
LINEAR STABILITY HYPOTHESIS AND RANK STRATIFICATION FOR NONLINEAR MODELS

Yaoyu Zhang^{1,4,*}, Zhongwang Zhang¹, Leyang Zhang², Zhiwei Bai¹, Tao Luo^{1,3}, Zhi-Qin John Xu^{1†}

¹ School of Mathematical Sciences, Institute of Natural Sciences, MOE-LSC
and Qing Yuan Research Institute, Shanghai Jiao Tong University

² Department of Mathematics, University of Illinois Urbana-Champaign

³ CMA-Shanghai

⁴ Shanghai Center for Brain Science and Brain-Inspired Technology

ABSTRACT

Models with nonlinear architectures/parameterizations such as deep neural networks (DNNs) are well known for their mysteriously good generalization performance at overparameterization. In this work, we tackle this mystery from a novel perspective focusing on the transition of the target recovery/fitting accuracy as a function of the training data size. We propose a rank stratification for general nonlinear models to uncover a model rank as an “effective size of parameters” for each function in the function space of the corresponding model. Moreover, we establish a linear stability theory proving that a target function almost surely becomes linearly stable when the training data size equals its model rank. Supported by our experiments, we propose a linear stability hypothesis that linearly stable functions are preferred by nonlinear training. By these results, model rank of a target function predicts a minimal training data size for its successful recovery. Specifically for the matrix factorization model and DNNs of fully-connected or convolutional architectures, our rank stratification shows that the model rank for specific target functions can be much lower than the size of model parameters. This result predicts the target recovery capability even at heavy overparameterization for these nonlinear models as demonstrated quantitatively by our experiments. Overall, our work provides a unified framework with quantitative prediction power to understand the mysterious target recovery behavior at overparameterization for general nonlinear models.

1 Introduction

How many data points are needed for a model to recover a target function is a basic yet fundamental problem for the theoretical understanding of model fitting. For example, in linear regression, a linear target function in general can be recovered when the training data size n is no less than the model parameter size m . Similarly, in band-limited signal recovery, we have the Nyquist-Shannon sampling theorem stating that a periodic signal with no higher frequency than f (with $m = 2f$ coefficients) can be exactly recovered from $n \geq 2f$ uniformly sampled points [1]. Above results indicate a phase transition of the target recovery accuracy as a function of the training data size at $n = m$. Then $n \geq m$ is often referred to as the overdetermined/underparameterized regime whereas $n < m$ is often referred to as the underdetermined/overparameterized regime. Traditional learning theory suggests that a model in the overparameterized regime is likely to overfit the data [2, 3], thus fails to explain why overparameterized deep neural networks often generalize well in practice [4, 5].

Motivated by the success of deep neural networks (DNNs), there emerges a trend to study general models with nonlinear architectures/parameterizations for target recovery, e.g., linear models with deep parameterization [6], deep matrix factorization models [7, 8, 9], deep linear networks [10, 11, 12]. It has been demonstrated that these nonlinear models are capable of recovering target functions even at heavy overparameterization. In this work, we refer to this phenomenon as the recovery mystery of nonlinear models. To understand this mystery, a popular approach is to analyze in detail the

*Corresponding author: zhyy.sjtu@sjtu.edu.cn.

†Corresponding author: xuzhiqin@sjtu.edu.cn.

training dynamics of these nonlinear models on a case-by-case basis in order to uncover the underlying implicit bias for each nonlinear model, such as the low-rank bias in deep matrix factorization models [7, 8, 9] and low-frequency bias in deep neural networks [13, 14, 15]. Following this approach, many works advance our understanding about the recovery mystery for certain nonlinear models. However, this approach encounters huge difficulty in quantitatively analyzing deep neural networks. In addition, it fails to provide a unified mechanism underlying the recovery mystery for general nonlinear models.

In this work, we take a novel approach to this recovery mystery. Specifically, we move the focus from the detailed training dynamics in previous implicit bias studies to a general macroscopic behavior—transition of the target recovery accuracy as a function of the training data size. Our study uncovers a new quasi-determined regime at overparameterization $n < m$ in which a given model is capable of recovering a given target function. This quasi-determined regime is determined by the linear stability for recovery of the target function (different from the numerical linear stability in Refs. [16, 17]). Supported by experiments, we propose a linear stability hypothesis that linearly stable interpolations are preferred by nonlinear training. Importantly, by proposing a rank stratification and establishing the linear stability theory, our work show that many long standing open problems in nonlinear model fitting reduce to the validity of this hypothesis. For example, the cause of target recovery at overparameterization, the effective size of parameters and the implicit bias for a nonlinear model, as well as the advantage of the general layer-based architecture and the superiority of the convolutional architecture specifically for neural networks. Specifically, for any nonlinear model, our rank stratification uncovers a model rank for each target function in the model function space, which quantifies the data size needed for its linear stability. Therefore, our rank stratification is a powerful tool to obtain quantitative predictions about the data size needed to recover a target function in any nonlinear model. These predictions are numerically demonstrated for matrix factorization models, two-layer tanh-NNs of a fully-connected or convolutional architecture, and remain to be demonstrated for various other models.

Our linear stability hypothesis, rank stratification and linear stability theory get inspiration from experiments, predict experiments and are supported by experiments. In Section 2, we show how we obtain the linear stability hypothesis from the experimental observation of a simple nonlinear model. To understand the condition of the linear stability, we propose a rank stratification for general nonlinear models in Section 3. Moreover, we demonstrate that the model rank of a target function obtained by rank stratification can exactly match with the transition of recovery in experiments, thus well serving as an “effective size of parameters” for this target function. In Section 4, we further establish the linear stability theory based on the rank stratification, which uncovers a quasi-determined regime at overparameterization with target recovery capability. In Section 5, we present the rank hierarchies obtained from rank stratification for NNs of different architectures. Our analysis quantifies the superiority of CNNs to fully-connected NNs for the CNN functions as further demonstrated by our experiments.

2 Linear stability hypothesis

Unlike linear regression, many nonlinear models (nonlinear in parameters) like DNNs are capable of accurately recovering certain target functions at overparameterization. As an example, we present the generalization accuracy of gradient descent training in recovering/fitting different target functions with various sample sizes for the following two models in Fig. 1. One is a simple nonlinear model $f_{\text{NL}}(\mathbf{x}; \boldsymbol{\theta}) = \theta_0 + \theta_1 x_1 + \theta_2 \theta_3 x_2$ with input $\mathbf{x} = [x_1, x_2]^T$ and parameter $\boldsymbol{\theta} = [\theta_0, \theta_1, \theta_2, \theta_3]^T$. The other is its linear counterpart $f_{\text{L}}(\mathbf{x}; \boldsymbol{\theta}) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$ with $\mathbf{x} = [x_1, x_2]^T$ and $\boldsymbol{\theta} = [\theta_0, \theta_1, \theta_2]^T$. These two models share the same model function space $\mathcal{F}_{\text{NL}} = \mathcal{F}_{\text{L}} = \{a_0 + a_1 x_1 + a_2 x_2 | a_0, a_1, a_2 \in \mathbb{R}\}$, however clearly differ in their target recovery performance. In Fig. 1(a), the linear model fails to recover all target functions with training data size less than 3 as predicted by the theory of linear regression. Surprisingly, as shown in Fig. 1(b), the nonlinear model \mathcal{F}_{NL} accurately recovers 1, x_1 and $1 + x_1$ with only 2 data points less than both its parameter size 4 and the dimension of function space 3. This experiment again confirms the long standing mystery that nonlinear models in general are capable of recovering certain target functions at overparameterization. Remark that, though model \mathcal{F}_{NL} is very simple, it is not easy to analyze its nonlinear training dynamics. In this situation, we take a novel approach to understand this recovery/generalization mystery by proposing the following question: *When is it possible to distinguish a target minimizer (based on certain local property) from infinitely many other global minimizers at overparameterization?* Note that a target minimizer is a global minimizer whose output function equals the target function.

To answer this question, we get inspiration from the following observation. With 2 training data points, model f_{NL} finds minimizers close to $[1, 0, 0, 0]^T$, $[0, 1, 0, 0]^T$, and $[1, 1, 0, 0]^T$ in fitting 1, x_1 , and $1 + x_1$, respectively, as shown in Table 1. By looking into the tangent function space $\mathcal{T}_{\boldsymbol{\theta}} = \text{span}\{\partial_{\theta_i} f(\cdot; \boldsymbol{\theta})\}_{i=1}^M = \text{span}\{1, x_1, \theta_3 x_2, \theta_2 x_2\}$ of all the global minimizers, we notice that $\mathcal{T}_{\boldsymbol{\theta}}$ is 2-d at $[1, 0, 0, 0]^T$, $[0, 1, 0, 0]^T$, and $[1, 1, 0, 0]^T$, whereas 3-d at all the other global minimizers. Remark that, given $n \geq r$ training data points, a global minimizer $\boldsymbol{\theta}^*$ with a r -d tangent function space possesses a special local property that the corresponding function $f(\cdot; \boldsymbol{\theta}^*)$ can be uniquely recovered in the

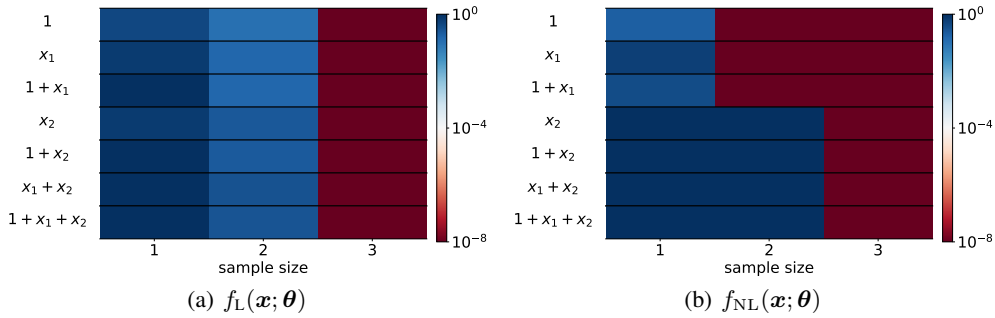


Figure 1: Average test error (color) vs. the number of samples (abscissa) for different target functions (ordinate). For all experiments, network parameters are initialized with a normal distribution with mean 0 and variance 10^{-8} , and trained with full-batch GD with a learning rate of 0.01. The training ends when the training error is less than 10^{-9} . Each test error is averaged over 100 trials with random initialization.

tangent function hyperplane $\tilde{\mathcal{T}}_{\theta^*} = f(\cdot; \theta^*) + \mathcal{T}_{\theta^*}$. For example, given 2 different training data points, target function $1 + x_1$ can be uniquely recovered in $\tilde{\mathcal{T}}_{\theta^*} = \{a_0 + a_1 x_1 | a_0, a_1 \in \mathbb{R}\}$ for $\theta^* = [1, 1, 0, 0]^T$. In this paper, we refer to this special local property as the linear stability for recovery or simply the linear stability (see Definition 1). Then the output functions of these linearly stable global minimizers are referred to as the linearly stable functions. Inspired by the above observation, we propose the following *linear stability hypothesis* for general nonlinear models that:

Linearly stable global minimizers are more likely to be selected by a nonlinear training process.

Clearly, linearly stable functions are more likely to be learned by our hypothesis. Remark that a sufficiently nonlinear training process is important for recovering a linearly stable target function in practice (an example will be shown in Fig. 3). By our hypothesis, the linear stability of a target function is the key for its successful recovery. This can be analyzed rigorously for general nonlinear models as detailed in the latter sections. In the following, we first propose a rank stratification for general models to uncover a minimal training data size needed for a target function to be linearly stable.

Table 1: Recovered parameter values with standard deviations for \mathcal{F}_{NL} over 100 trials for experiments of the first 3 rows in Fig. 1(b).

parameter	target function		
	1	x_1	$1 + x_1$
θ_0	$1.0 \times 10^0 \pm 9.8 \times 10^{-5}$	$4.4 \times 10^{-7} \pm 4.5 \times 10^{-5}$	$1.0 \times 10^0 \pm 2.7 \times 10^{-5}$
θ_1	$4.9 \times 10^{-5} \pm 3.0 \times 10^{-4}$	$1.0 \times 10^0 \pm 5.5 \times 10^{-5}$	$1.0 \times 10^0 \pm 7.1 \times 10^{-5}$
θ_2	$4.0 \times 10^{-4} \pm 1.8 \times 10^{-3}$	$1.2 \times 10^{-4} \pm 1.4 \times 10^{-3}$	$1.9 \times 10^{-4} \pm 1.4 \times 10^{-3}$
θ_3	$1.1 \times 10^{-4} \pm 1.9 \times 10^{-3}$	$6.0 \times 10^{-6} \pm 1.5 \times 10^{-3}$	$2.3 \times 10^4 \pm 1.4 \times 10^{-3}$

3 Rank stratification

For a linear model $\sum_{i=1}^m \theta_i \phi_i(\mathbf{x})$ with basis functions $\{\phi_i(\cdot)\}_{i=1}^m$, it is well known that any target function in its function space can be stably recovered when the size of training data $\{(\mathbf{x}_j, y_j)\}_{j=1}^n$ is no less than its effective size of parameters (or effective degrees of freedom) $\dim(\text{span}\{\phi_i(\cdot)\}_{i=1}^m)$. Remark that the above intuitive argument requires a mild assumption on data that $\text{rank}(\phi(\mathbf{X})) = \dim(\text{span}\{\phi_i(\cdot)\}_{i=1}^m)$, where $[\phi(\mathbf{X})]_{i,j} = \phi_i(\mathbf{x}_j)$ for $i \in [m], j \in [n]$. Therefore, for the linear model $\theta_0 + \theta_1 x_1 + \theta_2 x_2$ with 3 effective parameters, we observe target recovery at $n = 3$ as shown in Fig. 1. Above result for a linear model can be directly applied to understand the stability of recovery in the tangent function hyperplane of any parameter point for a nonlinear model. For any nonlinear model $f_{\theta} = f(\cdot; \theta)$, $\tilde{\mathcal{T}}_{\theta^*} = \{f(\cdot; \theta^*) + \mathbf{a}^T \nabla_{\theta} f(\cdot; \theta^*) | \mathbf{a} \in \mathbb{R}^M\}$ at $\theta^* \in \mathbb{R}^M$ has $\dim(\text{span}\{\partial_{\theta_i} f(\cdot; \theta^*)\}_{i=1}^M)$ effective parameters. In $\tilde{\mathcal{T}}_{\theta^*}$, $f(\cdot; \theta^*)$ in general can be stably recovered when $n \geq \dim(\text{span}\{\partial_{\theta_i} f(\cdot; \theta^*)\}_{i=1}^M)$. A special feature of many

nonlinear models is that the effective size of parameters changes over the parameter space. In this work, we formally define this effective size of parameters as the model rank of $\theta^* \in \mathbb{R}^M$ with respect to the model f_θ , i.e.,

$$\text{rank}_{f_\theta}(\theta^*) := \dim \left(\text{span} \left\{ \partial_{\theta_i} f(\cdot; \theta^*) \right\}_{i=1}^M \right). \quad (1)$$

This definition of model rank is consistent with the definition of rank in differential topology. Remark that, notation $\text{rank}(\cdot)$ without a subscript refers to the matrix rank by default in our work.

Understanding the distribution of the model rank over the parameter space is the first step for a linear stability analysis. As an example, for the nonlinear model $f_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 \theta_3 x_2$, the model rank at any point $\theta^* \in \mathbb{R}^4$ is adaptive as follows

$$\text{rank}_{f_\theta}(\theta^*) = \dim(\text{span}\{1, x_1, \theta_3^* x_2, \theta_2^* x_2\}) = \begin{cases} 2, & \theta_2^* = \theta_3^* = 0, \\ 3, & \text{others.} \end{cases} \quad (2)$$

To consider the linear stability for a target function f^* , one must note that $\mathcal{M}_{f^*} := \{\theta | f(\cdot; \theta) = f^*; \theta \in \mathbb{R}^M\}$ referred to as the target stratifold in this work is a disjoint union of manifolds with different dimensions and model ranks. For example, the target stratifold for $1 + x_1$ is $\{\theta | \theta_0 = 1, \theta_1 = 1, \theta_2 \theta_3 = 0\}$, on which rank-2 is attained only at $\theta = [1, 1, 0, 0]^T$ and rank-3 is attained elsewhere. When $n \geq 2$, target function $1 + x_1$ is stable for recovery at the tangent function hyperplane of $\theta = [1, 1, 0, 0]^T$ under a mild assumption. Then, our linear stability hypothesis predicts that $1 + x_1$ is likely to be recovered through training as numerically demonstrated in Fig. 1.

To quantify the minimal data size needed to recover a target function f^* in the model function space $\mathcal{F}_{f_\theta} := \{f(\cdot; \theta) | \theta \in \mathbb{R}^M\}$, we formally define its model rank as

$$\text{rank}_{f_\theta}(f^*) := \min_{\theta' \in \mathcal{M}_{f^*}} \text{rank}_{f_\theta}(\theta'), \quad (3)$$

with a slight misuse of the notion $\text{rank}_{f_\theta}(\cdot)$ to return the model rank for a function input. Then, the second step for a linear stability analysis is to stratify \mathcal{F}_{f_θ} into function sets of different model ranks from low to high, which forms a rank hierarchy. For the linear model $\theta_0 + \theta_1 x_1 + \theta_2 x_2$ with a constant rank, its whole function space is rank-3. For the nonlinear model $\theta_0 + \theta_1 x_1 + \theta_2 \theta_3 x_2$, the rank hierarchy is as follows,

$$\text{rank}_{f_\theta}(f^*) = \begin{cases} 2, & f^* \in \{a_0 + a_1 x_1 | a_0, a_1 \in \mathbb{R}\}, \\ 3, & f^* \in \{a_0 + a_1 x_1 + a_2 x_2 | a_2 \neq 0, a_0, a_1, a_2 \in \mathbb{R}\}. \end{cases} \quad (4)$$

This result shows that $1, x_1$ and $1 + x_1$ are rank-2, whereas all the other functions with nonzero coefficients in x_2 in Fig. 1 are rank-3. Then, our linear stability hypothesis predicts recovery with 2 data points for $1, x_1$ and $1 + x_1$, and 3 data points for other function through nonlinear training. Clearly, this prediction perfectly matches with the experimental results in Fig. 1.

In general, for a differentiable model f_θ with M parameters, the standard procedure of rank stratification is comprised of the following two steps: (1) stratify the parameter space into different model rank levels to obtain the rank hierarchy over the parameter space; (2) stratify the model function space into different model rank levels to obtain the rank hierarchy over the model function space. Remark that, the difficulty of rank stratification depends on the complexity of model architecture as shown in the following sections.

The rank stratification proposed above uncovers that different functions in the function space of a nonlinear model may have different effective sizes of parameters. Clearly, even when two models share the same model function space, different parameterization/architecture can lead to very different hierarchies as demonstrated by the above comparison between f_L and f_{NL} . By our linear stability hypothesis, this rank hierarchy indicates an implicit bias towards low model rank functions over the model function space as detailed later in Section 4. Overall, the proposed rank stratification is a powerful tool that could explicitly uncover an architecture-specific implicit bias of a nonlinear model. In the following subsection, we present the rank hierarchy obtained by the rank stratification in a table for a nonlinear matrix factorization model, and demonstrate the relation predicted by our hypothesis between the model rank of a target function and its experimental transition of the target recovery accuracy.

3.1 Matrix factorization model: rank hierarchy matches with the transition of target recovery

To demonstrate the power of the rank stratification for general nonlinear models, we consider in this section a practical nonlinear model of matrix factorization $f_\theta = \mathbf{A}\mathbf{B}$ with application in matrix completion. In Table 2, we present its rank hierarchy obtained through rank stratification (see Appendix Section A.1 for details). All elements in matrices \mathbf{A} and \mathbf{B} are trainable parameters. By Table 2, a target matrix f^* with matrix rank 1, 2, 3, or 4 possesses the model rank 7, 12, 15, or 16, respectively. It can be clearly seen from Fig. 2 that these model ranks exactly match with the transition

of the target recovery accuracy, i.e., the test error drops rapidly to almost 0 when the size of the observed entries equals the model rank of the target (marked by a yellow dashed line). This result further demonstrates the importance of rank stratification for understanding the target recovery behavior of nonlinear models, and supports the validity of the linear stability hypothesis.

model: $\mathbf{f}_\theta = \mathbf{A}\mathbf{B}$, $\theta = (\mathbf{A}, \mathbf{B})$, $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{d \times d}$	
$\text{rank}_{\mathbf{f}_\theta}(\mathbf{f}^*)$	\mathbf{f}^*
0	$\mathbf{0}_{d \times d}$
$2d - 1$	$\text{rank}(\mathbf{f}^*) = 1$
\vdots	\vdots
$2rd - r^2$	$\text{rank}(\mathbf{f}^*) = r$
\vdots	\vdots
d^2	$\text{rank}(\mathbf{f}^*) = d$

Table 2: Rank hierarchy for the matrix factorization model.

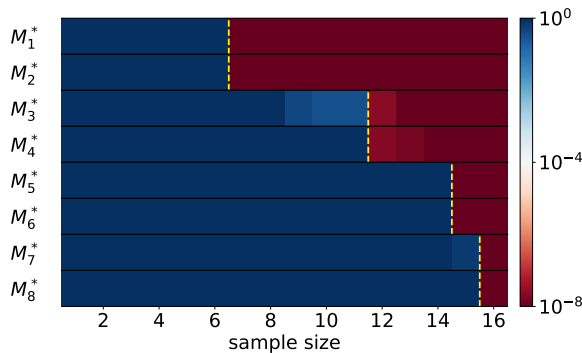


Figure 2: Average test error (color) vs. the number of samples (abscissa) for different target functions (ordinate). The yellow dashed line for each row indicates the transition when the size of the observed entries equals to the model rank of the target. Different rows represent different target matrices (see Appendix Section C for details). In particular, $\text{rank}(M_{2k-1}^*) = \text{rank}(M_{2k}^*) = k$ for $k = 1, 2, 3, 4$. For all experiments, the weights are initialized with a normal distribution with mean 0 and variance 10^{-8} , and trained with full-batch GD with a learning rate of 0.05. The training ends when the training error is less than 10^{-9} . Each test error is averaged over 50 trials with random initialization.

3.2 Quasi-determined regime

From above results, it is clear that there exists a new regime for nonlinear models at overparameterization where a target function can be successfully recovered. We name this regime the quasi-determined regime formally defined later in Definition 2. From experiments, *the quasi-determined regime covers a wide range of training data sizes from the model rank of the target function to the size of model parameters*. This adaptiveness to the target function is the key characteristic of the quasi-determined regime absent in conventional regime characterization. Remark that, the quasi-determined regime is specific to rank-adaptive models in which model ranks are non-constant over their function spaces. In a model with a constant model rank such as a linear model, the quasi-determined regime is empty because the model rank of any function in the model function space equals the dimension of the function space above which the fitting problem becomes determined/over-determined.

As shown in Fig. 3, target recovery in the quasi-determined regime is very different from that in the over-determined/underparameterized regime: (i) The success or accuracy of recovery depends on the initialization. One may need to tune the scale of initialization to a sufficiently small value, which leads to a highly nonlinear training dynamics, in order to achieve a good recovery accuracy. (ii) Increasing the size of training data above the model rank of the target function further increases the tolerance on the initialization scale and enhances the accuracy of recovery.

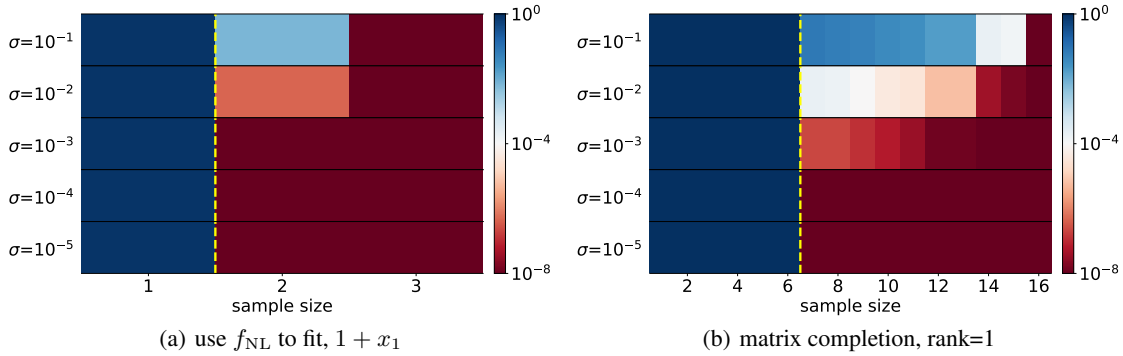


Figure 3: Average test error (color) vs. the number of samples (abscissa) over different sizes of standard deviation of initialization (ordinate). (a) Performance of fitting $1 + x$ by $f_{NL}(x, \theta)$. (b) Performance of completing a rank-1 matrix by the matrix factorization model. We use M_1^* in Fig. 2 as our target matrix, whose specific form is given in the appendix. The yellow dashed line for each row indicates the transition when the sample size equals to the model rank of the target. For all experiments, the network is initialized with a normal distribution with mean 0 and variance σ^2 , and trained with full-batch GD with a learning rate 0.05. Each test error is averaged over 50 trials with random initialization.

These two properties match with the widely-observed behavior of DNNs and other nonlinear models in practice that good hyperparameter tuning and a large training data size are two important factors for an accurate target recovery at overparameterization. Therefore, it is reasonable to believe that our quasi-determined regime is relevant to the training of general nonlinear models in practice. To understand the exact relation among our rank stratification, the quasi-determined regime and the linear stability hypothesis, we establish in the following the linear stability theory for recovery for general models.

4 Linear stability theory

In this section, we rigorously analyze the linear stability for general models to address when a function or a minimizer becomes linearly stable for recovery in the linearized function space, i.e., the tangent function hyperplane. By admitting the linear stability hypothesis, our results in the following yield quantitative predictions to the global target recovery behavior of nonlinear models. Remark that our linear stability hypothesis and analysis is inspired by the widely-used linear stability analysis in mathematics, which serves as a powerful tool to understand the first-order behavior of a nonlinear system. Also note that all the linear stability in our work refers to the linear stability for recovery in Definition 1. It is starkly different from the commonly considered numerical linear stability for neural networks originated from the numerical discretization of a continuous training dynamics [16, 17]. Our analysis begins with the following formal definition of the linear stability.

Definition 1 (linear stability for recovery). *Given any differentiable model f_θ with model function space \mathcal{F}_{f_θ} , loss function $\ell(\cdot, \cdot)$, and training data $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$,*

(i) *a parameter point $\theta^* \in \mathbb{R}^M$ is linearly stable if $f(\cdot; \theta^*)$ is the unique solution to*

$$\min_{f \in \tilde{\mathcal{T}}_{\theta^*}} \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i), y_i); \quad (5)$$

(ii) *a function $f^* \in \mathcal{F}_{f_\theta}$ is linearly stable if there exists a linearly stable parameter point θ' such that $f(\cdot; \theta') = f^*$.*

Without loss of generality, we consider $\ell(\cdot, \cdot)$ to be a continuously differentiable distance function by default and focus on studying the linear stability of the global minimizers attaining 0 loss as well as the corresponding interpolations, i.e., functions in \mathcal{F}_{f_θ} attaining 0 loss. By the linear stability hypothesis, above formal definition of the linear stability immediately gives us a novel regime with target recovery capability defined as follows.

Definition 2 (quasi-determined regime). *Using any model f_θ to fit a target function $f^* \in \mathcal{F}_{f_\theta}$ from data $S = \{(\mathbf{x}_i, f^*(\mathbf{x}_i))\}_{i=1}^n$ at overparameterization, the fitting problem is quasi-determined if f^* is linearly stable for recovery.*

In the following, we present our theory of linear stability beginning with a necessary and sufficient condition for linear stability.

Lemma 1 (linear stability condition, see Lemma 4 in Appendix for proof). *Given any differentiable model f_θ and training data $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$, a global minimizer θ^* satisfying $f(\mathbf{x}_i; \theta^*) = y_i$ for all $i \in [n]$ is linearly stable if and only if $\text{rank}_S(\theta^*) = \text{rank}_{f_\theta}(\theta^*)$, where the empirical model rank $\text{rank}_S(\theta^*) := \text{rank}(\nabla_{\theta} f(\mathbf{X}; \theta^*))$.*

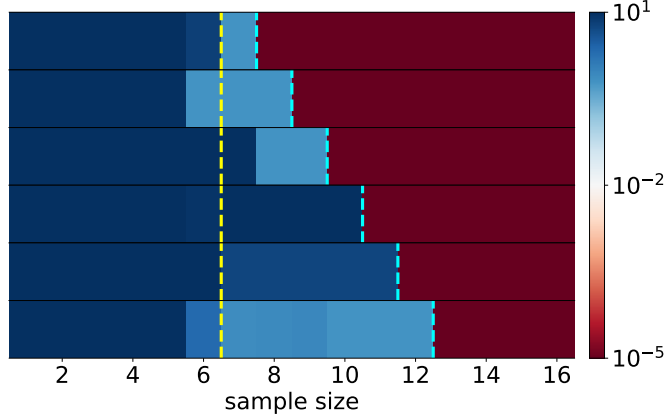


Figure 4: The relationship between the number of samples (abscissa) and the test error (color) for different sampling sequences with specifically designed orders (ordinate) to reconstruct a target matrix with matrix rank one. For a specific sampling sequence, $n_t = \min\{n | \text{rank}_{S_n}(\theta^*) = \text{rank}_{f_\theta}(\theta^*)\}$ is indicated by cyan dashed curve. The model rank of the target matrix is 7 indicated by the yellow dashed curve. For all experiments, the parameters are initialized by a normal distribution with mean 0 and variance 10^{-8} , and trained by the full-batch gradient descent with a learning rate 0.05. Each test error in the figure is averaged over 50 trials with random initialization.

Note that $\nabla_{\theta} f(\mathbf{X}; \theta^*) = [\nabla_{\theta} f(\mathbf{x}_1; \theta^*), \dots, \nabla_{\theta} f(\mathbf{x}_n; \theta^*)]$ with $\mathbf{X} := [\mathbf{x}_1, \dots, \mathbf{x}_n]$ is referred to as the empirical tangent matrix. By the above lemma, we can determine the linear stability of a target function by checking whether there exists a target minimizer satisfying this linear stability condition for the given training data. The intuition of Lemma 1 is as follows. In the tangent function hyperplane $\tilde{\mathcal{T}}_{\theta^*}$, $\text{rank}_S(\theta^*)$ quantifies the number of independent constraints from data, and $\text{rank}_{f_\theta}(\theta^*)$ quantifies the effective size of parameters. Therefore, the linearized problem Eq. (5) at a global minimizer θ^* becomes determined if and only if $\text{rank}_S(\theta^*) = \text{rank}_{f_\theta}(\theta^*)$. Lemma 1 implies cases in which a target function may not become linearly stable with $n = \text{rank}_{f_\theta}(\theta^*)$ training data points due to the lack of data independence. For these cases, our linear stability hypothesis predicts a transition of the target recovery accuracy later than $n = \text{rank}_{f_\theta}(\theta^*)$. We numerically verify this prediction by the following experiments. As shown in Fig. 2, the minimum sample size to recover a rank one matrix is 7. However, we can design a sample sequence $\{(i_1, j_1), (i_2, j_2), \dots\}$ with $S_n = \{((i_k, j_k), \mathbf{f}_{i_k j_k}^*)\}_{k=1}^n$ such that the minimum sample size to satisfy the linear stability condition in Lemma 1 is larger than 7, that is, $n_t = \min\{n | \text{rank}_{S_n}(\theta^*) = \text{rank}_{f_\theta}(\theta^*)\} > 7$. In Fig. 4, each row indicates a specially designed sample sequence with a different n_t indicated by the cyan dashed curve. Clearly, only when the linear stability condition is satisfied, i.e., sample size is no less than n_t , the test error drops rapidly to almost 0. All these experiments confirm that the experimental transition of the target recovery accuracy matches with the transition of the linear stability of the target. Again, they support the validity of our linear stability hypothesis.

Lemma 1 provides the exact condition about when a global minimizer becomes linearly stable. However, checking this condition for various global minimizers is a tedious job. In practice, it is important to have a more convenient and intuitive condition. For example, a model of m parameters can be recovered almost surely from $n \geq m$ data points for linear regression. In analogy, we find out that f^* becomes linearly stable with $n \geq \text{rank}_{f_\theta}(f^*)$ data points almost surely for a model analytic with respect to its parameters by the following theorem.

Theorem 1 (phase transition of linear stability for recovery, see Theorem 4 in Appendix for proof). *Given any analytic model f_θ , for any target function $f^* \in \mathcal{F}_{f_\theta}$ and n generic training data $S = \{(\mathbf{x}_i, f^*(\mathbf{x}_i))\}_{i=1}^n$,*

- (i) **Strictly under-determined regime:** *if $n < \text{rank}_{f_\theta}(f^*)$, then f^* is not linearly stable;*
- (ii) **Quasi-determined regime:** *if $n \geq \text{rank}_{f_\theta}(f^*)$, then f^* is linearly stable almost everywhere with respect to S .*

Above theorem proves a phase transition at $n = \text{rank}_{f_\theta}(f^*)$ in general for the linear stability of the target function f^* , which coincides with the experimentally observed phase transition in Figs. 1 and 2. By this theorem, we can conveniently refer to the fitting problem with training data sizes from $\text{rank}_{f_\theta}(f^*)$ to the parameter size M as the quasi-determined regime by default in our paper.

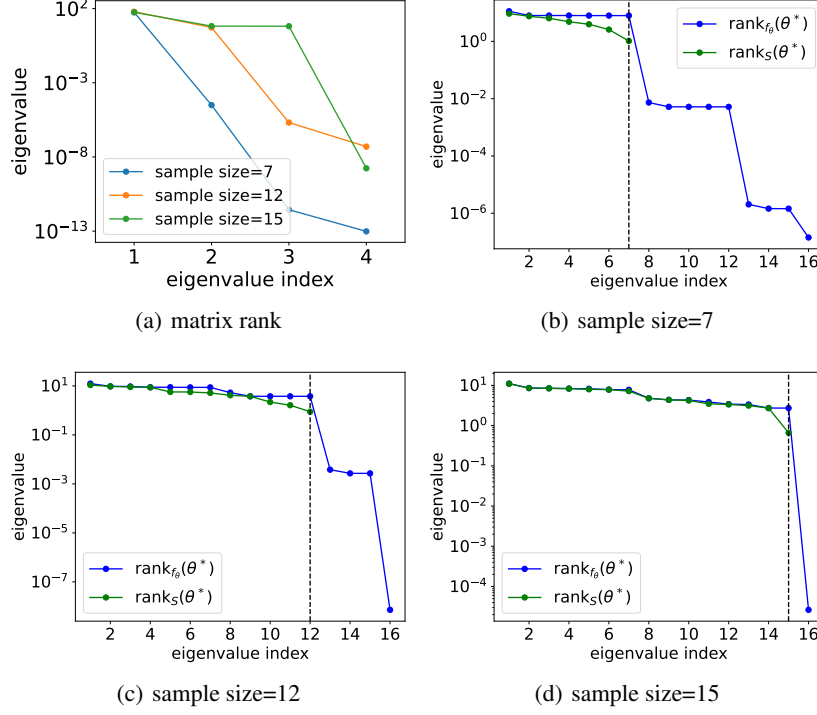


Figure 5: Learn a 4×4 full rank target matrix from different sample sizes. Given training samples of size $n = 7, 12, 15$, respectively (indicated by black dashed line in (b-d)), (a) average of the ordered eigenvalues of the recovered matrices are presented; (b,c,d) average of the ordered singular values of the empirical tangent matrix $[\partial_{\theta_s}(\mathbf{f}\theta^*)_{i_k j_k}]_{s \in [32], k \in [n]}$ vs. that of the tangent matrix $[\partial_{\theta_s}(\mathbf{f}\theta^*)_{i_k j_k}]_{s \in [32], k \in [16]}$ ($\{(i_k, j_k)\}_{k=1}^{16}$ takes all of the matrix indices) are presented. Here, θ^* is the recovered parameter vector at convergence. For all experiments, the parameters are initialized by a normal distribution with mean 0 and variance 10^{-8} , and trained by full-batch gradient descent with a learning rate 0.05. Each ordered eigenvalue in the figure is averaged over 50 ordered eigenvalues obtained from 50 trials of training with random initialization.

In general, the linear stability hypothesis indicates an implicit bias of nonlinear training towards linearly stable interpolations, which do not necessarily coincide with the target function. By Lemma 1, any interpolation with model rank higher than data size n is not linearly stable. Thus, we obtain the following corollary showing the implicit bias towards interpolations with lower model ranks by the linear stability hypothesis.

Corollary 1 (implicit bias of linear stability hypothesis, see Corollary 5 in Appendix for proof). *Given any model f_θ and training data $S = \{(x_i, y_i)\}_{i=1}^n$, if an interpolation $f' \in \mathcal{F}_{f_\theta}$ is linearly stable, then $\text{rank}_{f_\theta}(f') \leq n$.*

This bias towards interpolations with lower model ranks highlights the importance of the rank stratification in understanding the fitting behavior of a nonlinear model. When a rank hierarchy is obtained, we immediately obtain a quantitative understanding about the intrinsic preference of the nonlinear model. For example, the rank hierarchy in Table 2 shows that a target with a lower matrix rank has a lower model rank. Therefore the matrix factorization model is intrinsically biased towards a low matrix rank completion by our linear stability hypothesis. We can predict that, even in a strictly under-determined regime, a completion of the matrix with model rank no higher than the sample size is likely to be learned through nonlinear training. In the experiment shown in Fig. 5, we find that completions of matrix rank 1, 2 and 3 can be reliably learned from 7, 12 and 15 samples, respectively, given a full rank target matrix of size 4. In Fig. 5(b-d), we further observe that, at these learned parameter points, the empirical model rank equals the model rank. Therefore, despite the failure of recovering the full rank target matrix in these experiments, our linear stability hypothesis and its implicit bias Corollary 1 successfully predict the training behavior of the matrix factorization model.

5 Rank stratification for deep neural networks

By admitting the linear stability hypothesis and establishing the linear stability theory above, the rank hierarchy obtained by rank stratification becomes the key to understand the target recovery performance of a nonlinear model. In this section, we present our rank stratification results for DNNs with numerical demonstration of our theoretical predictions. Our results provide quantitative understandings to the following long-standing open problems: (i) the capability of target recovery for DNNs at overparameterization; (ii) the specialty or superiority of the DNN model in comparison to other models. (iii) the impact of architecture on the target recovery performance of DNNs. In the following, we first present the rank hierarchies for two-layer tanh-NNs with fully-connected or convolutional architectures and their numerical studies. Later, we provide a rank upper bound estimate with a partial rank hierarchy for general deep NNs. In the main text, we directly present the results of rank stratification and focus on their implications. All the theoretical details including the rank stratification, supporting theorems and proofs can be found in Appendix Section A.2.

5.1 Rank stratification for two-layer fully-connected NNs

In Table 3, we present the rank hierarchy for a two-layer fully-connected tanh-NN with m hidden neurons (see Appendix Section A.2.3 for details). Note that similar rank hierarchies can be obtained for two-layer NNs with other architectures and other common activation functions. From Table 3, it is clear that two-layer fully-connected tanh-NNs are rank-adaptive, i.e., different functions occupy different model rank levels, indicating that they are capable of recovering certain target functions at overparameterization. In addition, they possess the following special property—the model expressiveness can be (almost) arbitrarily increased without changing the existing rank hierarchy. For example, when the width of hidden layer increases from m to $m' > m$, the model function space is expanded while the function sets at all model rank levels no greater than $m(d+1)$ remaining unchanged. This is a profound property for a nonlinear model in the following sense. For conventional models, the improvement of expressiveness is in general at a cost of damaging the fitting performance over the original model function space. In linear regression or random feature models, adding a new independent variable or basis function improves the model expressiveness. However, one more data point is needed to recover all functions in the original model function space. Therefore, we always have a hard time to trade off between the model expressiveness and the data size needed for target recovery. However, for the two-layer tanh-NNs, the model rank as an effective size of parameters of any function in the model function space never grows no matter the increase of width m . We name this property the *free expressiveness property*. Remark that, this good property is also possessed by NNs with linear or polynomial activations, but their expressiveness cannot be improved to the extent of universal approximation. The practical implication of our result is that, when a two-layer tanh-NN is used for fitting, we do not need to trade the model expressiveness for a good fitting performance. One can simply use a wide NN with sufficient expressiveness even to fit relatively simple target functions without worrying about the generalization performance.

model: $f_{\theta}(\mathbf{x}) = \sum_{i=1}^m a_i \tanh(\mathbf{w}_i^T \mathbf{x}), \mathbf{x} \in \mathbb{R}^d, \theta = (a_i, \mathbf{w}_i)_{i=1}^m$	
$\text{rank}_{f_{\theta}}(f^*)$	f^*
0	0
$d+1$	$\mathcal{F}_1^{\text{NN}} \setminus \{0\} : \{a_1^* \sigma(\mathbf{w}_1^{*T} \mathbf{x}) a_1^* \neq 0, \mathbf{w}_1^* \neq \mathbf{0}\}$
\vdots	\vdots
$k(d+1)$	$\mathcal{F}_k^{\text{NN}} \setminus \mathcal{F}_{k-1}^{\text{NN}} : \{\sum_{i=1}^k a_i^* \sigma(\mathbf{w}_i^{*T} \mathbf{x}) a_i^* \neq 0, \mathbf{w}_i^* \neq \mathbf{0}, \mathbf{w}_i^* \neq \pm \mathbf{w}_j^* \text{ for any } i \neq j\}$
\vdots	\vdots
$m(d+1)$	$\mathcal{F}_m^{\text{NN}} \setminus \mathcal{F}_{m-1}^{\text{NN}} : \{\sum_{i=1}^m a_i^* \sigma(\mathbf{w}_i^{*T} \mathbf{x}) a_i^* \neq 0, \mathbf{w}_i^* \neq \mathbf{0}, \mathbf{w}_i^* \neq \pm \mathbf{w}_j^* \text{ for any } i \neq j\}$

Table 3: The rank hierarchy for two-layer fully-connected width- m tanh-NN.

5.2 Rank stratification for two-layer CNNs

In Table 4, a rank hierarchy is presented in the first two columns for the following simple two-layer tanh-CNN with weight sharing

$$f_{\theta}(\mathbf{x}) = \sum_{l=1}^m \sum_{i=1}^{d+1-s} a_{li} \tanh \left(\sum_{\alpha=1}^s x_{i+s-\alpha} K_{l;\alpha} \right), \quad \mathbf{x} = [x_1, \dots, x_d]^T \in \mathbb{R}^d. \quad (6)$$

In addition, to uncover the impact of model architecture on the rank hierarchy, we also estimate the model rank of functions in each function set listed in the first column for the corresponding CNN without weight sharing as well as the corresponding fully-connected NN. All the theoretical details as well as the results for CNNs with 2D image inputs can be found in Appendix Section A.2.4 and A.2.5. To illustrate the result in Table 4, we consider a target function generated by a two-layer width- k tanh-CNN with kernel size 3×3 and stride 1 on the MNIST dataset. Without loss of generality, we consider the case without ineffective neurons (whose output weight $a_{li} = 0$), i.e., $m_{\text{null}} = 0$. Given any $m \geq k$, the model rank of this target function is $685k$ in an m -kernel CNN, $6760k$ in an m -kernel CNN without weight sharing, and $530660k$ in a width- $26m$ fully-connected NN. It is clear that all these NNs possess the free expressiveness property, i.e., the model rank of this target function does not depend on m . However, comparing different architectures, the model rank of this target function in a CNN is almost three orders of magnitude less than that in a fully-connected NN. By our linear stability hypothesis and theory, this target function requires $\sim 10^3$ times more training data to be recovered by a fully-connect NN than by a CNN. Clearly, regarding the recovery of this target function, the CNN architecture is superior to the CNN architecture without weight sharing, and is far superior to the fully-connect architecture. Note that, the major gap of model rank is between the CNN without weight sharing ($6760k$) and the fully-connect NN ($530660k$). This two orders of magnitude gap of model rank highlights the importance of removing all the unnecessary connections in the design of an NN architecture.

f^*	CNN	CNN without weight sharing	Fully-connected NN
0	0	0	0
$\mathcal{F}_1^{\text{CNN}} \setminus \{0\}$	$d+1$	$(s+1)(d+1-s) - sm_{\text{null}}$	$(d+1)(d+1-s) - dm_{\text{null}}$
\vdots	\vdots	\vdots	\vdots
$\mathcal{F}_k^{\text{CNN}} \setminus \mathcal{F}_{k-1}^{\text{CNN}}$	$k(d+1)$	$k(s+1)(d+1-s) - sm_{\text{null}}$	$k(d+1)(d+1-s) - dm_{\text{null}}$
\vdots	\vdots	\vdots	\vdots
$\mathcal{F}_m^{\text{CNN}} \setminus \mathcal{F}_{m-1}^{\text{CNN}}$	$m(d+1)$	$m(s+1)(d+1-s) - sm_{\text{null}}$	$m(d+1)(d+1-s) - dm_{\text{null}}$

Table 4: The rank hierarchy for two-layer tanh-CNN with weight sharing in Eq. (6). For functions in each function set over the rank hierarchy, we also present their model rank in the corresponding CNN without weight sharing and the corresponding fully-connected NN. Here $m_{\text{null}} = |\{a_{li} | a_{li} = 0\}|$ is a variable counting the number of ineffective neurons in the target function. Note that, when a bias term is added for each hidden neuron (shared or not according to the architecture), the model rank of a function in $\mathcal{F}_k^{\text{CNN}} \setminus \mathcal{F}_{k-1}^{\text{CNN}}$ is $k(d+2)$, $k(s+2)(d+1-s) - sm_{\text{null}}$ and $k(d+2)(d+1-s) - sm_{\text{null}}$, respectively, for these three architectures.

5.3 Experimental demonstration of the linear stability hypothesis in two-layer NNs

In Fig. 6, we perform experiments to examine our linear stability hypothesis in two-layer tanh-NNs with different architectures. Specifically, we consider the following target function in our experiments:

$$f^*(\mathbf{x}) = \mathbf{W}^{*[2]} \tanh(\mathbf{W}^{*[1]} \mathbf{x}), \quad (7)$$

where $\mathbf{W}^{*[2]} = [1, 1, 1]$,

$$\mathbf{W}^{*[1]} = \begin{bmatrix} 0.6 & 0.8 & 1 & 0 & 0 \\ 0 & 0.6 & 0.8 & 1 & 0 \\ 0 & 0 & 0.6 & 0.8 & 1 \end{bmatrix}.$$

For the training dataset and the test dataset, we sample the input data from a standard normal distribution and obtain the output values through the target function. We use two-layer tanh-NNs (with a bias term for each hidden neuron) of various architectures and various kernels/widths to fit randomly sampled training datasets of various sizes from 1 to 63.

Note that, in a 1-kernel CNN with or without weight sharing or a width-3 fully-connected NN (labeled as 1x in Fig. 6(b-d)), the model rank of the target function equals the size of model parameters. In this situation, the CNN enables a significantly earlier transition of the target recovery accuracy than the other architectures as shown in Fig. 6(a). However, This result is trivial in the sense that all the recoveries happen at the conventional over-determined/underparameterized regime. In Fig. 6(b-d), we increase the kernels/widths of NNs by N times labeled by Nx for each architecture. Specifically for $N = 100$, the sizes of model parameters become 700, 1500 and 2100 respectively. The rank hierarchy in Table 4 gives rise to a constant model rank of 7, 15 and 21 (indicated by yellow dashed lines), respectively, regardless of the choice of N . In Fig. 6(b-d), we observe delayed transitions of the target recovery accuracy for $N > 1$, i.e., the test error drops to almost 0 at a sample size later than the model rank. However, it is easy to notice that the observed transition is far closer to the model rank than to the size of model parameters especially when N is large. We remark that various factors could contribute to a delayed transition of recovery in practice such as a suboptimal tuning of hyperparameters. In practice, it remains an important open problem to find an optimal training method and hyperparameters for NNs to enable a recovery of a target function as close as possible to its model rank.

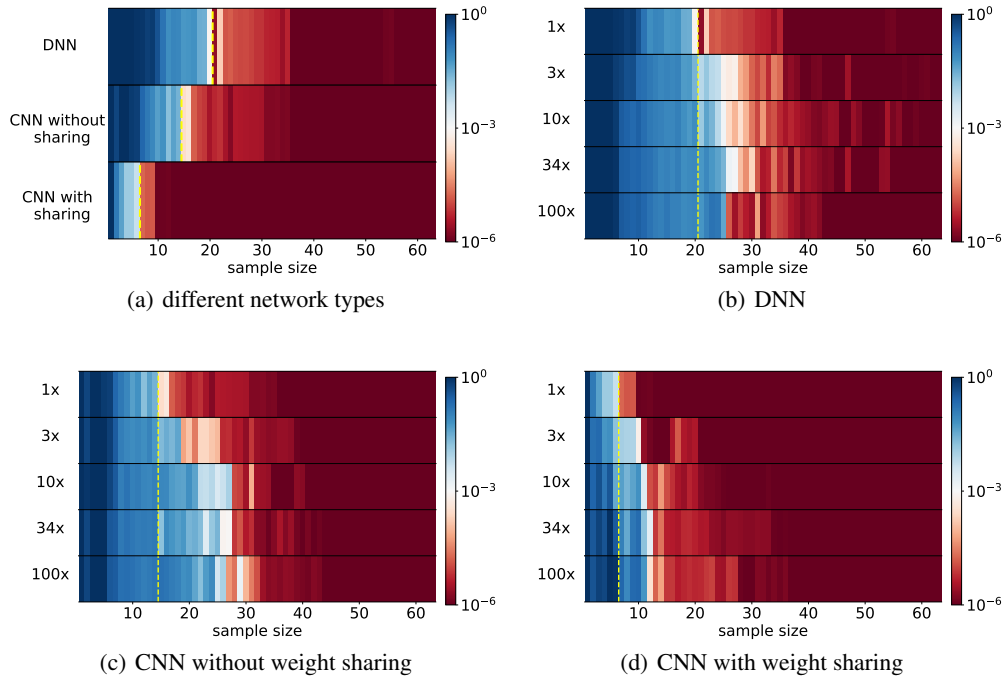


Figure 6: Average test error (color) for NNs of different architectures (ordinate) and sample sizes (abscissa) in fitting the target function Eq. (7). The yellow dashed line for each row indicates the model rank of the target in the corresponding NN. (a) Two-layer 1-kernel tanh-CNN vs. two-layer 1-kernel tanh-CNN without weight sharing vs. two-layer width-3 fully-connected tanh-NN. Note that these NNs are referred to as 1x for each architecture in (b-d). (b) Two-layer N -kernel tanh-CNN, (c) two-layer N -kernel tanh-CNN without weight sharing, and (d) two-layer width- $3N$ fully-connected tanh-NN labeled by Nx for $N = 1, 3, 10, 34, 100$. For all experiments, network parameters are initialized by a normal distribution with mean 0 and variance 10^{-20} , and trained by full-batch gradient descent with a fine-tuned learning rate.

5.4 Rank upper bound estimate for general deep NNs

For a general DNN model, its rank stratification is difficult. The model rank estimate of any given target function f^* in a DNN function space requires solving the following two challenging problems: (i) identifying the target stratifold \mathcal{M}_{f^*} in the parameter space; (ii) finding the minimal model rank over the target stratifold. With the help of the previously proposed critical embedding operators [18, 19], we can obtain a partial target stratifold and a rank upper bound estimate for general deep NNs shown in the following theorem (See Appendix Section A.2.1 for details).

Theorem 2 (rank upper bound estimate for DNNs, see Theorem 3 in Appendix for proof). *Given any NN with M_{wide} parameters, for any function $f^* \in \mathcal{F}_{\theta_{\text{narr}}}$ in the function space of a narrower NN with M_{narr} parameters, we have $\text{rank}_{f_{\theta_{\text{wide}}}}(f^*) \leq \text{rank}_{f_{\theta_{\text{narr}}}}(f^*) \leq M_{\text{narr}}$.*

Here narrower means no larger width in each hidden layer. Applying this theorem to a depth- L width- $\{m_i\}_{i=0}^L$ DNN, we obtain a partial rank hierarchy illustrated in Table 5. This partial rank hierarchy clearly shows that DNNs are rank-adaptive in general. Specifically, even in a very large DNN with M_{wide} parameters, there always exist families of functions with model ranks far less than M_{wide} , indicating the target recovery capability at heavy overparameterization. Importantly, Theorem 3 extends the free expressiveness property observed above for two-layer tanh-NNs to general DNNs. By our linear stability hypothesis and theory, this result indicates that one can simply use a wide NN with sufficient expressiveness without worrying about significant deterioration of the target recovery performance.

model: $f_{\theta}(x) = \mathbf{W}^{[L]}\sigma(\dots\sigma(\mathbf{W}^{[1]}x)\dots)$, $\mathbf{W}^{[l]} = \mathbb{R}^{m_l \times m_{l-1}}$, $m_L = 1$, $m_0 = d$	
f^*	upper bound of model rank
$\mathcal{F}_{\{1,1,\dots,1\}}$	$d + L - 1$
\vdots	\vdots
$\mathcal{F}_{\{m'_i\}_{i=1}^{L-1}}$, $1 \leq m'_i \leq m_i$	$dm'_1 + m'_1 m'_2 + \dots + m'_{L-2} m'_{L-1} + m'_{L-1}$
\vdots	\vdots
$\mathcal{F}_{\{m_i\}_{i=1}^{L-1}}$	$dm_1 + m_1 m_2 + \dots + m_{L-2} m_{L-1} + m_{L-1}$

Table 5: Partial rank hierarchy for general deep fully-connected NNs. $\mathcal{F}_{\{m_i\}_{i=1}^{L-1}}$ denotes the function space of an L -layer DNN with width- $\{m_i\}_{i=1}^{L-1}$ for hidden layers. For simplicity, we consider DNNs without bias terms.

6 Conclusions and discussion

In this work, we establish a framework to analyze quantitatively the mysterious target recovery behavior at overparameterization for general nonlinear models as illustrated in Fig. 7. We apply this framework to the matrix factorization model, two-layer tanh-NNs with a fully-connected or convolutional architecture, and successfully predict their target recovery behaviors even at heavy overparameterization. Remark that our framework relies on a linear stability hypothesis, which needs to be further verified. If this hypothesis is later systematically verified in experiments or even get proved theoretically in certain sense, the following five long standing open problems can be answered quantitatively as follows. Three problems are for general nonlinear models and two problems are specifically for DNNs.

- (1) *The cause of the target recovery at overparameterization for certain nonlinear models:* a rank-adaptive architecture/parameterization for the nonlinear model.
- (2) *The effective size of parameters for a nonlinear model:* the model rank quantifies the effective size of parameters for each function in the model function space. Remark that this problem had been proposed by Leo Breiman specifically for NNs almost three decades ago [4].
- (3) *The implicit bias of a nonlinear model through nonlinear training:* towards lower model rank interpolations.
- (4) *The advantage of the general layer-based architecture of neural networks:* free expressiveness, i.e., expressiveness can be arbitrarily improved through widening with almost no deterioration of the fitting performance.
- (5) *The superiority of CNNs to fully-connected NNs:* functions in the CNN function space in general possess far lower model ranks in CNNs than in the corresponding fully-connected NNs.

In certain sense, our theoretical framework reduces all these important problems to the validity of the linear stability hypothesis. Apart from the evidences provided in this work, the condensation phenomenon [20] during the nonlinear training of DNNs provides a rationale empirical evidence to support the linear stability hypothesis. For a two-layer ReLU NN, the condensation happens when input weights of hidden neurons (the input weights of a hidden neuron consist of all the weights from its input layer and its bias term) condense on isolated orientations. The rank of a condensed ReLU network is independent of, and far less than, the number of network parameters. For two-layer infinite width ReLU networks, [20] show that the condensation is a common feature of training networks in the nonlinear regime of the phase diagram (small initialization regime) for both synthetic and real datasets. Similar observations are made for three-layer ReLU NNs [21] and networks with different activation functions [22]. In addition, large learning rate [23] and dropout [24] can facilitate the condensation. Several works provides some preliminary theoretical support for different activations in the initial training stage [25, 26, 21]. Therefore, the condensation phenomenon suggests that nonlinear training of neural networks prefers low rank minimizers, which is consistent with the linear stability

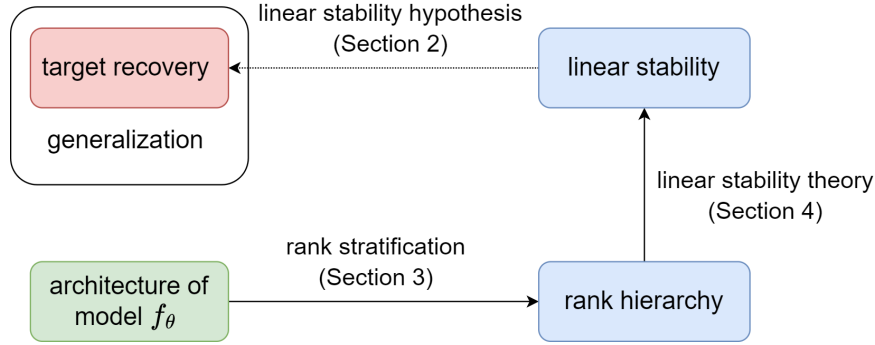


Figure 7: Illustration of the theoretical framework established in our work linking the architecture of a nonlinear model to its quantitative target recovery behavior at overparameterization.

hypothesis. In future works, we will look into details of this hypothesis, e.g., its requirement on the training dynamics for different models, through both experimental and theoretical means.

Acknowledgement

This work is sponsored by the National Key R&D Program of China Grant No. 2019YFA0709503 (Z. X.), the Shanghai Sailing Program, the Natural Science Foundation of Shanghai Grant No. 20ZR1429000 (Z. X.), the National Natural Science Foundation of China Grant No. 62002221 (Z. X.), the National Natural Science Foundation of China Grant No. 12101401 (T. L.), Shanghai Municipal Science and Technology Key Project No. 22JC1401500 (T. L.), the National Natural Science Foundation of China Grant No. 12101402 (Y. Z.), Shanghai Municipal of Science and Technology Project Grant No. 20JC1419500 (Y.Z.), the Lingang Laboratory Grant No.LG-QS-202202-08 (Y.Z.), Shanghai Municipal of Science and Technology Major Project No. 2021SHZDZX0102, and the HPC of School of Mathematical Sciences and the Student Innovation Center at Shanghai Jiao Tong University.

References

- [1] Claude E Shannon. Communication in the presence of noise. *Proceedings of the IEEE*, 72(9):1192–1201, 1984.
- [2] Vladimir Naumovich Vapnik. Adaptive and learning systems for signal processing communications, and control. *Statistical learning theory*, 1998.
- [3] Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- [4] Leo Breiman. Reflections after refereeing papers for nips. *The Mathematics of Generalization*, XX:11–15, 1995.
- [5] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. In *5th International Conference on Learning Representations, ICLR 2017*.
- [6] Blake Woodworth, Suriya Gunasekar, Jason D Lee, Edward Moroshko, Pedro Savarese, Itay Golan, Daniel Soudry, and Nathan Srebro. Kernel and rich regimes in overparametrized models. In *Conference on Learning Theory*, pages 3635–3673. PMLR, 2020.
- [7] Suriya Gunasekar, Blake E Woodworth, Srinadh Bhojanapalli, Behnam Neyshabur, and Nati Srebro. Implicit regularization in matrix factorization. *Advances in Neural Information Processing Systems*, 30, 2017.
- [8] Yuanzhi Li, Tengyu Ma, and Hongyang Zhang. Algorithmic regularization in over-parameterized matrix sensing and neural networks with quadratic activations. In *Conference On Learning Theory*, pages 2–47. PMLR, 2018.
- [9] Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix factorization. *Advances in Neural Information Processing Systems*, 32, 2019.
- [10] Andrew M Saxe, James L McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *arXiv preprint arXiv:1312.6120*, 2013.
- [11] Andrew K. Lampinen and Surya Ganguli. An analytic theory of generalization dynamics and transfer learning in deep linear networks. In *International Conference on Learning Representations*, 2019.

- [12] Suriya Gunasekar, Jason D Lee, Daniel Soudry, and Nati Srebro. Implicit bias of gradient descent on linear convolutional networks. *Advances in Neural Information Processing Systems*, 31, 2018.
- [13] Zhi-Qin J Xu, Yaoyu Zhang, and Yanyang Xiao. Training behavior of deep neural network in frequency domain. *International Conference on Neural Information Processing*, pages 264–274, 2019.
- [14] Zhi-Qin John Xu, Yaoyu Zhang, Tao Luo, Yanyang Xiao, and Zheng Ma. Frequency principle: Fourier analysis sheds light on deep neural networks. *Communications in Computational Physics*, 28(5):1746–1767, 2020.
- [15] Yaoyu Zhang, Tao Luo, Zheng Ma, and Zhi-Qin John Xu. A linear frequency principle model to understand the absence of overfitting in neural networks. *Chinese Physics Letters*, 38(3):038701, 2021.
- [16] Lei Wu, Chao Ma, and Weinan E. How sgd selects the global minima in over-parameterized learning: A dynamical stability perspective. *Advances in Neural Information Processing Systems*, 31, 2018.
- [17] Rotem Mulayoff, Tomer Michaeli, and Daniel Soudry. The implicit bias of minima stability: A view from function space. *Advances in Neural Information Processing Systems*, 34:17749–17761, 2021.
- [18] Yaoyu Zhang, Zhongwang Zhang, Tao Luo, and Zhi-Qin John Xu. Embedding principle of loss landscape of deep neural networks. *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, 2021.
- [19] Yaoyu Zhang, Yuqing Li, Zhongwang Zhang, Tao Luo, and Zhi-Qin John Xu. Embedding principle: a hierarchical structure of loss landscape of deep neural networks. *Journal of Machine Learning vol*, 1:1–45, 2022.
- [20] Tao Luo, Zhi-Qin John Xu, Zheng Ma, and Yaoyu Zhang. Phase diagram for two-layer relu neural networks at infinite-width limit. *J. Mach. Learn. Res.*, 22:71–1, 2021.
- [21] Hanxu Zhou, Qixuan Zhou, Tao Luo, Yaoyu Zhang, and Zhi-Qin John Xu. Towards understanding the condensation of neural networks at initial training. *Advances in Neural Information Processing Systems*, 2022.
- [22] Hanxu Zhou, Qixuan Zhou, Zhenyuan Jin, Tao Luo, Yaoyu Zhang, and Zhi-Qin John Xu. Empirical phase diagram for three-layer neural networks with infinite width. *Advances in Neural Information Processing Systems*, 2022.
- [23] Maksym Andriushchenko, Aditya Varre, Loucas Pillaud-Vivien, and Nicolas Flammarion. Sgd with large step sizes learns sparse features. *arXiv preprint arXiv:2210.05337*, 2022.
- [24] Zhongwang Zhang and Zhi-Qin John Xu. Implicit regularization of dropout. *arXiv preprint arXiv:2207.05952*, 2022.
- [25] Hartmut Maennel, Olivier Bousquet, and Sylvain Gelly. Gradient descent quantizes relu network features. *arXiv preprint arXiv:1803.08367*, 2018.
- [26] Franco Pellegrini and Giulio Biroli. An analytic theory of shallow networks dynamics for hinge loss classification. *Advances in Neural Information Processing Systems*, 33, 2020.
- [27] Kenji Fukumizu, Shoichiro Yamaguchi, Yoh-ichi Mototake, and Mirai Tanaka. Semi-flat minima and saddle points by embedding neural networks to overparameterization. *Advances in Neural Information Processing Systems*, 32:13868–13876, 2019.
- [28] Berfin Simsek, François Ged, Arthur Jacot, Francesco Spadaro, Clement Hongler, Wulfram Gerstner, and Johanni Brea. Geometry of the loss landscape in overparameterized neural networks: Symmetries and invariances. In *Proceedings of the 38th International Conference on Machine Learning*, pages 9722–9732. PMLR, 2021.
- [29] Zhiwei Bai, Tao Luo, Zhi-Qin John Xu, and Yaoyu Zhang. Embedding principle in depth for the loss landscape analysis of deep neural networks. *arXiv preprint arXiv:2205.13283*, 2022.

A Details of rank stratification

Definition 3 (model rank). *Given any differentiable (in parameters) model f_{θ} , the model rank for any $\theta^* \in \mathbb{R}^M$ is defined as*

$$\text{rank}_{f_{\theta}}(\theta^*) := \dim \left(\text{span} \left\{ \partial_{\theta_i} f(\cdot; \theta^*) \right\}_{i=1}^M \right), \quad (8)$$

where $\text{span} \left\{ \phi_i(\cdot) \right\}_{i=1}^M = \left\{ \sum_{i=1}^M a_i \phi_i(\cdot) \mid a_1, \dots, a_M \in \mathbb{R} \right\}$ and $\dim(\cdot)$ returns the dimension of a linear function space. Then the model rank for any function $f^* \in \mathcal{F}_{f_{\theta}}$ with model function space $\mathcal{F}_{f_{\theta}} := \{f(\cdot; \theta) \mid \theta \in \mathbb{R}^M\}$ is defined as

$$\text{rank}_{f_{\theta}}(f^*) := \min_{\theta' \in \mathcal{M}_{f^*}} \text{rank}_{f_{\theta}}(\theta'), \quad (9)$$

where the target stratifold $\mathcal{M}_{f^*} := \{\theta \mid f(\cdot; \theta) = f^*; \theta \in \mathbb{R}^M\}$.

Given a differentiable model f_θ , the standard procedure of rank stratification is comprised of the following two steps:
Step 1: Stratify the parameter space into different model rank levels to obtain the rank hierarchy over the parameter space;

Step 2: Stratify the model function space into different model rank levels to obtain the rank hierarchy over the model function space.

The difficulty of rank stratification depends on the complexity of model architecture. For example, this standard two-step rank stratification is straight-forward for $f_{\text{NL}} = \theta_0 + \theta_1 x_1 + \theta_2 \theta_3 x_2$ as illustrated in the main text. For a matrix factorization model, a linear algebra lemma is needed for its rank stratification. For DNN models, rank stratification is in general difficult. In our work, with the help of the previously discovered embedding principle and critical embedding operators, we obtain a complete rank hierarchy for two-layer tanh-NNs and a partial rank hierarchy for general multi-layer DNNs.

A.1 Matrix factorization

A.1.1 Theoretical preparation

Lemma 2 (linear algebra lemma). *Let \mathbf{A} and \mathbf{B} be two $(d \times d)$ matrices and $r_{\mathbf{A}} := \text{rank}(\mathbf{A}), r_{\mathbf{B}} := \text{rank}(\mathbf{B})$,*

$$\mathbf{\Gamma} = \begin{bmatrix} \mathbf{I} \otimes \mathbf{B} \\ \mathbf{A}^T \otimes \mathbf{I} \end{bmatrix},$$

where \mathbf{I} is the $(d \times d)$ identity matrix and \otimes is the Kronecker product. Then $\text{rank}(\mathbf{\Gamma}) = d^2 - (d - r_{\mathbf{A}})(d - r_{\mathbf{B}})$.

Proof. In order to compute the rank of $\mathbf{\Gamma}$, we consider the dimension of the null space $N(\mathbf{\Gamma})$ of $\mathbf{\Gamma}$ due to the relationship

$$\text{rank}(\mathbf{\Gamma}) + \dim(N(\mathbf{\Gamma})) = d^2.$$

We will show that $\dim(N(\mathbf{\Gamma})) = (d - r_{\mathbf{A}})(d - r_{\mathbf{B}})$, thus $\text{rank}(\mathbf{\Gamma}) = d^2 - (d - r_{\mathbf{A}})(d - r_{\mathbf{B}})$, as desired.

Let $n_{\mathbf{A}} = d - r_{\mathbf{A}}, n_{\mathbf{B}} = d - r_{\mathbf{B}}$, and suppose that $N(\mathbf{A}^T) = \text{span}\{\alpha_1, \alpha_2, \dots, \alpha_{n_{\mathbf{A}}}\}$, $N(\mathbf{B}) = \text{span}\{\beta_1, \beta_2, \dots, \beta_{n_{\mathbf{B}}}\}$ are the null spaces of \mathbf{A}^T and \mathbf{B} , respectively. Since for any $1 \leq i \leq n_{\mathbf{A}}, 1 \leq j \leq n_{\mathbf{B}}$

$$\begin{bmatrix} \mathbf{I} \otimes \mathbf{B} \\ \mathbf{A}^T \otimes \mathbf{I} \end{bmatrix} [\alpha_i \otimes \beta_j] = \begin{bmatrix} \mathbf{I} \alpha_i \otimes \mathbf{B} \beta_j \\ \mathbf{A}^T \alpha_i \otimes \mathbf{I} \beta_j \end{bmatrix} = \begin{bmatrix} \alpha_i \otimes \mathbf{0} \\ \mathbf{0} \otimes \beta_j \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix},$$

we have $N(\mathbf{A}^T) \otimes N(\mathbf{B}) \subseteq N(\mathbf{\Gamma})$.

On the other hand, let $\mathbf{0} \neq \mathbf{x} \in \mathbb{R}^{d^2}$ be in the null space $N(\mathbf{\Gamma})$ of $\mathbf{\Gamma}$, i.e. $\mathbf{\Gamma} \mathbf{x} = \mathbf{0}$. We have

$$\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix} = \begin{bmatrix} \mathbf{I} \otimes \mathbf{B} \\ \mathbf{A}^T \otimes \mathbf{I} \end{bmatrix} \mathbf{x} = \begin{bmatrix} (\mathbf{I} \otimes \mathbf{B}) \mathbf{x} \\ (\mathbf{A}^T \otimes \mathbf{I}) \mathbf{x} \end{bmatrix} = \begin{bmatrix} \text{vec}(\mathbf{B} \mathbf{X} \mathbf{I}^T) \\ \text{vec}(\mathbf{I} \mathbf{X} \mathbf{A}) \end{bmatrix} = \begin{bmatrix} \text{vec}(\mathbf{B} \mathbf{X}) \\ \text{vec}(\mathbf{X} \mathbf{A}) \end{bmatrix},$$

where \mathbf{X} is the inverse of the vectorization operator (formed by reshaping the vector $\mathbf{x} = [x_1, x_2, \dots, x_{d^2}]^T$), namely,

$$\mathbf{X} = \begin{bmatrix} x_1 & x_{d+1} & \cdots & x_{(d-1)d+1} \\ x_2 & x_{d+2} & \cdots & x_{(d-1)d+2} \\ \vdots & \vdots & \ddots & \vdots \\ x_d & x_{d+d} & \cdots & x_{d^2} \end{bmatrix}.$$

Therefore, we conclude $\mathbf{B} \mathbf{X} = \mathbf{0}$ and $\mathbf{A}^T \mathbf{X}^T = \mathbf{0}$. Note that the first equation $\mathbf{B} \mathbf{X} = \mathbf{0}$ implies each column of \mathbf{X} is a linear combination of $\{\beta_1, \beta_2, \dots, \beta_{n_{\mathbf{B}}}\}$. Thus, there exists $\mathbf{C}_{n_{\mathbf{B}} \times d}$ such that

$$\mathbf{X} = [\beta_1, \beta_2, \dots, \beta_{n_{\mathbf{B}}}] \mathbf{C}.$$

Since $\{\beta_1, \beta_2, \dots, \beta_{n_{\mathbf{B}}}\}$ is linearly independent, the second equation $\mathbf{A}^T \mathbf{X}^T = \mathbf{0}$ implies $\mathbf{A}^T \mathbf{C}^T = \mathbf{0}$. Thus, the i -th row \mathbf{C}_i of matrix \mathbf{C} satisfies $\mathbf{C}_i \in N(\mathbf{A}^T)$ for any $i \in [n_{\mathbf{B}}]$. By re-vectorizing $\mathbf{x} = \text{vec}(\mathbf{X}) = [x_1, x_2, \dots, x_{d^2}]^T$, we have

$$\mathbf{x} = \mathbf{C}_1 \otimes \beta_1 + \mathbf{C}_2 \otimes \beta_2 + \cdots + \mathbf{C}_{n_{\mathbf{B}}} \otimes \beta_{n_{\mathbf{B}}}.$$

Therefore, we conclude that $\mathbf{x} \in N(\mathbf{A}^T) \otimes N(\mathbf{B})$ and $N(\mathbf{\Gamma}) \subseteq N(\mathbf{A}^T) \otimes N(\mathbf{B})$.

Now we have $N(\mathbf{\Gamma}) = N(\mathbf{A}^T) \otimes N(\mathbf{B})$, by which

$$\dim(N(\mathbf{\Gamma})) = \dim(N(\mathbf{A}) \otimes N(\mathbf{B})) = (d - r_{\mathbf{A}})(d - r_{\mathbf{B}}),$$

and $\text{rank}(\mathbf{\Gamma}) = d^2 - (d - r_{\mathbf{A}})(d - r_{\mathbf{B}})$. □

A.1.2 Rank stratification

Matrix factorization model: $f_\theta = \mathbf{A}\mathbf{B}$ with $\theta = (\mathbf{A}, \mathbf{B})$, $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{d \times d}$.

Step 1: Stratify the parameter space into different model rank levels to obtain the rank hierarchy over the parameter space.

Given any parameter point $\mathbf{A} = [a_{ij}]_{i,j=1}^d \in \mathbb{R}^{d \times d}$, $\mathbf{B} = [b_{ij}]_{i,j=1}^d \in \mathbb{R}^{d \times d}$. Consider the tangent space

$$\text{span} \{ \mathbf{P}^{ij}, \mathbf{Q}^{ij} \}_{i,j=1}^d,$$

where $\mathbf{P}^{ij} = \frac{\partial \mathbf{f}_\theta}{\partial a_{ij}}$, $\mathbf{Q}^{ij} = \frac{\partial \mathbf{f}_\theta}{\partial b_{ij}}$ and $\text{rank}(\mathbf{P}^{ij}) = \text{rank}(\mathbf{Q}^{ij}) = 1$.

By vectorizing $\mathbf{P}_{d \times d}^{ij}$ and $\mathbf{Q}_{d \times d}^{ij}$, we get

$$\text{vec}(\mathbf{P}^{ij}) = [P_{11}^{ij}, \dots, P_{1d}^{ij}, P_{21}^{ij}, \dots, P_{2d}^{ij}, \dots, P_{d1}^{ij}, \dots, P_{dd}^{ij}]^T \in \mathbb{R}^{d^2},$$

$$\text{vec}(\mathbf{Q}^{ij}) = [Q_{11}^{ij}, \dots, Q_{1d}^{ij}, Q_{21}^{ij}, \dots, Q_{2d}^{ij}, \dots, Q_{d1}^{ij}, \dots, Q_{dd}^{ij}]^T \in \mathbb{R}^{d^2}.$$

Now we put these vectors into a matrix $\mathbf{\Gamma}_{2d^2 \times d^2}$, namely

$$\mathbf{\Gamma}_{2d^2 \times d^2} = [\text{vec}(\mathbf{Q}^{11}), \dots, \text{vec}(\mathbf{Q}^{dd}), \text{vec}(\mathbf{P}^{11}), \dots, \text{vec}(\mathbf{P}^{dd})]^T.$$

Clearly,

$$\text{rank}(\mathbf{\Gamma}) = \dim \left(\text{span} \{ \mathbf{P}^{ij}, \mathbf{Q}^{ij} \}_{i,j=1}^d \right).$$

Therefore, we only need to compute the rank of matrix $\mathbf{\Gamma}$.

By exploiting the Kronecker product of matrices, we are able to write $\mathbf{\Gamma}$ in a more concise form:

$$\mathbf{\Gamma} = \begin{bmatrix} \mathbf{B} & & & & \\ & \mathbf{B} & & & \\ & & \ddots & & \\ & & & \mathbf{B} & \\ a_{11}\mathbf{I} & a_{21}\mathbf{I} & \cdots & a_{d1}\mathbf{I} & \\ a_{12}\mathbf{I} & a_{22}\mathbf{I} & \cdots & a_{d2}\mathbf{I} & \\ \vdots & \vdots & \ddots & \vdots & \\ a_{1d}\mathbf{I} & a_{2d}\mathbf{I} & \cdots & a_{dd}\mathbf{I} & \end{bmatrix} = \begin{bmatrix} \mathbf{I} \otimes \mathbf{B} \\ \mathbf{A}^T \otimes \mathbf{I} \end{bmatrix}.$$

Let $r_A := \text{rank}(\mathbf{A})$, $r_B := \text{rank}(\mathbf{B})$. By Lemma 2, the rank of $\mathbf{\Gamma}$ is $d^2 - (d - r_A)(d - r_B)$. Therefore the matrix factorization model possesses the rank levels $\{d^2 - (d - r_1)(d - r_2) | r_1, r_2 \in [d]\}$ over its parameter space, each of which is occupied by $\{(\mathbf{A}, \mathbf{B}) | \text{rank}(\mathbf{A}) = r_1, \text{rank}(\mathbf{B}) = r_2\}$.

Step 2: Stratify the model function space into different model rank levels to obtain the rank hierarchy over the model function space.

Given any matrix $\mathbf{f}^* \in \mathbb{R}^{d \times d}$, let $r = \text{rank}(\mathbf{f}^*)$. By definition, the model rank of \mathbf{f}^* is the minimal model rank among all parameters recovering \mathbf{f}^* . Because

$$\text{rank}(\mathbf{A}\mathbf{B}) \leq \min\{\text{rank}(\mathbf{A}), \text{rank}(\mathbf{B})\},$$

any factorization $\mathbf{f}^* = \mathbf{A}^*\mathbf{B}^*$ satisfies $\text{rank}(\mathbf{A}^*) \geq r$ and $\text{rank}(\mathbf{B}^*) \geq r$. By the analysis in Step 1, we have $\text{rank}_{\mathbf{f}_\theta}(\theta^*) \geq d^2 - (d - r)^2 = 2rd - r^2$. By the singular value decomposition $\mathbf{f}^* = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$, $\mathbf{A}^* = \mathbf{U}\mathbf{\Sigma}^{\frac{1}{2}}$ and $\mathbf{B}^* = \mathbf{\Sigma}^{\frac{1}{2}}\mathbf{V}^T$ recover \mathbf{f}^* with $\text{rank}(\mathbf{A}^*) = \text{rank}(\mathbf{B}^*) = r$. Therefore, $\text{rank}_{\mathbf{f}_\theta}(\theta^*)$ attains its lower bound $2rd - r^2$, thus $\text{rank}_{\mathbf{f}_\theta}(\mathbf{f}^*) = 2rd - r^2$. Then, the matrix factorization model possesses the rank levels $\{2rd - r^2 | r \in [d]\}$ over its function space, each of which is occupied by $\{\mathbf{f} \in \mathbb{R}^{d \times d} | \text{rank}(\mathbf{f}) = r\}$ as illustrated in Table 6.

Remark that the above analysis serves as a proof of the following proposition.

Proposition 1 (rank hierarchy of a matrix factorization model). *Given a matrix factorization model $\mathbf{f}_\theta = \mathbf{A}\mathbf{B}$ with $\theta = (\mathbf{A}, \mathbf{B})$, $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{d \times d}$, for any matrix $\mathbf{f}^* \in \mathbb{R}^{d \times d}$, we have model rank*

$$\text{rank}_{\mathbf{f}_\theta}(\mathbf{f}^*) = 2rd - r^2,$$

where $r = \text{rank}(\mathbf{f}^*)$ is the matrix rank of \mathbf{f}^* .

Table 6: Rank hierarchy for the matrix factorization model.

model	$f_{\theta} = \mathbf{A}\mathbf{B}, \theta = (\mathbf{A}, \mathbf{B}), \mathbf{A}, \mathbf{B} \in \mathbb{R}^{d \times d}$	
$\text{rank}_{f_{\theta}}(\mathbf{f}^*)$	\mathbf{f}^*	$\arg \min_{\theta' \in \mathcal{M}_{f^*}} \text{rank}_{f_{\theta}}(\theta')$
0	$\mathbf{0}_{d \times d}$	$\mathbf{A} = \mathbf{B} = \mathbf{0}_{d \times d}$
$2d - 1$	$\text{rank}(\mathbf{f}^*) = 1$	$\text{rank}(\mathbf{A}) = \text{rank}(\mathbf{B}) = 1, \mathbf{A}\mathbf{B} = \mathbf{f}^*$
\vdots	\vdots	\vdots
$2rd - r^2$	$\text{rank}(\mathbf{f}^*) = r$	$\text{rank}(\mathbf{A}) = \text{rank}(\mathbf{B}) = r, \mathbf{A}\mathbf{B} = \mathbf{f}^*$
\vdots	\vdots	\vdots
d^2	$\text{rank}(\mathbf{f}^*) = d$	$\text{rank}(\mathbf{A}) = \text{rank}(\mathbf{B}) = d, \mathbf{A}\mathbf{B} = \mathbf{f}^*$

A.2 DNNs

A.2.1 Rank upper bound estimate via critical mappings

Rank stratification for DNNs is in general difficult. Luckily, the recent discovery of the embedding principle and the critical embedding operators provide powerful tools for the rank stratification [18, 19, 27, 28, 29]. In the following, we first provide a general definition of critical mappings, by which the previously proposed critical embeddings are special cases. Then, we prove Lemma 3, showing that uncovering critical mappings is an important means for obtaining an upper bound estimate of the model rank. This general result combined with the embedding principle directly provides a rank upper bound estimate for general deep NNs illustrated in Table. 5.

Definition 4 (critical mapping). *Given two differentiable models $f_{\theta_A} = f(\cdot; \theta_A)$ with $\theta_A \in \mathbb{R}^{M_A}$ and $g_{\theta_B} = g(\cdot; \theta_B)$ with $\theta_B \in \mathbb{R}^{M_B}$, $\mathcal{P} : \mathbb{R}^{M_A} \rightarrow \mathbb{R}^{M_B}$ is a critical mapping from model A to B if given any $\theta \in \mathbb{R}^{M_A}$, we have*

(i) *output preserving: $f_{\theta} = g_{\mathcal{P}(\theta)}$;*

(ii) *criticality preserving: for any data $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ and empirical risk function $R_S(\cdot)$, if $\nabla_{\theta} R_S(f_{\theta}) = \mathbf{0}$, then $\nabla_{\theta} R_S(g_{\mathcal{P}(\theta)}) = \mathbf{0}$.*

Lemma 3 (rank upper bound estimate). *Given two models $f_{\theta_A} = f(\cdot; \theta_A)$ with $\theta_A \in \mathbb{R}^{M_A}$ and $g_{\theta_B} = g(\cdot; \theta_B)$ with $\theta_B \in \mathbb{R}^{M_B}$, if there exists a critical mapping \mathcal{P} from model A to B, then $\text{rank}_g(f^*) \leq \text{rank}_f(f^*) \leq M_A$ for any $f^* \in \mathcal{F}_f$.*

Remark 1. *If $M_B \gg M_A$, this upper bound estimate is highly informative, indicating target recovery capability at heavy overparameterization for model B. Importantly, this lemma establishes the relation between our rank stratification and previous studies about the critical embedding for the DNN loss landscape analysis. As a result, the critical embedding intrinsic to the DNN architecture not only benefits optimization as studied in previous works, but also profoundly benefits the recovery/generalization performance.*

Proof. By Definition 2, for any $f^* \in \mathcal{F}_f$, there exists $\theta^* \in \mathbb{R}^{M_A}$ such that $\text{rank}_f(\theta^*) = \text{rank}_f(f^*)$. Then, $g_{\mathcal{P}(\theta^*)} = f_{\theta^*} = f^*$. Without loss of generality, we consider $R_S(h) = \frac{1}{2} \sum_{i=1}^n (h(\mathbf{x}_i) - y_i)^2$. Then $\nabla_{\theta} R_S(f_{\theta^*}) = \sum_{i=1}^n (y_i - f^*(\mathbf{x}_i)) \nabla_{\theta^*} f_{\theta^*}(\mathbf{x}_i)$ and $\nabla_{\mathcal{P}(\theta)} R_S(g_{\mathcal{P}(\theta^*)}) = \sum_{i=1}^n (y_i - f^*(\mathbf{x}_i)) \nabla_{\mathcal{P}(\theta)} g_{\mathcal{P}(\theta^*)}(\mathbf{x}_i)$. Because \mathcal{P} is criticality preserving for arbitrary data S , we have $\ker(\nabla_{\theta} f_{\theta^*}(\mathbf{X})) \subseteq \ker(\nabla_{\mathcal{P}(\theta)} g_{\mathcal{P}(\theta^*)}(\mathbf{X}))$ for any $\mathbf{X} := [\mathbf{x}_1, \dots, \mathbf{x}_n]$. Here, $\nabla_{\theta} f_{\theta^*}(\mathbf{X}) = [\nabla_{\theta} f_{\theta^*}(\mathbf{x}_1), \dots, \nabla_{\theta} f_{\theta^*}(\mathbf{x}_n)]$. Because $\text{rank}_S(\mathcal{P}(\theta^*)) + \dim(\ker(\nabla_{\mathcal{P}(\theta)} g_{\mathcal{P}(\theta^*)}(\mathbf{X}))) = \text{rank}_S(\theta^*) + \dim(\ker(\nabla_{\theta} f_{\theta^*})) = n$, we have $\text{rank}_S(\mathcal{P}(\theta^*)) \leq \text{rank}_S(\theta^*)$ for any data S (for the definition of $\text{rank}_S(\theta^*)$, see Appendix Section B). Taking the infinite data limit, we obtain $\text{rank}_g(\mathcal{P}(\theta^*)) \leq \text{rank}_f(\theta^*) \leq M_A$. Therefore, $\text{rank}_g(f^*) \leq \text{rank}_f(f^*) \leq M_A$ for any $f^* \in \mathcal{F}_f$. \square

Theorem 1 (Embedding Principle, Theorem 4.2 in Ref. [19]). *Given any NN and any K -neuron wider NN, there exists a K -step composition embedding \mathcal{P} satisfying that: For any given data S , loss function $\ell(\cdot, \cdot)$, activation function $\sigma(\cdot)$, given any critical point θ_{narr}^c of the narrower NN, $\theta_{\text{wide}}^c := \mathcal{P}(\theta_{\text{narr}}^c)$ is still a critical point of the K -neuron wider NN with the same output function, i.e., $\mathbf{f}_{\theta_{\text{narr}}^c} = \mathbf{f}_{\theta_{\text{wide}}^c}$.*

Here wider/narrower means no smaller/larger width in each hidden layer (see either of Refs. [18, 19] for a mathematical definition). As a direct consequence of Lemma 3 and Theorem 1, we obtain the following theorem.

Theorem 3 (rank upper bound estimate for DNNs, Theorem 2 in the main context). *Given any NN with M_{wide} parameters, for any function in the function space of a narrower NN with M_{narr} parameters $f^* \in \mathcal{F}_{\theta_{\text{narr}}}$, we have $\text{rank}_{f_{\theta_{\text{wide}}}}(f^*) \leq \text{rank}_{f_{\theta_{\text{narr}}}}(f^*) \leq M_{\text{narr}}$.*

This theorem gives rise to the partial rank hierarchy exhibited in Table. 5.

A.2.2 Theoretical preparation for two-layer NN rank stratification

Proposition 2. *Let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ be any analytic function such that $\sigma^{(n_j)}(0) \neq 0$ for an infinite sequence of distinct indices $\{n_j\}_{j=1}^{\infty}$. Given $d \in \mathbb{N}$ and m distinct weights $\mathbf{w}_1, \dots, \mathbf{w}_m \in \mathbb{R}^d \setminus \{\mathbf{0}\}$, such that $\mathbf{w}_k \neq \pm \mathbf{w}_j$ for all $1 \leq k < j \leq m$. Then $\{\sigma(\mathbf{w}_i^T \mathbf{x}), \sigma'(\mathbf{w}_i^T \mathbf{x})x_1, \dots, \sigma'(\mathbf{w}_i^T \mathbf{x})x_d\}_{i=1}^m$ is a linearly independent function set.*

Proof. For x sufficiently close to $0 \in \mathbb{R}$, we can write $\sigma(x) = \sum_{j=0}^{\infty} c_j x^j$, where $c_j = \sigma^{(j)}(0)/(j!)$. Then, $\sigma'(x) = \sum_{j=1}^{\infty} j c_j x^{j-1}$. Suppose that the set is not linearly independent. Choose not-all-zero constants $\{\alpha_i\}_{i=1}^m$ and $\{\beta_{i1}, \dots, \beta_{id}\}_{i=1}^m$ such that

$$\mathbf{x} \mapsto \sum_{i=1}^m \left(\alpha_i \sigma(\mathbf{w}_i^T \mathbf{x}) + \sum_{t=1}^d \beta_{it} \sigma'(\mathbf{w}_i^T \mathbf{x}) x_t \right)$$

is a zero map on \mathbb{R}^d , where x_t denotes the t -th component of input. For $k, j, i \in [d]$, define the sets

$$\begin{aligned} A_{k,j} &:= \{\mathbf{x} \in \mathbb{R}^d \mid \langle \mathbf{x}, \mathbf{w}_k \pm \mathbf{w}_j \rangle = 0\} \\ B_i &:= \{\mathbf{x} \in \mathbb{R}^d \mid \langle \mathbf{x}, \mathbf{w}_i \rangle = 0\}. \end{aligned}$$

Clearly, each $A_{k,j}$ is the union of two linear subspaces of dimension $(d-1)$, while each B_i is a possibly empty affine subspace of dimension $(d-1)$. Thus,

$$E := \left(\bigcup_{1 \leq k, j \leq m} A_{k,j} \right) \cup \left(\bigcup_{i=1}^d B_i \right)$$

has \mathcal{L}^d Lebesgue measure zero. Let $\mathbf{e} \in \mathbb{R}^d \setminus E$. Denote $p_i := \langle \mathbf{w}_i, \mathbf{e} \rangle$ for each $i \in [m]$. Since $p_i \neq p_j$ and $p_i + p_j \neq 0$ whenever $i \neq j$, we can, without loss of generality, assume that $|p_1| > |p_2| > \dots > |p_m| > 0$. For any sufficiently small ε and any i, t we have

$$\begin{aligned} \sigma(\mathbf{w}_i^T(\varepsilon \mathbf{e})) &= \sum_{j=0}^{\infty} (c_j p_i^j) \varepsilon^j, \\ \sigma'(\mathbf{w}_i^T(\varepsilon \mathbf{e})) (\varepsilon \mathbf{e})_t &= \varepsilon_t \sum_{j=1}^{\infty} (j c_j p_i^{j-1}) \varepsilon^j. \end{aligned}$$

Thus, for sufficiently small ε ,

$$\begin{aligned} \sum_{i=1}^m \left(\alpha_i \sigma(\mathbf{w}_i^T(\varepsilon \mathbf{e})) + \sum_{t=1}^d \beta_{it} \sigma'(\mathbf{w}_i^T(\varepsilon \mathbf{e})) (\varepsilon \mathbf{e})_t \right) &= \left(\sum_{i=1}^m \alpha_i \right) c_0 + \sum_{j=1}^{\infty} c_j \sum_{i=1}^m \left(\alpha_i + \frac{1}{p_i} \sum_{t=1}^d j \beta_{it} \varepsilon_t \right) p_i^j \varepsilon^j \\ &= 0. \end{aligned} \quad (10)$$

We have $c_j \sum_{i=1}^m \left(\alpha_i + \frac{1}{p_i} \sum_{t=1}^d j \beta_{it} \varepsilon_t \right) p_i^j = 0$ for all $j \in \mathbb{N}$. In particular, for any $j \geq 2$, since $n_j \geq 1$ and $c_{n_j} \neq 0$, we have $\sum_{i=1}^m \left(\alpha_i + \frac{1}{p_i} \sum_{t=1}^d n_j \beta_{it} \varepsilon_t \right) p_i^{n_j} = 0$, which yields

$$\alpha_1 + \frac{1}{p_1} \sum_{t=1}^d n_j \beta_{1t} \varepsilon_t = - \sum_{i=2}^m \left(\alpha_i + \frac{1}{p_i} \sum_{t=1}^d n_j \beta_{it} \varepsilon_t \right) \frac{p_i^{n_j}}{p_1^{n_j}}.$$

If $m = 1$, by taking limits $j \rightarrow \infty$, we have $\alpha_1 = \sum_{t=1}^d \beta_{1t} \varepsilon_t = 0$.

Otherwise, since $|p_1| > |p_i|$ for any $2 \leq i \leq m$, it follows that, by taking limits $j \rightarrow \infty$,

$$\lim_{j \rightarrow \infty} \left(\alpha_1 + \frac{1}{p_1} \sum_{t=1}^d n_j \beta_{1t} \varepsilon_t \right) = \lim_{j \rightarrow \infty} - \sum_{i=2}^m \left(\alpha_i + \frac{1}{p_i} \sum_{t=1}^d n_j \beta_{it} \varepsilon_t \right) \frac{p_i^{n_j}}{p_1^{n_j}} = 0.$$

Thus, we also have $\alpha_1 = \sum_{t=1}^d \beta_{1t} e_t = 0$. For $m > 2$, we may rewrite Eq. (10) as

$$\alpha_2 + \frac{1}{p_2} \sum_{t=1}^d n_j \beta_{2t} e_t = - \sum_{i=3}^m \left(\alpha_i + \frac{1}{p_i} \sum_{t=1}^d n_j \beta_{it} e_t \right) \frac{p_i^{n_j}}{p_2^{n_j}}$$

for each $j \geq 2$, and take limits as we do above to deduce that $\alpha_2 + \frac{1}{p_2} \sum_{t=1}^d n_j \beta_{2t} e_t = 0$. By repeating this procedure for at most m times, we conclude that $\alpha_i + \frac{1}{p_i} \sum_{t=1}^d n_j \beta_{it} e_t = 0$ for all $i \in [m]$. Then, $\alpha_i = \sum_{t=1}^d \beta_{it} e_t = 0$ for any $i \in [m]$. For each i , $\sum_{t=1}^d \beta_{it} e'_t$ is a linear function of e' on the open set $\mathbb{R}^d \setminus E$ which vanishes on a neighborhood of e , we must have $\alpha_i = \beta_{it} = 0$ for any $i \in [m]$, $t \in [d]$. Therefore, $\{\sigma(\mathbf{w}_i^T \mathbf{x}), \sigma'(\mathbf{w}_i^T \mathbf{x})x_1, \dots, \sigma'(\mathbf{w}_i^T \mathbf{x})x_d\}_{i=1}^m$ must be a linearly independent set. \square

Corollary 2 (model rank estimate for two-layer NNs). *Let $\sigma = \tanh$. Given $d \in \mathbb{N}$, weights $\mathbf{w}_1, \dots, \mathbf{w}_m \in \mathbb{R}^d$, $a_1, \dots, a_m \in \mathbb{R}$, we have*

$$\dim(\text{span}\{\sigma(\mathbf{w}_i^T \mathbf{x}), a_i \sigma'(\mathbf{w}_i^T \mathbf{x})x_1, \dots, a_i \sigma'(\mathbf{w}_i^T \mathbf{x})x_d\}_{i=1}^m) = m_{\mathbf{w}} + m_a d,$$

where $m_{\mathbf{w}} = \frac{1}{2} |\{\mathbf{w}_i, -\mathbf{w}_i | \mathbf{w}_i \neq \mathbf{0}, i \in [m]\}|$ indicating the number of independent neurons, $m_a = \frac{1}{2} |\{\mathbf{w}_i, -\mathbf{w}_i | \mathbf{w}_i \neq \mathbf{0}, a_i \neq 0, i \in [m]\}| + |\{\mathbf{w}_i | \mathbf{w}_i = \mathbf{0}, a_i \neq 0, i \in [m]\}|$ indicating the number of independent effective neurons. Here, $|\cdot|$ is the cardinality of a set, i.e., number of different elements in a set.

Proof. Note that $\sigma = \tanh$ is analytic and $\sigma^{(2n+1)}(0) \neq 0$ for all n . Because \tanh is an odd function, we have $\tanh(x) = -\tanh(-x)$ and $\tanh(0) = 0$. Therefore, given $\mathbf{w}_i, \mathbf{w}_j \neq \mathbf{0}$ with $\mathbf{w}_i = \pm \mathbf{w}_j$, $\text{span}\{\sigma(\mathbf{w}_i^T \mathbf{x})\} = \text{span}\{\sigma(\mathbf{w}_j^T \mathbf{x})\}$ and $\text{span}\{\sigma'(\mathbf{w}_i^T \mathbf{x})x_1, \dots, \sigma'(\mathbf{w}_i^T \mathbf{x})x_d\} = \text{span}\{\sigma'(\mathbf{w}_j^T \mathbf{x})x_1, \dots, \sigma'(\mathbf{w}_j^T \mathbf{x})x_d\}$. Since there are $m_{\mathbf{w}}$ different non-zero weights, by Proposition 2 we have

$$\dim(\text{span}\{\sigma(\mathbf{w}_i^T \mathbf{x})\}_{i=1}^m) = m_{\mathbf{w}}.$$

Furthermore, note that

$$\text{span}\{a_i \sigma'(\mathbf{w}_i^T \mathbf{x})x_1, \dots, a_i \sigma'(\mathbf{w}_i^T \mathbf{x})x_d\}_{i=1}^m = \text{span}\{\sigma'(\mathbf{w}_i^T \mathbf{x})x_1, \dots, \sigma'(\mathbf{w}_i^T \mathbf{x})x_d : a_i \neq 0, i \in [m]\}.$$

Thus, by Proposition 2,

$$\begin{aligned} & \dim(\text{span}\{\sigma'(\mathbf{w}_i^T \mathbf{x})x_1, \dots, \sigma'(\mathbf{w}_i^T \mathbf{x})x_d : \mathbf{w}_i \neq \mathbf{0}, a_i \neq 0, i \in [m]\}) \\ &= \frac{1}{2} |\{\mathbf{w}_i, -\mathbf{w}_i | \mathbf{w}_i \neq \mathbf{0}, a_i \neq 0, i \in [m]\}| \cdot d. \end{aligned}$$

Now suppose that $\mathbf{w}_j = \mathbf{0} \in \mathbb{R}^d$ for some $j \in [m]$. Since $\sigma'(\mathbf{w}^T \mathbf{x}) = \sigma'(0) \neq 0$ for all $x \in \mathbb{R}^d$,

$$\text{span}\{\sigma'(\mathbf{w}_j^T \mathbf{x})x_1, \dots, \sigma'(\mathbf{w}_j^T \mathbf{x})x_d\} = \text{span}\{x_1, \dots, x_d\}$$

which consists only of linear functions. By the nonlinearity of \tanh , we conclude that

$$\dim(\text{span}\{\sigma'(\mathbf{w}_i^T \mathbf{x})x_1, \dots, \sigma'(\mathbf{w}_i^T \mathbf{x})x_d\}) = m_a d$$

and thus $\dim(\text{span}\{\sigma(\mathbf{w}_i^T \mathbf{x}), \sigma'(\mathbf{w}_i^T \mathbf{x})x_1, \dots, \sigma'(\mathbf{w}_i^T \mathbf{x})x_d\}) = m_{\mathbf{w}} + m_a d$ as desired. \square

The above corollary directly gives rise to the following result.

Corollary 3. *Let $\sigma = \tanh$. Given distinct weights $\mathbf{w}_1, \dots, \mathbf{w}_m \in \mathbb{R}^d \setminus \{\mathbf{0}\}$ satisfying $\mathbf{w}_k \neq \pm \mathbf{w}_j$ for $k \neq j$, and $a_1, \dots, a_m \in \mathbb{R} \setminus \{0\}$, we have*

$$\dim(\text{span}\{\sigma(\mathbf{w}_i^T \mathbf{x}), a_i \sigma'(\mathbf{w}_i^T \mathbf{x})x_1, \dots, a_i \sigma'(\mathbf{w}_i^T \mathbf{x})x_d\}_{i=1}^m) = m(d+1).$$

Next we consider the estimate of the model rank for convolutional neural networks (CNNs) which are widely used in practice. Here we consider the case where the input has two-dimensional indices, which is the most general case for the image input. The following two propositions can be directly generalized to the model rank estimate of CNNs with an input of one index dimension in the main text.

Proposition 3 (model rank estimate for CNNs (with weight sharing)). *Given $m \in \mathbb{N}$, $d \in \mathbb{N}$ and $s \in [d]$. For any $l \in [m]$, let \mathbf{K}_l be a $(s \times s)$ matrix. Consider CNNs with stride = 1. For a \tanh -CNN f_{θ} with weight sharing,*

$$f_{\theta}(\mathbf{I}) = \sum_{l=1}^m \sum_{i,j=1}^{d+1-s} a_{lij} \tanh \left(\sum_{\alpha,\beta} I_{i+s-\alpha, j+s-\beta} K_{l;\alpha,\beta} \right), \mathbf{I} \in \mathbb{R}^{d \times d},$$

for $\mathbf{K}_l \neq \mathbf{0}, 1 \leq l \leq m$, its model rank at $\theta = (a_{lij}, \mathbf{K}_l)_{l,i,j}$ is $m_a s^2 + m_K (d+1-s)^2$, where $m_K = \frac{1}{2} |\{\mathbf{K}_l, -\mathbf{K}_l | l \in [m]\}|$ indicating the number of independent kernels, $m_a = \frac{1}{2} \sum_{\mathbf{K} \in \mathcal{K}} \dim(\text{span}\{a_{l,:,\cdot}\}_{l \in h(\mathbf{K})})$ indicating the number of independent effective neurons. Here $\mathcal{K} = \{\mathbf{K}_l, -\mathbf{K}_l | l \in [m]\}$, h is a function over \mathcal{K} s.t. for each $\mathbf{K} \in \mathcal{K}, h(\mathbf{K}) = \{l | l \in [m], \mathbf{K}_l = \pm \mathbf{K}\}$. $|\cdot|$ is the cardinality of a set, i.e., number of different elements in a set, and $a_{l,:,\cdot}$ denotes the $(d+1-s) \times (d+1-s)$ matrix whose entries are a_{lij} 's.

Remark. In presence of zero-kernels at a parameter point θ , i.e., $\mathbf{K}_l = \mathbf{0}$ for certain l 's, the model rank is obviously no less than that at a parameter point θ' obtained by replacing a_{lij} by 0 for these l 's and all i, j 's in θ . Because we always have $f(\cdot; \theta) = f(\cdot; \theta')$ and the model rank at θ' is always the same as that in a narrower NN with all zero-kernels removed, establishing the rank estimate for parameter points with nonzero-kernels is sufficient for the rank estimate over the mode function space.

Proof. We first consider the case in which $\mathbf{K}_l \pm \mathbf{K}_{l'} \neq \mathbf{0}$ for any distinct $l, l' \in [m]$. Let $\sigma = \tanh$. In this case the model rank is the dimension of the following function space (with respect to variable $\mathbf{I} \in \mathbb{R}^{d \times d}$)

$$\begin{aligned} & \text{span} \left\{ \frac{\partial f_{\theta}}{\partial a_{lij}}, \frac{\partial f_{\theta}}{\partial K_{l;\alpha,\beta}} \right\} \\ &= \text{span} \left\{ \sigma \left(\sum_{\alpha', \beta'} I_{i+s-\alpha', j+s-\beta'} K_{l;\alpha', \beta'} \right), \right. \\ & \quad \left. \sum_{i', j'=1}^{d+1-s} a_{li'j'} \sigma' \left(\sum_{\alpha', \beta'} I_{i'+s-\alpha', j'+s-\beta'} K_{l;\alpha', \beta'} \right) I_{i'+s-\alpha, j'+s-\beta} \right\}_{l,i,j,\alpha,\beta}, \end{aligned}$$

where $l \in [m]$ and $\alpha, \beta \in [s]$. Next, we prove by contradiction that the set of functions

$$\left\{ \sum_{i,j=1}^{d+1-s} a_{lij} \sigma' \left(\sum_{\alpha', \beta'} I_{i+s-\alpha', j+s-\beta'} K_{l;\alpha', \beta'} \right) I_{i+s-\alpha, j+s-\beta} \right\}_{l \in [m], \alpha, \beta \in [s]}$$

are linearly independent. If they are not linearly independent, there exist not all zero constants $\zeta_{l11}, \dots, \zeta_{lss}$ for $l \in [m]$, such that

$$\sum_{l=1}^m \sum_{\alpha, \beta=1}^s \zeta_{l\alpha\beta} \sum_{i,j=1}^{d+1-s} a_{lij} \sigma' \left(\sum_{\alpha', \beta'} I_{i+s-\alpha', j+s-\beta'} K_{l;\alpha', \beta'} \right) I_{i+s-\alpha, j+s-\beta} = 0,$$

which implies that the set of functions

$$\left\{ a_{lij} \sigma' \left(\sum_{\alpha', \beta'} I_{i+s-\alpha', j+s-\beta'} K_{l;\alpha', \beta'} \right) I_{i+s-\alpha, j+s-\beta} \right\}_{l,i,j,\alpha,\beta}$$

are linearly dependent, contradicting Proposition 2. Moreover, if $a_{lij} = 0$ for any $l \in [m]$ and all $i, j \in \{1, \dots, d+1-s\}$,

$$\sum_{i,j=1}^{d+1-s} a_{lij} \sigma' \left(\sum_{\alpha', \beta'} I_{i+s-\alpha', j+s-\beta'} K_{l;\alpha', \beta'} \right) I_{i+s-\alpha, j+s-\beta} = 0$$

for all $\alpha, \beta \in [s]$. Notice that two kernels with $\mathbf{K}_l = \pm \mathbf{K}_{l'}$ can be reduced to one while maintaining model rank if and only if the corresponding output weights $a_{l,:,\cdot}$ and $a_{l',:,\cdot}$ are linearly dependent. Then, similar to Corollary 2, we conclude that the model rank is $m_a s^2 + m_K (d+1-s)^2$. \square

Proposition 4 (model rank estimate for CNNs without weight sharing). *Given $m \in \mathbb{N}, d \in \mathbb{N}$ and $s \in [d]$. For any $l \in [m]$ and $i, j \in [d+1-s]$, let \mathbf{K}_{lij} be a $(s \times s)$ matrices. Consider CNNs with stride = 1. For a tanh CNN f_{θ} without weight sharing,*

$$f_{\theta}(\mathbf{I}) = \sum_{l=1}^m \sum_{i,j=1}^{d+1-s} a_{lij} \tanh \left(\sum_{\alpha,\beta} I_{i+s-\alpha, j+s-\beta} K_{lij;\alpha,\beta} \right), \mathbf{I} \in \mathbb{R}^{d \times d},$$

for $\mathbf{K}_{lij} \neq \mathbf{0}, l \in [m], i, j \in [d+1-s]$, its model rank at $\theta = (a_{lij}, \mathbf{K}_{lij})_{l,i,j}$ is $m_a s^2 + m_K$, where $m_K = \frac{1}{2} |\{p(\mathbf{K}_{lij}), -p(\mathbf{K}_{lij}) | l \in [m], i, j \in [d+1-s]\}|$ indicating the number of independent kernels,

$m_a = \frac{1}{2}|\{p(\mathbf{K}_{lij}), -p(\mathbf{K}_{lij})|l \in [m], i, j \in [d+1-s], a_{lij} \neq 0\}|$ indicating the number of independent effective neurons. Here p is the padding function over kernels, i.e., for each $(s \times s)$ kernel \mathbf{K}_{lij} , $p(\mathbf{K}_{lij}) \in \mathbb{R}^{d \times d}$ s.t. $p(\mathbf{K}_{lij})[i:i+s-1, j:j+s-1] = \mathbf{K}_{lij}$ and the other elements of $p(\mathbf{K}_{lij})$ are zero. $|\cdot|$ is the cardinality of a set, i.e., number of different elements in a set.

Remark. Similar to CNNs (with weight sharing), for CNNs without weight sharing, establishing the rank estimate for parameter points with nonzero-kernels is sufficient for the rank estimate over the mode function space.

Proof. Let $\sigma = \tanh$. The model rank is the dimension of the following function space

$$\begin{aligned} & \text{span} \left\{ \frac{\partial f_{\boldsymbol{\theta}}}{\partial a_{lij}}, \frac{\partial f_{\boldsymbol{\theta}}}{\partial K_{lij;\alpha,\beta}} \right\}_{l,i,j,\alpha,\beta} \\ &= \text{span} \left\{ \sigma \left(\sum_{\alpha',\beta'} I_{i+s-\alpha',j+s-\beta'} K_{lij;\alpha',\beta'} \right), \right. \\ & \quad \left. a_{lij}\sigma' \left(\sum_{\alpha',\beta'} I_{i+s-\alpha',j+s-\beta'} K_{lij;\alpha',\beta'} \right) I_{i+s-\alpha,j+s-\beta} \right\}_{l,i,j,\alpha,\beta}, \end{aligned}$$

where $l \in [m]$, $1 \leq i, j \leq d+1-s$, and $\alpha, \beta \in [s]$. Also note that if $a_{lij} = 0$ for some $l \in [m]$ and $i, j \in \{1, \dots, d+1-s\}$, then

$$a_{lij}\sigma' \left(\sum_{\alpha',\beta'} I_{i+s-\alpha',j+s-\beta'} K_{lij;\alpha',\beta'} \right) I_{i+s-\alpha,j+s-\beta} = 0$$

for all $\alpha, \beta \in [s]$. It follows from Proposition 2 that this space has dimension $m_a s^2 + m_K$. \square

Remark. In particular, for a target function $f^* \in \mathcal{F}_m^{\text{CNN}} \setminus \mathcal{F}_{m-1}^{\text{CNN}}$, the model rank of a equivalent CNNs model without weight sharing is $m_a s^2 + m_K = m(s^2 + 1)(d+1-s)^2 - s^2 m_{\text{null}}$, where $m_{\text{null}} = |\{a_{lij} | a_{lij} = 0\}|$ is a variable counting the number of ineffective neurons in the target function.

A.2.3 Two-layer fully-connected neural networks

Two-layer tanh-NN: $f_{\boldsymbol{\theta}}(\mathbf{x}) = \sum_{i=1}^m a_i \sigma(\mathbf{w}_i^T \mathbf{x})$, $\mathbf{x} \in \mathbb{R}^d$, $\boldsymbol{\theta} = (a_i, \mathbf{w}_i)_{i=1}^m$, $\sigma = \tanh$.

Step 1: Stratify the parameter space into different model rank levels to obtain the rank hierarchy over the parameter space.

Given any parameter point $\boldsymbol{\theta} = (a_i, \mathbf{w}_i)_{i=1}^m$. Consider the tangent space

$$\text{span} \left\{ \sigma(\mathbf{w}_i^T \mathbf{x}), a_i \sigma'(\mathbf{w}_i^T \mathbf{x}) x_1, \dots, a_i \sigma'(\mathbf{w}_i^T \mathbf{x}) x_d \right\}_{i=1}^m.$$

By Corollary 2, the dimension of tangent space is $m_a d + m_w$, where $m_w = \frac{1}{2}|\{\mathbf{w}_j, -\mathbf{w}_j | j \in [m], \mathbf{w}_j \neq \mathbf{0}\}|$ indicating the number of independent neurons, $m_a = \frac{1}{2}|\{\mathbf{w}_j, -\mathbf{w}_j | j \in [m], \mathbf{w}_j \neq \mathbf{0}, a_j \neq 0\}| + |\{\mathbf{w}_j | j \in [m], \mathbf{w}_j = \mathbf{0}, a_j \neq 0\}|$ indicating the number of independent effective neurons. Here, $|\cdot|$ is the cardinality of a set, i.e., number of different elements in a set.

Step 2: Stratify the model function space into different model rank levels to obtain the rank hierarchy over the model function space.

Given any target function f^* that can be recovered by a two-layer NN with width m . Without loss of generality, let $f^* = f_{\boldsymbol{\theta}^*} := \sum_{i=1}^k a_i^* \sigma(\mathbf{w}_i^{*T} \mathbf{x})$, $a_i^* \neq 0$, $\mathbf{w}_i^* \neq \mathbf{0}$, $\mathbf{w}_i^* \neq \pm \mathbf{w}_j^*$, $1 \leq k \leq m$, $\boldsymbol{\theta}^* = (a_i^*, \mathbf{w}_i^*)_{i=1}^k$. By Proposition 2, the set $\{\sigma(\mathbf{w}_i^T \mathbf{x}), a_i \sigma'(\mathbf{w}_i^T \mathbf{x}) x_1, \dots, a_i \sigma'(\mathbf{w}_i^T \mathbf{x}) x_d\}_{i=1}^k$ is linearly independent and $\text{rank}_{f_{\boldsymbol{\theta}}}(f_{\boldsymbol{\theta}^*}) = k(d+1)$. By definition, the model rank of f^* is the minimal model rank among all parameters recovering f^* . Suppose there exists a NN $f_{\boldsymbol{\theta}} = \sum_{i=1}^q a_i \sigma(\mathbf{w}_i^T \mathbf{x})$, $a_i \neq 0$, $\mathbf{w}_i \neq \mathbf{0}$, $\mathbf{w}_i \neq \pm \mathbf{w}_j$ for $i \neq j$ that can recover f^* , and $\{\sigma(\mathbf{w}_i^T \mathbf{x}), a_i \sigma'(\mathbf{w}_i^T \mathbf{x}) x_1, \dots, a_i \sigma'(\mathbf{w}_i^T \mathbf{x}) x_d\}_{i=1}^q = q(d+1) < k(d+1)$, then we have $q < k$. Since $\{\sigma(\mathbf{w}_i^T \mathbf{x})\}_{i=1}^k$ is linearly independent and $\dim(\text{span}\{\sigma(\mathbf{w}_i^T \mathbf{x})\}_{i=1}^q) \leq q < k$, this contradicts $f_{\boldsymbol{\theta}}$ being able to recover f^* . Therefore, $\text{rank}_{f_{\boldsymbol{\theta}}}(f_{\boldsymbol{\theta}^*}) \geq k(d+1)$ and $\text{rank}_{f_{\boldsymbol{\theta}}}(f_{\boldsymbol{\theta}^*})$ attains its lower bound $k(d+1)$ at $\boldsymbol{\theta}^*$. Thus $\text{rank}_{f_{\boldsymbol{\theta}}}(f^*) = k(d+1)$. Then, the two-layer NN model possesses the rank levels $\{k(d+1) | k \in [m]\}$ over its function space, each of which is occupied by $\{f : \mathbb{R}^d \rightarrow \mathbb{R} | f = f_{\boldsymbol{\theta}} := \sum_{i=1}^k a_i \sigma(\mathbf{w}_i^T \mathbf{x}), a_i \neq 0, \mathbf{w}_i \neq \mathbf{0}, \mathbf{w}_i \neq \pm \mathbf{w}_j\}$ as illustrated in Table 3 in the main text.

Remark that the above analysis serves as a proof of the following proposition.

Proposition 5 (rank hierarchy of two-layer tanh-NN). *Given a two-layer NN*

$$f_{\boldsymbol{\theta}}(\mathbf{x}) = \sum_{i=1}^m a_i \tanh(\mathbf{w}_i^T \mathbf{x}), \mathbf{x} \in \mathbb{R}^d, \boldsymbol{\theta} = (a_i, \mathbf{w}_i)_{i=1}^m,$$

for any target function f^* that can be recovered by a two-layer NN with width m , say $f^* = f_{\boldsymbol{\theta}^*} := \sum_{i=1}^k a_i^* \sigma(\mathbf{w}_i^{*T} \mathbf{x})$, $a_i^* \neq 0$, $\mathbf{w}_i^* \neq \mathbf{0}$, $\mathbf{w}_i^* \neq \pm \mathbf{w}_j^*$, $1 \leq k \leq m$, $\boldsymbol{\theta}^* = (a_i^*, \mathbf{w}_i^*)_{i=1}^k$, we have model rank

$$\text{rank}_{f_{\boldsymbol{\theta}}}(f^*) = k(d+1).$$

A.2.4 Two-layer convolutional neural networks

Two-layer tanh-CNN with weight sharing. Consider the 2-layer width- m tanh convolution neural networks without weight sharing. Given $m \in \mathbb{N}$, $d \in \mathbb{N}$ and $s \in [d]$. For any $l \in [m]$ and $i, j \in \{1, \dots, d+1-s\}$, let \mathbf{K}_l be a $(s \times s)$ convolutional kernel. For a 2-d input I , consider tanh-CNNs with stride = 1:

$$f_{\boldsymbol{\theta}}(I) = \sum_{l=1}^m \sum_{i,j=1}^{d+1-s} a_{lij} \sigma \left(\sum_{\alpha,\beta} I_{i+s-\alpha, j+s-\beta} K_{l;\alpha,\beta} \right), \quad I \in \mathbb{R}^{d \times d}, \sigma = \tanh$$

Step 1: Stratify the parameter space into different model rank levels to obtain the rank hierarchy over the parameter space.

Given any parameter point $\boldsymbol{\theta} = (a_{lij}, \mathbf{K}_l)_{l,i,j} (\mathbf{K}_l \neq \mathbf{0})$. Consider the tangent space

$$\text{span} \left\{ \sigma \left(\sum_{\alpha',\beta'} I_{i+s-\alpha', j+s-\beta'} K_{l;\alpha',\beta'} \right), \sum_{i',j'=s+1}^{d+1} a_{i'j'l} \sigma' \left(\sum_{\alpha',\beta'} I_{i'+s-\alpha', j'+s-\beta'} K_{l;\alpha',\beta'} \right) I_{i'+s-\alpha, j'+s-\beta} \right\}_{l,i,j,\alpha,\beta}.$$

By Proposition 3, the dimension of tangent space is $m_a s^2 + m_K (d+1-s)^2$, where $m_K = \frac{1}{2} |\{\mathbf{K}_l, -\mathbf{K}_l | l \in [m]\}|$ indicating the number of independent kernels, $m_a = \frac{1}{2} \sum_{\mathbf{K} \in \mathcal{K}} \dim(\text{span}\{a_{l,:,\cdot}\}_{l \in h(\mathbf{K})})$ indicating the number of independent effective neurons. Here $\mathcal{K} = \{\mathbf{K}_l, -\mathbf{K}_l | l \in [m]\}$ and h is a function over \mathcal{K} s.t. for each $\mathbf{K} \in \mathcal{K}$, $h(\mathbf{K}) = \{l | l \in [m], \mathbf{K}_l = \pm \mathbf{K}\}$. $|\cdot|$ is the cardinality of a set, i.e., number of different elements in a set.

Step 2: Stratify the model function space into different model rank levels to obtain the rank hierarchy over the model function space.

Given any target function f^* that can be recovered by a two-layer NN with width m . Without loss of generality, let $f^* = f_{\boldsymbol{\theta}^*} := \sum_{l=1}^k \sum_{i,j=1}^{d+1-s} a_{lij}^* \sigma \left(\sum_{\alpha,\beta} I_{i+s-\alpha, j+s-\beta} K_{l;\alpha,\beta}^* \right)$, $\mathbf{K}_l^* \neq \mathbf{0}$, $\mathbf{K}_l^* \neq \pm \mathbf{K}_{l'}^*$ for any $l \neq l', \forall l, \exists a_{lij}^* \neq 0$, $1 \leq k \leq m$, $\boldsymbol{\theta}^* = (a_{lij}^*, \mathbf{K}_l^*)_{l,i,j}$. By Proposition 2, $\text{rank}_{f_{\boldsymbol{\theta}}(\boldsymbol{\theta}^*)} = k(s^2 + (d+1-s)^2)$. By definition, the model rank of f^* is the minimal model rank among all parameters recovering f^* . Suppose there exists a NN $f_{\boldsymbol{\theta}} = \sum_{l=1}^q \sum_{i,j=1}^{d+1-s} a_{lij} \sigma \left(\sum_{\alpha,\beta} I_{i+s-\alpha, j+s-\beta} K_{l;\alpha,\beta} \right)$, $\mathbf{K}_l^* \neq \mathbf{0}$, $\mathbf{K}_l^* \neq \pm \mathbf{K}_{l'}^*$ for any $l \neq l', \forall l, \exists a_{lij}^* \neq 0$ that can recover f^* and the dimension of tangent space $q(s^2 + (d+1-s)^2) < k(s^2 + (d+1-s)^2)$, then we have $q < k$. Since $\{\sigma \left(\sum_{\alpha,\beta} I_{i+s-\alpha, j+s-\beta} K_{l;\alpha,\beta}^* \right)\}_{l=1}^k$ is linearly independent and $\dim(\text{span}\{\sigma \left(\sum_{\alpha,\beta} I_{i+s-\alpha, j+s-\beta} K_{l;\alpha,\beta}^* \right)\}_{l=1}^q) \leq q < k$, this contradicts $f_{\boldsymbol{\theta}}$ being able to recover f^* . Therefore, $\text{rank}_{f_{\boldsymbol{\theta}}(\boldsymbol{\theta}^*)} \geq k(s^2 + (d+1-s)^2)$ and $\text{rank}_{f_{\boldsymbol{\theta}}(\boldsymbol{\theta}^*)}$ attains its lower bound $k(s^2 + (d+1-s)^2)$ at $\boldsymbol{\theta}^*$. Thus $\text{rank}_{f_{\boldsymbol{\theta}}}(f^*) = k(s^2 + (d+1-s)^2)$. Then, the two-layer CNN model possesses the rank levels $\{k(s^2 + (d+1-s)^2) | k \in [m]\}$ over its function space, each of which is occupied by $\{f : \mathbb{R}^d \rightarrow \mathbb{R} | f = f_{\boldsymbol{\theta}} := \sum_{l=1}^k \sum_{i,j=1}^{d+1-s} a_{lij} \sigma \left(\sum_{\alpha,\beta} I_{i+s-\alpha, j+s-\beta} K_{l;\alpha,\beta} \right)\}$ as illustrated in Table 7.

A.2.5 Details of architecture comparison

On the basis of the above rank hierarchy for two-layer CNNs with weight sharing, we further exhibit in Table. 8 the model rank in other architectures such as CNNs without weight sharing and fully-connected NNs illustrated in Fig. 8. Remark that the total size of hidden neurons $m(d+1-s)^2$ is fixed over different architectures. The model of CNN without weight sharing is introduced below.

Table 7: The rank hierarchy of two-layer width- m tanh-CNN (with weight sharing).

model: $f_{\theta}(I) = \sum_{l=1}^m \sum_{i,j=1}^{d+1-s} a_{lij} \sigma \left(\sum_{\alpha,\beta} I_{i+s-\alpha,j+s-\beta} K_{l;\alpha,\beta} \right)$, $I \in \mathbb{R}^{d \times d}$, $\theta = (a_{lij}, \mathbf{K}_l)_{l,i,j}$	
$\text{rank}_{f_{\theta}}(f^*)$	f^*
0	0
$s^2 + (d+1-s)^2$	$\mathcal{F}_1^{\text{CNN}} \setminus \{0\} : \left\{ \sum_{i,j=1}^{d+1-s} a_{ij}^* \sigma \left(\sum_{\alpha,\beta} I_{i+s-\alpha,j+s-\beta} K_{1;\alpha,\beta}^* \right) \mid \mathbf{K}_1^* \neq \mathbf{0}, \exists i, j, a_{ij}^* \neq 0 \right\}$
\vdots	\vdots
$k(s^2 + (d+1-s)^2)$	$\mathcal{F}_k^{\text{CNN}} \setminus \mathcal{F}_{k-1}^{\text{CNN}} : \left\{ \sum_{l=1}^k \sum_{i,j=1}^{d+1-s} a_{lij}^* \sigma \left(\sum_{\alpha,\beta} I_{i+s-\alpha,j+s-\beta} K_{l;\alpha,\beta}^* \right) \mid \mathbf{K}_l^* \neq \mathbf{0}, \mathbf{K}_l^* \neq \pm \mathbf{K}_{l'}^* \text{ for any } l \neq l', \forall l, \exists a_{lij}^* \neq 0 \right\}$
\vdots	\vdots
$m(s^2 + (d+1-s)^2)$	$\mathcal{F}_m^{\text{CNN}} \setminus \mathcal{F}_{m-1}^{\text{CNN}} : \left\{ \sum_{l=1}^m \sum_{i,j=1}^{d+1-s} a_{lij}^* \sigma \left(\sum_{\alpha,\beta} I_{i+s-\alpha,j+s-\beta} K_{l;\alpha,\beta}^* \right) \mid \mathbf{K}_l^* \neq \mathbf{0}, \mathbf{K}_l^* \neq \pm \mathbf{K}_{l'}^* \text{ for any } l \neq l', \forall l, \exists a_{lij}^* \neq 0 \right\}$

Two-layer tanh-CNN without weight sharing. Consider the 2-layer width- m tanh convolution neural networks without weight sharing. Given $m \in \mathbb{N}$, $d \in \mathbb{N}$ and $s \in [d]$. For any $l \in [m]$ and $i, j \in \{1, \dots, d+1-s\}$, let \mathbf{K}_{lij} be a $(s \times s)$ convolutional kernel. For a 2-d input I , consider tanh-CNNs with stride = 1:

$$f_{\theta}(I) = \sum_{l=1}^m \sum_{i,j=1}^{d+1-s} a_{lij} \sigma \left(\sum_{\alpha,\beta} I_{i+s-\alpha,j+s-\beta} K_{lij;\alpha,\beta} \right), \quad I \in \mathbb{R}^{d \times d}, \sigma = \tanh$$

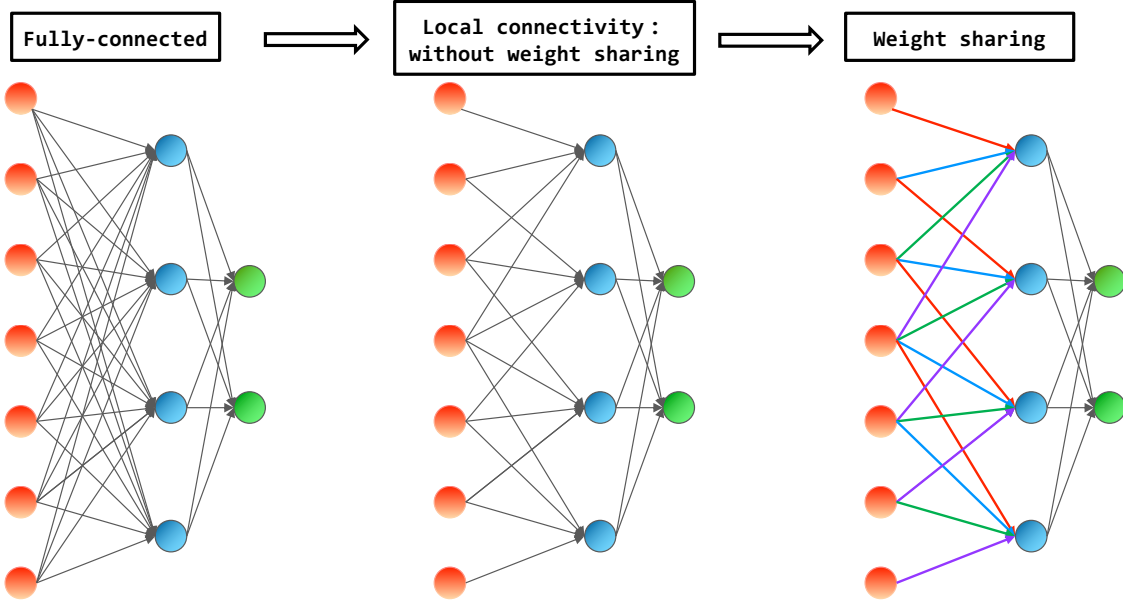


Figure 8: Illustration of architectures from fully-connected NN to CNN for comparison.

By Table. 8, for a common data of $d = 28$, suppose the target function can be recovered by a CNN model with width m , then the model rank for different NN architectures varies a lot from $685m$ (CNN with weight sharing) to $6760m$ (CNN without weight sharing) to $530660m$ (fully-connected NN), indicating a huge difference in their target recovery/generalization performance with limit training data.

Table 8: The rank hierarchy for two-layer tanh-CNN with m -kernels of size $s \times s$ and stride 1. The input $\mathbf{x} \in \mathbb{R}^{d \times d}$. For functions in each function set over the rank hierarchy, we also present their model rank in the corresponding CNN without weight sharing and the corresponding DNN. Here $m_{\text{null}} = |\{a_{lij} | a_{lij} = 0\}|$ counts the number of zero output weights in the target function.

f^*	CNN	CNN without weight sharing	Fully-connected NN
0	0	0	0
$\mathcal{F}_1^{\text{CNN}} \setminus \{0\}$	$s^2 + (d+1-s)^2$	$(s^2+1)(d+1-s)^2 - s^2 m_{\text{null}}$	$(d^2+1)(d+1-s)^2 - d^2 m_{\text{null}}$
\vdots	\vdots	\vdots	\vdots
$\mathcal{F}_k^{\text{CNN}} \setminus \mathcal{F}_{k-1}^{\text{CNN}}$	$k(s^2 + (d+1-s)^2)$	$k(s^2+1)(d+1-s)^2 - s^2 m_{\text{null}}$	$k(d^2+1)(d+1-s)^2 - d^2 m_{\text{null}}$
\vdots	\vdots	\vdots	\vdots
$\mathcal{F}_m^{\text{CNN}} \setminus \mathcal{F}_{m-1}^{\text{CNN}}$	$m(s^2 + (d+1-s)^2)$	$m(s^2+1)(d+1-s)^2 - s^2 m_{\text{null}}$	$m(d^2+1)(d+1-s)^2 - d^2 m_{\text{null}}$

B Details of linear stability theory

Definition 5 (linear stability for recovery). *Given any differentiable model f_{θ} with model function space $\mathcal{F}_{f_{\theta}}$, loss function $\ell(\cdot, \cdot)$, and training data $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$,*

(i) *a parameter point $\theta^* \in \mathbb{R}^M$ is linearly stable if $f(\cdot; \theta^*)$ is the unique solution to*

$$\min_{f \in \tilde{\mathcal{P}}_{\theta^*}} \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i), y_i), \quad (11)$$

where the tangent function hyperplane $\tilde{\mathcal{P}}_{\theta^*} := f(\cdot; \theta^*) + \mathcal{P}_{\theta^*} = \{f(\cdot; \theta^*) + \mathbf{a}^T \nabla_{\theta} f(\cdot; \theta^*) | \mathbf{a} \in \mathbb{R}^M\}$;

(ii) *a function $f^* \in \mathcal{F}_{f_{\theta}}$ is linearly stable if there exists a linearly stable parameter point θ' such that $f(\cdot; \theta') = f^*$.*

Definition 6 (empirical tangent matrix and empirical model rank). *Given any differentiable model f_{θ} and training data $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$, at any parameter point θ^* , $\nabla_{\theta} f(\mathbf{X}; \theta^*) = [\nabla_{\theta} f(\mathbf{x}_1; \theta^*), \dots, \nabla_{\theta} f(\mathbf{x}_n; \theta^*)]$ is referred to as the empirical tangent matrix. Then the empirical model rank is defined as follows*

$$\text{rank}_S(\theta^*) = \text{rank}(\nabla_{\theta} f(\mathbf{X}; \theta^*)).$$

Assumption 1. *The loss function $\ell : \mathbb{R} \times \mathbb{R} \rightarrow [0, \infty)$ is a continuously differentiable function satisfying $\ell(x, y) = 0$ if and only if $x = y$.*

Lemma 4 (linear stability condition for recovery). *Given any differentiable model f_{θ} and training data $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$, a global minimizer $\theta^* \in \mathbb{R}^M$ satisfying $f(\mathbf{x}_i; \theta^*) = y_i$ for all $i \in [n]$ is linearly stable if and only if $\text{rank}_S(\theta^*) = \text{rank}_{f_{\theta}}(\theta^*)$.*

Proof. Let

$$\tilde{R}_S(\mathbf{a}) = \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i, \theta^*) + \mathbf{a}^T \nabla_{\theta} f(\mathbf{x}_i; \theta^*), y_i).$$

Because $\ell(f(\mathbf{x}_i, \theta^*), y_i) = 0$ for all $i \in [n]$, it follows that \mathbf{a} is a global minimum of \tilde{R}_S if and only if $\mathbf{a} \in \ker(\nabla_{\theta} f(\mathbf{X}; \theta^*)^T) = \{\nabla_{\theta} f(\mathbf{X}; \theta^*)^T \mathbf{a} = \mathbf{0} | \mathbf{a} \in \mathbb{R}^M\}$. Now if θ^* is linearly stable, because $f(\cdot, \theta^*)$ is the unique solution, for any $\mathbf{a} \in \mathbb{R}^M$ such that $\tilde{R}_S(\mathbf{a}) = 0$, we must have $\mathbf{a} \in \ker(\nabla_{\theta} f(\cdot; \theta^*)^T) = \{\nabla_{\theta} f(\cdot; \theta^*)^T \mathbf{a} = 0 | \mathbf{a} \in \mathbb{R}^M\}$, thus $\ker(\nabla_{\theta} f(\mathbf{X}; \theta^*)^T) \subseteq \ker(\nabla_{\theta} f(\cdot; \theta^*)^T)$. But since $\ker(\nabla_{\theta} f(\cdot; \theta^*)^T) \subseteq \ker(\nabla_{\theta} f(\mathbf{X}; \theta^*)^T)$ based on the fact that zero function attains 0 at any data point, we have $\ker(\nabla_{\theta} f(\cdot; \theta^*)^T) = \ker(\nabla_{\theta} f(\mathbf{X}; \theta^*)^T)$. Because $\text{rank}_S(\theta^*) + \dim(\ker(\nabla_{\theta} f(\mathbf{X}; \theta^*)^T)) = \text{rank}_{f_{\theta}}(\theta^*) + \dim(\ker(\nabla_{\theta} f(\cdot; \theta^*)^T)) \equiv M$, we obtain $\text{rank}_S(\theta^*) = \text{rank}_{f_{\theta}}(\theta^*)$.

Conversely, if $\text{rank}_S(\theta^*) = \text{rank}_{f_{\theta}}(\theta^*)$, then $\ker(\nabla_{\theta} f(\cdot; \theta^*)^T) = \ker(\nabla_{\theta} f(\mathbf{X}; \theta^*)^T)$. Thus, for any $\mathbf{a} \in \mathbb{R}^M$ with $\tilde{R}_S(\mathbf{a}) = 0$, we have $\theta \in \ker(\nabla_{\theta} f(\mathbf{X}; \theta^*)^T) = \ker(\nabla_{\theta} f(\cdot; \theta^*)^T)$. Therefore, $f(\cdot, \theta^*)$ is the unique solution in its tangent function hyperplane, i.e., θ^* is linearly stable. \square

Lemma 5. Given m linearly independent analytic functions $\phi_1(\mathbf{x}), \dots, \phi_m(\mathbf{x})$ with $\phi_i : \mathbb{R}^d \rightarrow \mathbb{R}$ for all $i \in [m]$, $\text{rank}(\Phi(\mathbf{X})) = m$ a.e. with respect to $\mathcal{L}^{d \times m}$ Lebesgue measure, where

$$\Phi(\mathbf{X}) := \begin{bmatrix} \phi_1(\mathbf{x}_1) & \dots & \phi_m(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ \phi_1(\mathbf{x}_m) & \dots & \phi_m(\mathbf{x}_m) \end{bmatrix}.$$

Proof. Clearly, $\det(\Phi(\cdot)) : \mathbb{R}^{d \times m} \rightarrow \mathbb{R}$ is an analytic function over $\mathbb{R}^{d \times m}$. In addition, because $\{\phi_i\}_{i=1}^m$ are linearly independent, there exists $\mathbf{X} \in \mathbb{R}^{d \times m}$ such that $\det(\Phi(\cdot)) \neq 0$, i.e., $\det(\Phi(\cdot))$ is a non-zero analytic function. Therefore, $\text{rank}(\Phi(\mathbf{X})) = m$ a.e. with respect to $\mathcal{L}^{d \times m}$ Lebesgue measure. \square

Corollary 4. Given m linearly independent analytic functions $\phi_1(\mathbf{x}), \dots, \phi_m(\mathbf{x})$ with $\phi_i : \mathbb{R}^d \rightarrow \mathbb{R}$ for all $i \in [m]$ and $\dim(\text{span}(\{\phi_i(\cdot)\}_{i=1}^m)) = r$, $\text{rank}(\Phi(\mathbf{X})) = \min\{n, r\}$ a.e. with respect to $\mathcal{L}^{d \times m}$ Lebesgue measure.

Proof. It is obvious that $\text{rank}(\Phi(\mathbf{X})) \leq \min\{n, r\}$. For $n \leq r$, we can always pick n independent functions from $\{\phi_i(\cdot)\}_{i=1}^m$. By Lemma 5, $\Phi(\mathbf{X})$ has a rank- n submatrix of $\Phi(\mathbf{X})$ a.e. with respect to Lebesgue measure. For $n > r$, we have that the submatrix of the first r rows of $\Phi(\mathbf{X})$ has rank r a.e. by Lemma 5. Therefore, $\text{rank}(\Phi(\mathbf{X})) = \min\{n, r\}$ a.e. with respect to $\mathcal{L}^{d \times m}$ Lebesgue measure. \square

Theorem 4 (phase transition of linear stability for recovery). Given any analytic model f_θ , for any target function $f^* \in \mathcal{F}_{f_\theta}$ and n generic training data $S = \{(\mathbf{x}_i, f^*(\mathbf{x}_i))\}_{i=1}^n$,

- (i) **Strictly under-determined regime:** if $n < \text{rank}_{f_\theta}(f^*)$, then f^* is not linearly stable;
- (ii) **Quasi-determined regime:** if $n \geq \text{rank}_{f_\theta}(f^*)$, then f^* is linearly stable almost everywhere with respect to S .

Proof. (i) For any $\theta^* \in \mathcal{M}_{f^*}$, we have $\text{rank}_S(\theta^*) \leq n < \text{rank}_{f_\theta}(f^*) \leq \text{rank}_{f_\theta}(\theta^*)$. Therefore the linear stability condition cannot be satisfied, i.e., f^* is not linearly stable.

(ii) Given any $\theta^* \in \mathcal{M}_{f^*}$ with $\text{rank}(\theta^*) = \text{rank}_{f_\theta}(f^*)$, by Corollary 4, $\text{rank}(\nabla_{\theta} f(\mathbf{X}; \theta^*)^T) = \min\{n, \text{rank}_{f_\theta}(\theta^*)\}$ almost everywhere. Because $n \geq \text{rank}_{f_\theta}(f^*)$, we have $\text{rank}_S(\theta^*) = \text{rank}_{f_\theta}(\theta^*)$ almost everywhere. By the linear stability condition Lemma 1, f^* is linearly stable almost everywhere. \square

Corollary 5 (implicit bias of linear stability hypothesis). Given any model f_θ and training dataset $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$, if an interpolation $f' \in \mathcal{F}_{f_\theta}$ is linearly stable, then $\text{rank}_{f_\theta}(f') \leq n$.

Proof. By the linear stability condition Lemma 1, there exists $\theta' \in \mathcal{M}_{f'}$ such that $\text{rank}_S(\theta') = \text{rank}_{f_\theta}(\theta')$. Because $\text{rank}_S(\theta') \leq n$, we have $\text{rank}_{f_\theta}(\theta') \leq n$. Then $\text{rank}_{f_\theta}(f') \leq \text{rank}_{f_\theta}(\theta') \leq n$. \square

C Details of experiments

For Fig. 2, the target matrices we use are as follows:

$$\begin{aligned} \mathbf{M}_1^* &= \begin{bmatrix} 1 & 0.3 & 0.7 & -0.4 \\ 2 & 0.6 & 1.4 & -0.8 \\ 4 & 1.2 & 2.8 & -1.6 \\ 7 & 2.1 & 4.9 & -2.8 \end{bmatrix}, & \mathbf{M}_2^* &= \begin{bmatrix} 4 & 0.6 & 1.8 & 0.8 \\ 6 & 0.9 & 2.7 & 1.2 \\ 8 & 1.2 & 3.6 & 1.6 \\ 18 & 2.7 & 8.1 & 3.6 \end{bmatrix}, \\ \mathbf{M}_3^* &= \begin{bmatrix} -1.8 & 2.4 & 7.7 & -5.3 \\ 0.4 & 1.8 & 5.4 & -3.6 \\ 3.2 & 1.8 & 4.8 & -3. \\ 6.6 & 2.4 & 5.9 & -3.5 \end{bmatrix}, & \mathbf{M}_4^* &= \begin{bmatrix} 7.6 & 3.3 & 19.8 & -7.3 \\ 7.6 & 2.1 & 10.7 & -2.4 \\ 8.8 & 1.8 & 7.6 & -0.2 \\ 19.2 & 3.6 & 14.1 & 0.9 \end{bmatrix}, \\ \mathbf{M}_5^* &= \begin{bmatrix} -1.8 & 2.4 & 7.7 & -5.3 \\ 0.4 & 1.8 & 5.4 & -3.6 \\ 3.2 & 1.8 & 4.8 & -3 \\ 6.6 & 2.4 & 5.9 & -3.5 \end{bmatrix}, & \mathbf{M}_6^* &= \begin{bmatrix} 8.5 & 9.3 & 22.5 & -6.1 \\ 8.2 & 6.1 & 12.5 & -1.6 \\ 11.5 & 19.8 & 15.7 & 3.4 \\ 20.4 & 11.6 & 17.7 & 2.5 \end{bmatrix}, \\ \mathbf{M}_7^* &= \begin{bmatrix} 3.6 & -1.2 & 8.1 & -3.5 \\ 8.1 & -3.5 & 3.6 & -1.2 \\ 9.1 & -1.7 & 11.4 & -0.6 \\ 11.4 & -0.6 & 9.1 & -1.7 \end{bmatrix}, & \mathbf{M}_8^* &= \begin{bmatrix} 12.1 & 17.3 & 24.1 & -4.9 \\ 16.3 & 24.1 & 16.1 & 1.1 \\ 14.2 & 25.8 & 16.9 & 4.3 \\ 22.2 & 15.6 & 18.5 & 3.1 \end{bmatrix}. \end{aligned}$$

For Fig. 3, the target matrix we use is M_1^* as defined above.

For Fig. 4, the target matrix we use is M_2^* as defined above. The sampling sequences we use are listed as follows for each row in Fig. 4, respectively:

$\{(3, 1), (4, 3), (2, 1), (1, 3), (2, 4), (4, 1), (1, 1), (1, 2), (4, 2), (4, 4), (3, 2), (3, 4), (3, 3), (2, 2), (2, 3), (1, 4)\}$,
 $\{(3, 4), (2, 1), (2, 3), (4, 3), (4, 1), (4, 4), (1, 1), (3, 3), (1, 2), (1, 4), (1, 3), (2, 4), (3, 2), (2, 2), (3, 1), (4, 2)\}$,
 $\{(2, 4), (3, 3), (3, 1), (4, 4), (4, 3), (3, 4), (1, 3), (1, 4), (2, 3), (3, 3), (1, 1), (1, 2), (4, 2), (2, 2), (2, 1), (4, 1)\}$,
 $\{(4, 4), (2, 3), (4, 2), (1, 2), (1, 4), (3, 2), (4, 1), (3, 1), (1, 1), (3, 4), (1, 3), (2, 2), (2, 4), (2, 1), (3, 3), (4, 3)\}$,
 $\{(2, 4), (3, 4), (4, 1), (1, 2), (2, 2), (4, 4), (1, 1), (3, 1), (3, 2), (4, 2), (2, 1), (1, 3), (4, 3), (3, 3), (2, 3), (1, 4)\}$,
 $\{(4, 3), (4, 4), (2, 1), (3, 4), (3, 3), (3, 1), (2, 3), (1, 1), (4, 1), (2, 4), (1, 4), (1, 3), (1, 2), (2, 2), (3, 2), (4, 2)\}$.

For Fig. 5, the target matrix we use is M_2^* as defined above. For three sets of experiments with 7, 12, 15 training samples, we sample the following sets of indices of the target matrix, respectively:

$\{(1, 1), (1, 2), (1, 3), (1, 4), (2, 2), (3, 3), (4, 4)\}$,
 $\{(1, 1), (1, 2), (1, 3), (1, 4), (2, 1), (2, 2), (2, 3), (2, 4), (3, 3), (3, 4), (4, 3), (4, 4)\}$,
 $\{(1, 1), (1, 2), (1, 3), (1, 4), (2, 1), (2, 2), (2, 3), (2, 4), (3, 1), (3, 2), (3, 3), (3, 4), (4, 2), (4, 3), (4, 4)\}$.

For Fig. 6, we consider the following target function:

$$f_{\theta}(\mathbf{x}) = \mathbf{W}^{[2]}\sigma(\mathbf{W}^{[1]}\mathbf{x}),$$

$$\text{where } \mathbf{W}^{[1]} = \begin{bmatrix} 0.6 & 0.8 & 1 & 0 & 0 \\ 0 & 0.6 & 0.8 & 1 & 0 \\ 0 & 0 & 0.6 & 0.8 & 1 \end{bmatrix}, \mathbf{W}^{[2]} = [1, 1, 1].$$

For the training dataset and the test dataset, we construct the input data through the standard normal distribution and obtain the output values from the target function. The size of the training dataset varies whereas the size of the test dataset is fixed to 1000. The learning rate for the experiments in each setup is fine-tuned from 0.05 to 0.5 for a better generalization performance.