

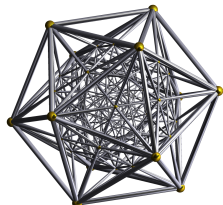
ICASSP 2023 Short Course

Learning Nonlinear and Deep Representations from High-Dimensional Data: From Theory to Practice

Lecture 6: Low-dimensional Representations in Deep Networks II

Sam Buchanan, Yi Ma, Qing Qu
Atlas Wang, John Wright, Yuqian Zhang, Zhihui Zhu

June 8, 2023



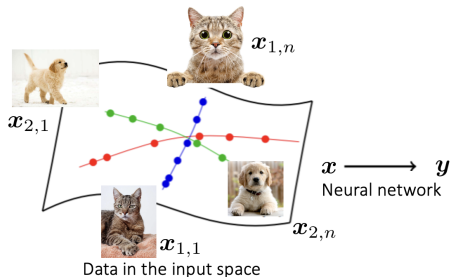
Outline

- 1 Neural Collapse, Transfer Learning, & Intermediate Layers
- 2 Law of Parsimony in Gradient Dynamics
- 3 Progressive Feature Separation in Deep Neural Networks
- 4 Efficient Deep Matrix Completion
- 5 Conclusion

Neural Collapse in Classification

Labels: $k = 1, \dots, K$

- $K = 10$ classes (MNIST, CIFAR10, etc)
- $K = 1000$ classes (ImageNet)



Cat	Dog	...	Truck
[[[
1	0		0
0	1		0
⋮	⋮		⋮
0	0		1
]]]

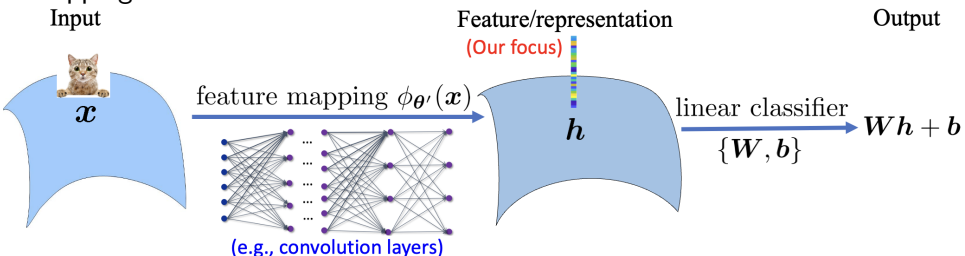
One-hot labeling vectors in \mathbb{R}^K

Assume balanced dataset where each class has n training samples

- If not, we can use data augmentation to make them balanced

Deep Neural Network Classifiers

A deep neural network classifier often contains two parts: a feature mapping and a linear classifier



- Output: $f(x; \theta) = \mathbf{W} \phi_{\theta'}(x) + \mathbf{b}$ with $\theta = (\theta', \mathbf{W}, \mathbf{b})$.
- Training problem:

$$\min_{\theta', \mathbf{W}, \mathbf{b}} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \underbrace{\mathcal{L}_{\text{CE}}(\mathbf{W} \phi_{\theta'}(\mathbf{x}_{k,i}) + \mathbf{b}, \mathbf{y}_k)}_{\text{cross-entropy (CE) loss}} + \lambda \underbrace{\|(\theta', \mathbf{W}, \mathbf{b})\|_F^2}_{\text{weight decay}}$$

Neural Collapse in Classification

Prevalence of neural collapse during the terminal phase of deep learning training

 Vardan Papyan,  X. Y. Han, and David L. Donoho

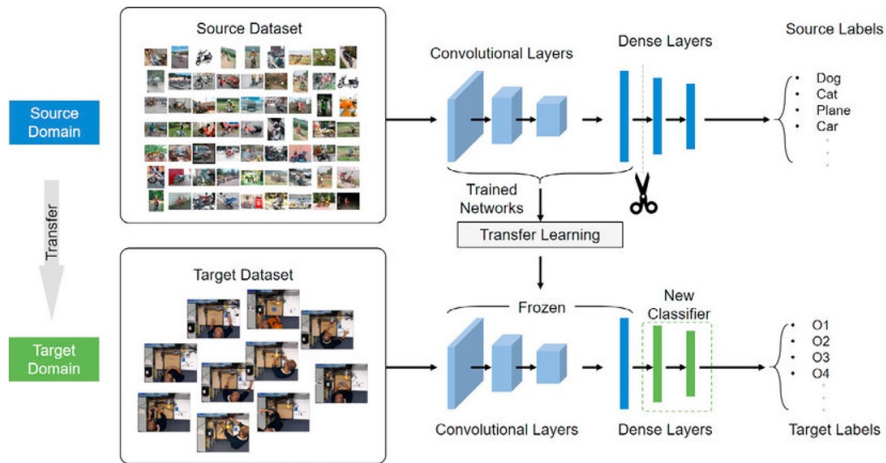
[+ See all authors and affiliations](#)

PNAS October 6, 2020 117 (40) 24652-24663; first published September 21, 2020;
<https://doi.org/10.1073/pnas.2015509117>

Contributed by David L. Donoho, August 18, 2020 (sent for review July 22, 2020; reviewed by Helmut Boelsckei and Stéphane Mallat)

- Reveals common outcome of learned features and classifiers across a variety of architectures and dataset
- Precise mathematical structure within the features and classifier

Implications of Neural Collapse in Transfer Learning?



Efficient Fine-tuning of Pre-trained Models

- **Full model fine-tuning** (use pre-trained model as an initialization)
 - Expensive & prone to overfitting
- **Linear probing** with penultimate layer features (use pre-trained model as a feature extractor)
 - cheap but worse performance

Efficient Fine-tuning of Pre-trained Models

- **Full model fine-tuning** (use pre-trained model as an initialization)
 - Expensive & prone to overfitting
- **Linear probing** with penultimate layer features (use pre-trained model as a feature extractor)
 - cheap but worse performance



Message: We can design simple and efficient fine-tuning approaches of pre-trained models via Neural Collapse.

Measure of Variability Collapse on Downstream Data

One key metric for Neural Collapse:

$$\mathcal{NC}_1 = \text{trace}(\Sigma_W \Sigma_B^\dagger).$$

- **Within-class covariance:**

$$\Sigma_W = \frac{1}{nK} \sum_{k=1}^K \sum_{i=1}^n (\mathbf{h}_{k,i} - \bar{\mathbf{h}}_k) (\mathbf{h}_{k,i} - \bar{\mathbf{h}}_k)^\top$$

- **Between-class covariance:**

$$\Sigma_B = \frac{1}{K} \sum_{k=1}^K (\bar{\mathbf{h}}_k - \mathbf{h}_G) (\bar{\mathbf{h}}_k - \mathbf{h}_G)^\top$$

More Variability Collapse, Better Transfer Performance

Neural Collapse of pre-trained models evaluated on **downstream data**:

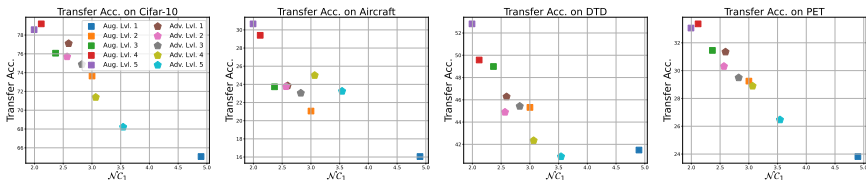


Figure: Transfer accuracy and \mathcal{NCC}_1 of Cifar-100 pre-trained models on different downstream tasks.

More Variability Collapse, Better Transfer Performance

Neural Collapse of pre-trained models evaluated on **downstream data**:

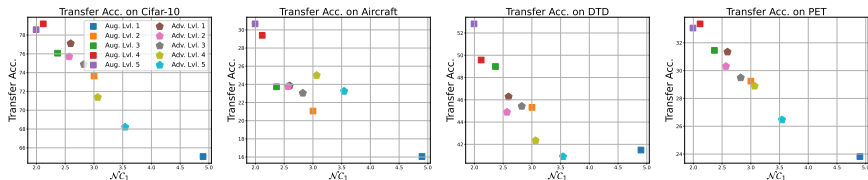


Figure: Transfer accuracy and \mathcal{NCC}_1 of Cifar-100 pre-trained models on different downstream tasks.

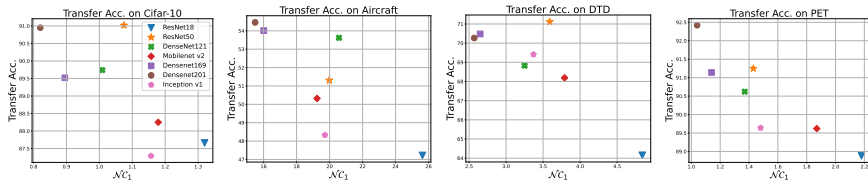


Figure: Transfer accuracy and \mathcal{NCC}_1 of public ImageNet-1k pre-trained models on different downstream tasks.

Parameter-Efficient Fine-Tuning via Neural Collapse

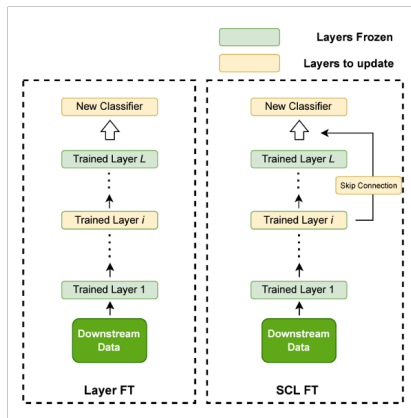
Observation: Better transfer performance can be achieved by making the last-layer features more collapsed on the downstream data.

Parameter-Efficient Fine-Tuning via Neural Collapse

Observation: Better transfer performance can be achieved by making the last-layer features more collapsed on the downstream data.

Skip connection layer fine-tuning (SCL-FT):

- Only fine-tuning one key intermediate layer.
- Improving feature collapse via skip connections.

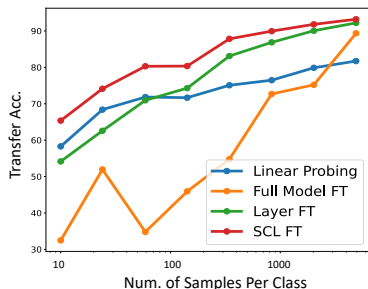


Better Performance with Significantly Fewer Parameters

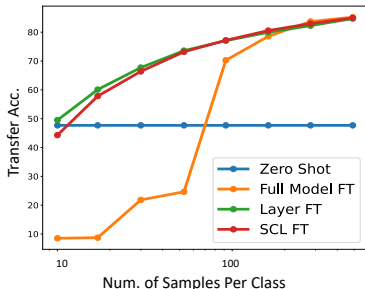
Backbone	ResNet50					ViT-B			CLIP		
Dataset	Cifar-10	Cifar-100	Aircraft	DTD	PET	Aircraft	DTD	PET	Aircraft	DTD	PET
Transfer accuracy											
Linear Probe / Zero Shot	85.33	65.47	43.23	68.46	89.26	43.65	73.88	92.23	12.87	32.34	39.66
Layer FT	94.04	77.47	70.27	67.66	89.40	65.83	77.13	92.94	67.63	79.47	91.09
SCL FT	94.94	78.32	70.72	72.87	91.69	65.80	77.34	93.13	66.58	79.04	90.02
Full Model FT	85.51	78.88	80.77	76.12	73.24	64.66	76.54	93.02	59.11	72.82	84.44
<i>$\mathcal{N}C_1$ evaluated on the penultimate layer feature h^{L-1}</i>											
Linear Probe / Zero Shot	1.84	18.36	20.36	3.52	1.45	17.91	1.99	0.66	17.47	2.96	3.77
Layer FT	0.28	3.22	3.37	1.68	0.68	6.98	1.62	0.44	1.30	0.42	0.32
SCL FT	0.22	2.61	1.02	0.64	0.39	7.48	1.33	0.35	1.65	0.61	0.46
Full Model FT	0.17	0.15	0.61	0.31	0.28	3.78	1.11	0.21	0.49	0.17	0.18
Percentage of parameters fine-tuned											
Linear Probe / Zero Shot	0.09%	0.86%	0.86%	0.41%	0.32%	0.09%	0.04%	0.03%	0.0%	0.0%	0.0%
Layer FT	6.52%	7.24%	7.24%	6.82%	6.73%	8.18%	8.14%	8.13%	8.16%	8.11%	8.10%
SCL FT	6.52%	7.24%	7.24%	6.82%	6.73%	8.18%	8.14%	8.13%	8.16%	8.11%	8.10%
Full Model FT	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%

Table: Transfer learning results for linear probing / zero-shot, layer FT, SCL FT and full model FT on downstream datasets. We use publicly available ResNet50, ViT-B and CLIP models.

Less Overfitting Compared to Full Model FT



(a) ResNet18, Cifar-10



(b) CLIP, Cifar-100

Figure: Transfer accuracy for different fine-tuning methods with varying size of downstream training dataset.

We fine-tune ImageNet-1k pre-trained ResNet18 models and CLIP using subsets of the Cifar-10 and Cifar-100 downstream datasets, respectively with varying sizes.

Parameter-Efficient Fine-Tuning via Neural Collapse

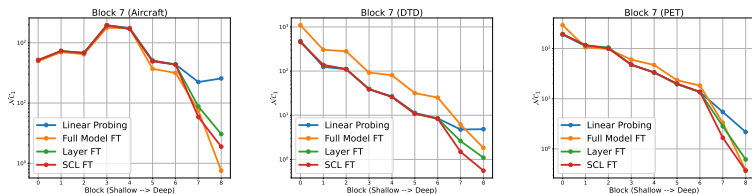


Figure: Layer-wise $\mathcal{N}C_1$ for different fine-tuning methods on ResNet18.

Parameter-Efficient Fine-Tuning via Neural Collapse

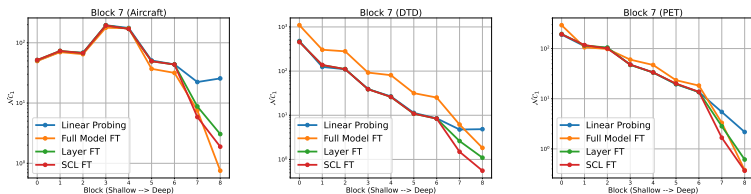


Figure: Layer-wise \mathcal{NC}_1 for different fine-tuning methods on ResNet18.

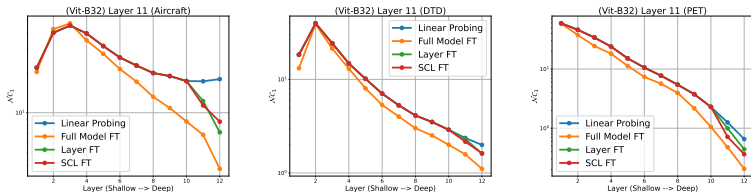


Figure: Layer-wise \mathcal{NC}_1 for different fine-tuning methods on Vision Transformer ViT-B32.

Representation Structures within the Intermediate Layers?

Is there any structure in the representations of intermediate layers?¹

$$\mathcal{NC}_1^l = \text{trace}(\Sigma_W^l \Sigma_B^{l\dagger}).$$

- Between-class covariance:

$$\Sigma_B^l = \frac{1}{K} \sum_{k=1}^K \left(\bar{z}_k^l - z_G^l \right) \left(\bar{z}_k^l - z_G^l \right)^\top$$

- Within-class covariance:

$$\Sigma_W^l = \frac{1}{nK} \sum_{k=1}^K \sum_{i=1}^n \left(z_{k,i}^l - \bar{z}_k^l \right) \left(z_{k,i}^l - \bar{z}_k^l \right)^\top$$

¹V. Pappas, Traces of class/cross-class structure pervade deep learning spectra, JMLR, 2021. He & Su, A Law of Progressive Separation for Deep Learning, 2022.

Representation Structures within the Intermediate Layers?

Is there any structure in the representations of intermediate layers?¹

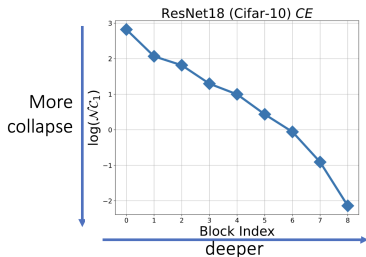
$$\mathcal{NC}_1^l = \text{trace}(\Sigma_W^l \Sigma_B^{l\dagger}).$$

- Between-class covariance:

$$\Sigma_B^l = \frac{1}{K} \sum_{k=1}^K \left(\bar{z}_k^l - z_G^l \right) \left(\bar{z}_k^l - z_G^l \right)^\top$$

- Within-class covariance:

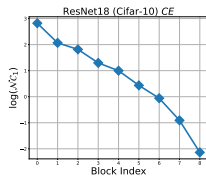
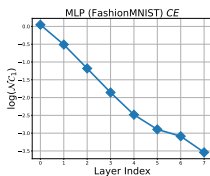
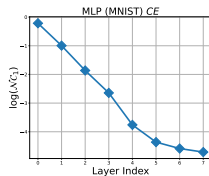
$$\Sigma_W^l = \frac{1}{nK} \sum_{k=1}^K \sum_{i=1}^n \left(z_{k,i}^l - \bar{z}_k^l \right) \left(z_{k,i}^l - \bar{z}_k^l \right)^\top$$



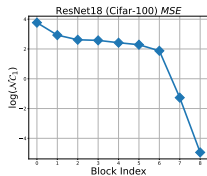
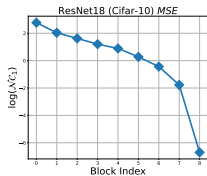
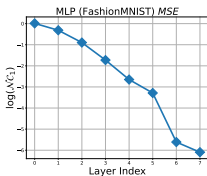
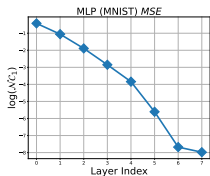
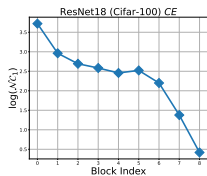
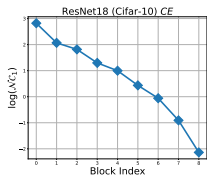
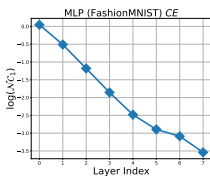
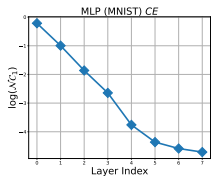
Effects of Depth: Creating **progressive** separation and concentration of data from shallow to deep layers!

¹V. Pappas, Traces of class/cross-class structure pervade deep learning spectra, JMLR, 2021. He & Su, A Law of Progressive Separation for Deep Learning, 2022.

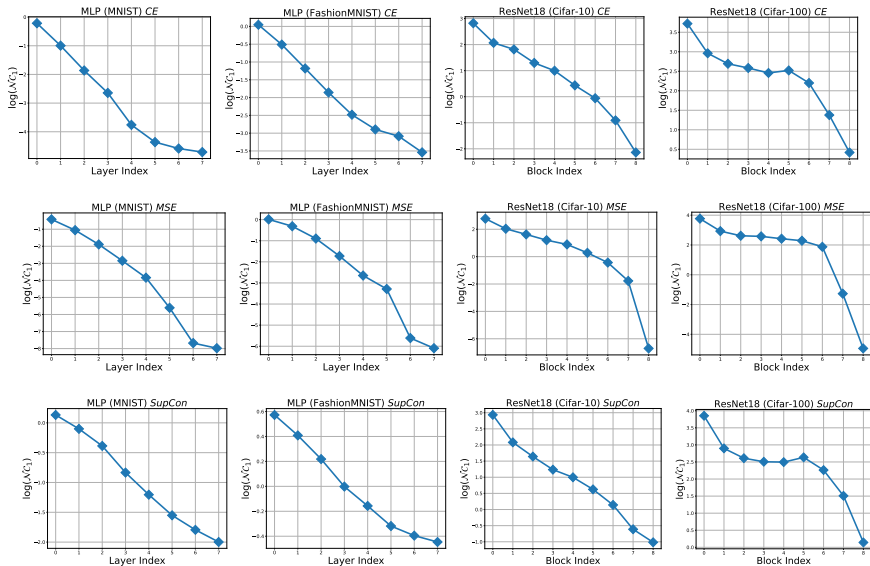
Prevalence of Progressive Separation in Deep Learning



Prevalence of Progressive Separation in Deep Learning



Prevalence of Progressive Separation in Deep Learning



Towards Understanding Progressive Collapse?

Given linearly separable data \mathbf{X} , we measure a certain metric of data separation of each layer's output for linear and nonlinear deep networks:

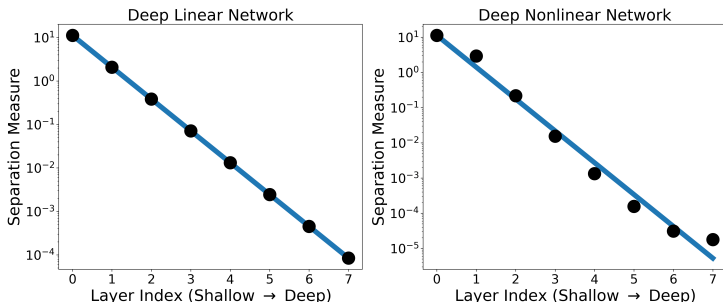
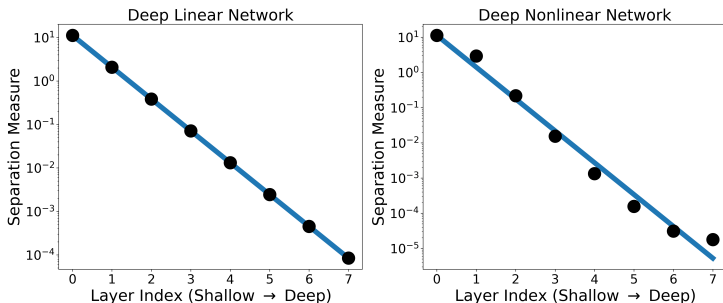


Figure: Progressive collapse with linear decay² on deep linear and nonlinear networks. The x -axis denotes the layer index and the y -axis denotes the separation measure.

²He & Su, A Law of Progressive Separation for Deep Learning, 2022.

Towards Understanding Progressive Collapse?

Given linearly separable data \mathbf{X} , we measure a certain metric of data separation of each layer's output on linear and nonlinear deep networks:



Even training overparameterized deep linear networks could have **infinite many solutions**, why do such benign structures happen?

Outline

- 1 Neural Collapse, Transfer Learning, & Intermediate Layers
- 2 Law of Parsimony in Gradient Dynamics**
- 3 Progressive Feature Separation in Deep Neural Networks
- 4 Efficient Deep Matrix Completion
- 5 Conclusion

Main Message

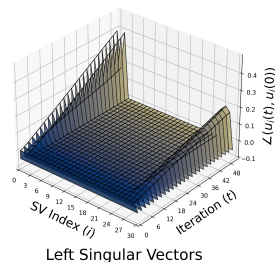
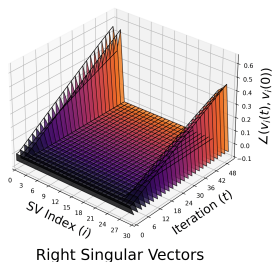
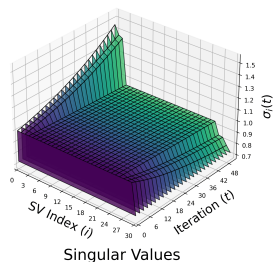
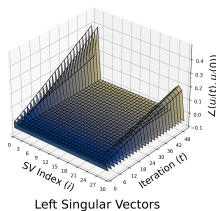
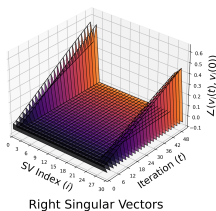
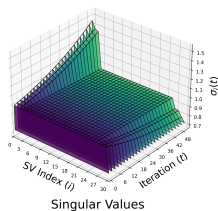


Figure: Evolution of SVD of the weight matrix $W_1(t) = U_1(t)\Sigma_1(t)V_1(t)^\top$.

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

Main Message



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

Setup on Deep Linear Networks

- **Training data** $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N \subset \mathbb{R}^{d_x} \times \mathbb{R}^{d_y}$ with

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

Setup on Deep Linear Networks

- **Training data** $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N \subset \mathbb{R}^{d_x} \times \mathbb{R}^{d_y}$ with

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

- **Deep linear network (DLN):**

$$f_{\Theta}(\mathbf{x}) := \mathbf{W}_L \cdots \mathbf{W}_1 \mathbf{x} = \mathbf{W}_{L:1} \mathbf{x},$$

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

Setup on Deep Linear Networks

- **Training data** $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N \subset \mathbb{R}^{d_x} \times \mathbb{R}^{d_y}$ with

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

- **Deep linear network (DLN):**

$$f_{\Theta}(\mathbf{x}) := \mathbf{W}_L \cdots \mathbf{W}_1 \mathbf{x} = \mathbf{W}_{L:1} \mathbf{x},$$

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

- **Loss function:**

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

Training DLNs via Gradient Descent (GD)

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

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

depending on the size of \mathbf{W}_l .

Training DLNs via Gradient Descent (GD)

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

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

depending on the size of \mathbf{W}_l .

- **Learning dynamics of GD.** We update all weights via GD for $t = 1, 2, \dots$ as

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

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

Training DLNs via Gradient Descent (GD)

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

- The weight matrices are square except the last layer, i.e., $d_x = d_1 = d_2 = \dots = d_{L-1} = d$ for some $d \in \mathbb{N}_+$.
- The input data is *whitened* in the sense that $\mathbf{X}\mathbf{X}^\top = \mathbf{I}_{d_x}$.³
- The cross correlation matrix $\mathbf{Y}\mathbf{X}^\top$ has certain low-dimensional structures (e.g., low-rank or wide matrix).

³For any full rank $\mathbf{X} \in \mathbb{R}^{d_x \times N}$ with $N \geq d_x$, whitened data can always be obtained with a data pre-processing step such as preconditioning.

Training DLNs via Gradient Descent (GD)

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

- The weight matrices are square except the last layer, i.e., $d_x = d_1 = d_2 = \dots = d_{L-1} = d$ for some $d \in \mathbb{N}_+$.
- The input data is *whitened* in the sense that $\mathbf{X}\mathbf{X}^\top = \mathbf{I}_{d_x}$.³
- The cross correlation matrix $\mathbf{Y}\mathbf{X}^\top$ has certain low-dimensional structures (e.g., low-rank or wide matrix).

Under the simplified settings, would GD possess any parsimonious structures during training?

³For any full rank $\mathbf{X} \in \mathbb{R}^{d_x \times N}$ with $N \geq d_x$, whitened data can always be obtained with a data pre-processing step such as preconditioning.

The Evolution of Singular Spaces in GD Iterates for DLNs

We train a $L = 3$ layer DLN with $d_x = d_y = 30$ and $r := \text{rank}(\mathbf{Y}) = 3$.

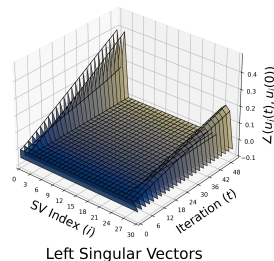
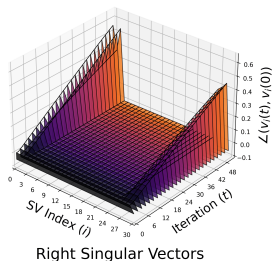
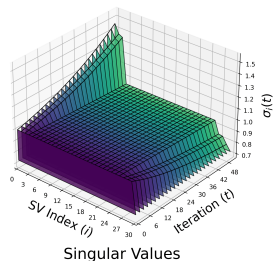


Figure: Evolution of SVD of the weight matrix $\mathbf{W}_1(t) = \mathbf{U}_1(t)\mathbf{\Sigma}_1(t)\mathbf{V}_1(t)^\top$.

The Evolution of Singular Spaces in GD Iterates for DLNs

We train a $L = 3$ layer DLN with $d_x = d_y = 30$ and $r := \text{rank}(\mathbf{Y}) = 3$.

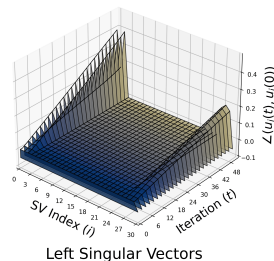
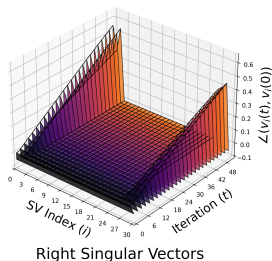
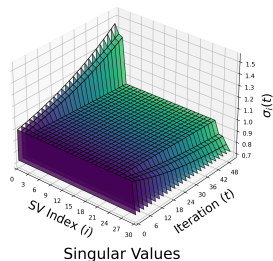


Figure: Evolution of SVD of the weight matrix $\mathbf{W}_1(t) = \mathbf{U}_1(t)\mathbf{\Sigma}_1(t)\mathbf{V}_1(t)^\top$.

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

The Evolution of Singular Spaces in GD Iterates for DLNs

We train a $L = 3$ layer DLN with $d_x = d_y = 30$ and $r := \text{rank}(\mathbf{Y}) = 3$.

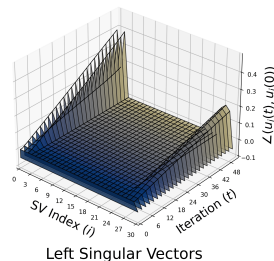
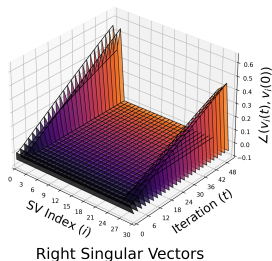
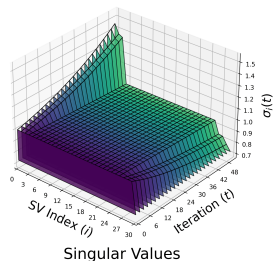


Figure: Evolution of SVD of the weight matrix $\mathbf{W}_1(t) = \mathbf{U}_1(t)\mathbf{\Sigma}_1(t)\mathbf{V}_1(t)^\top$.

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

The Evolution of Singular Spaces in GD Iterates for DLNs

We train a $L = 3$ layer DLN with $d_x = d_y = 30$ and $r := \text{rank}(\mathbf{Y}) = 3$.

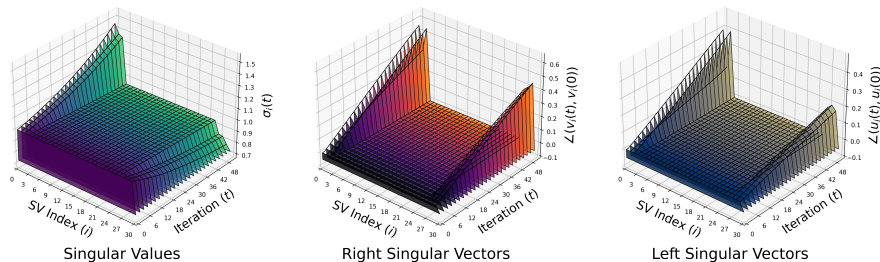
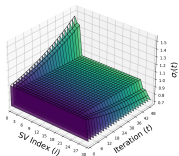


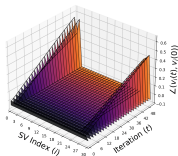
Figure: Evolution of SVD of the weight matrix $\mathbf{W}_1(t) = \mathbf{U}_1(t)\mathbf{\Sigma}_1(t)\mathbf{V}_1(t)^\top$.

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

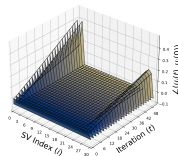
Layer 1



Singular Values

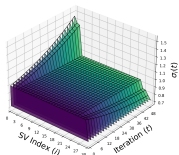


Right Singular Vectors

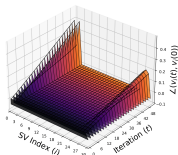


Left Singular Vectors

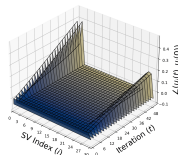
Layer 2



Singular Values

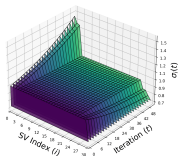


Right Singular Vectors

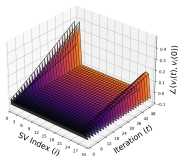


Left Singular Vectors

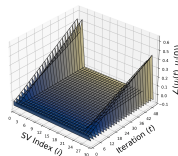
Layer 3



Singular Values



Right Singular Vectors



Left Singular Vectors

The Evolution of Singular Spaces in GD Iterates for DLNs

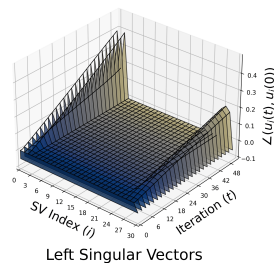
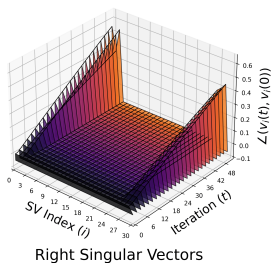
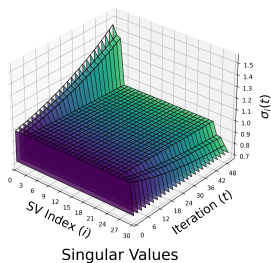


Figure: Evolution of SVD of the weight matrix $W_1(t) = U_1(t)\Sigma_1(t)V_1(t)^\top$.

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

The Law of Parsimony in GD

Theorem (Yaras et al.'23)

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

- **Case 1:** Suppose $\mathbf{Y}\mathbf{X}^\top \in \mathbb{R}^{d_y \times d_x}$ with $d_y = r$ and $m := d_x - 2d_y > 0$. Then $\exists \{\mathbf{U}_l\}_{l=1}^L \subseteq \mathcal{O}^d$ and $\{\mathbf{V}_l\}_{l=1}^L \subseteq \mathcal{O}^d$ satisfying $\mathbf{V}_{l+1} = \mathbf{U}_l$ for all $l \in [L-1]$, such that $\mathbf{W}_l(t)$ admits the following decomposition

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

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

The Law of Parsimony in GD

Theorem (Yaras et al.'23)

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

- **Case 1:** Suppose $\mathbf{Y}\mathbf{X}^\top \in \mathbb{R}^{d_y \times d_x}$ with $d_y = r$ and $m := d_x - 2d_y > 0$. Then $\exists \{\mathbf{U}_l\}_{l=1}^L \subseteq \mathcal{O}^d$ and $\{\mathbf{V}_l\}_{l=1}^L \subseteq \mathcal{O}^d$ satisfying $\mathbf{V}_{l+1} = \mathbf{U}_l$ for all $l \in [L-1]$, such that $\mathbf{W}_l(t)$ admits the following decomposition

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

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

- **Case 2:** Suppose $\mathbf{Y}\mathbf{X}^\top \in \mathbb{R}^{d_y \times d_x}$ is of rank $r \in \mathbb{N}_+$ with $d_y = d_x$, and $m = d_x - 2r > 0$. Similar results hold with different $\rho(t)$.

The Law of Parsimony in GD

- **Dynamics of singular values and vectors of weight matrices.**

Let $\widetilde{\mathbf{W}}_l(t) = \widetilde{\mathbf{U}}_l(t)\widetilde{\boldsymbol{\Sigma}}_l(t)\widetilde{\mathbf{V}}_l^\top(t)$, we can rewrite our decomposition as

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

⁴M. Huh et al. The Low-Rank Simplicity Bias in Deep Networks, TMLR'23.

The Law of Parsimony in GD

- **Dynamics of singular values and vectors of weight matrices.**

Let $\widetilde{\mathbf{W}}_l(t) = \widetilde{\mathbf{U}}_l(t)\widetilde{\boldsymbol{\Sigma}}_l(t)\widetilde{\mathbf{V}}_l^\top(t)$, we can rewrite our decomposition as

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

- **Invariance of subspaces in the weights.** Both $\mathbf{U}_{l,2}$ and $\mathbf{V}_{l,2}$ of size $d - 2r$ are unchanged throughout training. The learning process occurs **only** within an **invariant subspace** of dimension $2r$!

⁴M. Huh et al. The Low-Rank Simplicity Bias in Deep Networks, TMLR'23.

The Law of Parsimony in GD

- **Dynamics of singular values and vectors of weight matrices.**

Let $\tilde{\mathbf{W}}_l(t) = \tilde{\mathbf{U}}_l(t)\tilde{\mathbf{\Sigma}}_l(t)\tilde{\mathbf{V}}_l^\top(t)$, we can rewrite our decomposition as

$$\mathbf{W}_l(t) = \begin{bmatrix} \mathbf{U}_{l,1}\tilde{\mathbf{U}}_l(t) & \mathbf{U}_{l,2} \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{\Sigma}}_l(t) & \mathbf{0} \\ \mathbf{0} & \rho(t)\mathbf{I}_m \end{bmatrix} \begin{bmatrix} \mathbf{V}_{l,1}\tilde{\mathbf{V}}_l(t) & \mathbf{V}_{l,2} \end{bmatrix}^\top,$$

- **Invariance of subspaces in the weights.** Both $\mathbf{U}_{l,2}$ and $\mathbf{V}_{l,2}$ of size $d - 2r$ are unchanged throughout training. The learning process occurs **only** within an **invariant subspace** of dimension $2r$!
- **Implicit low-rank bias.**⁴ As $\lim_{\varepsilon \rightarrow 0} \rho(t) = 0$ for all $t \geq 0$, all the weights $\mathbf{W}_l(t)$ and the end-to-end matrix $\mathbf{W}_{L:1}(t)$ are inherently low-rank (e.g., at most rank $2r$).

⁴M. Huh et al. The Low-Rank Simplicity Bias in Deep Networks, TMLR'23.

The Evolution of Singular Spaces in More Generic Settings

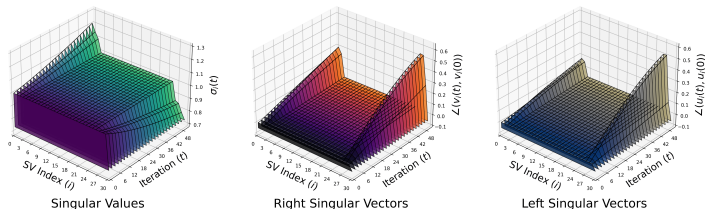


Figure: Evolution of SVD of weight matrices without whitened data.

The Evolution of Singular Spaces in More Generic Settings

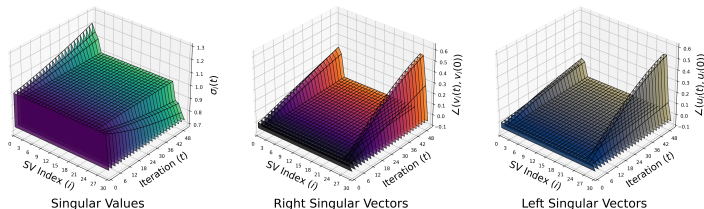


Figure: Evolution of SVD of weight matrices without whitened data.

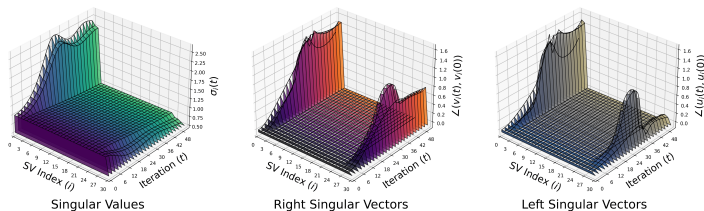
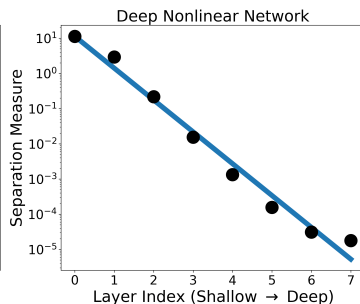
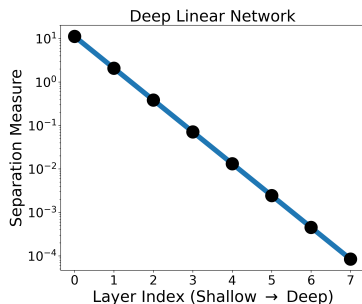


Figure: Evolution of SVD of weight matrices with momentum.

Outline

- 1 Neural Collapse, Transfer Learning, & Intermediate Layers
- 2 Law of Parsimony in Gradient Dynamics
- 3 Progressive Feature Separation in Deep Neural Networks
- 4 Efficient Deep Matrix Completion
- 5 Conclusion

Main Message



The law of parsimony in GD explains progressive feature separation.

Problem Setup: Train DLNs for Classification Problems

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

$$\mathbf{z}_{k,i}^l := \mathbf{W}_l \dots \mathbf{W}_1 \mathbf{x}_{k,i} = \mathbf{W}_{l:1} \mathbf{x}_{k,i}, \quad \forall l \in [L],$$

Problem Setup: Train DLNs for Classification Problems

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

$$\mathbf{z}_{k,i}^l := \mathbf{W}_l \dots \mathbf{W}_1 \mathbf{x}_{k,i} = \mathbf{W}_{l:1} \mathbf{x}_{k,i}, \quad \forall l \in [L],$$

- **Measure of Data Separation:** we replace $\mathcal{NC}_1^l = \text{trace}(\mathbf{\Sigma}_W^l \mathbf{\Sigma}_B^{l\dagger})$ with a simpler measure

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

Problem Setup: Train DLNs for Classification Problems

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

$$\mathbf{z}_{k,i}^l := \mathbf{W}_l \dots \mathbf{W}_1 \mathbf{x}_{k,i} = \mathbf{W}_{l:1} \mathbf{x}_{k,i}, \quad \forall l \in [L],$$

- **Measure of Data Separation:** we replace $\mathcal{NC}_1^l = \text{trace}(\mathbf{\Sigma}_W^l \mathbf{\Sigma}_B^{l\dagger})$ with a simpler measure

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

$$\mathbf{\Sigma}_W^l = \frac{1}{nK} \sum_{k=1}^K \sum_{i=1}^n \left(\mathbf{z}_{k,i}^l - \bar{\mathbf{z}}_k^l \right) \left(\mathbf{z}_{k,i}^l - \bar{\mathbf{z}}_k^l \right)^\top, \quad \mathbf{\Sigma}_B^l = \frac{1}{K} \sum_{k=1}^K \left(\bar{\mathbf{z}}_k^l - \mathbf{z}_G^l \right)$$

Progressive Feature Separation with Linear Decay Rate

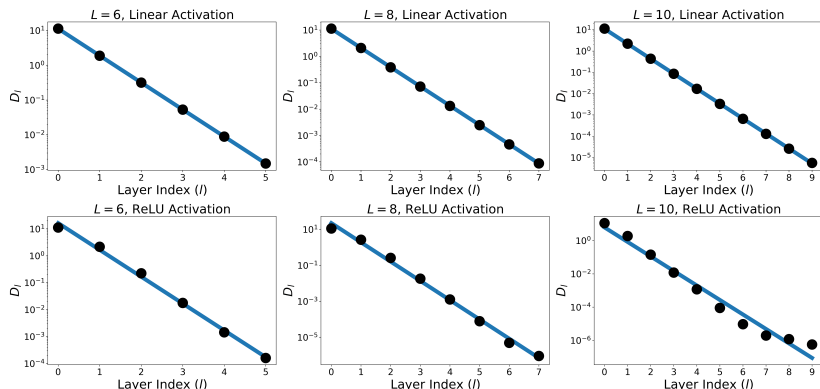


Figure: Linear decay of feature separation in trained deep networks.

Progressive Feature Separation with Linear Decay Rate

Theorem (Wang et al.'23)

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

- **Global Optimality:** $\mathbf{W}_{L:1}\mathbf{X} = \mathbf{Y}$.
- **Balancedness:** For all weights

$$\mathbf{W}_{l+1}^\top \mathbf{W}_{l+1} = \mathbf{W}_l \mathbf{W}_l^\top, \forall l \in [L-2],$$

$$\|\mathbf{W}_L^\top \mathbf{W}_L - \mathbf{W}_{L-1} \mathbf{W}_{L-1}^\top\|_F \leq \varepsilon^2 \sqrt{d-K}.$$

Progressive Feature Separation with Linear Decay Rate

Theorem (Wang et al.'23)

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

- **Global Optimality:** $\mathbf{W}_{L:1}\mathbf{X} = \mathbf{Y}$.
- **Balancedness:** For all weights

$$\mathbf{W}_{l+1}^\top \mathbf{W}_{l+1} = \mathbf{W}_l \mathbf{W}_l^\top, \forall l \in [L-2],$$

$$\|\mathbf{W}_L^\top \mathbf{W}_L - \mathbf{W}_{L-1} \mathbf{W}_{L-1}^\top\|_F \leq \varepsilon^2 \sqrt{d - K}.$$

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

- **Global Optimality:** $\mathbf{W}_{L:1}\mathbf{X} = \mathbf{Y}$.
- **Balancedness:** For all weights

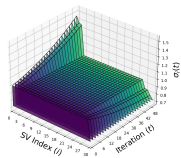
$$\begin{aligned} \mathbf{W}_{l+1}^\top \mathbf{W}_{l+1} &= \mathbf{W}_l \mathbf{W}_l^\top, \forall l \in [L-2], \\ \|\mathbf{W}_L^\top \mathbf{W}_L - \mathbf{W}_{L-1} \mathbf{W}_{L-1}^\top\|_F &\leq \varepsilon^2 \sqrt{d-2K}. \end{aligned}$$

- **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 $\sigma_i(\mathbf{W}_l) = \varepsilon, \forall i \in \mathcal{A}$.

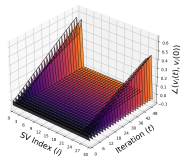
Then, it holds for all $l = 0, 1, \dots, L-2$ that

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

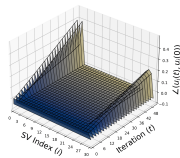
Layer 1



Singular Values

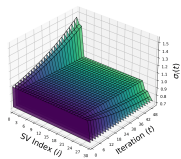


Right Singular Vectors

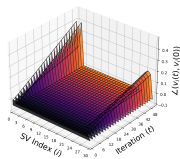


Left Singular Vectors

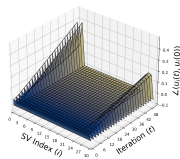
Layer 2



Singular Values

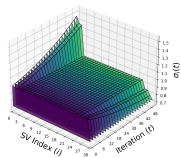


Right Singular Vectors

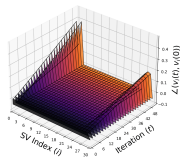


Left Singular Vectors

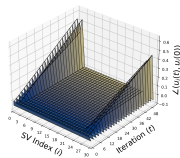
Layer 3



Singular Values



Right Singular Vectors



Left Singular Vectors

Effects of Initialization Scale ε

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

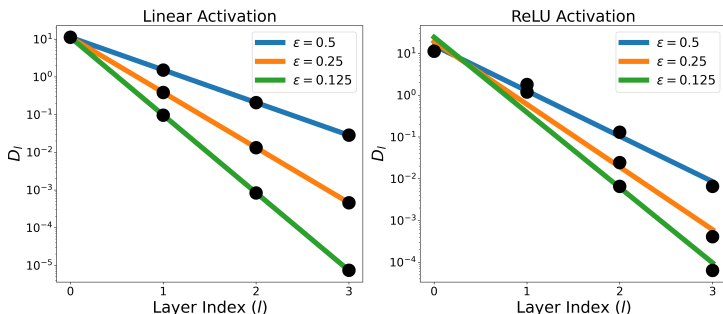


Figure: Linear decay of feature separation measure D_l in trained deep networks with varying initialization scale ε .

Tradeoffs Between Decay Rate and Convergence

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

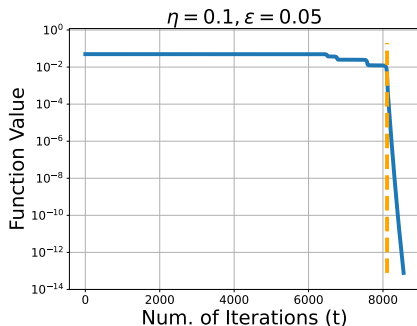
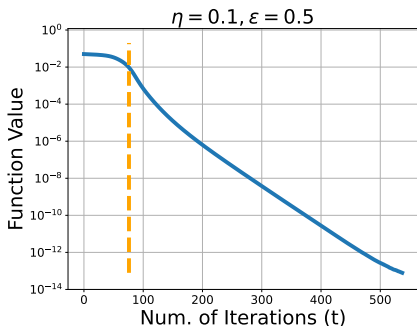


Figure: The dynamics of GD for DLNs with learning rate $\eta = 0.1$.

Is the Orthogonal Initialization Critical?

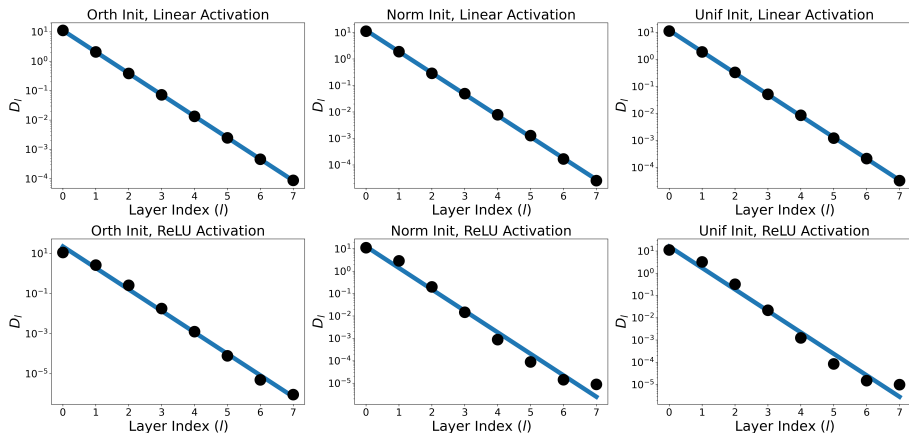


Figure: Linear decay of feature separation in trained DLNs with different initialization types (left to right: Orth., Norm, Unif).

Outline

- 1 Neural Collapse, Transfer Learning, & Intermediate Layers
- 2 Law of Parsimony in Gradient Dynamics
- 3 Progressive Feature Separation in Deep Neural Networks
- 4 Efficient Deep Matrix Completion**
- 5 Conclusion

Main Message

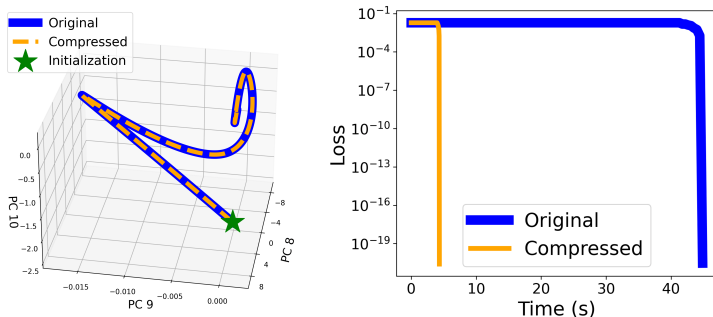


Figure: Efficient training of deep linear networks.

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

Problem Setup for Deep Matrix Completion

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

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

Problem Setup for Deep Matrix Completion

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

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

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

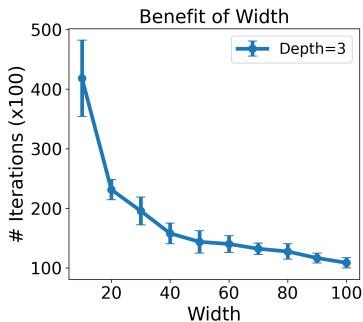
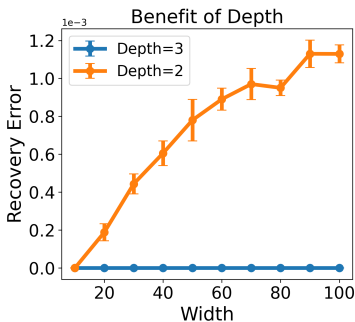
Problem Setup for Deep Matrix Completion

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

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

- 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.

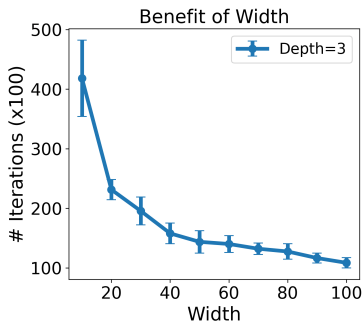
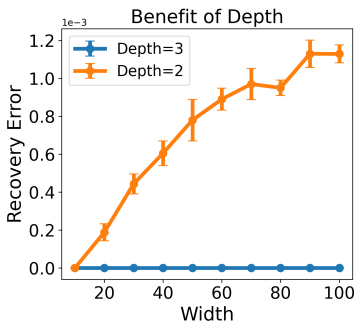
Why Deep Matrix Factorization and Overparameterization?



- **Benefits of Depth (Left):** Improved sample complexity⁵ and less prone to overfitting.

⁵Arora, S., Cohen, N., Hu, W., & Luo, Y. (2019). Implicit regularization in deep matrix factorization. *Advances in Neural Information Processing Systems*, 32.

Why Deep Matrix Factorization and Overparameterization?



- **Benefits of Depth (Left):** Improved sample complexity⁵ and less prone to overfitting.
- **Benefits of Width (Right):** Increasing the width of the network results in accelerated convergence in terms of iterations.

⁵Arora, S., Cohen, N., Hu, W., & Luo, Y. (2019). Implicit regularization in deep matrix factorization. *Advances in Neural Information Processing Systems*, 32.

Overparameterization: A Double Edged Sword

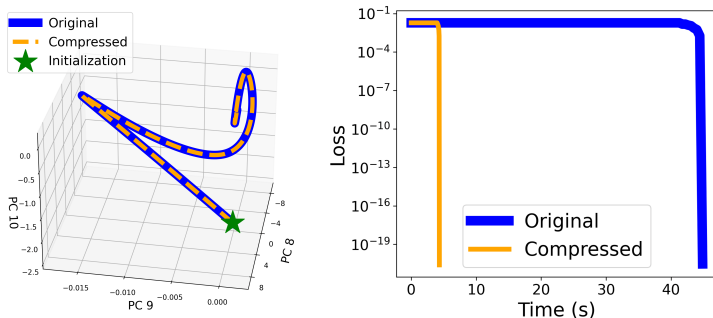


Figure: Efficient training of deep linear networks.

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

How to Achieve the Best of Two Worlds?

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

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

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

How to Achieve the Best of Two Worlds?

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

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

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 $\mathbf{W}_{L:1}$:

$$\begin{aligned} \mathbf{W}_{L:1}(t) &= [\mathbf{U}_{L,1} \quad \mathbf{U}_{L,2}] \begin{bmatrix} \widetilde{\mathbf{W}}_{L:1}(t) & \mathbf{0} \\ \mathbf{0} & \rho^L(t) \mathbf{I}_m \end{bmatrix} \begin{bmatrix} \mathbf{V}_{1,1}^\top \\ \mathbf{V}_{1,2}^\top \end{bmatrix} \\ &= \mathbf{U}_{L,1} \widetilde{\mathbf{W}}_{L:1}(t) \mathbf{V}_{1,1}^\top + \rho^L(t) \mathbf{U}_{L,2} \mathbf{V}_{1,2}^\top, \end{aligned}$$

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

How to Achieve the Best of Two Worlds?

- **The effects of small initialization ε and depth L :**

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

How to Achieve the Best of Two Worlds?

- **The effects of small initialization ε and depth L :**

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

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

The Simple Case: Deep Matrix Factorization

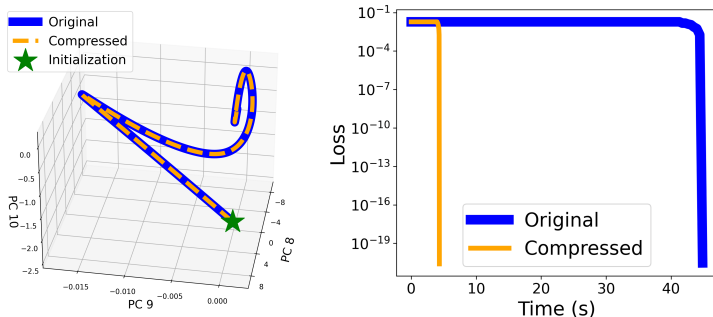
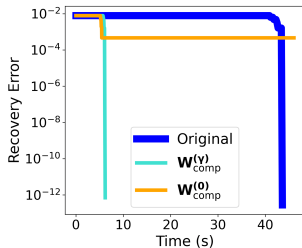
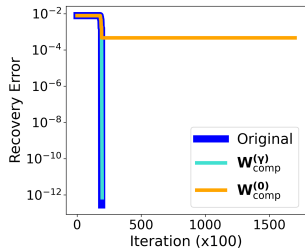
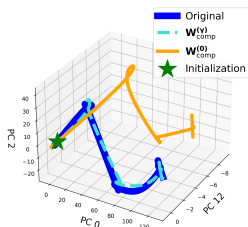


Figure: Efficient training of deep linear networks.

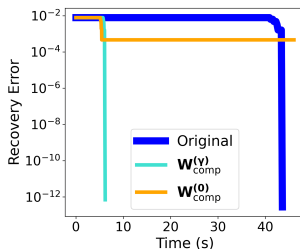
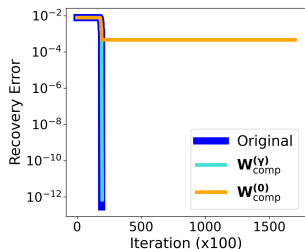
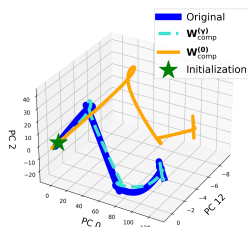
Comparison on the number of parameters: original network Ld^2
vs. compressed network $L\hat{r}^2$.

From Deep Matrix Factorization to Completion?



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

From Deep Matrix Factorization to Completion?

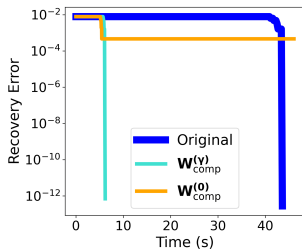
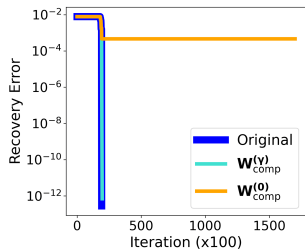
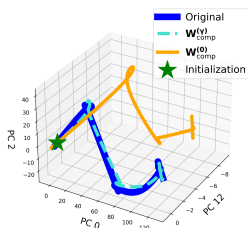


- However, directly applying our approach from deep matrix factorization to completion does not work well...
- This is due to the fact that the law of parsimony in GD:

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

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

From Deep Matrix Factorization to Completion?

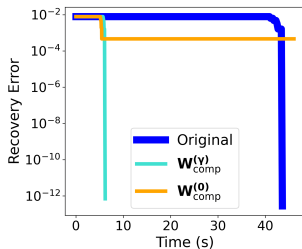
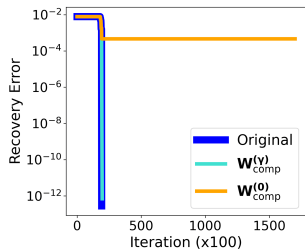
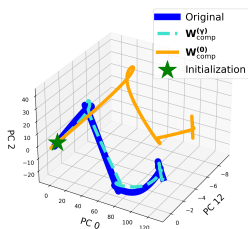


- **Remedy:** update both $V_{1,1}(t)$ and $U_{L,1}(t)$ factors via GD with a **discrepant** learning rate γ in the “compressed network”:⁶

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

⁶This is done simultaneously with the GD updates on the subnetwork $\widetilde{\mathbf{W}}_{L:1}(t)$, which uses the original learning rate η .

From Deep Matrix Factorization to Completion?



- **Remedy:** update both $\mathbf{V}_{1,1}(t)$ and $\mathbf{U}_{L,1}(t)$ factors via GD with a **discrepant** learning rate $\gamma\eta$ in the “compressed network”:⁶

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

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

⁶This is done simultaneously with the GD updates on the subnetwork $\widetilde{\mathbf{W}}_{L:1}(t)$, which uses the original learning rate η .

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 $d \gg \hat{r}$?

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 $d \gg \hat{r}$?

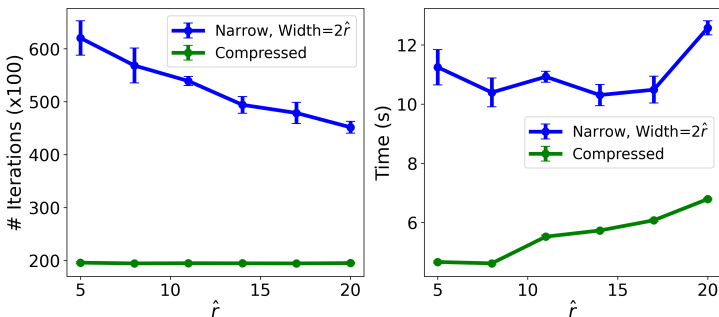


Figure: Efficiency of compressed networks vs. narrow network.

Compressed Networks vs. Narrow Networks?

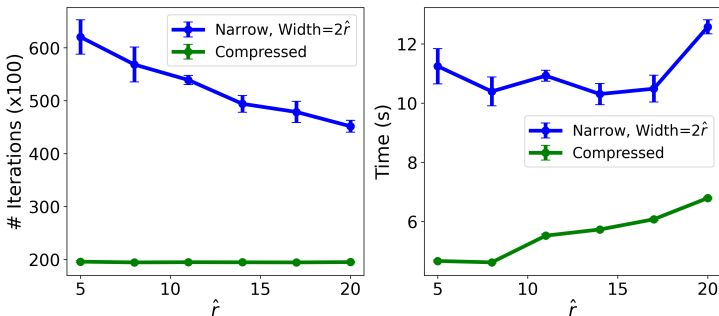


Figure: Efficiency of compressed networks vs. narrow network.

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

Outline

- 1 Neural Collapse, Transfer Learning, & Intermediate Layers
- 2 Law of Parsimony in Gradient Dynamics
- 3 Progressive Feature Separation in Deep Neural Networks
- 4 Efficient Deep Matrix Completion
- 5 Conclusion**

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 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.
- 3 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.
- 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.

Conclusion and Coming Attractions

Learning common deep networks for low-dim structure

- **Low-dimensional features:** understand low-dim. features (sparse and neural collapse (NC)) learned in deep classifiers trained with one-hot labeling based losses
- **Law-of-parsimony in GD:** efficient network compression & training, and understanding intermediate layers of deep networks

Next lecture: Theory and practice of classifying and representing data with nonlinear low-dimensional structures.

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.**

