# Connectivity Shapes Implicit Regularization in Matrix Factorization Models for Matrix Completion

Zhiwei Bai

Shanghai Jiao Tong University

School of Mathematical Sciences & Institute of Natural Sciences
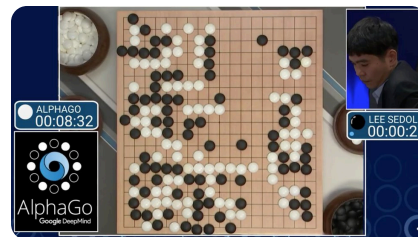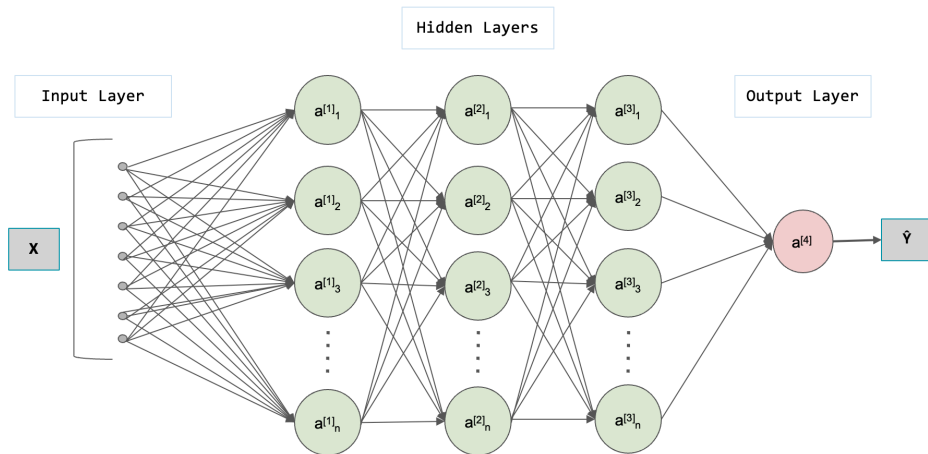
September 26, 2024

# 目录
# CONTENTS

# 1. Introduction and Motivation

# Background: DNNs as Function Approximator

- **Deep Neural Networks (DNNs)** have achieved remarkable success in various fields.



- **DNNs as Function Approximator**

$$f\left( \vphantom{\text{cat}} \right) = \text{"Cat"}$$

$$f\left( \vphantom{\text{wave}} \right) = \text{"How are you"}$$

$$f\left( \vphantom{\text{go}} \right) = \text{"5-5"}_{\text{(next move)}}$$

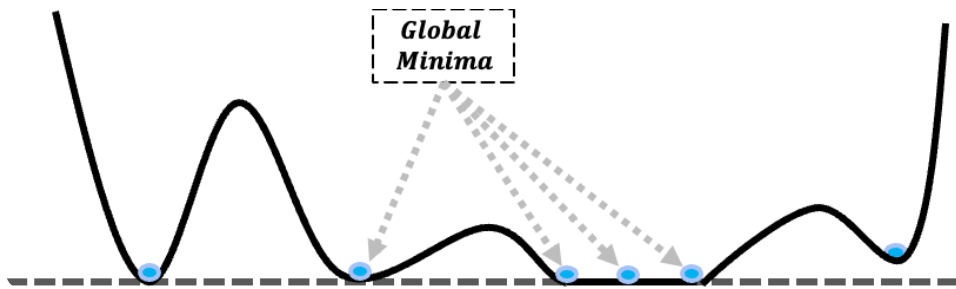- **Key Structure: Composition of Functions Layer by Layer**

$$\boldsymbol{f}_{\boldsymbol{\theta}}^{[l]}(\boldsymbol{x}) = \sigma(\boldsymbol{W}^{[l]} \boldsymbol{f}_{\boldsymbol{\theta}}^{[l-1]}(\boldsymbol{x}) + \boldsymbol{b}^{[l]}), l = 1, 2, \cdots, L-1.$$

$$\boldsymbol{f}_{\boldsymbol{\theta}}^{[l]}(\boldsymbol{X}) = \sum_{i=1}^{h} \text{softmax}_{\text{row}}\left( \frac{\boldsymbol{X}\boldsymbol{W}_{Q_i}\boldsymbol{W}_{K_i}^{\top}\boldsymbol{X}^{\top}}{\sqrt{d_k}} \right) \boldsymbol{X}\boldsymbol{W}_{V_i}\boldsymbol{W}_{O_i}$$

# Background: How to Understand the Learning Behavior?

- **Theory: Understanding the learning behavior**

- **DNNs: Overparameterization**

$$\#param >> \#data$$



**Global Minima**

🔍 **Q: Which Global Minimum is learned?**

- **Mathematical Formulation**

  - Empirical risk:

  $$R_S(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^{n} \ell(\boldsymbol{f}(\boldsymbol{x}_i; \boldsymbol{\theta}), \boldsymbol{y}_i)$$

  - Model: $\boldsymbol{f}(\boldsymbol{x}; \boldsymbol{\theta})$
  - Data: $S = \{(\boldsymbol{x}_i, \boldsymbol{y}_i)\}_{i=1}^{n}$
  - Loss function: $\ell(\cdot, \cdot)$
  - Learning dynamics: $\dot{\boldsymbol{\theta}} = -\nabla R_S(\boldsymbol{\theta})$ with $\boldsymbol{\theta}_0 \sim N(\boldsymbol{0}, \sigma^2)$

🔍 **How to analyze the learning dynamics?**

# Background: the Generalization Mystery

- ## DNNs' capacity is very large



(a) learning curves  (b) convergence slowdown  (c) generalization error growth

Sufficiently large for memorizing the entire random dataset

💡 Q: Is explicit regularization necessary?

[Zhang et al.] Understanding deep learning requires rethinking generalization. ICLR 2017 (Best Paper)↱

- ## DNNs generalize well without explicit regularization

| model | # params | random crop | weight decay | train accuracy | test accuracy | |
|---|---|---|---|---|---|---|
| Inception | 1,649,402 | yes | yes | 100.0 | 89.05 | Not Sufficient |
| | | yes | no | 100.0 | 89.31 | |
| | | no | yes | 100.0 | 86.03 | |
| | | no | no | 100.0 | 85.75 | |
| (fitting random labels) | | no | no | 100.0 | 9.78 | |
| Inception w/o BatchNorm | 1,649,402 | no | yes | 100.0 | 83.00 | |
| | | no | no | 100.0 | 82.00 | |
| (fitting random labels) | | no | no | 100.0 | 10.12 | |
| Alexnet | 1,387,786 | yes | yes | 99.90 | 81.22 | Not Necessary |
| | | yes | no | 99.82 | 79.66 | |
| | | no | yes | 100.0 | 77.36 | |
| | | no | no | 100.0 | 76.07 | |
| (fitting random labels) | | no | no | 99.82 | 9.86 | |
| MLP 3x512 | 1,735,178 | no | yes | 100.0 | 53.35 | |
| | | no | no | 100.0 | 52.39 | |
| (fitting random labels) | | no | no | 100.0 | 10.48 | |
| MLP 1x512 | 1,209,866 | no | yes | 99.80 | 50.39 | |
| | | no | no | 100.0 | 50.51 | |
| (fitting random labels) | | no | no | 99.34 | 10.61 | |

💡 Explicit regularization may improve generalization performance, but is neither necessary nor sufficient

# The Generalization Mystery $\implies$ Implicit Regularization

- **Matrix Completion**

$$\begin{bmatrix} 1 & 2 & 3 \\ \star & 4 & \star \\ \star & \star & 9 \end{bmatrix} \implies \begin{bmatrix} 1 & 2 & 3 \\ {\color{red}0} & 4 & {\color{red}0} \\ {\color{red}0} & {\color{red}0} & 9 \end{bmatrix}$$

- **Non–Factorization Model (Overparameterization):**

$$\boldsymbol{f_\theta} = \boldsymbol{W} \in \mathbb{R}^{d \times d}, \boldsymbol{\theta} = \mathrm{vec}(\boldsymbol{W}) \in \mathbb{R}^{d^2}$$

- **Linear** w.r.t. $\boldsymbol{\theta}$

- **Convex Optimization**

$$R_S(\boldsymbol{\theta}) = \frac{1}{n} \sum_{k=1}^n ((\boldsymbol{f_\theta})_{i_k j_k} - \boldsymbol{M}_{i_k j_k})^2$$

- **Implicit Regularization** $(\dot{\boldsymbol{\theta}} = -\nabla R_S(\boldsymbol{\theta}))$

$$\min_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|_2 = \|\boldsymbol{W} - \boldsymbol{W}_0\|_F$$

- **Matrix Completion**

$$\begin{bmatrix} 1 & 2 & 3 \\ \star & 4 & \star \\ \star & \star & 9 \end{bmatrix} \implies \begin{bmatrix} 1 & 2 & 3 \\ {\color{red}2} & 4 & {\color{red}6} \\ {\color{red}3} & {\color{red}6} & 9 \end{bmatrix}$$

- **Matrix Factorization Model (Composition Structure, Overparameterization)**

$$\boldsymbol{f_\theta} = \boldsymbol{A}\boldsymbol{B} \in \mathbb{R}^{d \times d}, \boldsymbol{A}, \boldsymbol{B} \in \mathbb{R}^{d \times d}$$

- **Non–Linear** w.r.t. $\boldsymbol{\theta}$

- **Non–Convex Optimization**

$$R_S(\boldsymbol{\theta}) = \frac{1}{n} \sum_{k=1}^n ((\boldsymbol{f_\theta})_{i_k j_k} - \boldsymbol{M}_{i_k j_k})^2$$

- **Implicit Regularization** $(\dot{\boldsymbol{\theta}} = -\nabla R_S(\boldsymbol{\theta}))$

$$????????????$$

# Recent Works on Implicit Regularization in Matrix Factorization

[Gunasekar et al.] 2017 NeurIPS

[Jin et al.] 2023 ICML

[Arora et al.] 2019 NeurIPS

## 💡 Nuclear Norm Minimization

**Theorem 1.** In the case where the **observation matrices** $\{\mathcal{A}_i\}_{i=1}^{m}$ **commute**, the symmetrical matrix factorization model $f_\theta = UU^\top$ finds the minimal nuclear norm solution.

## 💡 Nuclear Norm Minimization

**Theorem 2.** In the case where the **observation matrices** $\{\mathcal{A}_i\}_{i=1}^{m}$ **commute**, the asymmetrical matrix factorization model $f_\theta = AB$ finds the minimal nuclear norm solution.

## 💡 Rank Minimization

**Theorem 3.** In the case where the **observation matrices** $\{\mathcal{A}_i\}_{i=1}^{m}$ **satisfy the RIP condition**, the symmetrical matrix factorization model $f_\theta = UU^\top$ finds the minimal rank solution.

Are these characterizations sufficient? Do they describe the whole picture of matrix factorization models?

# Examples

- Observation Matrices Commute:

$$E_{ij}E_{mn} = \delta_{jm}E_{in} = E_{mn}E_{ij} = \delta_{ni}E_{mj}$$

$$\implies \begin{bmatrix} \textcolor{red}{\times} & \star & \star & \textcolor{red}{\checkmark} \\ \textcolor{red}{\times} & \star & \star & \star \\ \textcolor{red}{\times} & \star & \star & \star \\ \textcolor{red}{\times} & \textcolor{red}{\times} & \textcolor{red}{\times} & \textcolor{red}{\times} \end{bmatrix}$$

- Counterexample:

$$\begin{bmatrix} 0 & 1 \\ 2 & \star \end{bmatrix} \stackrel{\textcolor{red}{GD}}{\implies} \begin{bmatrix} 0 & 1 \\ 2 & \textcolor{red}{0} \end{bmatrix}, \begin{bmatrix} 1 & 2 \\ \star & 3 \end{bmatrix} \stackrel{\textcolor{red}{GD}}{\implies} \begin{bmatrix} 1 & 2 \\ \textcolor{red}{1.5} & 3 \end{bmatrix}$$

- GD still learned the minimal nuclear norm solution although the observation matrices do not commute

- Restricted Isometry Property (RIP):
  The measurement operator $\mathcal{A}$ satisfies the $(\delta, r)$ RIP if

$$(1 - \delta)\|\boldsymbol{Z}\|_{\mathrm{F}}^2 \leq \|\mathcal{A}(\boldsymbol{Z})\|_2^2 \leq (1 + \delta)\|\boldsymbol{Z}\|_{\mathrm{F}}^2$$

  for all $\boldsymbol{Z} \in \mathbb{R}^{d \times d}$ with $\mathrm{rank}(\boldsymbol{Z}) \leq r$

- Counterexample:

$$\begin{bmatrix} 1 & 2 \\ 3 & \star \end{bmatrix} \stackrel{\textcolor{red}{GD}}{\implies} \begin{bmatrix} 1 & 2 \\ 3 & \textcolor{red}{6} \end{bmatrix}, \begin{bmatrix} 1 & 2 \\ 10 & \star \end{bmatrix} \stackrel{\textcolor{red}{GD}}{\implies} \begin{bmatrix} 1 & 2 \\ 10 & \textcolor{red}{20} \end{bmatrix}$$

- GD still learned the minimal rank solution although the observation matrices do not satisfy the $(\delta, r)$ RIP condition

How to construct a unified understanding of when, how, and why they achieve different implicit regularization effects?

# Empirical Observations

- ## The connectivity of observed data affects the implicit regularization



- ## Low rank bias in connected case
- ## Low nuclear norm bias in certain disconnected case

# 2. Connectivity Affects Implicit Regularization

# Definition of Connectivity

## Observation matrix $P$

$$P_{ij} = \begin{cases} 1, & M_{ij} \text{ is observed and non-zero} \\ 0, & \text{otherwise} \end{cases}$$

## Associated Observation Graph $G_M$

**Definition 1 (Associated Observation Graph).** The associated observation graph $G_M$ is the bipartite graph with adjacency matrix $\begin{bmatrix} \mathbf{0} & P^\top \\ P & \mathbf{0} \end{bmatrix}$, with isolated vertices removed.

## Connectivity

**Definition 2. Connected:** $G_M$ is connected; **Disconnected:** $G_M$ is disconnected
The connected components of $M$ are defined as the connected components of $G_M$.

# Examples of Connectivity

📖 **Disconnectivity with Complete Bipartite Components**

**Definition 3. Disconnectivity with Complete Bipartite Components:** Graph $G_M$ is disconnected and each connected component forms a complete bipartite subgraph.
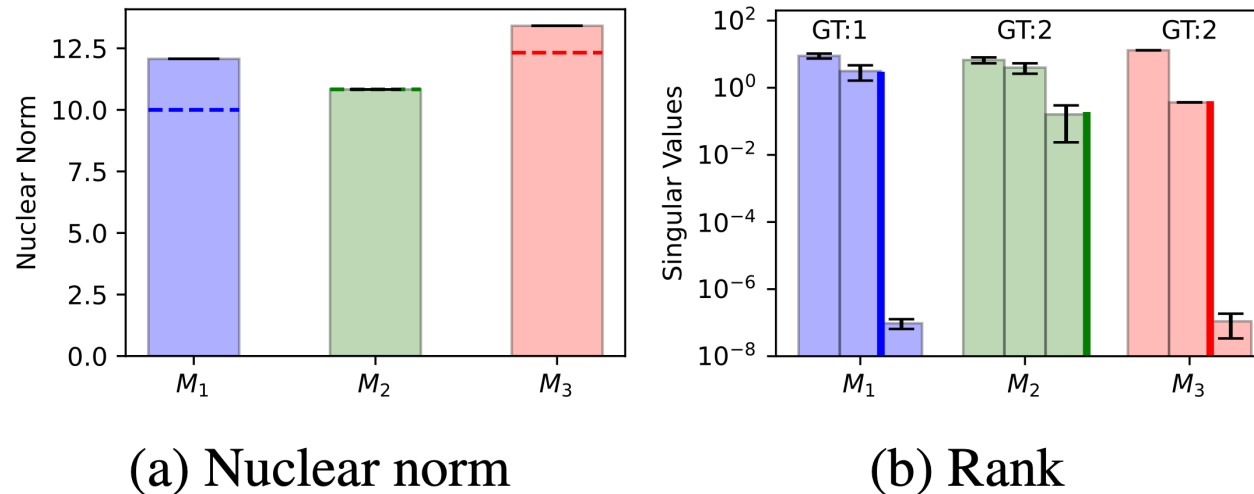
- **Disconnected**

$$M_1 = \begin{bmatrix} 1 & 2 & \star \\ 3 & \star & \star \\ \star & \star & 5 \end{bmatrix}$$



- **Disconnected (complete bipartite components)**

$$M_2 = \begin{bmatrix} 1 & 2 & \star \\ 3 & 4 & \star \\ \star & \star & 5 \end{bmatrix}$$



- **Connected**

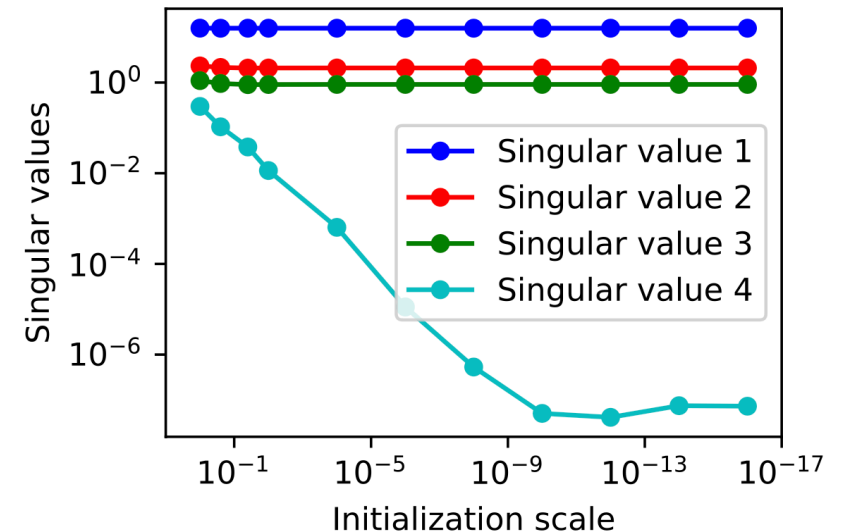$$M_3 = \begin{bmatrix} 1 & 2 & \star \\ 3 & 4 & \star \\ 6 & \star & 5 \end{bmatrix}$$

# Connectivity Affects Implicit Regularization

- **Disconnected**

$$M_1 = \begin{bmatrix} 1 & 2 & \star \\ 3 & \star & \star \\ \star & \star & 5 \end{bmatrix}$$

- **Disconnected (complete bipartite components)**

$$M_2 = \begin{bmatrix} 1 & 2 & \star \\ 3 & 4 & \star \\ \star & \star & 5 \end{bmatrix}$$

- **Connected**

$$M_3 = \begin{bmatrix} 1 & 2 & \star \\ 3 & 4 & \star \\ 6 & \star & 5 \end{bmatrix}$$

- 



(a) Nuclear norm          (b) Rank

# Connected Case—Initialization Matters

- **Matrix Completion**

$$\begin{bmatrix} 4 & 0.6 & 1.8 & 0.8 \\ 6 & 0.9 & 2.7 & \star \\ 8 & 2.2 & 2.6 & 1.6 \\ 8 & 2.7 & 5.1 & 3.6 \end{bmatrix}$$

- **Matrix Factorization Model** $f_\theta = AB$



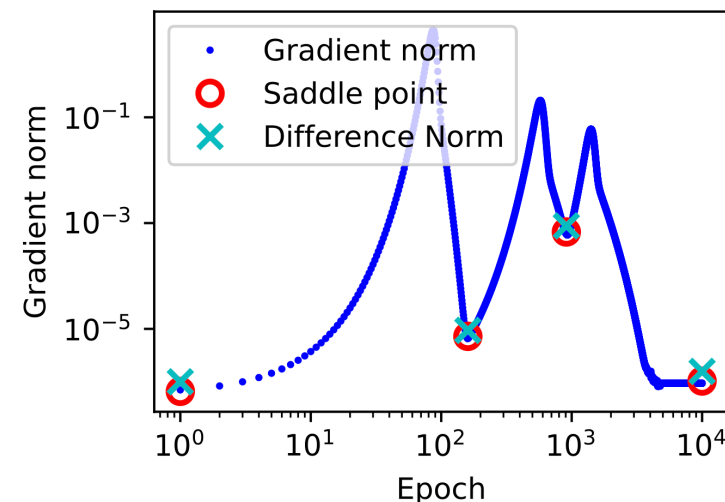- Large initialization: rank–4
- Small initialization: rank–3

💡 Learning lowest–rank solution in infinitesimal initialization

# Connected Case—Traversing Progressive Optima
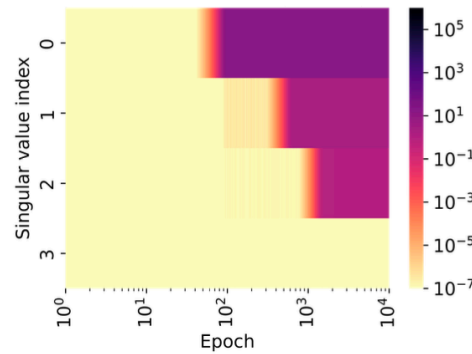
- **Training Loss at Small Initialization**



- **Gradient Norm during Training**



- Training Loss: stepwise decline
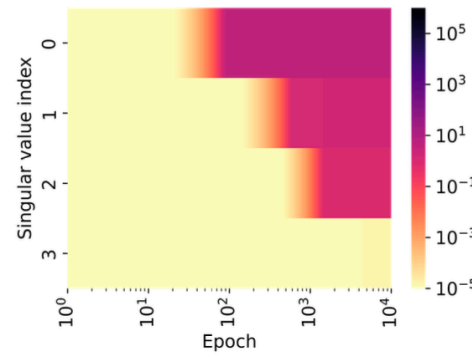- Saddle Points: Experience optimal approximation of each rank

💡 Traversing progressive optima at each rank

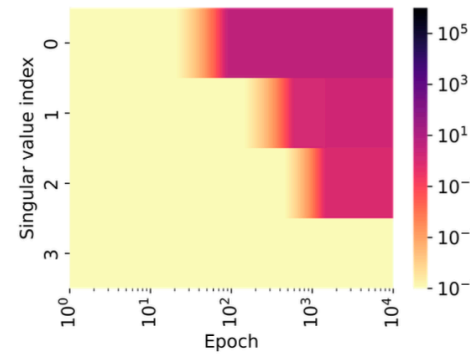# Connected Case—Alignment of $\mathrm{row}(\boldsymbol{A})$ and $\mathrm{col}(\boldsymbol{B})$

- **Evolution of Singular Values**



(e) Singular values of $\boldsymbol{W}$  (f) Singular values of $\boldsymbol{A}$  (g) Singular values of $\boldsymbol{B}$  (h) Singular values of $\boldsymbol{W}_{\mathrm{aug}}$

- **Rank increases step by step**

- $\mathrm{rank}(\boldsymbol{A}) = \mathrm{rank}\left(\boldsymbol{B}^{\top}\right) = \mathrm{rank}\left(\boldsymbol{W}_{\mathrm{aug}}\right)$, where $\boldsymbol{W}_{\mathrm{aug}} = \begin{bmatrix} \boldsymbol{A} \\ \boldsymbol{B}^{\top} \end{bmatrix}$

- $\implies \mathrm{row}(\boldsymbol{A}) = \mathrm{col}(\boldsymbol{B})$, which induces an **invariant manifold** in theoretical analysis

# Disconnected Case—Alignment of $\mathrm{row}(A)$ and $\mathrm{col}(B)$

- Evolution of Singular Values



$$\begin{bmatrix} 1 & \star & 3 \\ \star & 5 & \star \\ 3 & \star & 9 \end{bmatrix}$$

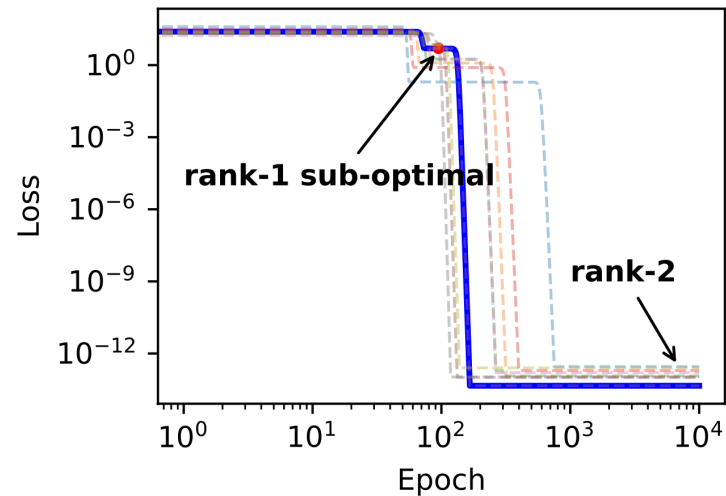(a) Matrix completion  (b) Singular values of $A$  (c) Singular values of $B$  (d) Singular values of $W_{\mathrm{aug}}$

- Alignment of the row space of $A$ and the column space of $B$:
$$\mathrm{row}(A) = \mathrm{col}(B)$$

- Lowest–rank solution is not learned ($\mathrm{rank}$-$2$) in disconnected case!

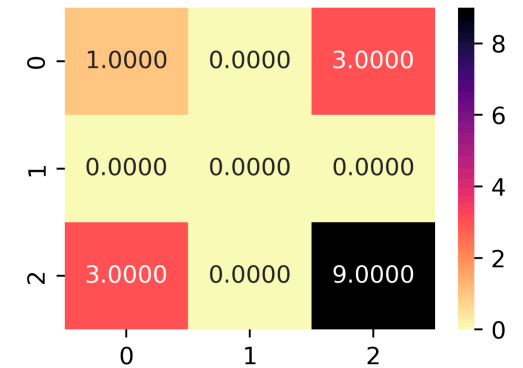- Lowest nuclear norm solution is learned in this disconnected case

# Disconnected Case—Learn Sub-optimal Saddle Point

- **Training Loss at Small Initialization**



- **Rank–1 Sub–optimal**

$$
\begin{bmatrix}
1 & \star & 3 \\
\star & 5 & \star \\
3 & \star & 9
\end{bmatrix}
$$



- **Dynamics: decouple into two independent systems in the disconnected case**

$$
\begin{cases}
\dot{\boldsymbol{a}}_i = -\frac{2}{n} \sum_{j \in I_i} (\boldsymbol{a}_i \cdot \boldsymbol{b}_{\cdot,j} - \boldsymbol{M}_{ij}) \boldsymbol{b}_{\cdot,j}^\top, i \in \{1,3\} \\
\dot{\boldsymbol{b}}_{\cdot,j} = -\frac{2}{n} \sum_{i \in I_j} (\boldsymbol{a}_i \cdot \boldsymbol{b}_{\cdot,j} - \boldsymbol{M}_{ij}) \boldsymbol{a}_i^\top, j \in \{1,3\}
\end{cases}
\qquad
\begin{cases}
\dot{\boldsymbol{a}}_2 = -\frac{2}{n} (\boldsymbol{a}_2 \cdot \boldsymbol{b}_{\cdot,2} - \boldsymbol{M}_{22}) \boldsymbol{b}_{\cdot,2}^\top \\
\dot{\boldsymbol{b}}_{\cdot,2} = -\frac{2}{n} (\boldsymbol{a}_2 \cdot \boldsymbol{b}_{\cdot,2} - \boldsymbol{M}_{ij}) \boldsymbol{a}_2^\top
\end{cases}
$$

# 3. Training Dynamics Analysis

# Hierarchical Intrinsic Invariant Manifold

## Hierarchical Intrinsic Invariant Manifold (HIIM)

**Proposition 1 (Hierarchical Intrinsic Invariant Manifold (HIIM)).** Let $f_\theta = AB$ be a matrix factorization model and $\{\alpha_1, \cdots, \alpha_k\}$ be $k$ linearly independent vectors. Define the manifold $\Omega_k$ as

$$\Omega_k := \Omega_k(\alpha_1, \cdots, \alpha_k) = \{\theta = (A, B) \mid \mathrm{row}(A) = \mathrm{col}(B) = \mathrm{span}\{\alpha_1, \cdots, \alpha_k\}\}$$

The manifold $\Omega_k$ possesses the following properties:

**(1) Invariance under Gradient Flow:** Given data $S$ and the gradient flow dynamics $\dot{\theta} = -\nabla R_S(\theta)$, if the initial point $\theta_0 \in \Omega_k$, then $\theta(t) \in \Omega_k$ for all $t \geq 0$.

**(2) Intrinsic Property:** $\Omega_k$ is a data–independent invariant manifold, meaning that for any data $S$, $\Omega_k$ remains invariant under the gradient flow dynamics.

**(3) Hierarchical Structure:** The manifolds $\Omega_k$ form a hierarchy: $\Omega_0 \subsetneq \Omega_1 \subsetneq \cdots \subsetneq \Omega_{k-1} \subsetneq \Omega_k$.

# Disconnected Case: Intrinsic Sub-$\Omega_k$ Invariant Manifold

✒ **Intrinsic Sub–$\Omega_k$ Invariant Manifold**

**Proposition 2 (Intrinsic Sub–$\Omega_k$ Invariant Manifold).** Let $f_\theta = AB$ be a matrix factorization model, $M$ be an incomplete matrix and $\Omega_k$ be an invariant manifold defined in Prop. 1. If $M$ is disconnected with $m$ connected components, then there exist $m$ sub– $\Omega_k$ manifolds $\omega_k$ such that $\omega_k \subsetneqq \Omega_k$, each possessing the following properties:

**(1) Invariance under Gradient Flow:** Given data $S$ and the gradient flow dynamics $\dot{\theta} = -\nabla R_S(\theta)$, if the initial point $\theta_0 \in \omega_k$, then $\theta(t) \in \omega_k$ for all $t \geq 0$.
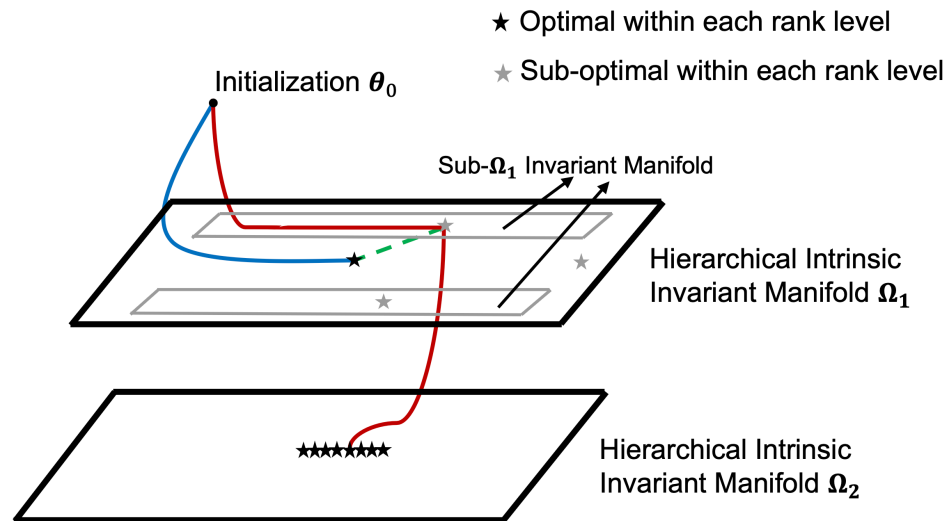
**(2) Intrinsic Property:** $\omega_k$ is a data–value–independent invariant manifold, meaning that for a fixed sampling pattern in $M$ and any observed values $S, \omega_k$ remains invariant under the gradient flow.

**(3) Strict Subset Relation:** The output set $\{f_\theta \mid \theta \in \omega_k\}$ is a proper subset of $\{f_\theta \mid \theta \in \Omega_k\}$, namely, $\{f_\theta \mid \theta \in \omega_k\} \subsetneqq \{f_\theta \mid \theta \in \Omega_k\}$

# Intuitive Illustration

- **Illustration of Training Trajectories**



★ Optimal within each rank level

☆ Sub-optimal within each rank level

Initialization $\boldsymbol{\theta}_0$

Sub-$\boldsymbol{\Omega}_1$ Invariant Manifold

Hierarchical Intrinsic Invariant Manifold $\boldsymbol{\Omega}_1$

Hierarchical Intrinsic Invariant Manifold $\boldsymbol{\Omega}_2$

Blue line represents the trajectory converging to the lowest–rank solution. Red line represents the actual trajectory experienced by the model

- Connected case: Model traverses with invariant manifold $\Omega_k$
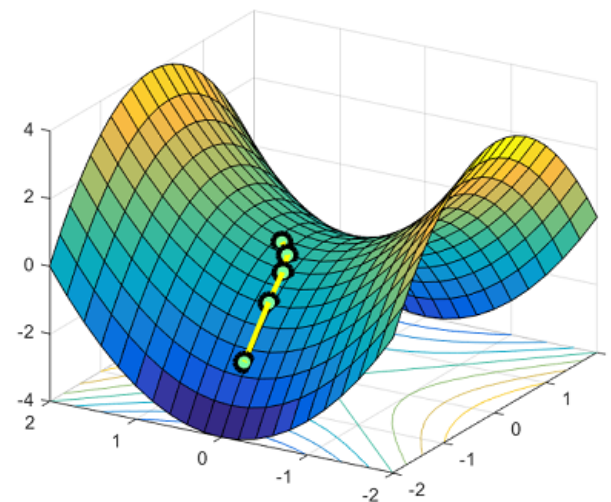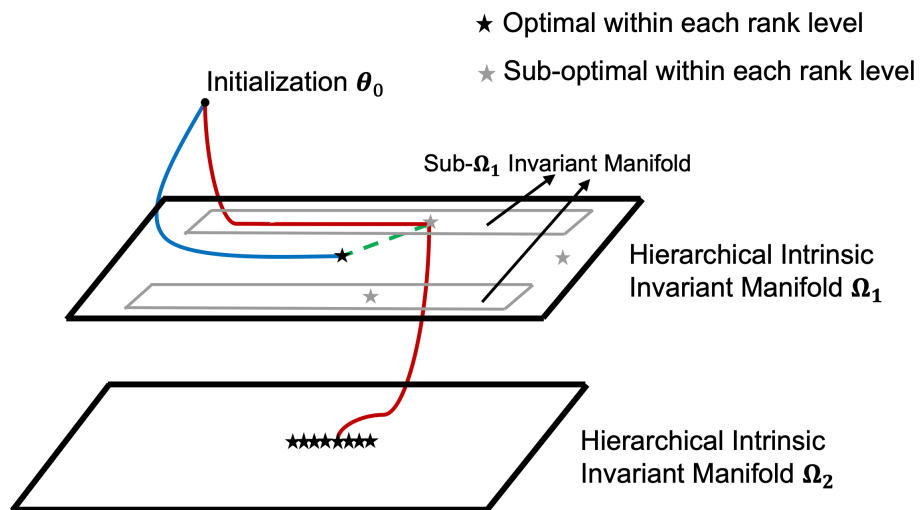
- Disconnected case:

  - Sub–$\Omega_k$ invariant manifold emerges

  - Each sub–$\Omega_k$ induces a sub–optimal saddle point

  - Sub–optima prevent the model from learning the lowest–rank solution

# Loss Landscape does not Contain any Local Minima

> ✒ **Loss Landscape**
>
> **Theorem 1 (Loss Landscape).** Given any data $S$, the critical points of $R_S(\boldsymbol{\theta})$ are either strict saddle points or global minima.

- **Gradient descent easily escapes saddle points**

# Assumptions for Encountered Critical Points

## 🔥 Assumption 1 Top Eigenvalue

**Assumption 1 (Top Eigenvalue).** Let $\delta M = (A_c B_c - M)_{S_x}$ be the residual matrix at the critical point $\theta_c = (A_c, B_c)$. Assume that the top singular value of the matrix $\delta M$ is unique.

## 🔥 Assumption 2 Second–order Stationary Point

**Assumption 2 (Second–order Stationary Point).** Let $\Omega$ be an $\Omega_k$ invariant manifold or sub– $\Omega_k$ invariant manifold defined in Prop. 1 or 2. Assume $\theta_c$ is a second–order stationary point within $\Omega$, i.e., $\nabla R_S(\theta_c) = 0$ and $\theta^\top \nabla^2 R_S(\theta_c)\theta \geq 0$ for all $\theta \in \Omega$.

# Characterization of Training Dynamics

**Theorem 2 (Transition to the Next Rank–level Invariant Manifold).** Consider the dynamics $\dot{\boldsymbol{\theta}} = -\nabla R_S(\boldsymbol{\theta})$. Let $\varphi\left(\boldsymbol{\theta}_0, t\right)$ denote the value of $\boldsymbol{\theta}(t)$ when $\boldsymbol{\theta}(0) = \boldsymbol{\theta}_0$. Let $\boldsymbol{\Omega}$ be an $\boldsymbol{\Omega}_k$ or sub– $\boldsymbol{\Omega}_k$ invariant manifold. Let $\boldsymbol{\theta}_c \in \boldsymbol{\Omega}$ be a critical point satisfying Assump. 1 and 2. Then, for randomly selected $\boldsymbol{\theta}_0$, with probability 1 with respect to $\boldsymbol{\theta}_0$, the limit

$$\tilde{\varphi}\left(\boldsymbol{\theta}_c, t\right) := \lim_{\alpha \to 0} \varphi\left(\boldsymbol{\theta}_c + \alpha \boldsymbol{\theta}_0, t + \tfrac{1}{\lambda_1}\log \tfrac{1}{\alpha}\right)$$

exists and falls into an invariant manifold $\boldsymbol{\Omega}_{k+1}$. Here $\lambda_1$ is the top eigenvalue of $-\nabla^2 R_S\left(\boldsymbol{\theta}_c\right)$.

# Proof Sketch

- Linear Approximation near critical point $\boldsymbol{\theta}_c$:
  $$\frac{d\boldsymbol{\theta}}{dt} \approx H(\boldsymbol{\theta}_0 - \boldsymbol{\theta}_c).$$
- Solution $\boldsymbol{\theta}(t) = e^{tH}(\boldsymbol{\theta}_0 - \boldsymbol{\theta}_c) + \boldsymbol{\theta}_c$, specifically
  $$\boldsymbol{\theta}(t) = \sum_{i=1}^{s} \sum_{j=1}^{l_i} e^{\lambda_i t} \langle \boldsymbol{\theta}_0 - \boldsymbol{\theta}_c, q_{ij} \rangle q_{ij} + \boldsymbol{\theta}_c$$
- Dominant eigenvalue trajectory:
  $$\boldsymbol{\theta}(t_0) = \sum_{j=1}^{l_1} e^{\lambda_1 t_0} \langle \boldsymbol{\theta}_0 - \boldsymbol{\theta}_c, q_{1j} \rangle q_{1j} + O(e^{\lambda_2 t_0})$$
- The first principal component $\sum_{j=1}^{l_1} e^{\lambda_1 t_0} \langle \boldsymbol{\theta}_0 - \boldsymbol{\theta}_c, q_{1j} \rangle q_{1j}$ corresponds to an $\Omega_1$ invariant manifold under Assump. 1 and 2
- Escaping $\boldsymbol{\theta}_c$ increases rank by 1, **entering** $\Omega_{k+1}$

- **Escape from the top eigen–direction**



Alignment of $\mathrm{row}(\boldsymbol{A})$ and $\mathrm{col}(\boldsymbol{B})$

# 4. Implicit Regularization Analysis

# Minimum Rank Regularization

## Minimum Rank Regularization

**Theorem 3 (Minimum Rank).** Consider the dynamics $\dot{\boldsymbol{\theta}} = -\nabla R_S(\boldsymbol{\theta})$, where $\boldsymbol{\theta}(t) = (\boldsymbol{A}(t), \boldsymbol{B}(t))$, and denote $\boldsymbol{W}_t = \boldsymbol{A}(t)\boldsymbol{B}(t)$. Assume $\boldsymbol{W}_t$ achieves an optimal within each invariant manifold $\boldsymbol{\Omega}_k$. For a full rank initialization $\boldsymbol{W}_0$, if the limit $\widehat{\boldsymbol{W}} = \lim_{\alpha \to 0} \boldsymbol{W}_\infty (\alpha \boldsymbol{W}_0)$ exists and is a global optimum with $\widehat{\boldsymbol{W}}_{ij} = \boldsymbol{M}_{ij}$ for all $(i,j) \in S_{\boldsymbol{x}}$, then

$$\widehat{\boldsymbol{W}} \in \operatorname{argmin}_{\boldsymbol{W}} \operatorname{rank}(\boldsymbol{W}) \quad \text{s.t.} \quad \boldsymbol{W}_{ij} = \boldsymbol{M}_{ij}, \forall (i,j) \in S_{\boldsymbol{x}}$$

- In connected case, experiments provide strong evidence that model achieves an optimal within each invariant manifold $\boldsymbol{\Omega}_k$

# Minimum Nuclear Norm Regularization

- In disconnected case, the minimum nuclear norm may still serve as a characterization

> ✒ **Minimum Nuclear Norm Regularization**
>
> **Theorem 4 (Minimum Nuclear Norm Guarantee).** Consider the dynamics $\dot{\boldsymbol{\theta}} = -\nabla R_S(\boldsymbol{\theta})$, where $\boldsymbol{\theta}(t) = (\boldsymbol{A}(t), \boldsymbol{B}(t))$, and let $\boldsymbol{W}_t = \boldsymbol{A}(t)\boldsymbol{B}(t)$. If the observation graph associated with the incomplete matrix $\boldsymbol{M}$ is **disconnected with complete bipartite components**, and if for a full rank initialization $\boldsymbol{W}_0$, the limit $\widehat{\boldsymbol{W}} = \lim_{\alpha \to 0} \boldsymbol{W}_\infty (\alpha \boldsymbol{W}_0)$ exists and is a global optimum with $\widehat{\boldsymbol{W}}_{ij} = \boldsymbol{M}_{ij}$ for all $(i,j) \in S_{\boldsymbol{x}}$, then
>
> $$\widehat{\boldsymbol{W}} \in \operatorname{argmin}_{\boldsymbol{W}} \|\boldsymbol{W}\|_* \quad \text{s.t.} \quad \boldsymbol{W}_{ij} = \boldsymbol{M}_{ij}, \forall (i,j) \in S_{\boldsymbol{x}}$$

# 5. Discussion and Conclusion

# Generalize to Neural Networks: From Linear to Nonlinear

- **Matrix Factorization:**
$$f_{\boldsymbol{\theta}} = \boldsymbol{AB}$$

- **Linear** w.r.t input $\boldsymbol{x}$

- **Implicit Bias: Low rank**

- **Neural Networks:**
$$f_{\boldsymbol{\theta}}(\boldsymbol{x}) = \sum_{i=1}^{m} a_i \sigma(\boldsymbol{w}_i^\top \boldsymbol{x})$$

- **Non–Linear** w.r.t input $\boldsymbol{x}$

- **Implicit Bias: ???????**

- **Model Rank for Non–linear Models:**
$$\operatorname{rank}_{f_{\boldsymbol{\theta}}}(\boldsymbol{\theta}^*) := \dim \left( \operatorname{span} \left\{ \partial_{\theta_i} f\left(\cdot; \boldsymbol{\theta}^*\right) \right\}_{i=1}^{M} \right)$$

- 💡 **Experiments: Non–linear models has low model rank bias**

[Zhang et al.] Yaoyu Zhang*, Zhongwang Zhang, Leyang Zhang, *Zhiwei Bai*, Tao Luo, Zhi–Qin John Xu. Optimistic estimate uncovers the potential of nonlinear models. arXiv preprint arXiv: 2307.08921, 2023.

# Generalize to Transformer Architecture

- Matrix Factorization Model is a Component of the Transformer Architecture

$$Y = \sum_{i=1}^{h} \text{softmax}_{\text{row}} \left( \frac{X W_{Q_i} W_{K_i}^{\top} X^{\top}}{\sqrt{d_k}} \right) X W_{V_i} W_{O_i}$$
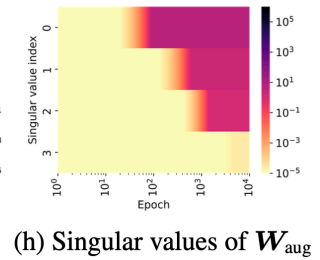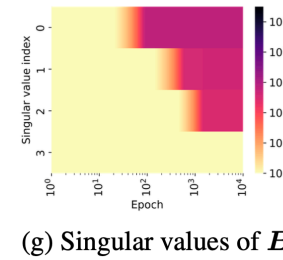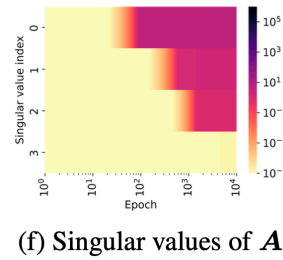
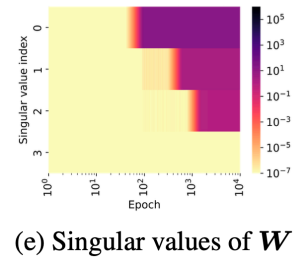- Low-rank (model rank) behavior in the Attention Module



(a) Singular values of $\boldsymbol{W_Q}$ (b) Singular values of $\boldsymbol{W_K}$ (c) Singular values of $\boldsymbol{W_V}$ (d) Singular values of $\boldsymbol{W_O}$
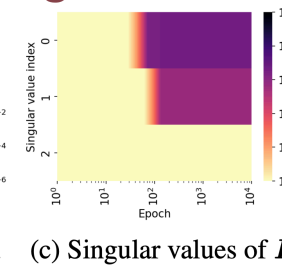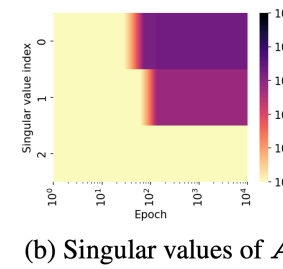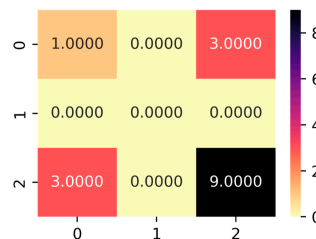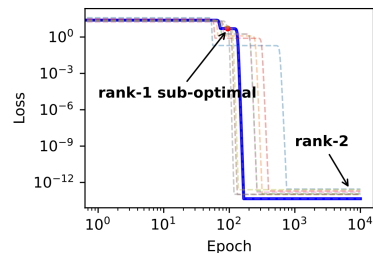
# Take Home Messages

- **Implicit Regularization** of **Overparameterized** models $\implies$ **Generalization**

- **Connected Case: Hierarchical Invariant Manifold Traversal**; Model achieves optima within each invariant manifold $\implies$ **Minimum Rank Regularization**



$$\begin{bmatrix} 4 & 0.6 & 1.8 & 0.8 \\ 6 & 0.9 & 2.7 & \star \\ 8 & 2.2 & 2.6 & 1.6 \\ 8 & 2.7 & 5.1 & 3.6 \end{bmatrix}$$

(e) Singular values of $\boldsymbol{W}$  (f) Singular values of $\boldsymbol{A}$  (g) Singular values of $\boldsymbol{B}$  (h) Singular values of $\boldsymbol{W}_{\text{aug}}$

- **Disconnected Case: Sub-optima emerges** $\implies$ **preventing low rank; Disconnected with complete bipartite components** $\implies$ **Minimum Nuclear Norm Regularization**



$$\begin{bmatrix} 1 & \star & 3 \\ \star & 5 & \star \\ 3 & \star & 9 \end{bmatrix}$$

(b) Singular values of $\boldsymbol{A}$  (c) Singular values of $\boldsymbol{B}$  (d) Singular values of $\boldsymbol{W}_{\text{aug}}$

# Thanks!

**Mail:** bai299@sjtu.edu.cn          **WeChat:** bai299O_O          **Web:** https://ZhiweiBai.github.io/