d_x = d_1 = d_2 = \cdots = d_{L-1} = d$ for some $d \in \mathbb{N}_+$.
- The input data is *whitened* in the sense that $oldsymbol{X}oldsymbol{X}^ op=oldsymbol{I}_{d_x}.^1$
- The cross correlation matrix YX[⊤] has certain *low-dimensional* structures (e.g., low-rank or wide matrix).

Throughout training of deep networks, the gradient descent leads to certain parsimonious structures in the weight matrices.

¹For any full rank $X \in \mathbb{R}^{d_x \times N}$ with $N \ge d_x$, whitened data can always be obtained with a data pre-processing step such as preconditioning.



Figure: Evolution of SVD of the weight matrix $W_1(t) = U_1(t)\Sigma_1(t)V_1(t)^{\top}$.



Figure: Evolution of SVD of the weight matrix $W_1(t) = U_1(t)\Sigma_1(t)V_1(t)^{\top}$.

• Left: the evolution of singular values $\sigma_{1i}(t)$ throughout training $t \ge 0$;

▲ □ ▶ ▲ □ ▶ ▲ □ ▶



Figure: Evolution of SVD of the weight matrix $W_1(t) = U_1(t)\Sigma_1(t)V_1(t)^{\top}$.

- Left: the evolution of singular values $\sigma_{1i}(t)$ throughout training $t \ge 0$;
- Middle: the evolution of $\angle(\boldsymbol{v}_{1i}(t), \boldsymbol{v}_{1i}(0))$ throughout training $t \ge 0$;

< □ > < □ > < □ > < □ > < □ > < □ >



Figure: Evolution of SVD of the weight matrix $W_1(t) = U_1(t)\Sigma_1(t)V_1(t)^{\top}$.

- Left: the evolution of singular values $\sigma_{1i}(t)$ throughout training $t \ge 0$;
- Middle: the evolution of $\angle(\boldsymbol{v}_{1i}(t), \boldsymbol{v}_{1i}(0))$ throughout training $t \ge 0$;
- **Right:** the evolution of $\angle(\boldsymbol{u}_{1i}(t), \boldsymbol{u}_{1i}(0))$ throughout training $t \ge 0$.





Layer 2



Layer 3



Qing Qu (EECS, University of Michigan)

Law of Parsimony in G

э

8/34

The Evolution of Singular Spaces in GD Iterates for DLNs



Figure: Evolution of SVD of the weight matrix $W_1(t) = U_1(t)\Sigma_1(t)V_1(t)^{\top}$.

The GD learning process takes place only within a **minimal invariant subspace** of each weight matrix, while the remaining singular subspaces stay **unaffected** throughout training.

Qing Qu (EECS, University of Michigan)

Law of Parsimony in GD

June 30, 2023

< ロ > < 同 > < 回 > < 回 >

Theorem (Yaras et al.'23)

Suppose we train an *L*-layer DLN $f_{\Theta}(\cdot)$ on (X, Y) via GD, the iterates $\{W_l(t)\}_{l=1}^L$ for all $t \ge 0$ satisfy the following:

• Case 1: Suppose $YX^{\top} \in \mathbb{R}^{d_y \times d_x}$ is of rank $r \in \mathbb{N}_+$ with $d_y = d_x$, and $m = d_x - 2r > 0$. Then $\exists \{U_l\}_{l=1}^L \subseteq \mathcal{O}^d$ and $\{V_l\}_{l=1}^L \subseteq \mathcal{O}^d$ satisfying $V_{l+1} = U_l$ for all $l \in [L-1]$, such that $W_l(t)$ admits the following decomposition

$$\boldsymbol{W}_{l}(t) = \boldsymbol{U}_{l} \begin{bmatrix} \widetilde{\boldsymbol{W}}_{l}(t) & \boldsymbol{0} \\ \boldsymbol{0} & \rho(t)\boldsymbol{I}_{m} \end{bmatrix} \boldsymbol{V}_{l}^{\top}, \quad \forall l \in [L-1], \ t \geq 0,$$

where $\widetilde{W}_l(t) \in \mathbb{R}^{2r \times 2r}$ for all $l \in [L-1]$ with $\widetilde{W}_l(0) = \varepsilon I_{2r}$.

イロト 不得下 イヨト イヨト 二日

Theorem (Yaras et al.'23)

Suppose we train an *L*-layer DLN $f_{\Theta}(\cdot)$ on (X, Y) via GD, the iterates $\{W_l(t)\}_{l=1}^L$ for all $t \ge 0$ satisfy the following:

• Case 1: Suppose $YX^{\top} \in \mathbb{R}^{d_y \times d_x}$ is of rank $r \in \mathbb{N}_+$ with $d_y = d_x$, and $m = d_x - 2r > 0$. Then $\exists \{U_l\}_{l=1}^L \subseteq \mathcal{O}^d$ and $\{V_l\}_{l=1}^L \subseteq \mathcal{O}^d$ satisfying $V_{l+1} = U_l$ for all $l \in [L-1]$, such that $W_l(t)$ admits the following decomposition

$$\boldsymbol{W}_{l}(t) = \boldsymbol{U}_{l} \begin{bmatrix} \widetilde{\boldsymbol{W}}_{l}(t) & \boldsymbol{0} \\ \boldsymbol{0} & \rho(t)\boldsymbol{I}_{m} \end{bmatrix} \boldsymbol{V}_{l}^{\top}, \quad \forall l \in [L-1], \ t \geq 0,$$

where $\widetilde{W}_l(t) \in \mathbb{R}^{2r \times 2r}$ for all $l \in [L-1]$ with $\widetilde{W}_l(0) = \varepsilon I_{2r}$.

• Case 2: Suppose $YX^{\top} \in \mathbb{R}^{d_y \times d_x}$ with $d_y = r$ and $m := d_x - 2d_y > 0$. Similar results hold with different $\rho(t)$.

イロト 不得 トイヨト イヨト

э

• Dynamics of singular values and vectors of weight matrices. Let $\widetilde{W}_l(t) = \widetilde{U}_l(t)\widetilde{\Sigma}_l(t)\widetilde{V}_l^{\top}(t)$, we can rewrite our decomposition as

$$\boldsymbol{W}_{l}(t) = \begin{bmatrix} \boldsymbol{U}_{l,1} \widetilde{\boldsymbol{U}}_{l}(t) & \boldsymbol{U}_{l,2} \end{bmatrix} \begin{bmatrix} \widetilde{\boldsymbol{\Sigma}}_{l}(t) & \boldsymbol{0} \\ \boldsymbol{0} & \rho(t) \boldsymbol{I}_{m} \end{bmatrix} \begin{bmatrix} \boldsymbol{V}_{l,1} \widetilde{\boldsymbol{V}}_{l}(t) & \boldsymbol{V}_{l,2} \end{bmatrix}^{\top},$$

²M. Huh et al. The Low-Rank Simplicity Bias in Deep Networks, TMLR'23. https://minyoungg.github.io/overparam/

Qing Qu (EECS, University of Michigan)

Law of Parsimony in GD

June 30, 2023

11/34

• Dynamics of singular values and vectors of weight matrices. Let $\widetilde{W}_l(t) = \widetilde{U}_l(t)\widetilde{\Sigma}_l(t)\widetilde{V}_l^{\top}(t)$, we can rewrite our decomposition as

$$\boldsymbol{W}_{l}(t) = \begin{bmatrix} \boldsymbol{U}_{l,1} \widetilde{\boldsymbol{U}}_{l}(t) & \boldsymbol{U}_{l,2} \end{bmatrix} \begin{bmatrix} \widetilde{\boldsymbol{\Sigma}}_{l}(t) & \boldsymbol{0} \\ \boldsymbol{0} & \rho(t) \boldsymbol{I}_{m} \end{bmatrix} \begin{bmatrix} \boldsymbol{V}_{l,1} \widetilde{\boldsymbol{V}}_{l}(t) & \boldsymbol{V}_{l,2} \end{bmatrix}^{\top},$$

• Invariance of subspaces in the weights. Both $U_{l,2}$ and $V_{l,2}$ of size d-2r are unchanged throughout training. The learning process occurs only within an invariant subspace of dimension 2r!

²M. Huh et al. The Low-Rank Simplicity Bias in Deep Networks, TMLR'23. https://minyoungg.github.io/overparam/

Qing Qu (EECS, University of Michigan)

Law of Parsimony in GD

• Dynamics of singular values and vectors of weight matrices. Let $\widetilde{W}_l(t) = \widetilde{U}_l(t)\widetilde{\Sigma}_l(t)\widetilde{V}_l^{\top}(t)$, we can rewrite our decomposition as

$$\boldsymbol{W}_{l}(t) = \begin{bmatrix} \boldsymbol{U}_{l,1} \widetilde{\boldsymbol{U}}_{l}(t) & \boldsymbol{U}_{l,2} \end{bmatrix} \begin{bmatrix} \widetilde{\boldsymbol{\Sigma}}_{l}(t) & \boldsymbol{0} \\ \boldsymbol{0} & \rho(t) \boldsymbol{I}_{m} \end{bmatrix} \begin{bmatrix} \boldsymbol{V}_{l,1} \widetilde{\boldsymbol{V}}_{l}(t) & \boldsymbol{V}_{l,2} \end{bmatrix}^{\top},$$

- Invariance of subspaces in the weights. Both U_{l,2} and V_{l,2} of size d - 2r are unchanged throughout training. The learning process occurs only within an invariant subspace of dimension 2r!
- Implicit low-rank bias.² As lim_{ε→0} ρ(t) = 0 for all t ≥ 0, all the weights W_l(t) and the end-to-end matrix W_{L:1}(t) are inherently low-rank (e.g., at most rank 2r).

²M. Huh et al. The Low-Rank Simplicity Bias in Deep Networks, TMLR'23. https://minyoungg.github.io/overparam/

Qing Qu (EECS, University of Michigan)

Law of Parsimony in GD

The Evolution of Singular Spaces in More Generic Settings



Figure: Evolution of SVD of weight matrices without whitened data.

The Evolution of Singular Spaces in More Generic Settings



Figure: Evolution of SVD of weight matrices without whitened data.



Figure: Evolution of SVD of weight matrices_with_momentum. =

Qing Qu (EECS, University of Michigan)

Law of Parsimonv in GE

June 30, 2023

12/34

Outline

1 Law of Parsimony in Gradient Dynamics

2 Efficient Deep Matrix Completion via Network Compression

3 Understanding Hierarchical Representations in Deep Neural Networks

4 Conclusion

Main Message



Figure: Efficient training of deep linear networks.

The law of parsimony in GD leads to efficient network compression.

< 1 k

Problem Setup for Deep Matrix Completion

Consider recovering $\Phi \in \mathbb{R}^{d \times d}$ with $r := \operatorname{rank}(\Phi) \ll d$ with minimum number of observation encoded by $\Omega \in \{0, 1\}^{d \times d}$:

$$\min_{\boldsymbol{\Theta}} \ell_{\mathrm{mc}}(\boldsymbol{\Theta}) := \frac{1}{2} \| \boldsymbol{\Omega} \odot (\boldsymbol{W}_{L:1} - \boldsymbol{\Phi}) \|_{F}^{2}.$$

Problem Setup for Deep Matrix Completion

Consider recovering $\Phi \in \mathbb{R}^{d \times d}$ with $r := \operatorname{rank}(\Phi) \ll d$ with minimum number of observation encoded by $\Omega \in \{0, 1\}^{d \times d}$:

$$\min_{\boldsymbol{\Theta}} \ell_{\mathrm{mc}}(\boldsymbol{\Theta}) := \frac{1}{2} \| \boldsymbol{\Omega} \odot (\boldsymbol{W}_{L:1} - \boldsymbol{\Phi}) \|_{F}^{2}.$$

• If full observation $\Omega = \mathbf{1}_d \mathbf{1}_d^\top$ available, the problem simplifies to deep matrix factorization.

Problem Setup for Deep Matrix Completion

Consider recovering $\Phi \in \mathbb{R}^{d \times d}$ with $r := \operatorname{rank}(\Phi) \ll d$ with minimum number of observation encoded by $\Omega \in \{0, 1\}^{d \times d}$:

$$\min_{\boldsymbol{\Theta}} \ell_{\mathrm{mc}}(\boldsymbol{\Theta}) := \frac{1}{2} \| \boldsymbol{\Omega} \odot (\boldsymbol{W}_{L:1} - \boldsymbol{\Phi}) \|_F^2.$$

- If full observation $\Omega = \mathbf{1}_d \mathbf{1}_d^\top$ available, the problem simplifies to deep matrix factorization.
- If the network depth L = 2, it reduces to the Burer-Monteiro factorization formulation.

く 目 ト く ヨ ト く ヨ ト

Why Deep Matrix Factorization and Overparameterization?



- Benefits of Depth (Left): Improved sample complexity³ and less prone to overfitting.
- Benefits of Width (Right): Increasing the width of the network results in accelerated convergence in terms of iterations.

³Arora, S., Cohen, N., Hu, W., & Luo, Y. (2019). Implicit regularization in deep matrix factorization. Advances in Neural Information Processing Systems, 32

Qing Qu (EECS, University of Michigan)

Law of Parsimony in GE

Overparameterization: A Double Edged Sword



Figure: Efficient training of deep linear networks.

Cons: Increasing the depth and width of the network leads to much **more parameters**. Could be **expensive to optimize!**

• **Deep matrix factorization.** As a starting point, consider the simple deep matrix factorization setting:

$$\min_{\boldsymbol{\Theta}} \ \frac{1}{2} \| \boldsymbol{W}_{L:1} - \boldsymbol{\Phi} \|_F^2,$$

with $\Omega = \mathbf{1}_d \mathbf{1}_d^{\top}$. We optimize the problem via GD from ε -scale orthogonal initialization.

< 回 > < 回 > < 回 >

• **Deep matrix factorization.** As a starting point, consider the simple deep matrix factorization setting:

$$\min_{\boldsymbol{\Theta}} \ \frac{1}{2} \| \boldsymbol{W}_{L:1} - \boldsymbol{\Phi} \|_F^2,$$

with $\Omega = \mathbf{1}_d \mathbf{1}_d^{\top}$. We optimize the problem via GD from ε -scale orthogonal initialization.

• Law of parsimony in GD for the end-to-end matrix $W_{L:1}$:

$$\begin{split} \boldsymbol{W}_{L:1}(t) &= \begin{bmatrix} \boldsymbol{U}_{L,1} & \boldsymbol{U}_{L,2} \end{bmatrix} \begin{bmatrix} \widetilde{\boldsymbol{W}}_{L:1}(t) & \boldsymbol{0} \\ \boldsymbol{0} & \rho^{L}(t)\boldsymbol{I}_{m} \end{bmatrix} \begin{bmatrix} \boldsymbol{V}_{1,1}^{\top} \\ \boldsymbol{V}_{1,2}^{\top} \end{bmatrix} \\ &= \boldsymbol{U}_{L,1}\widetilde{\boldsymbol{W}}_{L:1}(t)\boldsymbol{V}_{1,1}^{\top} + \rho^{L}(t)\boldsymbol{U}_{L,2}\boldsymbol{V}_{1,2}^{\top}, \end{split}$$

where we overestimate the rank $\hat{r} > r$ and let $m = d - 2\hat{r}$.

• The effects of small initialization ε and depth L:

$$\begin{aligned} \boldsymbol{W}_{L:1}(t) &= \boldsymbol{U}_{L,1} \widetilde{\boldsymbol{W}}_{L:1}(t) \boldsymbol{V}_{1,1}^{\top} + \rho^{L}(t) \boldsymbol{U}_{L,2} \boldsymbol{V}_{1,2}^{\top} \\ &\approx \boldsymbol{U}_{L,1} \widetilde{\boldsymbol{W}}_{L:1}(t) \boldsymbol{V}_{1,1}^{\top}, \quad \forall t \geq 0, \end{aligned}$$

< □ > < 同 > < 回 > < 回 > < 回 >

э

• The effects of small initialization ε and depth L:

$$\begin{aligned} \boldsymbol{W}_{L:1}(t) &= \boldsymbol{U}_{L,1} \widetilde{\boldsymbol{W}}_{L:1}(t) \boldsymbol{V}_{1,1}^{\top} + \boldsymbol{\rho}^{L}(t) \boldsymbol{U}_{L,2} \boldsymbol{V}_{1,2}^{\top} \\ &\approx \boldsymbol{U}_{L,1} \widetilde{\boldsymbol{W}}_{L:1}(t) \boldsymbol{V}_{1,1}^{\top}, \quad \forall t \geq 0, \end{aligned}$$

Claim: With small initialization, running GD on the original weights $\{W_l\}_{l=1}^L \subseteq \mathbb{R}^{d \times d}$ is **almost equivalent** to running GD on the compressed weights $\{\widetilde{W}_l\}_{l=1}^L \subseteq \mathbb{R}^{2\widehat{r} \times 2\widehat{r}}$.

The Simple Case: Deep Matrix Factorization



Figure: Efficient training of deep linear networks.

Comparison on the number of parameters: original network Ld^2 vs. compressed network $L\hat{r}^2$.



• However, directly applying our approach from deep matrix factorization to completion does not work well...



- However, directly applying our approach from deep matrix factorization to completion does not work well...
- This is due to the fact that the law of parsimony in GD:

$$\boldsymbol{W}_{L:1}(t) \approx \boldsymbol{U}_{L,1} \widetilde{\boldsymbol{W}}_{L:1}(t) \boldsymbol{V}_{1,1}^{\top}, \quad \forall t \ge 0,$$

does NOT hold, because $\Omega\odot\Phi$ is not low-rank for arbitrary $\Omega.$



Remedy: update both V_{1,1}(t) and U_{L,1}(t) factors via GD with a discrepant learning rate γη in the "compressed network":⁴

$$\boldsymbol{W}_{\text{comp}}^{(\gamma)}(t) := \boldsymbol{U}_{L,1}(t) \widetilde{\boldsymbol{W}}_{L:1}(t) \boldsymbol{V}_{1,1}^{\top}(t).$$

⁴This is done simultaneously with the GD updates on the subnetwork $\widetilde{W}_{L:1}(t)$, which uses the original learning rate η .

Qing Qu (EECS, University of Michigan)



Remedy: update both V_{1,1}(t) and U_{L,1}(t) factors via GD with a discrepant learning rate γη in the "compressed network":⁴

$$\boldsymbol{W}_{\text{comp}}^{(\gamma)}(t) := \boldsymbol{U}_{L,1}(t) \widetilde{\boldsymbol{W}}_{L:1}(t) \boldsymbol{V}_{1,1}^{\top}(t).$$

• **Complexity:** original network $O(Ld^2)$ vs compressed network O(Ld).

⁴This is done simultaneously with the GD updates on the subnetwork $\widetilde{W}_{L:1}(t)$, which uses the original learning rate η .

Qing Qu (EECS, University of Michigan)

Outline

1 Law of Parsimony in Gradient Dynamics

2 Efficient Deep Matrix Completion via Network Compression

3 Understanding Hierarchical Representations in Deep Neural Networks

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへで

4 Conclusion

Networks

Main Message



For classification problem, the law of parsimony in GD explains progressive feature separation in deep linear networks.

Problem Setup: Train DLNs for Classification Problems

- Balanced Training Data: $\{(x_{k,i}, y_k)\}_{i \in [n], k \in [K]}$ for K-class classification: $x_{k,i} \in \mathbb{R}^d$ is the *i*-th sample in the *k*-th class, $y_k \in \mathbb{R}^K$ is an one-hot label.
- Feature in the *l*-th Layer of DLN:

$$oldsymbol{z}_{k,i}^l := oldsymbol{W}_l \dots oldsymbol{W}_1 oldsymbol{x}_{k,i} = oldsymbol{W}_{l:1} oldsymbol{x}_{k,i}, \; orall l \in [L],$$

Problem Setup: Train DLNs for Classification Problems

- Balanced Training Data: $\{(x_{k,i}, y_k)\}_{i \in [n], k \in [K]}$ for K-class classification: $x_{k,i} \in \mathbb{R}^d$ is the *i*-th sample in the *k*-th class, $y_k \in \mathbb{R}^K$ is an one-hot label.
- Feature in the *l*-th Layer of DLN:

$$oldsymbol{z}_{k,i}^l := oldsymbol{W}_l \dots oldsymbol{W}_1 oldsymbol{x}_{k,i} = oldsymbol{W}_{l:1} oldsymbol{x}_{k,i}, \; orall l \in [L],$$

• Measure of Data Separation: To characterize the network's capability to separate data across layers, we use a metric (He & Su. 2022, Tirer et al. (2022))

$$D_l \ := \ \mathsf{trace}(\boldsymbol{\Sigma}_W^l)/\mathsf{trace}(\boldsymbol{\Sigma}_B^l),$$

$$\boldsymbol{\Sigma}_{W}^{l} = \frac{1}{nK} \sum_{k=1}^{K} \sum_{i=1}^{n} \left(\boldsymbol{z}_{k,i}^{l} - \bar{\boldsymbol{z}}_{k}^{l} \right) \left(\boldsymbol{z}_{k,i} - \bar{\boldsymbol{z}}_{k}^{l} \right)^{\top}, \ \boldsymbol{\Sigma}_{B}^{l} = \frac{1}{K} \sum_{k=1}^{K} \left(\bar{\boldsymbol{z}}_{k}^{l} - \boldsymbol{z}_{G}^{l} \right)^{(k-1)} \left(\bar{\boldsymbol{z}}_{k,i}^{l} - \bar{\boldsymbol{z}}_{K}^{l} \right)^{(k-1)}$$

Networks

Progressive Feature Separation with Linear Decay Rate



Figure: Linear decay of feature separation in trained deep networks.

< □ > < □ > < □ > < □ > < □ > < □ >

Progressive Feature Separation with Linear Decay Rate

Theorem (Wang et al.'23)

Suppose we train a *L*-layer DLN with parameters $\Theta = \{W_l\}_{l=1}^L$ via GD with orthogonal initialization of ε -scaling, where input $X \in \mathbb{R}^{d \times N}$ is orthogonal and square and $d_l = d > 2K$. If Θ satisfies the following:

- Global Optimality: $W_{L:1}X = Y$.
- Balancedness: For all weights

$$\boldsymbol{W}_{l+1}^{\top} \boldsymbol{W}_{l+1} = \boldsymbol{W}_{l} \boldsymbol{W}_{l}^{\top}, \forall l \in [L-2], \\ \|\boldsymbol{W}_{L}^{\top} \boldsymbol{W}_{L} - \boldsymbol{W}_{L-1} \boldsymbol{W}_{L-1}^{\top}\|_{F} \leq \varepsilon^{2} \sqrt{d-K}.$$

Networks

Progressive Feature Separation with Linear Decay Rate

Theorem (Wang et al.'23)

Suppose we train a *L*-layer DLN with parameters $\Theta = \{W_l\}_{l=1}^L$ via GD with orthogonal initialization of ε -scaling, where input $X \in \mathbb{R}^{d \times N}$ is orthogonal and square and $d_l = d > 2K$. If Θ satisfies the following:

- Global Optimality: $W_{L:1}X = Y$.
- Balancedness: For all weights

$$\boldsymbol{W}_{l+1}^{\top} \boldsymbol{W}_{l+1} = \boldsymbol{W}_{l} \boldsymbol{W}_{l}^{\top}, \forall l \in [L-2], \\ \|\boldsymbol{W}_{L}^{\top} \boldsymbol{W}_{L} - \boldsymbol{W}_{L-1} \boldsymbol{W}_{L-1}^{\top}\|_{F} \leq \varepsilon^{2} \sqrt{d-K}.$$

• Unchanged Spectrum: There exists an index set $\mathcal{A} \subseteq [d]$ with $|\mathcal{A}| = d - 2K$ such that for all $l \in [L-1]$ that $\sigma_i(\mathbf{W}_l) = \varepsilon, \ \forall i \in \mathcal{A}.$

Progressive Feature Separation with Linear Decay Rate

Theorem (Wang et al.'23)

Suppose we train a *L*-layer DLN with parameters $\Theta = \{W_l\}_{l=1}^L$ via GD with orthogonal initialization of ε -scaling, where input $X \in \mathbb{R}^{d \times N}$ is orthogonal and square and $d_l = d > 2K$. If Θ satisfies the following:

- Global Optimality: $W_{L:1}X = Y$.
- Balancedness: For all weights

$$\boldsymbol{W}_{l+1}^{\top} \boldsymbol{W}_{l+1} = \boldsymbol{W}_{l} \boldsymbol{W}_{l}^{\top}, \forall l \in [L-2], \\ \|\boldsymbol{W}_{L}^{\top} \boldsymbol{W}_{L} - \boldsymbol{W}_{L-1} \boldsymbol{W}_{L-1}^{\top}\|_{F} \leq \varepsilon^{2} \sqrt{d-K}.$$

• Unchanged Spectrum: There exists an index set $\mathcal{A} \subseteq [d]$ with $|\mathcal{A}| = d - 2K$ such that for all $l \in [L - 1]$ that $\sigma_i(\mathbf{W}_l) = \varepsilon$, $\forall i \in \mathcal{A}$. Then, it holds for all l = 0, 1, ..., L - 2 that

$$D_{l+1}/D_l \le 2(\sqrt{K}+1)\varepsilon^2.$$

Networks





Layer 2





Layer 3



Qing Qu (EECS, University of Michigan)

Law of Parsimony in GE

June 30, 2023

26 / 34

3

Effects of Initialization Scale ε

As predicted by our theory, the decay ratio critically depends on the scale of initialization ε :



Figure: Linear decay of feature separation measure D_l in trained deep networks with varying initialization scale ε .

< A > <

Networks

Is the Orthogonal Initialization Critical?



Figure: Linear decay of feature separation in trained DLNs with different initialization types (left to right: Orth., Norm, Unif).

- 4 回 ト 4 ヨ ト 4 ヨ ト

Outline

1 Law of Parsimony in Gradient Dynamics

2 Efficient Deep Matrix Completion via Network Compression

3 Understanding Hierarchical Representations in Deep Neural Networks





References

- 1 Yaras, C.*, Wang, P.*, Hu, W., Zhu, Z., Balzano, L., Qu, Q. (2023). The Law of Parsimony in Gradient Descent for Learning Deep Linear Networks. arXiv preprint arXiv:2306.01154.
- 2 Wang, P., Yaras, C., Li, X., Hu, W., Zhu, Z., Balzano, L., Qu, Q. (2023). Understanding Hierarchical Representation Learning in Deep Networks via Neural Collapse. Working paper.
- 3 Li, X., Liu S., Zhou, J., Lu, X., Fernandez-Granda, C., Zhu, Z., Qu, Q. (2023) Principled and Efficient Transfer Learning of Deep Models via Neural Collapse, arXiv preprint arXiv:2212.12206.
- 4 Zhu, Z., Ding, T., Zhou, J., Li, X., You, C., Sulam, J., Qu, Q. (2021). A geometric analysis of neural collapse with unconstrained features. Advances in Neural Information Processing Systems, 34, 29820-29834.

3

Conclusion

The GD learning process takes place only within a **minimal invariant subspace** of each weight matrix, while the remaining singular subspaces stay **unaffected** throughout training.

- Efficient training via network compression.
- Understanding representations in deep networks.

Thank You! Questions?

Call for Papers

- IEEE JSTSP Special Issue on Seeking Low-dimensionality in Deep Neural Networks (SLowDNN) Manuscript Due: Nov. 30, 2023.
- Conference on Parsimony and Learning (CPAL) January 2024, Hongkong, Manuscript Due: **Aug. 28, 2023**.





Tradeoffs Between Decay Rate and Convergence

However, there is trade-off between decay rate ε and training speed of GD:



Figure: The dynamics of GD for DLNs with learning rate $\eta = 0.1$.

< ∃⇒

Compressed Networks vs. Narrow Networks?

Question: Does law of parsimony imply that optimizing a narrow network of the same width $2\hat{r}$ would perform just as efficiently as the compressed network with a true width of $d \gg \hat{r}$?



Figure: Efficiency of compressed networks vs. narrow network.

∃ →

< □ > < A >

Compressed Networks vs. Narrow Networks?



Figure: Efficiency of compressed networks vs. narrow network.

Answer: No! Over-parameterized networks are "easier" to train.