
Rank Diminishing in Deep Neural Networks

Ruili Feng¹, Kecheng Zheng^{2,1}, Yukun Huang¹, Deli Zhao^{2,3},
Michael Jordan⁴, Zheng-Jun Zha^{1*}

¹University of Science and Technology of China, Hefei, China

²Ant Research, ³Alibaba Group, Hangzhou, China

⁴University of California, Berkeley

ruilifengustc@gmail.com, {zkcys001, kevinh}@mail.ustc.edu.cn,
zhaodeli@gmail.com, jordan@cs.berkeley.edu, zhazj@ustc.edu.cn.

Abstract

The rank of neural networks measures information flowing across layers. It is an instance of a key structural condition that applies across broad domains of machine learning. In particular, the assumption of low-rank feature representations leads to algorithmic developments in many architectures. For neural networks, however, the intrinsic mechanism that yields low-rank structures remains vague and unclear. To fill this gap, we perform a rigorous study on the behavior of network rank, focusing particularly on the notion of rank deficiency. We theoretically establish a universal monotonic decreasing property of network rank from the basic rules of differential and algebraic composition, and uncover rank deficiency of network blocks and deep function couplings. By virtue of our numerical tools, we provide the first empirical analysis of the per-layer behavior of network rank in practical settings, *i.e.*, ResNets, deep MLPs, and Transformers on ImageNet. These empirical results are in direct accord with our theory. Furthermore, we reveal a novel phenomenon of independence deficit caused by the rank deficiency of deep networks, where classification confidence of a given category can be linearly decided by the confidence of a handful of other categories. The theoretical results of this work, together with the empirical findings, may advance understanding of the inherent principles of deep neural networks. Code to detect the rank behavior of networks can be found in <https://github.com/RuiLiFeng/Rank-Diminishing-in-Deep-Neural-Networks>.

1 Introduction

In mathematics, the rank of a smooth function measures the volume of independent information captured by the function [25]. Deep neural networks are highly smooth functions, thus the rank of a network has long been an essential concept in machine learning that underlies many tasks such as information compression [55, 65, 40, 63, 56], network pruning [35, 64, 6, 29, 10], data mining [7, 28, 11, 66, 21, 32], computer vision [68, 67, 34, 30, 32, 70, 69, 17, 18], and natural language processing [9, 31, 8, 12]. Numerous methods are either designed to utilize the mathematical property of network ranks, or are derived from an assumption that low-rank structures are to be preferred.

Yet a rigorous investigation to the behavior of rank of general networks, combining both theoretical and empirical arguments, is still absent in current research, weakening our confidence in the being able to predict performance. To the best of our knowledge, there are only a few previous works discussing the rank behavior of specific network architectures, like attention blocks [15] and BatchNorms [13, 5]

*Corresponding author.

in pure MLP structures. The empirical validation of those methods are also limited to shallow networks, specific architectures, or merely the final layers of deep networks, leaving the global behavior of general deep neural networks mysterious due to prohibitive space-time complexity for measuring them. Rigorous work on network rank that combines both strong theoretical and empirical evidence would have significant implications.

In this paper, we make several contributions towards this challenging goal. We find that the two essential ingredients of deep learning, the chain rules of differential operators and matrix multiplications, are enough to establish a universal principle—that network rank decreases monotonically with the depth of networks. Two factors further enhance the speed of decreasing: a) the explicit rank deficiency of many frequently used network modules, and b) an intrinsic potential of spectrum centralization enforced by the nature of coupling of massive composite functions. To empirically validate our theory, we design numerical tools to efficiently and economically examine the rank behavior of deep neural networks. This is a non-trivial task, as rank is very sensitive to noise and perturbation, and computing ranks of large networks is computationally prohibitive in time and space. Finally, we uncover an interesting phenomenon of independence deficit in multi-class classification networks. We find that many classes do not have their own unique representations in the classification network, and some highly irrelevant classes can decide the outputs of others. This independence deficit can significantly deteriorate the performance of networks in generalized data domains where each class demands a unique representation. In conclusion, the results of this work, together with the numerical tools we invent, may advance understanding of intrinsic properties of deep neural networks, and provide foundations for a broad study of low-dimensional structures in machine learning.

2 Preliminaries

Settings We consider the general deep neural network with L layers. It is a smooth vector-valued function $F : \mathbb{R}^n \rightarrow \mathbb{R}^d$, where \mathbb{R}^n and \mathbb{R}^d are the ambient space of inputs and outputs, respectively. Deep neural networks are coupling of multiple layers, thus we write F as

$$F = f^L \circ f^{L-1} \circ \dots \circ f^1. \quad (1)$$

For simplicity, we further write the k -th sub-network² of F as

$$F_k = f^k \circ \dots \circ f^1, \quad (2)$$

and we use $\mathcal{F}_k = F_k(\mathcal{X})$ to denote the feature space of the k -th sub-network on the data domain \mathcal{X} . We are more interested in the behavior of network rank in the feature spaces rather than scalar outputs (which trivially have rank 1). Thus, for classification or regression networks that output a scalar value, we will consider $F = F_L$ as the transformation from the input space to the final feature space instead. Thus, we always have $n \gg 1$ and $d \gg 1$. For example, for ResNet-50 [23] architecture on ImageNet, we only consider the network slice from the inputs to the last feature layer with 2,048 units.

Rank of Function The rank of a function $f = (f_1, \dots, f_d)^T : \mathbb{R}^n \rightarrow \mathbb{R}^d$ refers to the highest rank of its Jacobian matrix J_f over its input domain \mathcal{X} , which is defined as

$$\text{Rank}(f) = \text{Rank}(J_f) = \max_{x \in \mathcal{X}} \text{Rank}((\partial f_i(x)/\partial x_j)_{n \times d}). \quad (3)$$

It is well-known that the region of non-highest rank is a zero-measure set on the feature manifold by Sard’s theorem [25], so we are safe to ignore them in the definition of the function ranks. The rank of a function represents the volume of information captured by it in the output [25]. That is why it is so important to investigate the behavior of neural networks and many practical applications. Theoretically, by the rank theorem and Sard’s theorem of manifolds [25], we can know that rank of the function equals the intrinsic dimension of its output feature space, as in the following lemma.³

Lemma 1. *Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}^d$ is smooth almost everywhere. Let $\text{Rank}(f) = r$. If data domain \mathcal{X} is a manifold embedded in \mathbb{R}^n and $\phi : \mathcal{U} \rightarrow \mathcal{O}$ is a smooth bijective parameterization from an open subset $\mathcal{U} \subset \mathbb{R}^s$ to $\mathcal{O} \subset \mathcal{X}$, then we have $\dim(f(\mathcal{X})) = \text{Rank}(J_{f \circ \phi}) \leq r$. Thus, the rank of function f gives an upper bound for the intrinsic dimension $\dim(f(\mathcal{X}))$ of the output space.*

²In this paper, sub-network means network slice from the input to some intermediate feature layer; layer network means an independent component of the network, without skip connections from the outside to it, like bottleneck layer of ResNet-50.

³Due to space limitation, all the related proofs are attached in the Appendix.

It is worth mentioning that the intrinsic dimension $\dim(\mathbf{f}(\mathcal{X}))$ of the feature space is usually hard to measure, so the rank of the network gives an operational estimate of it.

3 Numerical Tools

Validating the rank behavior of deep neural networks is a challenging task because it involves operations of high complexity on large-scale non-sparse matrices, which is infeasible both in time and space. Computing the full Jacobian representation of sub-networks of ResNet-50, for example, consumes over 150G GPU memory and several days at a single input point. In accuracy, this is even more challenging as rank is very sensitive to small perturbations. The numerical accuracy of float32, $1.19e - 7$ [45], cannot be trivially neglected in computing matrix ranks. Thus, in this section we establish some numerical tools for validating our subsequent arguments, and provide rigorous theoretical support for them.

3.1 Numerical Rank: Stable Alternative to Rank

The rank of large matrices is known to be unstable: it varies significantly under even small noise perturbations [48]. Matrices perturbed by even small Gaussian noises are almost surely of full rank, regardless of the true rank of the original matrix. Thus in practice we have to use an alternative: we count the number of singular values larger than some given threshold ϵ as the numerical rank of the matrix. Let $\mathbf{W} \in \mathbb{R}^{n \times d}$ be a given matrix. Its numerical rank with tolerance ϵ is

$$\text{Rank}_\epsilon(\mathbf{W}) = \#\{i \in \mathbb{N}_+ : i \leq \min\{n, d\}, \sigma_i \geq \epsilon \|\mathbf{W}\|_2\}, \quad (4)$$

where $\|\mathbf{W}\|_2$ is the ℓ_2 norm (spectral norm) of matrix W , $\sigma_i, i = 1, \dots, \min\{n, d\}$ are its singular values, and $\#$ is the counting measurement for finite sets. We can prove that the numerical rank is stable under small perturbations. Based on Weyl inequalities [57], we have the following theorem.

Theorem 1. *For any given matrix \mathbf{W} , almost every tolerance $\epsilon > 0$, and any perturbation matrix \mathbf{D} , there exists a positive constant $\delta_{\max}(\epsilon)$ such that $\forall \delta \in [0, \delta_{\max}(\epsilon))$, $\text{Rank}_\epsilon(\mathbf{W} + \delta \mathbf{D}) = \text{Rank}_\epsilon(\mathbf{W})$. If \mathbf{W} is a low-rank matrix without random perturbations, then there is a ϵ_{\max} such that for any $\epsilon < \epsilon_{\max}$, $\text{Rank}_\epsilon(\mathbf{W} + \delta \mathbf{D}) = \text{Rank}_\epsilon(\mathbf{W}) = \text{Rank}(\mathbf{W})$ for all $\delta \in [0, \delta_{\max}(\epsilon))$.*

This property of the numerical rank metric makes it a suitable tool for investigating the rank behavior of neural networks. Possible small noises can be filtered out in Jacobian matrices of networks by using numerical rank. It is worth mentioning that random matrices no longer have full rank almost surely under the numerical rank. Instead their rank distribution can be inferred from the well-known Marcenko–Pastur distribution [38] of random matrices. So under numerical rank, low-rank matrices will be commonly seen. In this paper, we always use the numerical rank when measuring ranks.

3.2 Partial Rank of the Jacobian: Estimating Lower Bound of Lost Rank in Deep Networks

To enable the validation of trend of the network ranks, we propose to compute only the rank of sub-matrices of the Jacobian as an alternative. Those sub-matrices are also the Jacobian matrices with respect to a fixed small patch of inputs. Rigorously, given a function \mathbf{f} and its Jacobian \mathbf{J}_f , we denote partial rank of the Jacobian as the rank of a sub-matrix of the Jacobian that consists of the j_1 -th, j_2 -th, ..., j_K -th column of the original Jacobian

$$\text{PartialRank}(\mathbf{J}_f) = \text{Rank}(\text{Sub}(\mathbf{J}_f, j_1, \dots, j_K)) = \text{Rank}((\partial \mathbf{f}_i / \partial \mathbf{x}_{j_k})_{d \times K}), \quad (5)$$

where $1 \leq j_1 < \dots < j_K \leq n$. The partial rank can be efficiently computed and the variance of partial ranks of adjacent sub-networks gives a lower bound on the variance of their ranks. We demonstrate this in appendix F.

3.3 Classification Dimension: Estimating Final Feature Dimension

Measuring the intrinsic dimension of feature manifolds is known to be intractable. So we turn to an approximation procedure. For most classification networks, a linear regression over the final feature manifold decides the final network prediction and accuracy. So we can estimate the intrinsic dimension as the minimum number of principal components in the final feature space to preserve a high classification accuracy. Let $\text{cls} : \mathbb{R}^d \rightarrow \mathbb{R}^c$ be the classification predictions based on the final

Networks	ResNet-18	ResNet-50	GluMixer-24	ResMLP-S24	Swin-T	ViT-T
ClsDim	149	131	199	196	344	109
Ambient Dim.	512	2048	384	384	768	192

Table 1: Classification dimensions (with respect to 95% classification performance of the ambient feature space \mathbb{R}^d) and ambient dimensions of the final feature manifolds of different networks. All networks have low intrinsic dimensions for final features.

feature representation $\mathbf{F}(\mathbf{x})$, Pro_k be the operator that project a vector to the subspace spanned by the top- k PCA components of the final feature representations, and $\mathbf{1}_{\text{cond}}$ the indicator for condition. The classification dimension is then defined as

$$\text{ClsDim}(\mathbf{F}(\mathcal{X})) = \min_k \{k : \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathbb{P}_{\mathcal{X}, \mathcal{Y}}} [\mathbf{1}_{\text{Cls}(\text{Pro}_k(\mathbf{F}(\mathbf{x}))) = \mathbf{y}}] \geq 1 - \epsilon\}, \quad (6)$$

which is the minimum dimensionality needed to reconstruct the classification accuracy of the whole model.

4 Simple Principle of Rank Diminishing for General Networks

The specific designs of neural networks are vast and diverse, but most of them share two fundamental ingredients of deep networks, *i.e.*, *the chain rule and matrix multiplication*. This provides us a chance to analyze the behavior of deep networks that is architecture-agnostic. So in the theoretical aspect, we will focus on how these two intrinsic structures endow impetus of rank diminishing to neural networks.

A most straightforward observation comes from the basic rule of matrix multiplication that for any two matrices \mathbf{A} and \mathbf{B} , we have $\text{Rank}(\mathbf{AB}) \leq \min\{\text{Rank}(\mathbf{A}), \text{Rank}(\mathbf{B})\}$ [26]. Taking this into the chain rule of differential of $\mathbf{J}_{\mathbf{F}} = \mathbf{J}_{f^L} \mathbf{J}_{f^{L-1}} \dots \mathbf{J}_{f^1}$, we then have $\text{Rank}(\mathbf{J}_{\mathbf{F}_k}) = \text{Rank}(\mathbf{J}_{f^k \circ \mathbf{F}_{k-1}}) = \text{Rank}(\mathbf{J}_{f^k} \mathbf{J}_{\mathbf{F}_{k-1}}) \leq \text{Rank}(\mathbf{J}_{\mathbf{F}_{k-1}})$, $k = 2, \dots, L$, which is Eq. (7). We can then get the following principle of rank diminishing.

Theorem 2 (Principle of Rank Diminishing). *Suppose that each layer $f_i, i = 1, \dots, L$ of network \mathbf{F} is almost everywhere smooth⁴ and data domain \mathcal{X} is a manifold, then both the rank of sub-networks and intrinsic dimension of feature manifolds decrease monotonically by depth:*

$$\text{Rank}(\mathbf{f}_1) \geq \text{Rank}(\mathbf{f}_2 \circ \mathbf{f}_1) \geq \dots \geq \text{Rank}(\mathbf{f}_{L-1} \circ \dots \circ \mathbf{f}_1) \geq \text{Rank}(\mathbf{F}_L), \quad (7)$$

$$\dim(\mathcal{X}) \geq \dim(\mathcal{F}_1) \geq \dim(\mathcal{F}_2) \geq \dots \geq \dim(\mathcal{F}_L). \quad (8)$$

This principle describes the behavior of generic neural networks with almost everywhere smooth components, which exhibits the monotonic decreasing (but not strictly) of network ranks and intrinsic dimensionality of feature manifolds.

A flaw of the above principle is that it does not tell whether the rank must decrease. So we need further analysis for the chance of strictly decreasing. A most direct reason to support the strictly decreasing comes from the structural impetus of rank deficiency of numerous network components. Frequently used operations like pooling, downsampling, and dense layer can loose ranks considerably as they explicitly decrease the ambient dimensions of feature representations, or have low rank weight matrices [39]. Specifically, we can have a global criterion for whether a network component will lose ranks as follows.

Theorem 3 (Structural Impetus of Strictly Decreasing). ⁵ *Roughly speaking, if almost everywhere on the input feature manifold, there is a direction such that moving along this direction keeps the output invariant, then the intrinsic dimension of the output feature manifold will be strictly lower than that of the input. The maximum number of independent such directions gives a lower bound on the number of lost intrinsic dimensions.*

For example, the ReLU activation gives the same output for all inputs that only differs from each other in the negative parts. Thus the feature manifold after a ReLU activation will lose the dimension

⁴It means having arbitrary order gradients except for a zero measure set in the input domain

⁵The rigorous version is given in the Appendix.

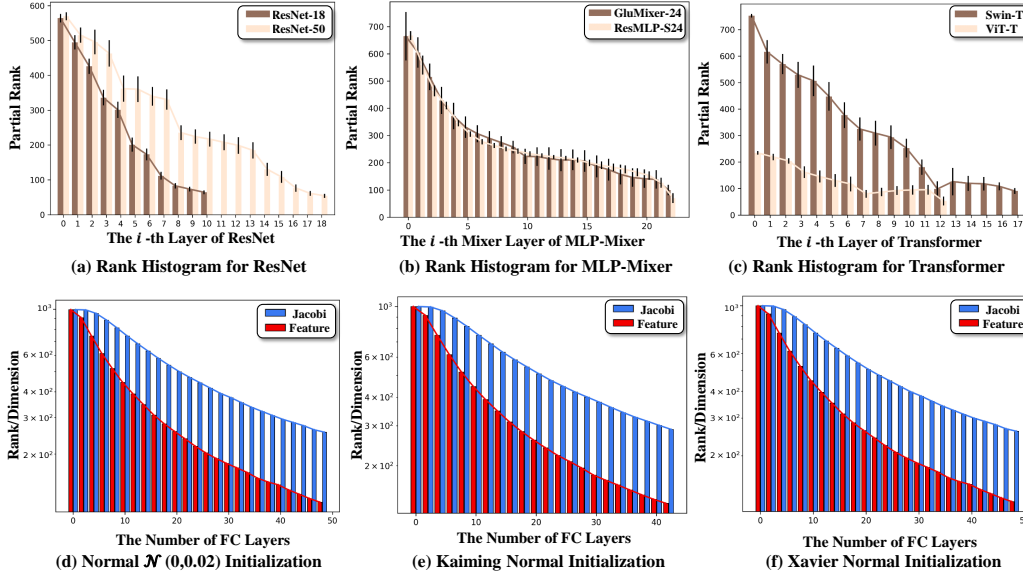


Figure 1: Partial rank of Jacobian matrices of CNN, MLP, and Transformer architecture networks for different layers on ImageNet (top row); rank of Jacobian matrices and feature dimensions of linear MLP network following conditions of Theorem 5 (bottom row). All the models show a similar trend of fast decreasing of ranks as predicted by Theorems 4 and 5.

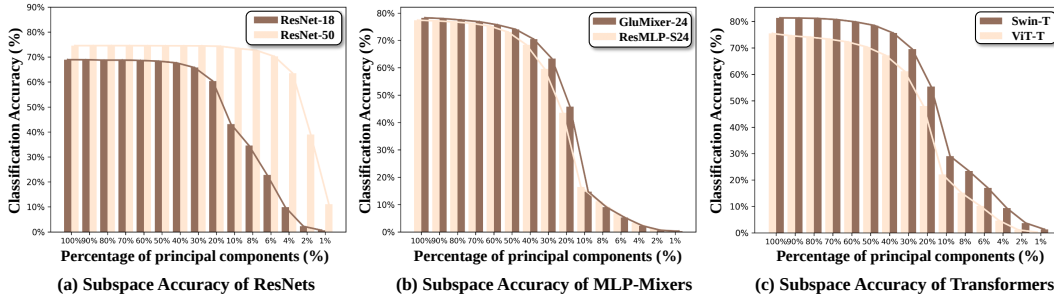


Figure 2: Classification Accuracy (top-1) of using subspaces spanned by top- k % eigenvectors (principal components) of the final feature manifolds. For all networks a small percentage (see Tab. 1) of eigenvectors is enough to reproduce the classification accuracy of the whole network, indicating a low intrinsic dimension of final feature manifolds. **Note that the x -axes are non-linear.**

for those negative parts. For general linear layers (*e.g.*, pure convolution and dense layers), the feature dimensions that belong to the orthogonal complementary spaces of the weights will be lost after applying the linear transformations, as the weight matrices are inactive to changes in their orthogonal complementary spaces.

5 Limiting Behavior of Network Ranks for Infinitely Deep Networks

From a theoretical perspective, it is of great interest to consider the limiting behavior—the rank of infinitely deep neural networks. However, it is infeasible to directly consider this problem and give any rigorous analysis for general cases—it is impossible to exhaust the tedious discussion of how specific structures can influence network ranks. Thus we need to concentrate on a simplified but still representative math model of this problem.

Theory Setup and Its Necessity In brief, we will assume in this section that the Jacobi matrices are random matrices independently following some distributions in a fixed Euclidean space. We then consider the limiting behavior of the singular values of series $\{\mathbf{J}_{F_n}\}_{n=1}^{\infty}$. Modeling the Jacobi

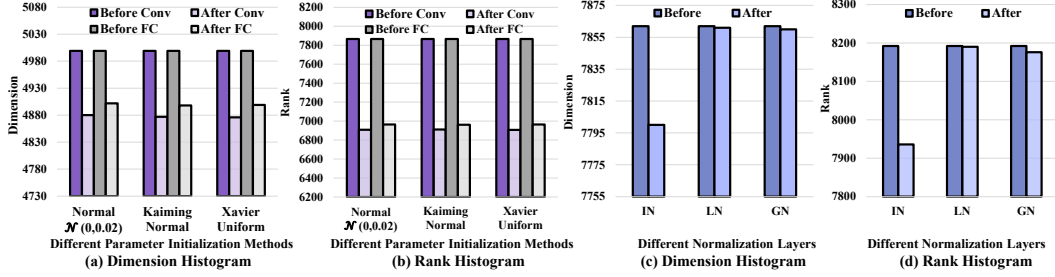


Figure 3: PCA dimension of feature spaces and rank of Jacobian matrix for commonly seen network components under standard Gaussian inputs and randomized weights. Convolution and FC layers tend to lose rank considerably; normalization layers, like InstanceNorm (IN) [54], LayerNorm (LN) [2], and GroupNorm (GN) [60], lose rank modestly. But none can preserve rank.

matrices with random matrices allows us to neglect the influence of specific architectures and reveal the intrinsic property using large number law of probability. While limiting them in the same space is a common practice in analyzing infinitely deep networks [50, 42, 46, 3, 36], as it offers the feasibility to define limits of variables. For example, if the sizes of Jacobi matrices change from time to time, some singular values may suddenly disappear or appear in the series $\{\mathbf{J}_{\mathbf{F}_n}\}_{n=1}^{\infty}$, thus it is unable to define limits for them.

After this simplification, the problem is still representative for understanding what happens in reality. The transformer, recurrent neural networks, and some very deep CNN architectures with constant width [61] can be viewed as examples of this setting. Another interesting setting is to assume that the Jacobi matrices follow the most simple distribution—Gaussian distributions. This can lead to a much stronger yet intuitive result on network ranks. So we are most interested in the limiting behavior when the Jacobi matrices follow weak regularized distributions and the Gaussian distributions. These two cases then lead to Theorems 4 and 5 correspondingly.

Theorem 4. *Let the network be $\mathbf{F} = \mathbf{f}^L \circ \dots \circ \mathbf{f}^1$, and all the ambient dimensions of feature manifolds be the same as the ambient dimension of inputs, i.e., $\mathbf{f}^k : \mathbb{R}^n \rightarrow \mathbb{R}^n, k = 1, \dots, L$. Suppose the Jacobian matrix of each layer \mathbf{f}_i independently follows some distribution μ , and $\mathbb{E}_{\mu}[\max\{\log \|\mathbf{J}_{\mathbf{f}^k}\|_2, 0\}] < \infty$. Let σ_k denote the k -th largest singular value of $\mathbf{J}_{\mathbf{F}}$. Then there is an integer $r < n$ and positive constants μ_r, \dots, μ_n that only depend on μ such that we have for μ -almost everywhere,*

$$\frac{\sigma_k}{\|\mathbf{J}_{\mathbf{F}}\|_2} \sim \exp(-L\mu_k) \rightarrow 0, k = r, \dots, n, \text{ as } L \rightarrow \infty, \quad (9)$$

meaning that for any tolerance $\epsilon > 0$, $\text{Rank}_{\epsilon}(\mathbf{F})$ drops below $r + 1$ as $L \rightarrow \infty$.

Theorem 5. *Let the network be $\mathbf{F} = \mathbf{f}^L \circ \dots \circ \mathbf{f}^1$, and all the ambient dimensions of feature manifolds be the same as the ambient dimension of inputs, i.e., $\mathbf{f}^k : \mathbb{R}^n \rightarrow \mathbb{R}^n, k = 1, \dots, L$. Suppose that $\mathbf{J}_{\mathbf{f}^i}$ independently follows the standard Gaussian distribution. Let σ_k denote the k -th largest singular value of $\mathbf{J}_{\mathbf{F}}$. Then almost surely*

$$\lim_{L \rightarrow \infty} \left(\frac{\sigma_k}{\|\mathbf{J}_{\mathbf{F}}\|_2} \right)^{\frac{1}{L}} = \exp \frac{1}{2} \left(\psi \left(\frac{n-k+1}{2} \right) - \psi \left(\frac{n}{2} \right) \right) < 1, k = 2, \dots, n, \quad (10)$$

where $\psi = \Gamma/\Gamma'$ and Γ is the Gamma function. That means for a large L and any tolerance ϵ , $\text{Rank}_{\epsilon}(\mathbf{F})$ drops to 1 exponentially with speed nC^L , where $C < 1$ is a positive constant that only depends on n .

The proofs of these two theorems rely on the advances of random matrix theory, especially the study of Lyapunov exponents of linear co-cycle systems [51, 59, 19, 47, 33] (see appendix A for detail). Both of the two theorems reveal that infinitely deep networks have an intrinsic intention to drop ranks rapidly and centralize energy in the largest singular vectors. This intention is purely aroused by the two essential ingredients, chain rule and matrix multiplication, as here we omit the consideration of specific network designs, but focus on the large number law of coupling of massive smooth functions.

These two theorems point out that the ratio of non-largest singular value divided by the largest one will converge to zero with an exponential speed. As a result, the numerical rank of the network will decrease quickly after the layer is deep enough. While this is not the hard rank we use in Theorem 2, we can still use it to infer the low rank structure of feature representations. Let $\mathbf{q}_1, \dots, \mathbf{q}_n$ be the singular vectors of \mathbf{J}_F corresponding to the singular values in Theorem 4, and $\alpha_i = \langle \Delta \mathbf{x} / \|\Delta \mathbf{x}\|_2, \mathbf{q}_i \rangle$, then when the depth $L \rightarrow \infty$ and $\Delta \mathbf{x}$ is small, we have $\mathbf{F}(\mathbf{x} + \Delta \mathbf{x}) = \mathbf{F}(\mathbf{x}) + \mathbf{J}_F \Delta \mathbf{x} + o(\|\Delta \mathbf{x}\|_2) = \sigma_1 \sum_{i=1}^r \alpha_i \mathbf{q}_i + \sigma_1 \sum_{j=2}^n \frac{\sigma_j}{\sigma_1} \alpha_j \mathbf{q}_j + o(\|\Delta \mathbf{x}\|_2) \approx \sigma_1 \sum_{i=1}^r \alpha_i \mathbf{q}_i$, which means the neighborhood of data \mathbf{x} is mapped into an r -dimensional subspace spanned by the singular vectors of the largest singular values.

This can also explain why layer-wise regularization techniques are helpful for training deep networks, as those regularization techniques re-normalize the singular value distributions of deep networks, easing the trend of the ratio $\frac{\sigma_i}{\sigma_1}$ becoming infinitesimal. Specifically, if the data distribution near \mathbf{x} is assumed to be standard Gaussian, then the feature distribution is $\mathcal{N}(\mathbf{F}(\mathbf{x}), \mathbf{J}_F^T \mathbf{J}_F)$, which has a very low rank covariance matrix due to $\frac{\sigma_i}{\sigma_1}, i \geq r$ is tiny. Then applying Batch Normalization will try to pull this feature distribution back to Gaussian with identity covariance, thus ease the trend of feature dimension collapse. This has also been discussed in related work of the Batch Normalization [13, 5]. We will discuss some other techniques to remiss the rank diminishing in the Appendix.

These two theorems are also connected with the famous gradient explosion issue of deep neural networks [4, 44], where the largest singular value of the Jacobian matrix tends to infinity when the layer gets deeper. This issue could be viewed as a special case of Theorem 5 that investigates the behavior of all singular values of deep neural networks. The behavior of network ranks in fact manipulates the well-known gradient explosion issue. Rigorously, we have the following conclusion.

Corollary 1. *Under the condition of Theorem 5, then almost surely gradient explosion happens at an exponential speed, i.e., $\log \|\mathbf{J}_F\|_2 = \log \sigma_1 \sim \frac{L}{2}(\log 2 + \psi(n/2)) \rightarrow \infty$ when L is large.*

6 Validating the Rank Behavior

6.1 Validating Rank Diminishing by Depth

In this section, we numerically validate our theory in three types of architectures of benchmark deep neural networks, CNNs, MLPs, and Transformers, in the ImageNet [14] data domain. Information of those networks is listed in Tab. A2. For validating the tendency of network rank of Jacobian matrices, we use the numerical rank of sub-matrices of Jacobian on the central $16 \times 16 \times 3$ image patch of input images. We report the results of other choices of patches in the Appendix. Details of the experiment setup can be found in appendix E.

Rank Diminishing of Jacobians As is discussed in Sec. 3.2 and Lemma 2, the partial rank of the Jacobian is a powerful weapon for us to detect the behavior of huge Jacobian matrices, which are infeasible to compute in practice. The decent value of partial ranks of adjacent sub-networks provides a lower bound to that of full ranks of them. Fig. 1 (a,b,c) report the partial rank of Jacobian matrices of three types of architectures, where we can find consistent diminishing of partial ranks in each layer, indicating a larger rank losing for the full rank of Jacobian matrices.

Implicit Impetus Theorem 5 gives an exponential speed of rank decent by layers. We find that it corresponds well with practice. We investigate this exponential law in a toy network of MLP-50, which is composed of 50 dense layers, each with 1,000 hidden units. The MLP-50 network takes Gaussian noise vectors of \mathbb{R}^{1000} as inputs, and returns a prediction of 1,000 categories. As all the feature manifolds are linear subspaces in this case, their intrinsic dimensions can be directly measured by the numerical rank of their covariance matrices. We report the full rank of Jacobian matrices and intrinsic dimensions of feature manifolds under three different randomly chosen weights in Fig. 1 (d,e,f). Due to the digital accuracy of float32, we stop calculation in each setting when the absolute values of elements of the matrices are lower than $1.19e - 7$. We can find standard curves of exponential laws in all cases for both ranks of Jacobian and intrinsic dimensions of features.

Structural Impetus We validate the structural impetus in Fig. 3. To give an estimation for general cases, here we use Gaussian noises with the size of $128 \times 8 \times 8$ as inputs, and randomize weights of the network components to be validated. We plug those components into a simple fully-connected (FC)

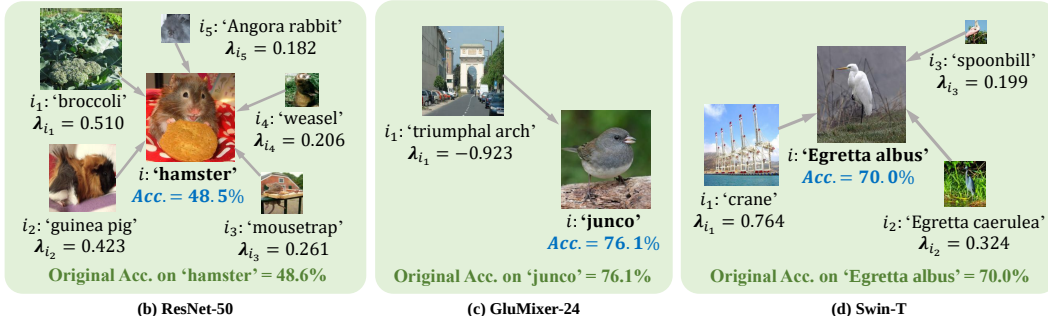
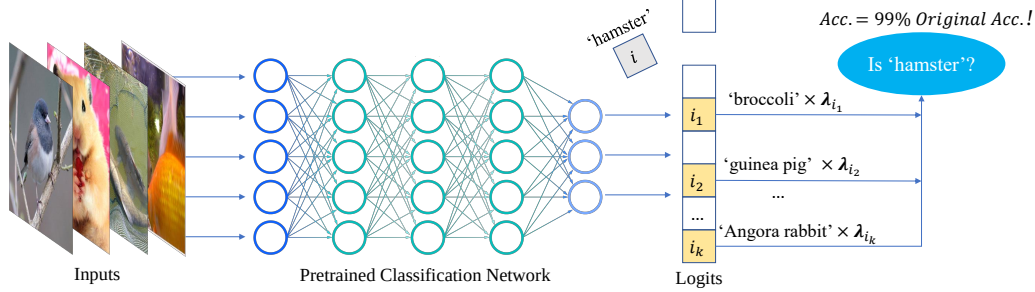


Figure 4: Independence deficit. Classification confidence of some ImageNet categories are linearly decided by a few other categories with fixed coefficients in the whole data domain. We illustrate this phenomenon in (a). Here we present some results from ResNet-50, GluMixer-24, and Swin-T. In the (b,c,d) we illustrate the categories of i_1, \dots, i_k (in the surrounding) to linearly decide category i (in the center) and their corresponding weights $\lambda_{i_1}, \dots, \lambda_{i_k}$. The classification accuracy on the validation set of using Eq. (12), instead of the true logits, to predict the label is reported in blue (if tested on positive samples only, the accuracy rates are 98%, 90%, 82% for cases in (b,c,d) correspondingly). For comparison, the original accuracy for the corresponding categories are reported in green. We can find that 1) a few other categories can decide the confidence of the target category i ; 2) some very irrelevant categories contribute the largest weights. For example in (c), the logits of class ‘junco’ is the negative of ‘triumphal arch’. Both of them indicate a rather drastic competition of different categories for independent representations in final features due to the tight rank budgets.

layer of 8,192 hidden units. As the structure is simple, we directly measure the intrinsic dimension of feature spaces and the full rank of Jacobian matrices before and after the features pass the network components to be measured. The dimension is determined by the number of PCA eigenvalues [27, 58] larger than $1.19e - 7 \times N \times \sigma_{\max}$, where N is the number of PCA eigenvalues, and σ_{\max} is the largest PCA eigenvalue. The batch size is set to 5,000. We find that the convolution (the kernel size is 3×3) and FC layers (the weight size is 8,192) tend to lose rank considerably, while different normalization layers also lose rank modestly. But none of them can preserve rank invariant.

6.2 Validating Low-Rank Terminal Spaces

This section seeks evidence to support a low-rank terminal feature space of deep networks, which can then support the significant diminishing of ranks in the previous layers. The evidence is consisted of a numerical validation that investigates the classification dimension of the terminal feature layer, together with a semantic validation that reveals the independence deficit of deep networks.

Numerical Validation To get an estimation of how many dimensions remain in the final feature representation, we measure the classification dimension in Fig. 2 and Tab. 1. We report the classification accuracy produced by projecting final feature representations to its top $k\%$ eigenvectors in Fig. 2. We choose a threshold of ϵ such that this procedure can reproduce 95% of the original accuracy of the network. The corresponding ClsDim is reported in Tab. 1. As discussed in Sec. 3.3, this gives an estimation of the intrinsic dimension of the final feature manifold. We can find a universal low-rank structure for all types of networks.

Semantic Validation We want to show that there are only a few independent representations to decide the classification scores for all the 1,000 categories of ImageNet. Specifically, can we predict the outputs of the network for some categories based on the outputs for a few other categories, as illustrated in Fig. 4 (a)? And if we can, will those categories be strongly connected to each other? A surprising fact is that, we can find many counter examples of irrelevant categories dominating the network outputs for given categories regarding various network architectures. This interesting phenomenon indicates a rather drastic competing in the final feature layer for the tight rank budgets of all categories, which yields non-realistic dependencies of different categories.

To find the dependencies of categories in final features, we can solve the following Lasso problem [52],

$$\boldsymbol{\lambda}^* = \arg \min_{\boldsymbol{\lambda}_i = -1} \mathbb{E}_{\boldsymbol{x}} [\|\boldsymbol{\lambda}^T \mathbf{W} \mathbf{F}(\boldsymbol{x})\|_2^2] + \eta \|\boldsymbol{\lambda}\|_1, \quad (11)$$

where $\mathbf{F}(\boldsymbol{x}) \in \mathbb{R}^{1000}$ is the slice of network from inputs to the final feature representation, \boldsymbol{x} is the sample from ImageNet \mathcal{X} , and \mathbf{W} is the final dense layer. The solution $\boldsymbol{\lambda}^*$ will be a sparse vector, with k non-zero elements $\lambda_{i_1} \geq \lambda_{i_2} \geq \dots \geq \lambda_{i_k}$, $k \ll 1000$. We can then get

$$\text{logits}(\boldsymbol{x}, i) \approx \lambda_{i_1} \text{logits}(\boldsymbol{x}, i_1) + \dots + \lambda_{i_k} \text{logits}(\boldsymbol{x}, i_k), i \notin \{i_1, \dots, i_k\}, k \ll 1000, \forall \boldsymbol{x} \in \mathcal{X}, \quad (12)$$

where $\text{logits}(\boldsymbol{x}, i_j)$, $j = 1, \dots, k$ is the logits of network for category i_j , *i.e.*, $\text{logits}(\boldsymbol{x}, i_j) = \mathbf{W}_{i_j} \mathbf{F}(\boldsymbol{x})$. It is easy to see that outputs for category i are linearly decided by outputs for i_1, \dots, i_k and are dominated by outputs for i_1 .

Fig. 4 demonstrates the solutions of Eq. (12) for three different categories in ImageNet with $\eta = 20$, and network architectures ResNet-50, GluMixer-24, and Swin-T. The results are surprising. It shows that many categories of the network predictions are in fact ‘redundant’, as they are purely decided by the predictions of the other categories with simple linear coefficients. In this case, the entanglement of different categories cannot be avoided, thus the network may perform poorly under domain shift. An even more surprising finding is that, some very irrelevant categories hold the largest weights when deciding the predictions of the redundant categories, which means that the networks just neglect the unique representations of those categories in training and yield over-fitting when predicting them.

7 Related Work

Previous studies of rank deficiency in deep neural networks follow two parallel clues. One is the study of rank behavior in specific neural network architectures. Dong et al. [15] studies deep networks consisting of pure self-attention networks, and proves that they converge exponentially to a rank-1 matrix under the assumption of globally bounded weight matrices. Daneshmand et al. [13] studies the effect of BatchNorm on MLPs and shows that BatchNorm can prevent drastic diminishing of network ranks in some small networks and datasets. Both of those works avoid directly validating the behavior of network ranks in intermediate layers due to the lacking of efficient numerical tools. An independent clue is the study of implicit self-regularization, which finds that weight matrices tend to lose ranks after training. Martin and Mahoney [39] studies this phenomenon in infinitely-wide, over-parametric neural networks with tools from random matrix theory. Arora et al. [1] studies this phenomenon in deep matrix decomposition. Those works focus on the theoretical behavior of rank of weight matrices induced by the training instead of network ranks.

8 Conclusion

This paper studies the rank behavior of deep neural networks. In contrast to previous work, we focus on directly validating rank behavior with deep neural networks of diverse benchmarks and various settings for real scenarios. We first formalize the analysis and measurement of network ranks. Then under the proposed numerical tools and theoretical analysis, we demonstrate the universal rank diminishing of deep neural networks from both empirical and theoretical perspectives. We further support the rank-deficient structure of networks by revealing the independence deficit phenomenon, where network predictions for a category can be linearly decided by a few other, even irrelevant categories. The results of this work may advance understanding of the behavior of fundamental network architectures and provide intuition for a wide range of work pertaining to network ranks.

Acknowledgement

The authors thank Prof. Dangzheng Liu of the University of Science and Technology of China for discussion. This work was supported by National Key R&D Program of China under Grant 2020AAA0105702, National Natural Science Foundation of China (NSFC) under Grants 62225207 and U19B2038, and the University Synergy Innovation Program of Anhui Province under Grants GXXT-2019-025.

References

- [1] S. Arora, N. Cohen, W. Hu, and Y. Luo. Implicit regularization in deep matrix factorization. *Adv. Neural Inform. Process. Syst.*, 32, 2019.
- [2] J. L. Ba, J. R. Kiros, and G. E. Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- [3] S. Bai, J. Z. Kolter, and V. Koltun. Deep equilibrium models. In *Adv. Neural Inform. Process. Syst.*, 2019.
- [4] Y. Bengio, P. Simard, and P. Frasconi. Learning long-term dependencies with gradient descent is difficult. *IEEE Trans. Neural Netw.*, 5(2):157–166, 1994.
- [5] N. Bjorck, C. P. Gomes, B. Selman, and K. Q. Weinberger. Understanding batch normalization. *Adv. Neural Inform. Process. Syst.*, 31, 2018.
- [6] C. Blakeney, Y. Yan, and Z. Zong. Is pruning compression?: Investigating pruning via network layer similarity. In *IEEE Winter Conf. Appl. Comput. Vis.*, pages 914–922, 2020.
- [7] X. Cai, C. Ding, F. Nie, and H. Huang. On the equivalent of low-rank linear regressions and linear discriminant analysis based regressions. In *Proc. ACM SIGKDD Int. Conf. Knowl. Discov. Data Min.*, pages 1124–1132, 2013.
- [8] B. Chen, T. Dao, E. Winsor, Z. Song, A. Rudra, and C. Ré. Scatterbrain: Unifying sparse and low-rank attention. *Adv. Neural Inform. Process. Syst.*, 34, 2021.
- [9] P. Chen, S. Si, Y. Li, C. Chelba, and C.-J. Hsieh. Groupreduce: Block-wise low-rank approximation for neural language model shrinking. *Adv. Neural Inform. Process. Syst.*, 31, 2018.
- [10] P. Chen, H.-F. Yu, I. Dhillon, and C.-J. Hsieh. Drone: Data-aware low-rank compression for large NLP models. *Adv. Neural Inform. Process. Syst.*, 34:29321–29334, 2021.
- [11] Y. Cheng, L. Yin, and Y. Yu. Lorslim: Low rank sparse linear methods for top-n recommendations. In *IEEE Int. Conf. on Data Min.*, pages 90–99. IEEE, 2014.
- [12] J. Chiu, Y. Deng, and A. Rush. Low-rank constraints for fast inference in structured models. *Adv. Neural Inform. Process. Syst.*, 34, 2021.
- [13] H. Daneshmand, J. Kohler, F. Bach, T. Hofmann, and A. Lucchi. Batch normalization provably avoids ranks collapse for randomly initialised deep networks. *Adv. Neural Inform. Process. Syst.*, 33:18387–18398, 2020.
- [14] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A large-scale hierarchical image database. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 248–255. Ieee, 2009.
- [15] Y. Dong, J.-B. Cordonnier, and A. Loukas. Attention is not all you need: Pure attention loses rank doubly exponentially with depth. In *Int. Conf. Mach. Learn.*, pages 2793–2803. PMLR, 2021.
- [16] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *Int. Conf. Learn. Represent.*, 2020.
- [17] R. Feng, Z. Lin, J. Zhu, D. Zhao, J. Zhou, and Z.-J. Zha. Uncertainty principles of encoding gans. In *Int. Conf. Mach. Learn.*, pages 3240–3251. PMLR, 2021.
- [18] R. Feng, D. Zhao, and Z.-J. Zha. Understanding noise injection in gans. In *Int. Conf. Mach. Learn.*, pages 3284–3293. PMLR, 2021.
- [19] H. Furstenberg and H. Kesten. Products of random matrices. *Ann. Math. Stat.*, 31(2):457–469, 1960.
- [20] H. Furstenberg and H. Kesten. Random matrix products and measures on projective spaces. *Israel J. Math.*, 46:12–32, 1983.

- [21] D. Goldfarb and Z. Qin. Robust low-rank tensor recovery: Models and algorithms. *SIAM J. Matrix Anal. Appl.*, 35(1):225–253, 2014.
- [22] B. Hanin and D. Rolnick. How to start training: The effect of initialization and architecture. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Adv. Neural Inform. Process. Syst.*, volume 31. Curran Associates, Inc., 2018. URL <https://proceedings.neurips.cc/paper/2018/file/d81f9c1be2e08964bf9f24b15f0e4900-Paper.pdf>.
- [23] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 770–778, 2016.
- [24] D. Hendrycks and K. Gimpel. Gaussian error linear units (GELUs). *arXiv preprint arXiv:1606.08415*, 2016.
- [25] M. W. Hirsch. *Differential topology*, volume 33. Springer Science & Business Media, 2012.
- [26] K. Hoffman. *Linear algebra*. Englewood Cliffs, NJ, Prentice-Hall, 1971.
- [27] H. Hotelling. Relations between two sets of variates. *Biometrika*, 28(3/4):321–377, 1936.
- [28] C.-J. Hsieh, K.-Y. Chiang, and I. S. Dhillon. Low rank modeling of signed networks. In *Proc. ACM SIGKDD Int. Conf. Knowl. Discov. Data Min.*, pages 507–515, 2012.
- [29] Y.-C. Hsu, T. Hua, S. Chang, Q. Lou, Y. Shen, and H. Jin. Language model compression with weighted low-rank factorization. In *Int. Conf. Learn. Represent.*, 2021.
- [30] L. Jing, L. Yang, J. Yu, and M. K. Ng. Semi-supervised low-rank mapping learning for multi-label classification. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1483–1491, 2015.
- [31] R. Karimi Mahabadi, J. Henderson, and S. Ruder. Compacter: Efficient low-rank hypercomplex adapter layers. *Adv. Neural Inform. Process. Syst.*, 34, 2021.
- [32] M. Kheirandishfard, F. Zohrizadeh, and F. Kamangar. Deep low-rank subspace clustering. In *IEEE Conf. Comput. Vis. Pattern Recog. Worksh.*, pages 864–865, 2020.
- [33] J. F. C. Kingman. Subadditive ergodic theory. *Ann. Probab.*, pages 883–899, 1973.
- [34] S. Kong and C. Fowlkes. Low-rank bilinear pooling for fine-grained classification. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 365–374, 2017.
- [35] M. Lin, R. Ji, Y. Wang, Y. Zhang, B. Zhang, Y. Tian, and L. Shao. Hrank: Filter pruning using high-rank feature map. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1529–1538, 2020.
- [36] J. Liu, K. Kawaguchi, B. Hooi, Y. Wang, and X. Xiao. Eignn: Efficient infinite-depth graph neural networks. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan, editors, *Adv. Neural Inform. Process. Syst.*, volume 34, pages 18762–18773. Curran Associates, Inc., 2021. URL <https://proceedings.neurips.cc/paper/2021/file/9bd5ee6fe55aaeb673025dbcb8f939c1-Paper.pdf>.
- [37] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 10012–10022, 2021.
- [38] V. A. Marčenko and L. A. Pastur. Distribution of eigenvalues for some sets of random matrices. *Math. USSR Sb.*, 1(4):457, 1967.
- [39] C. H. Martin and M. W. Mahoney. Implicit self-regularization in deep neural networks: Evidence from random matrix theory and implications for learning. *J. Mach. Learn. Res.*, 22(165):1–73, 2021.
- [40] J. Mu, R. Xiong, X. Fan, D. Liu, F. Wu, and W. Gao. Graph-based non-convex low-rank regularization for image compression artifact reduction. *IEEE Trans. Image Process.*, 29:5374–5385, 2020.
- [41] V. Nair and G. E. Hinton. Rectified linear units improve restricted Boltzmann machines. In *Int. Conf. Mach. Learn.*, pages 807–814. PMLR, 2010.
- [42] A. Nazaret and D. Blei. Variational inference for infinitely deep neural networks. In *Int. Conf. Mach. Learn.*, pages 16447–16461. PMLR, 2022.
- [43] C. M. Newman. The distribution of Lyapunov exponents: Exact results for random matrices. *Commun. Math. Phys.*, 103(1):121–126, 1986.

- [44] R. Pascanu, T. Mikolov, and Y. Bengio. On the difficulty of training recurrent neural networks. In *Int. Conf. Mach. Learn.*, pages 1310–1318. PMLR, 2013.
- [45] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Adv. Neural Inform. Process. Syst.*, 32, 2019.
- [46] S. Peluchetti and S. Favaro. Doubly infinite residual neural networks: a diffusion process approach. *J. Mach. Learn. Res.*, 22(175):1–48, 2021. URL <http://jmlr.org/papers/v22/20-706.html>.
- [47] Y. B. Pesin. Characteristic Lyapunov exponents and smooth ergodic theory. *Russ. Math. Surv.*, 32(4):55, 1977.
- [48] A. Quarteroni, R. Sacco, and F. Saleri. *Numerical mathematics*, volume 37. Springer Science & Business Media, 2010.
- [49] N. Shazeer. GLU variants improve transformer. *arXiv preprint arXiv:2002.05202*, 2020.
- [50] S. Sonoda and N. Murata. Transport analysis of infinitely deep neural network. *J. Mach. Learn. Res.*, 20(1):31–82, 2019.
- [51] R. Temam. *Infinite-dimensional dynamical systems in mechanics and physics*, volume 68. Springer Science & Business Media, 2012.
- [52] R. Tibshirani. Regression shrinkage and selection via the lasso. *J. R. Stat. Soc. Series B Stat. Methodol.*, 58(1):267–288, 1996.
- [53] H. Touvron, P. Bojanowski, M. Caron, M. Cord, A. El-Nouby, E. Grave, G. Izacard, A. Joulin, G. Synnaeve, J. Verbeek, et al. ResMLP: Feedforward networks for image classification with data-efficient training. *arXiv preprint arXiv:2105.03404*, 2021.
- [54] D. Ulyanov, A. Vedaldi, and V. Lempitsky. Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*, 2016.
- [55] T. Vogels, S. P. Karimireddy, and M. Jaggi. Practical low-rank communication compression in decentralized deep learning. *Adv. Neural Inform. Process. Syst.*, 33:14171–14181, 2020.
- [56] H. Wang and N. Ahuja. Rank-r approximation of tensors using image-as-matrix representation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, volume 2, pages 346–353. IEEE, 2005.
- [57] H. Weyl. Das asymptotische verteilungsgesetz der eigenwerte linearer partieller differentialgleichungen. *Math. Ann.*, 71:441–479, 1912.
- [58] S. Wold, K. Esbensen, and P. Geladi. Principal component analysis. *Chemom. Intell. Lab. Syst.*, 2(1-3): 37–52, 1987.
- [59] A. Wolf, J. B. Swift, H. L. Swinney, and J. A. Vastano. Determining Lyapunov exponents from a time series. *Physica D*, 16(3):285–317, 1985.
- [60] Y. Wu and K. He. Group normalization. In *Eur. Conf. Comput. Vis.*, pages 3–19, 2018.
- [61] L. Xiao, Y. Bahri, J. Sohl-Dickstein, S. Schoenholz, and J. Pennington. Dynamical isometry and a mean field theory of CNNs: How to train 10,000-layer vanilla convolutional neural networks. In J. Dy and A. Krause, editors, *Int. Conf. Mach. Learn.*, volume 80 of *Proceedings of Machine Learning Research*, pages 5393–5402. PMLR, 10–15 Jul 2018.
- [62] L. Xiao, Y. Bahri, J. Sohl-Dickstein, S. Schoenholz, and J. Pennington. Dynamical isometry and a mean field theory of cnns: How to train 10,000-layer vanilla convolutional neural networks. In *Int. Conf. Mach. Learn.*, pages 5393–5402. PMLR, 2018.
- [63] J. Ye. Generalized low rank approximations of matrices. *Mach. Learn.*, 61(1):167–191, 2005.
- [64] X. Yu, T. Liu, X. Wang, and D. Tao. On compressing deep models by low rank and sparse decomposition. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 7370–7379, 2017.
- [65] Z. Zha, X. Yuan, B. Wen, J. Zhou, J. Zhang, and C. Zhu. From rank estimation to rank approximation: Rank residual constraint for image restoration. *IEEE Trans. Image Process.*, 29:3254–3269, 2019.
- [66] M. Zhan, S. Cao, B. Qian, S. Chang, and J. Wei. Low-rank sparse feature selection for patient similarity learning. In *IEEE Int. Conf. on Data Min.*, pages 1335–1340. IEEE, 2016.

- [67] T. Zhang, B. Ghanem, S. Liu, C. Xu, and N. Ahuja. Low-rank sparse coding for image classification. In *Int. Conf. Comput. Vis.*, pages 281–288, 2013.
- [68] Y. Zhang, Z. Jiang, and L. S. Davis. Learning structured low-rank representations for image classification. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 676–683, 2013.
- [69] J. Zhu, R. Feng, Y. Shen, D. Zhao, Z.-J. Zha, J. Zhou, and Q. Chen. Low-rank subspaces in gans. *Adv. Neural Inform. Process. Syst.*, 34:16648–16658, 2021.
- [70] H. Zou, T. Hastie, and R. Tibshirani. Sparse principal component analysis. *J. Comput. Graph. Stat.*, 15(2): 265–286, 2006.

Appendix

Contents

A Proofs	13
A.1 Proof to Lemma 1	13
A.2 Proof to Theorem 1	14
A.3 Proof to Lemma 2	15
A.4 Proof to Theorem 2	15
A.5 Proof to Theorem 3	15
A.6 Proof to Theorem 4	16
A.6.1 Lyapunov exponents are limits of logarithms of subspace spectral norm divided by layer depth L	17
A.6.2 Singular value distributions of Jacobian matrices of deep function coupling	18
A.6.3 Existence of Lyapunov exponents for Jacobian matrices of deep function coupling	20
A.7 Proof to Theorem 5	20
A.8 Proof to Corollary 1	20
B Influences of Structures	20
C Orthogonal Initialization	22
D Code	22
E Experiments Setup	23
F Partial Rank of the Jacobian: Estimating Lower Bound of Lost Rank in Deep Networks	23
F.1 Partial Rank of Jacobians under Different Input Patches	23
G Estimating Dimension Diminishing in Features	23

A Proofs

A.1 Proof to Lemma 1

This Lemma is the direct result of the *rank theorem of manifolds*.

Theorem 6 (Rank Theorem [25]). *Suppose $\mathbf{f} : \mathcal{M} \rightarrow \mathcal{N}$ is a smooth function from m -dimensional manifold \mathcal{M} to n -dimensional manifold \mathcal{N} , and $\text{Rank}_{\mathcal{M},\mathcal{N}}(\mathbf{J}_{\mathbf{f}}) = r$. Then for each $\mathbf{x} \in \mathcal{M}$, there exists a smooth chart $(\mathcal{U}, \mathbf{m})$ around \mathbf{x} and a smooth chart $(\mathcal{V}, \mathbf{n})$ around $\mathbf{f}(\mathbf{x})$, such that*

$$\mathbf{n} \circ \mathbf{f} \circ \mathbf{m}^{-1} : \mathbf{m}(\mathcal{U}) \subset \mathbb{R}^m \rightarrow \mathbf{n}(\mathcal{V}) \subset \mathbb{R}^n \quad (\text{A13})$$

is given by $\mathbf{n} \circ \mathbf{f} \circ \mathbf{m}^{-1}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m) = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m, 0, \dots, 0)$.

The rank of a function $\text{Rank}_{\mathcal{M},\mathcal{N}}$ defined on the manifold is the rank under local chart systems of the input manifold \mathcal{M} and output manifold \mathcal{N} . Let ϕ be the chart for point $\mathbf{x} \in \mathcal{O} \subset \mathcal{M}$, and the identity map be the chart for point $\mathbf{f}(\mathbf{x}) \in \mathbb{R}^d$. Then it is easy to find that

$$\text{Rank}(\mathbf{I}^{-1} \circ \mathbf{f} \circ \phi) = \text{Rank}(\mathbf{f} \circ \phi) = \text{Rank}_{\mathcal{M},\mathcal{N}}(\mathbf{f}). \quad (\text{A14})$$

Then by Theorem 6, we know that

$$\dim(\mathbf{f}(\mathcal{X})) = \text{Rank}_{\mathcal{M},\mathcal{N}}(\mathbf{f}) = \text{Rank}(\mathbf{f} \circ \phi) \leq \text{Rank}(\mathbf{f}) = r. \quad (\text{A15})$$

The last equal sign comes from the rank inequality of matrix multiplication $\text{Rank}(\mathbf{A}\mathbf{B}) \leq \min\{\text{Rank}(\mathbf{A}), \text{Rank}(\mathbf{B})\}$, which we will discuss later.

A.2 Proof to Theorem 1

Proof to this Theorem needs Weyl's inequalities [57] for singular values of sum of matrices.

Theorem 7 (Weyl's inequalities). *Let \mathbf{A}, \mathbf{B} be $p \times n$ complex matrices, $\sigma_i(\cdot)$ be the i -th largest singular value of the matrix. Then*

$$|\sigma_i(\mathbf{A} + \mathbf{B}) - \sigma_i(\mathbf{B})| \leq \sigma_1(\mathbf{B}), 1 \leq i \leq p, n. \quad (\text{A16})$$

Let p be the number of singular values of \mathbf{W} and \mathbf{D} . By this theorem, we have

$$\sigma_i(\mathbf{W}) - \delta\sigma_1(\mathbf{D}) \leq \sigma_i(\mathbf{W} + \delta\mathbf{D}) \leq \sigma_i(\mathbf{W}) + \delta\sigma_1(\mathbf{D}), i = 1, \dots, p. \quad (\text{A17})$$

To measure the numerical rank, we need to estimate the relative quantities of singular values, which are,

$$\frac{\sigma_i(\mathbf{W}) - \delta\sigma_1(\mathbf{D})}{\sigma_1(\mathbf{W}) + \delta\sigma_1(\mathbf{D})} \leq \frac{\sigma_i(\mathbf{W} + \delta\mathbf{D})}{\sigma_1(\mathbf{W} + \delta\mathbf{D})} \leq \frac{\sigma_i(\mathbf{W}) + \delta\sigma_1(\mathbf{D})}{\sigma_1(\mathbf{W}) - \delta\sigma_1(\mathbf{D})}, i = 2, \dots, p. \quad (\text{A18})$$

Now assume that ϵ does not belong to the following set (which is a zero measure set in \mathbb{R}_+)

$$\Sigma_{\mathbf{W}} = \left\{ \frac{\sigma_i(\mathbf{W})}{\sigma_1(\mathbf{W})} : i = 2, \dots, p \right\}, \quad (\text{A19})$$

and $\text{Rank}_{\epsilon}(\mathbf{W}) = r$. We know that

$$\frac{\sigma_i(\mathbf{W})}{\sigma_1(\mathbf{W})} > \epsilon, i = 2, \dots, r; \frac{\sigma_i(\mathbf{W})}{\sigma_1(\mathbf{W})} < \epsilon, i = r + 1, \dots, p. \quad (\text{A20})$$

Thus, we have that $\forall \delta < \delta_{\max}$,

$$\frac{\sigma_i(\mathbf{W} + \delta\mathbf{D})}{\sigma_1(\mathbf{W} + \delta\mathbf{D})} \geq \frac{\sigma_i(\mathbf{W}) - \delta\sigma_1(\mathbf{D})}{\sigma_1(\mathbf{W}) + \delta\sigma_1(\mathbf{D})} > \epsilon, i = 2, \dots, r, \quad (\text{A21})$$

$$\frac{\sigma_i(\mathbf{W})}{\sigma_1(\mathbf{W})} \leq \frac{\sigma_i(\mathbf{W}) + \delta\sigma_1(\mathbf{D})}{\sigma_1(\mathbf{W}) - \delta\sigma_1(\mathbf{D})} < \epsilon, i = r + 1, \dots, p, \quad (\text{A22})$$

provided that

$$\delta_{\max} = \min \left\{ \frac{1}{\sigma_1(\mathbf{D})} \left(\frac{\sigma_r(\mathbf{W}) + \sigma_1(\mathbf{W})}{\epsilon + 1} - \sigma_1(\mathbf{W}) \right), \frac{\sigma_r(\mathbf{W})}{2\sigma_1(\mathbf{D})}, \frac{1}{\sigma_1(\mathbf{D})} \left(\sigma_1(\mathbf{W}) - \frac{\sigma_{r+1}(\mathbf{W}) + \sigma_1(\mathbf{W})}{\epsilon + 1} \right), \frac{\sigma_1(\mathbf{W})}{2\sigma_1(\mathbf{D})} \right\}. \quad (\text{A23})$$

Thus we can conclude

$$\text{Rank}_{\epsilon}(\mathbf{W} + \delta\mathbf{D}) = \text{Rank}_{\epsilon}(\mathbf{W}), \forall \delta \in [0, \delta_{\max}). \quad (\text{A24})$$

When $\text{Rank}(\mathbf{W}) = r < p$, it is then easy to see if $\epsilon < \frac{\sigma_r(\mathbf{W})}{\sigma_1(\mathbf{W})}$, we have

$$\frac{\sigma_i(\mathbf{W})}{\sigma_1(\mathbf{W})} > \epsilon, i = 2, \dots, r; \frac{\sigma_i(\mathbf{W})}{\sigma_1(\mathbf{W})} = 0 < \epsilon, i = r + 1, \dots, p. \quad (\text{A25})$$

Thus setting $\epsilon_{\max} = \frac{\sigma_r(\mathbf{W})}{\sigma_1(\mathbf{W})}$, we can always have δ_{\max} acquired by Eq. (A23), such that

$$\text{Rank}_{\epsilon}(\mathbf{W}) = \text{Rank}(\mathbf{W}) = \text{Rank}_{\epsilon}(\mathbf{W} + \delta\mathbf{D}), \forall \delta \in [0, \delta_{\max}). \quad (\text{A26})$$

A.3 Proof to Lemma 2

Let \mathbf{A}_i be the i -th column of matrix \mathbf{A} . Given two matrices $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{B} \in \mathbb{R}^{n \times d}$, we have

$$\mathbf{AB} = (\mathbf{AB}_1, \dots, \mathbf{AB}_d). \quad (\text{A27})$$

Thus for any $1 \leq i_1 < \dots < i_K \leq d$,

$$(\mathbf{AB})_{i_1, \dots, i_K} = (\mathbf{AB}_{i_1}, \dots, \mathbf{AB}_{i_K}) = \mathbf{A}(\mathbf{B})_{i_1, \dots, i_K}. \quad (\text{A28})$$

By the rank theorem [26] of matrices, we have

$$\text{Rank}((\mathbf{AB})_{i_1, \dots, i_K}) = \text{Rank}((\mathbf{B})_{i_1, \dots, i_K}) - \dim(\text{Ker}(\mathbf{A}) \cap \text{Im}((\mathbf{B})_{i_1, \dots, i_K})), \quad (\text{A29})$$

and

$$\text{Rank}(\mathbf{AB}) = \text{Rank}(\mathbf{B}) - \dim(\text{Ker}(\mathbf{A}) \cap \text{Im}(\mathbf{B})). \quad (\text{A30})$$

As $(\mathbf{B})_{i_1, \dots, i_K} \subset \mathbf{B}$, it is straightforward to get that

$$\text{Im}((\mathbf{B})_{i_1, \dots, i_K}) \subset \text{Im}(\mathbf{B}). \quad (\text{A31})$$

Thus

$$\text{Ker}(\mathbf{A}) \cap \text{Im}((\mathbf{B})_{i_1, \dots, i_K}) \subset \text{Ker}(\mathbf{A}) \cap \text{Im}(\mathbf{B}). \quad (\text{A32})$$

Then we have

$$\begin{aligned} \text{Rank}(\mathbf{B}) - \text{Rank}(\mathbf{AB}) &= \dim(\text{Ker}(\mathbf{A}) \cap \text{Im}(\mathbf{B})) \geq \dim(\text{Ker}(\mathbf{A}) \cap \text{Im}((\mathbf{B})_{i_1, \dots, i_K})) \\ &= \text{Rank}((\mathbf{B})_{i_1, \dots, i_K}) - \text{Rank}((\mathbf{AB})_{i_1, \dots, i_K}) \geq 0. \end{aligned} \quad (\text{A33})$$

Note that $\text{Rank}(\mathbf{f}_2 \circ \mathbf{f}_1) = \text{Rank}(\mathbf{J}_{\mathbf{f}_2} \mathbf{J}_{\mathbf{f}_1})$. Then we complete the proof.

A.4 Proof to Theorem 2

The key to this principle is the rank theorem of matrices [26], which is

$$\text{Rank}(\mathbf{AB}) = \text{Rank}(\mathbf{B}) - \dim(\text{Ker}(\mathbf{A}) \cap \text{Im}(\mathbf{B})). \quad (\text{A34})$$

Note that $\text{Rank}(\mathbf{AB}) = \text{Rank}(\mathbf{B}^T \mathbf{A}^T)$ and $\text{Rank}(\mathbf{A}^T) = \text{Rank}(\mathbf{A})$. Then we have

$$\begin{aligned} \text{Rank}(\mathbf{AB}) &= \text{Rank}(\mathbf{B}^T \mathbf{A}^T) = \text{Rank}(\mathbf{A}^T) - \dim(\text{Ker}(\mathbf{B}^T) \cap \text{Im}(\mathbf{A}^T)) \\ &= \text{Rank}(\mathbf{A}) - \dim(\text{Ker}(\mathbf{B}^T) \cap \text{Im}(\mathbf{A}^T)). \end{aligned} \quad (\text{A35})$$

The dimension of a linear subspace will at least be zero, thus the above equations suggest

$$\text{Rank}(\mathbf{AB}) \leq \text{Rank}(\mathbf{B}), \text{Rank}(\mathbf{AB}) \leq \text{Rank}(\mathbf{A}). \quad (\text{A36})$$

Applying this argument to the chain rule of differentials then yields the conclusion. Further using Lemma 1 gives the diminishing of intrinsic dimensions of feature manifolds.

A.5 Proof to Theorem 3

We first give the rigorous version of this theorem as follows.

Theorem 8. *Let e_x be the exponential map from a small neighborhood \mathcal{U}_x of point x on the input feature manifold to its tangent space at x , and $\mathbf{v}_x = e_x(x)$. Let \mathcal{X} be the input manifold, $s = \dim(\mathcal{X})$, \mathbf{f}^i be the layer network, $r = \text{Rank}(\mathbf{f}^i)$, and $\mathbf{f}^i(\mathcal{X})$ be the output manifold. If for almost everywhere on the input feature manifold, there is a unit vector $\mathbf{v} \in e_x(\mathcal{U}_x)$, such that the layer network \mathbf{f}^i satisfies*

$$\lim_{t \rightarrow 0} \frac{\|\mathbf{f}^i \circ e_x^{-1}(\mathbf{v}_x) - \mathbf{f}^i \circ e_x^{-1}(\mathbf{v}_x + t\mathbf{v})\|_2}{t} = 0, \quad (\text{A37})$$

then $\dim(\mathbf{f}^i(\mathcal{X})) < s$. If the number of such independent \mathbf{v} in $e_x(\mathcal{U}_x)$ is k , then $\dim(\mathbf{f}^i(\mathcal{X})) \leq s - k$.

Now we prove this theorem.

Note that Eq. (A37) implies

$$\mathbf{J}_{e_x^{-1}} \mathbf{v} \in \text{Ker}(\mathbf{J}_{f^i}). \quad (\text{A38})$$

As it is also easy to see

$$\mathbf{J}_{e_x^{-1}} \mathbf{v} \in \text{Im}(\mathbf{J}_{e_x^{-1}}), \quad (\text{A39})$$

we can conclude

$$\mathbf{0} \neq \mathbf{J}_{e_x^{-1}} \mathbf{v} \in \text{Ker}(\mathbf{J}_{f^i}) \cap \text{Im}(\mathbf{J}_{e_x^{-1}}), \quad (\text{A40})$$

where $\mathbf{0} \neq \mathbf{J}_{e_x^{-1}} \mathbf{v}$ comes from the full rank property of exponential map and its inverse. Thus we have

$$\dim(\text{Ker}(\mathbf{J}_{f^i}) \cap \text{Im}(\mathbf{J}_{e_x^{-1}})) \geq 1. \quad (\text{A41})$$

Specifically, if linearly independent $\mathbf{v}_1, \dots, \mathbf{v}_k$ satisfy Eq. (A37), we can conclude

$$\mathbf{0} \neq \mathbf{J}_{e_x^{-1}} \mathbf{v}_i \in \text{Ker}(\mathbf{J}_{f^i}) \cap \text{Im}(\mathbf{J}_{e_x^{-1}}), i = 1, \dots, k. \quad (\text{A42})$$

As $\mathbf{J}_{e_x^{-1}}$ has full rank due to the property of exponential map, we know that $\mathbf{J}_{e_x^{-1}} \mathbf{v}_i, i = 1, \dots, k$ are linearly independent. Then

$$\dim(\text{Ker}(\mathbf{J}_{f^i}) \cap \text{Im}(\mathbf{J}_{e_x^{-1}})) \geq k. \quad (\text{A43})$$

Thus the rank theorem of matrices [26] reads

$$\text{Rank}(\mathbf{J}_{f^i \circ e_x^{-1}}) = \text{Rank}(\mathbf{J}_{f^i} \mathbf{J}_{e_x^{-1}}) = \text{Rank}(\mathbf{J}_{e_x^{-1}}) - \dim(\text{Ker}(\mathbf{J}_{f^i}) \cap \text{Im}(\mathbf{J}_{e_x^{-1}})) \quad (\text{A44})$$

$$= s - \dim(\text{Ker}(\mathbf{J}_{f^i}) \cap \text{Im}(\mathbf{J}_{e_x^{-1}})) \leq s - k. \quad (\text{A45})$$

Combining this result with Theorem 6 proves our result.

A.6 Proof to Theorem 4

The proof to this theorem relies on the existence of Lyapunov exponents of dynamic systems. Given a linearized dynamic system

$$\dot{\mathbf{v}}(t) = \mathbf{X}_t \mathbf{v}, \quad \mathbf{v}(0) = \mathbf{v}_0 \in \mathbb{R}^n, \quad (\text{A46})$$

its (largest) Lyapunov exponent is defined as

$$\lambda = \limsup_{t \rightarrow \infty} \frac{1}{t} \|\mathbf{v}\|_2. \quad (\text{A47})$$

Further, for a sequence of subspace $\mathcal{L}_h \subset \mathcal{L}_{h-1} \subset \dots \subset \mathcal{L}_1 \subset \mathcal{L}_0 = \mathbb{R}^n$, we can define the corresponding Lyapunov exponents of all those subspaces as

$$\lambda_i = \lim_{t \rightarrow \infty} \frac{1}{t} \log \|\mathbf{v}\|_2, \quad i = 1, \dots, h+1, \quad \mathbf{v}_0 \in \mathcal{L}_{i-1} \setminus \mathcal{L}_i, \quad (\text{A48})$$

and we have

$$\lambda = \lambda_1 > \lambda_2 > \dots > \lambda_h. \quad (\text{A49})$$

It may be surprising to find that such Lyapunov exponents exist, as \mathbf{v}_0 can traverse the entire subspace $\mathcal{L}_{i-1} \setminus \mathcal{L}_i$. We will demonstrate the existence of the Lyapunov exponents for our case later in Sec. A.6.3, which is the classical results from the Furstenberg-Kesten theorem [19] and multiplicative ergodic theorem [47]. Before that, we will first assume the existence of those Lyapunov exponents for simplicity of analysis.

Now consider the case of function couplings

$$\mathbf{F} = \mathbf{f}^L \circ \dots \circ \mathbf{f}^2 \circ \mathbf{f}^1, \quad (\text{A50})$$

which has the Jacobian matrix

$$\mathbf{J}_{\mathbf{F}} = \mathbf{J}_{\mathbf{f}^L} \mathbf{J}_{\mathbf{f}^{L-1}} \dots \mathbf{J}_{\mathbf{f}^2} \mathbf{J}_{\mathbf{f}^1}. \quad (\text{A51})$$

Apparently, the following dynamic system induces the Jacobi matrix of \mathbf{F} ,

$$\dot{\mathbf{v}}(t) = \mathbf{J}_{\mathbf{f}^t} \mathbf{v}, \quad \mathbf{v}(0) = \mathbf{v}_0 \in \mathbb{R}^n, \quad t = 1, \dots, L. \quad (\text{A52})$$

Thus its Lyapunov exponents are given by

$$\lambda_i = \lim_{L \rightarrow \infty} \frac{1}{L} \log \|\mathbf{J}_{\mathbf{F}} \mathbf{v}_0\|_2, \quad i = 1, \dots, h+1, \quad \mathbf{v}_0 \in \mathcal{L}_{i-1} \setminus \mathcal{L}_i, \quad (\text{A53})$$

for a chain of subspaces $\{0\} = \mathcal{L}_{h+1} \subset \mathcal{L}_h \subset \mathcal{L}_{h-1} \subset \dots \subset \mathcal{L}_1 \subset \mathcal{L}_0 = \mathbb{R}^n$.

A.6.1 Lyapunov exponents are limits of logarithms of subspace spectral norm divided by layer depth L

We first demonstrate that the Lyapunov exponents are limits of logarithm of the spectral norm of \mathbf{F} on $\mathcal{L}_{i-1} \setminus \mathcal{L}_i$ divided by layer depth L when $L \rightarrow \infty$, for $i = 1, \dots, r$.

It is easy to see

$$\frac{1}{L} \log \|\mathbf{v}_0\|_2 \|\mathbf{J}_F \frac{\mathbf{v}_0}{\|\mathbf{v}_0\|_2}\|_2 = \frac{1}{L} (\log \|\mathbf{v}_0\|_2 + \log \|\mathbf{J}_F \frac{\mathbf{v}_0}{\|\mathbf{v}_0\|_2}\|_2). \quad (\text{A54})$$

When $L \rightarrow \infty$, $\frac{1}{L} \log \|\mathbf{v}_0\|_2 \rightarrow 0$ for any \mathbf{v}_0 , we have

$$\lambda_i = \lim_{L \rightarrow \infty} \frac{1}{L} \log \|\mathbf{J}_F \mathbf{v}_0\|_2, \quad \|\mathbf{v}_0\|_2 = 1, \quad \mathbf{v}_0 \in \mathcal{L}_{i-1} \setminus \mathcal{L}_i. \quad (\text{A55})$$

Let

$$\lambda_i^L = \sup_{\|\mathbf{v}_0\|_2=1, \mathbf{v}_0 \in \mathcal{L}_{i-1} \setminus \mathcal{L}_i} \frac{1}{L} \log \|\mathbf{J}_F \mathbf{v}_0\|_2. \quad (\text{A56})$$

Note that

$$\sup_{\|\mathbf{v}_0\|_2=1, \mathbf{v}_0 \in \mathcal{L}_{i-1} \setminus \mathcal{L}_i} \frac{1}{L} \log \|\mathbf{J}_F \mathbf{v}_0\|_2 = \frac{1}{L} \log \sup_{\|\mathbf{v}_0\|_2=1, \mathbf{v}_0 \in \mathcal{L}_{i-1} \setminus \mathcal{L}_i} \|\mathbf{J}_F \mathbf{v}_0\|_2, \quad (\text{A57})$$

and

$$\sup_{\|\mathbf{v}_0\|_2=1, \mathbf{v}_0 \in \mathcal{L}_{i-1} \setminus \mathcal{L}_i} \|\mathbf{J}_F \mathbf{v}_0\|_2 = \|\mathbf{J}_F\|_{2,i}, \quad (\text{A58})$$

where $\|\cdot\|_{2,i}$ denote the spectral norm of a linear operator constrained on $\mathcal{L}_{i-1} \setminus \mathcal{L}_i$. Then we have

$$\lambda_i^L = \frac{1}{L} \log \|\mathbf{J}_F\|_{2,i}. \quad (\text{A59})$$

Let $\mathbf{e}_{i_1}, \dots, \mathbf{e}_{i_k}$ be a set of standard orthogonal basis of i_k dimensional subspace $\mathcal{L}_{i-1} \setminus \mathcal{L}_i$. If the Lyapunov exponents exist, by Eq. (A55) we have for any $\epsilon > 0$, there is $N \in \mathbb{N}$ such that for all $L > N$,

$$\lambda_i - \frac{\epsilon}{2} \leq \frac{1}{L} \log \|\mathbf{J}_F \mathbf{e}_j\|_2 \leq \lambda_i + \frac{\epsilon}{2}, \quad j = 1, \dots, i_k. \quad (\text{A60})$$

Let $\mathbf{v} = \sum_{j=1}^{i_k} \alpha_j \mathbf{e}_j \in \mathcal{L}_{i-1} \setminus \mathcal{L}_i$, where $\alpha_j, j = 1, \dots, i_k, \sum_{j=1}^{i_k} \alpha_j^2 = 1$ is the coordinate of unit vector \mathbf{v} under the basis $\mathbf{e}_j, j = 1, \dots, i_k$. Assume that $\|\mathbf{J}_F \mathbf{e}_1\|_2 \geq \|\mathbf{J}_F \mathbf{e}_j\|_2, j = 2, \dots, i_k$. We

then have $\frac{|\alpha_j| \|\mathbf{J}_F \mathbf{e}_j\|_2}{\sum_{j=1}^{i_k} |\alpha_j| \|\mathbf{J}_F \mathbf{e}_1\|_2} \leq 1, j = 1, \dots, i_k$, and

$$\frac{1}{L} \log \|\mathbf{J}_F \mathbf{v}\|_2 \leq \frac{1}{L} \log \sum_{j=1}^{i_k} |\alpha_j| \|\mathbf{J}_F \mathbf{e}_j\|_2 \quad (\text{A61})$$

$$= \frac{1}{L} \log \left(\frac{|\alpha_1| \|\mathbf{J}_F \mathbf{e}_1\|_2}{\sum_{j=1}^{i_k} |\alpha_j| \|\mathbf{J}_F \mathbf{e}_1\|_2} + \sum_{j=2}^{i_k} \frac{|\alpha_j| \|\mathbf{J}_F \mathbf{e}_j\|_2}{\sum_{j=1}^{i_k} |\alpha_j| \|\mathbf{J}_F \mathbf{e}_1\|_2} \right) \left(\sum_{j=1}^{i_k} |\alpha_j| \|\mathbf{J}_F \mathbf{e}_1\|_2 \right) \quad (\text{A62})$$

$$= \frac{1}{L} \log \left(\frac{|\alpha_1| \|\mathbf{J}_F \mathbf{e}_1\|_2}{\sum_{j=1}^{i_k} |\alpha_j| \|\mathbf{J}_F \mathbf{e}_1\|_2} + \sum_{j=2}^{i_k} \frac{|\alpha_j| \|\mathbf{J}_F \mathbf{e}_j\|_2}{\sum_{j=1}^{i_k} |\alpha_j| \|\mathbf{J}_F \mathbf{e}_1\|_2} \right) + \frac{1}{L} \log \sum_{j=1}^{i_k} |\alpha_j| \|\mathbf{J}_F \mathbf{e}_1\|_2 \quad (\text{A63})$$

$$\leq \frac{1}{L} \log \left(1 + \sum_{j=2}^{i_k} \frac{|\alpha_j| \|\mathbf{J}_F \mathbf{e}_j\|_2}{\sum_{j=1}^{i_k} |\alpha_j| \|\mathbf{J}_F \mathbf{e}_1\|_2} \right) + \frac{1}{L} \log \sum_{j=1}^{i_k} |\alpha_j| \|\mathbf{J}_F \mathbf{e}_1\|_2 \quad (\text{A64})$$

$$\leq \frac{1}{L} \sum_{j=2}^{i_k} \frac{|\alpha_j| \|\mathbf{J}_F \mathbf{e}_j\|_2}{\sum_{j=1}^{i_k} |\alpha_j| \|\mathbf{J}_F \mathbf{e}_1\|_2} + \frac{1}{L} \log \sum_{j=1}^{i_k} |\alpha_j| \|\mathbf{J}_F \mathbf{e}_1\|_2 \quad (\text{A65})$$

$$\leq \frac{1}{L} (i_k - 1) + \frac{1}{L} \log \sum_{j=1}^{i_k} |\alpha_j| + \frac{1}{L} \log \|\mathbf{J}_F \mathbf{e}_1\|_2 \quad (\text{A66})$$

$$\leq \frac{1}{L} (i_k - 1) + \frac{1}{2L} \log(1^2 + \dots + 1^2) (\alpha_1^2 + \dots + \alpha_{i_k}^2) + \frac{1}{L} \log \|\mathbf{J}_F \mathbf{e}_1\|_2 \quad (\text{A67})$$

$$= \frac{1}{L} (i_k - 1) + \frac{1}{2L} \log i_k + \frac{1}{L} \log \|\mathbf{J}_F \mathbf{e}_1\|_2 \quad (\text{A68})$$

$$\leq \frac{1}{L} (i_k - 1) + \frac{1}{2L} \log i_k + \lambda_i + \frac{\epsilon}{2}. \quad (\text{A69})$$

Thus, if we set $N_0 = \max\{N, \frac{2i_k - 2 + \log i_k}{\epsilon}\}$, then when $L > N_0$ we have $\frac{1}{L} (i_k - 1) + \frac{1}{2L} \log i_k < \frac{\epsilon}{2}$ and

$$\frac{1}{L} \log \|\mathbf{J}_F \mathbf{v}\|_2 \leq \lambda_i + \epsilon, \forall \mathbf{v} \in \mathcal{L}_{i-1} \setminus \mathcal{L}_i, \|\mathbf{v}\|_2 = 1. \quad (\text{A70})$$

Combining Eqs. (A60) and (A70), we have for any $\epsilon > 0$, there is $N_0 \in \mathbb{N}$, such that when $L > N_0$, we always have

$$\lambda_i - \frac{\epsilon}{2} \leq \lambda_i^L = \sup_{\|\mathbf{v}_0\|_2=1, \mathbf{v}_0 \in \mathcal{L}_{i-1} \setminus \mathcal{L}_i} \frac{1}{L} \log \|\mathbf{J}_F \mathbf{v}_0\|_2 = \frac{1}{L} \log \|\mathbf{J}_F\|_{2,i} \leq \lambda_i + \epsilon. \quad (\text{A71})$$

Thus, if the Lyapunov exponents exist, *i.e.*, the existence of limits of Eq. (A53), we have

$$\lambda_i = \lim_{L \rightarrow \infty} \frac{1}{L} \log \|\mathbf{J}_F\|_{2,i} = \lim_{L \rightarrow \infty} \lambda_i^L. \quad (\text{A72})$$

A.6.2 Singular value distributions of Jacobian matrices of deep function coupling

In Sec. A.6.1 we have proved that the Lyapunov exponents (if they exist) are limits of logarithms of subspace spectral norms divided by L . Here we use this property to prove the deficiency of numerical ranks, *i.e.*, Eq. (9).

We first introduce the Courant-Fischer min-max theorem [26] of singular values.

Theorem 9 (Courant-Fischer Min-max Theorem). *Let \mathbf{A} be a $d \times n$ complex matrix and $\sigma_i(\mathbf{A})$ denote its i -th largest singular value, $i = 1, \dots, \min\{d, n\}$. Then we have*

$$\sigma_i(\mathbf{A}) = \sup_{\dim(\mathcal{V})=i} \inf_{\mathbf{v} \in \mathcal{V}, \|\mathbf{v}\|_2=1} \|\mathbf{A}\mathbf{v}\|_2, \quad (\text{A73})$$

$$\sigma_i(\mathbf{A}) = \inf_{\dim(\mathcal{V})=n-i+1} \sup_{\mathbf{v} \in \mathcal{V}, \|\mathbf{v}\|_2=1} \|\mathbf{A}\mathbf{v}\|_2, \quad (\text{A74})$$

$$(\text{A75})$$

where \mathcal{V} traverses subspaces of \mathbb{R}^n .

This theorem also serves as one of the definitions to singular values.

Now assume that $\dim(\mathcal{L}_0 \setminus \mathcal{L}_1) = r$, and we consider only the case of $d = n$ for simplicity. For any $\epsilon > 0$, we have $N \in \mathbb{N}$ such that when $L > N$,

$$\inf_{\mathbf{v} \in \mathcal{L}_0 \setminus \mathcal{L}_1, \|\mathbf{v}\|_2=1} \|\mathbf{J}_F \mathbf{v}\|_2 \geq \exp L(\lambda_1 - \epsilon) \quad (\text{A76})$$

$$(\text{A77})$$

due to Eq. (A55). As $\dim(\mathcal{L}_0 \setminus \mathcal{L}_1) = r$, by Theorem 9, we have

$$\sigma_1(\mathbf{J}_F) \geq \cdots \geq \sigma_r(\mathbf{J}_F) \geq \exp L(\lambda_1 - \epsilon). \quad (\text{A78})$$

For $\mathbf{v} \in \mathcal{L}_0 = \mathbb{R}^n$ and $\|\mathbf{v}\|_2 = 1$, as $\mathcal{L}_0 = \mathcal{L}_0 \setminus \mathcal{L}_1 \oplus \cdots \oplus \mathcal{L}_h \setminus \mathcal{L}_{h+1}$ (\oplus denotes direct sum of linear quotient subspaces in the Banach space), there is $\mathbf{v}_i \in \mathcal{L}_{i-1} \setminus \mathcal{L}_i, i = 1, \dots, h+1$, such that

$$\|\mathbf{v}_1\|_2^2 + \cdots + \|\mathbf{v}_{h+1}\|_2^2 = 1 \quad (\text{A79})$$

and

$$\mathbf{v} = \mathbf{v}_1 + \cdots + \mathbf{v}_{h+1}. \quad (\text{A80})$$

Then by the conclusion of Sec. A.6.1, we have

$$\limsup_{L \rightarrow \infty} \|\mathbf{J}_F \mathbf{v}\|_2 \leq \limsup_{L \rightarrow \infty} \|\mathbf{J}_F \mathbf{v}_1\|_2 + \cdots + \limsup_{L \rightarrow \infty} \|\mathbf{J}_F \mathbf{v}_{h+1}\|_2 \quad (\text{A81})$$

$$\leq \|\mathbf{v}_1\|_2 \lim_{L \rightarrow \infty} \|\mathbf{J}_F\|_{2,2} + \cdots + \|\mathbf{v}_{h+1}\|_2 \lim_{L \rightarrow \infty} \|\mathbf{J}_F\|_{2,h+1} \leq \sum_{i=1}^{h+1} \|\mathbf{v}_i\|_2 \lambda_i \leq \lambda_1. \quad (\text{A82})$$

Thus there is $N_1 \in \mathbb{N}$, such that when $L > N_1$, we have

$$\sup_{\mathbf{v} \in \mathcal{L}_0, \|\mathbf{v}\|_2=1} \|\mathbf{J}_F \mathbf{v}\|_2 \leq \lambda_1 + \epsilon. \quad (\text{A83})$$

As $\dim(\mathcal{L}_1) = n - 1 + 1$, by Theorem 9, we have when $L > N_1$,

$$\sigma_r(\mathbf{J}_F) \leq \cdots \leq \sigma_1(\mathbf{J}_F) \leq \exp L(\lambda_1 + \epsilon). \quad (\text{A84})$$

In conclusion, when $L > N_0 = \max\{N, N_1\}$, we have

$$\exp L(\lambda_1 - \epsilon) \leq \sigma_1(\mathbf{J}_F) \leq \exp L(\lambda_1 + \epsilon). \quad (\text{A85})$$

Thus when $L \rightarrow \infty$, we have

$$\sigma_1(\mathbf{J}_F) \sim \exp L\lambda_1. \quad (\text{A86})$$

Using the same argument for $\sigma_2(\mathbf{J}_F), \dots, \sigma_n(\mathbf{J}_F)$, we can find that if let $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \cdots \geq \hat{\lambda}_n$ be the Lyapunov exponents counting repetitions, *i.e.*,

$$\hat{\lambda}_k = \lambda_i, \text{ if } \sum_{j=1}^{i-1} \dim(\mathcal{L}_{j-1} \setminus \mathcal{L}_j) < k \leq \sum_{j=1}^i \dim(\mathcal{L}_{j-1} \setminus \mathcal{L}_j), i = 1, \dots, h+1, \quad (\text{A87})$$

then

$$\sigma_i(\mathbf{J}_F) \sim \exp L\hat{\lambda}_i. \quad (\text{A88})$$

Note that

$$\hat{\lambda}_1 = \cdots = \hat{\lambda}_r = \lambda_1, \hat{\lambda}_i \leq \lambda_2 < \lambda_1, i = r+1, \dots, n. \quad (\text{A89})$$

Thus we have

$$\frac{\sigma_i(\mathbf{J}_F)}{\sigma_1(\mathbf{J}_F)} \sim \exp L(\hat{\lambda}_i - \lambda_1) \rightarrow 0, i = r+1, \dots, n. \quad (\text{A90})$$

As a consequence, $\text{Rank}_\epsilon(\mathbf{F}) \leq r$ for any $\epsilon > 0$ when $L \rightarrow \infty$.

A.6.3 Existence of Lyapunov exponents for Jacobian matrices of deep function coupling

In above analysis, we have proven Theorem 4 under the existence of Lyapunov exponents. In this section, we introduce the classical result of multiplicative ergodic theorem in the specific domain of random matrices, which is proposed by Furstenberg and Kesten [19, 20].

Theorem 10 (Multiplicative Ergodic Theorem (Theorem 3.9 of [20])). *Let μ be a probability measure on all invertible matrices of $\mathbb{R}^{n \times n}$ which satisfies*

$$\mathbb{E}_\mu[\max\{\log \|\mathbf{J}_{f^k}^{\pm 1}\|_2, 0\}] < \infty, k = 1, \dots, L. \quad (\text{A91})$$

If each \mathbf{J}_{f^k} independently follows μ , then we have a chain of subspaces $\{0\} = \mathcal{L}_{h+1} \subset \mathcal{L}_h \subset \dots \subset \mathcal{L}_1 \subset \mathcal{L}_0 = \mathbb{R}^n$ and corresponding positive real constants $\lambda_1 > \lambda_2 > \dots > \lambda_{h+1}$ such that almost surely

$$\lambda_i = \lim_{t \rightarrow \infty} \frac{1}{t} \log \|\mathbf{J}_{\mathbf{F}}^t \mathbf{v}\|_2, \forall \mathbf{v} \in \mathcal{L}_{i-1} \setminus \mathcal{L}_i, i = 1, \dots, h+1, \quad (\text{A92})$$

which means the existence of the Lyapunov exponents.

Combining this theorem and the arguments above, we can finally prove Theorem 4.

A.7 Proof to Theorem 5

This theorem can be deduced from the Lyapunov components of Ginibre matrices (polynomial ensemble of square matrices sampled *i.i.d* from standard Gaussian).

Theorem 11 (Exact Lyapunov Exponent Distribution for Ginibre Matrices [43]). *If μ in Theorem 10 is standard Gaussian, then $h+1 = n$, and*

$$\lambda_i = \log \left(2 + \psi \left(\frac{n-i+1}{2} \right) \right), i = 1, \dots, n. \quad (\text{A93})$$

Combining this theorem with Theorem 4 can directly yield our result.

A.8 Proof to Corollary 1

This theorem is the direct result of Theorems 4 and 11. Note that it is easy to get

$$\lambda_1 = \lim_{L \rightarrow \infty} \frac{1}{L} \log \|\mathbf{J}_{\mathbf{F}}\|_2 \quad (\text{A94})$$

for standard Gaussian μ .

B Influences of Structures

Skip Connection Skip Connection is the most direct method to solve rank diminishing at the initialization period. In our formulation, the definition of a layer network requires it to accept inputs purely from its predecessor layer as

$$\mathbf{x}^i = \mathbf{f}^i(\mathbf{x}^{i-1}), \mathbf{x}^{i-1} = \mathbf{f}^{i-1}(\mathbf{x}^{i-2}). \quad (\text{A95})$$

However, when we add a skip connection from its ancestor layer $\mathbf{f}^s, s < i-1$, we have

$$\mathbf{x}^i = \mathbf{f}^i(\mathbf{x}^{i-1}, \mathbf{x}^s), \mathbf{x}^{i-1} = \mathbf{f}^{i-1}(\mathbf{x}^{i-2}), \mathbf{x}^{i-2} = \mathbf{f}^{i-2} \circ \dots \circ \mathbf{f}^s(\mathbf{x}^{s-1}), \mathbf{x}^s = \mathbf{f}^s(\mathbf{x}^{s-1}). \quad (\text{A96})$$

It actually makes the coupling of layers

$$\hat{\mathbf{f}}^s = \left(\begin{array}{c} \mathbf{f}^{i-1} \circ \dots \circ \mathbf{f}^s \\ \mathbf{f}^s \end{array} \right), \quad (\text{A97})$$

the true predecessor layer to \mathbf{f}^i , as

$$\mathbf{x}^i = \mathbf{f}^i(\hat{\mathbf{x}}), \hat{\mathbf{x}} = \hat{\mathbf{f}}^s(\mathbf{x}^{s-1}). \quad (\text{A98})$$

Thus the true layer depth is cut down by $i-s$, remaining $L-(i-s)$ layers. Skip connection is usually used with the residual network. This structure can ease rank diminishing inside the layer $\hat{\mathbf{f}}^s$, which we will discuss later. Overall, skip connection shortens the length of the chain of Jacobian matrices, thus restraining rank diminishing.

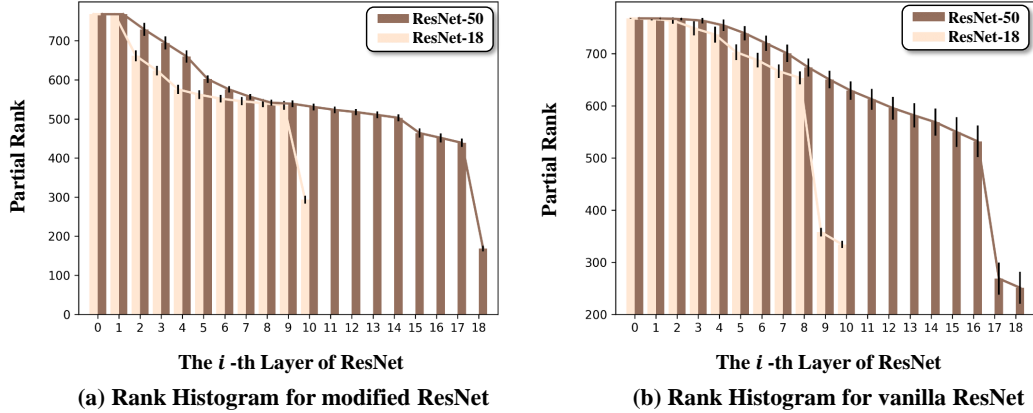


Figure A5: The partial rank of Jacobian matrices and perturbed PCA dimensions at the i -th layer of the modified ResNet-18 and ResNet-50 on ImageNet. We remove the operation of downsampling so that the feature dimension of the modified ResNet does not change (e.g., $H \times W \times C = 802, 816$ for ResNet-50). (a) The results of the modified ResNet-18 and ResNet-50. (b) The results of the initial ResNet-18 and ResNet-50. All the results are measured before training and using random initialization.

BatchNorm Some previous works [13, 5] discuss the role of BatchNorm in restraining rank diminishing. They show that BatchNorm may slow down the speed of rank diminishing in neural networks in some specific cases.

Residual Network Residual Network is another useful tool to restrain rank diminishing at the initialization period. The residual network r has the form

$$\mathbf{x}^o = r(\mathbf{x}^i) = \mathbf{x}^i + \text{Res}(\mathbf{x}^i), \quad (\text{A99})$$

where \mathbf{x}^o and \mathbf{x}^i are the output feature and input feature, respectively. Usually the residual term $\text{Res}(\mathbf{x}^i)$ is small compared with the input \mathbf{x}^i at the initialization period, as is pointed out by related works [22]. Assume that

$$\|\mathbf{J}_{\text{Res}}\|_2 < \epsilon, \quad (\text{A100})$$

where ϵ is very small. Then we have

$$\mathbf{J}_r = \mathbf{I} + \mathbf{J}_{\text{Res}} \quad (\text{A101})$$

is a diagonally dominant matrix, thus it has full rank. This means its kernel space $\text{Ker}(\mathbf{J}_r) = \{0\}$ is a zero dimension space. Thus by the Rank Theorem, for any predecessor layer \mathbf{f} , $r \circ \mathbf{f}$ will not lose rank as

$$\begin{aligned} \text{Rank}(r \circ \mathbf{f}) &= \text{Rank}(\mathbf{J}_r \mathbf{J}_f) = \text{Rank}(\mathbf{J}_f) - \dim(\text{Ker}(\mathbf{J}_r) \cap \text{Im}(\mathbf{J}_f)) \\ &= \text{Rank}(\mathbf{J}_f) = \text{Rank}(\mathbf{f}). \end{aligned} \quad (\text{A102})$$

However, in the well-trained ResNet18 and ResNet50 networks, we still observe considerable diminishing of ranks in Fig. 1. The reason for this phenomenon could be that, during training the magnitude of the residual term $\text{Res}(\mathbf{x})$ becomes large. Then the argument above no longer stand and hence the rank becomes lower. Taking the 16-th layer of the ResNet50 as example, we find that the rank of this layer drops from 530 to 119 after training, while the relative magnitude $\frac{\|\text{Res}(\mathbf{x})\|_2}{\|\mathbf{x}\|_2}$ increases from 0.5127 to 0.9557 after training. This may explain why the residual connection is less effective in preventing network ranks after training.

Influences of the Pooling Layers and Width To better validate the rank behavior of deep neural networks (e.g., ResNet), we remove the operation of downsampling in the ResNet so that the feature dimension (e.g., $x \in \mathbb{R}^{H \times W \times C=802,816}$ for ResNet-50) will not change. This modified ResNet can exclude the effect of pooling layers and changes of layer width. As shown in Fig.A5, we show the partial rank of Jacobian matrices and perturbed PCA dimensions at the i -th layer of the modified ResNet-18 and ResNet-50 on ImageNet. We can find that the curves of partial ranks and

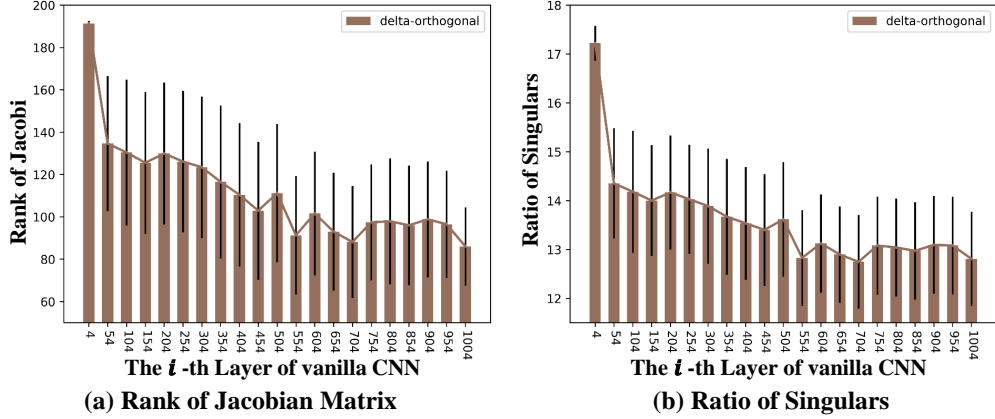


Figure A6: The rank of Jacobi matrices and perturbed PCA dimensions at the i -th layer of the 10,000 CNN initialized by Delta-Orthogonal on CIFAR-10.

perturbed PCA dimensions share a similar and consistent trend of decreasing as the initial networks. The consistent behavior of partial ranks and perturbed PCA dimensions also shows a monotonic decreasing property of network ranks. Thus the overall trend of rank diminishing seems to be independent of the pooling layers and changes of width. On the other hand, the partial rank witnesses a considerable drop near the terminal layer after applying pooling layers, which means it does have an effect on the network ranks.

C Orthogonal Initialization

Previous work [62] has demonstrated that under carefully designed initialization rule, we can train very deep (1,000 layer) plain CNNs. It is then curious to investigate that whether the phenomenon of rank diminishing happens in this case. To this end, we empirically measure the numerical rank of Jacobi matrices and perturbed PCA dimensions of an extremely deep CNN initialized by the Delta-Orthogonal method [62]. While computing the Jacobi matrices of very deep networks is infeasible in time, we only compute the 4-1,004 layer of the network. We omit the first 4 layers as there are downsampling architectures. We measure two metrics of each layer:

1. the numerical rank of the Jacobi matrix;
2. the ratio between sum of small singular values and large singular values $\frac{\sum_{i=11}^{3072} \sigma_i}{\sum_{i=1}^{10} \sigma_i}$, where 3,072 is the width of this network, and $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{3072}$ are singular values of the Jacobi matrix.

As shown in FigA6, a generally decreasing trend of the Jacobi rank and singular ratio can be observed. This result is consistent with Theorem 4. Although orthogonal initialization can suppress rank decrease to a certain extent, the impact of low rank cannot be ignored when the network is deep enough.

D Code

Algorithm A1 provides the pseudo-code of partial rank of the Jacobian. The implementation of the Algorithm A1 can refer to the 'rank_jacobian.py' python file.

Algorithm A2 provides the pseudo-code of perturbed PCA dimension of feature spaces. The implementation of the Algorithm A2 can refer to the 'rank_perturb.py' python file.

Algorithm A3 provides the pseudo-code of the classification dimension. The implementation of the Algorithm A3 can refer to the 'run_cls_dim.py' python file.

Algorithm A4 provides the pseudo-code of independence deficit. The implementation of the Algorithm A4 can refer to the 'run_deficit.py' python file.

Arch.	Network	Activ.	#Param.	Main Block	#Layer	Top-1 Acc.
ResNets	ResNet-18 [23]	ReLU [41]	11.7M	Bottleneck	11	69.8%
	ResNet-50 [23]	ReLU [41]	25.6M	Bottleneck	19	76.1%
MLP-like	GluMixer-24 [49]	SiLU [24]	25.0M	Mixer-Block	24	78.1%
	ResMLP-S24 [53]	GELU [24]	30.0M	Mixer-Block	24	79.4%
Transformer	ViT-T [16]	GELU [24]	5.7M	ViT-Block	13	75.5%
	Swin-T [37]	GELU [24]	29.0M	Swin-Block	18	81.3%

Table A2: Information of networks used in empirical validations. All pretrained on ImageNet.

E Experiments Setup

Information of those networks used in validations is listed in Tab. A2. When measuring rank, we set $\epsilon = \text{eps} \times N$, where eps is the digital accuracy of float32 (*i.e.*, $1.19e - 7$) and N is the number of singular values of the matrix to measure. This threshold represents the minimum digital accuracy of numerical rank we can capture in data stored as float32. All the experiments are conducted on the validation set of ImageNet and NVIDIA A100-SXM-80G GPUs.

F Partial Rank of the Jacobian: Estimating Lower Bound of Lost Rank in Deep Networks

To enable the validation of trend of the network ranks, we propose to compute only the rank of sub-matrices of the Jacobian as an alternative. Those sub-matrices are also the Jacobian matrices with respect to a fixed small patch of inputs. Rigorously, given a function \mathbf{f} and its Jacobian $\mathbf{J}_{\mathbf{f}}$, we denote partial rank of the Jacobian as the rank of a sub-matrix of the Jacobian that consists of the j_1 -th, j_2 -th, ..., j_K -th column of the original Jacobian

$$\text{PartialRank}(\mathbf{J}_{\mathbf{f}}) = \text{Rank}(\text{Sub}(\mathbf{J}_{\mathbf{f}}, j_1, \dots, j_K)) = \text{Rank}((\partial \mathbf{f}_i / \partial \mathbf{x}_{j_k})_{d \times K}), \quad (\text{A103})$$

where $1 \leq j_1 < \dots < j_K \leq n$. We can efficiently compute sub-matrix of the Jacobian by zero padding to small patches of input images. For any data point $\mathbf{x} \in \mathbb{R}^n$, let $\text{Sub}(\mathbf{x}, j_1, \dots, j_K) = (\mathbf{x}_{j_1}, \dots, \mathbf{x}_{j_K})^T \in \mathbb{R}^K$, and ψ pad $\text{Sub}(\mathbf{x}, j_1, \dots, j_K)$ to the spatial size of \mathbf{x} with zeros: $\psi(\text{Sub}(\mathbf{x}, j_1, \dots, j_K)) = (0, \dots, 0, \mathbf{x}_{j_1}, 0, \dots, \mathbf{x}_{j_K}, 0, \dots, 0)^T \in \mathbb{R}^n$ with $\psi(\text{Sub}(\mathbf{x}, j_1, \dots, j_K))_{j_k} = \mathbf{x}_{j_k}$, $k = 1, \dots, K$. We then have $\mathbf{J}_{\mathbf{f} \circ \psi} = \text{Sub}(\mathbf{J}_{\mathbf{f}}, j_1, \dots, j_K)$. As K can be very small compared with n , computing $\mathbf{J}_{\mathbf{f} \circ \psi}$ can be very cheap in time and space. The partial rank of Jacobian matrices of the network layers measures information captured among the spatial footprint j_1, \dots, j_K of the original input. They inherit the order relation of the rank of full Jacobian matrices. Thus we can validate the rank diminishing of network Jacobian matrices through the partial rank.

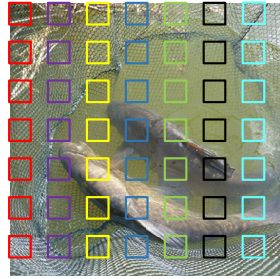
Lemma 2. For differentiable $\mathbf{f}_1, \mathbf{f}_2$, $|\text{Rank}(\mathbf{f}_1) - \text{Rank}(\mathbf{f}_2 \circ \mathbf{f}_1)| \geq |\text{Rank}(\text{Sub}(\mathbf{f}_1, j_1, \dots, j_K)) - \text{Rank}(\text{Sub}(\mathbf{f}_2 \circ \mathbf{f}_1, j_1, \dots, j_K))|, \forall 1 \leq K \leq n, 1 \leq j_1, \dots, j_K \leq n$. Thus variance of partial ranks of adjacent sub-networks gives a lower bound on the variance of their ranks.

F.1 Partial Rank of Jacobians under Different Input Patches

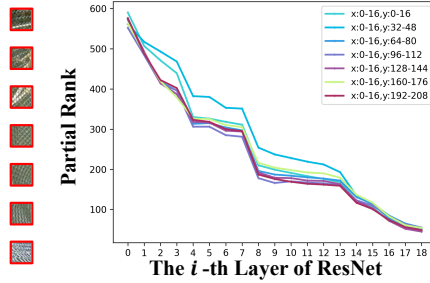
In Fig. A7 we report partial ranks of different input image patches (marked with colored boxes in Fig. A7(a)) for the layers of ResNet-50 on ImageNet. We can find that the curves of partial ranks share a similar and consistent trend among different input patches. Thus, picking one patch, for example, the central patch of $16 \times 16 \times 3$ pixels we use in Sec. 6.1, could be enough to demonstrate the overall behavior of network ranks. The consistent behavior of all those partial ranks also shows that partial rank is a good tool to investigate network ranks.

G Estimating Dimension Diminishing in Features

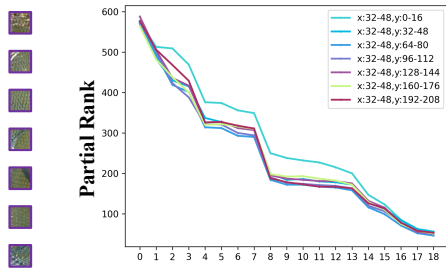
Measuring the intrinsic dimension of feature manifolds is known to be hard. However, we manage to give a rough estimation to the dimension dropped by different layer networks. To do this, we use a new metric called the Perturbed PCA Dimension. It measures the expectation of PCA dimension of small local neighborhoods over the feature manifold.



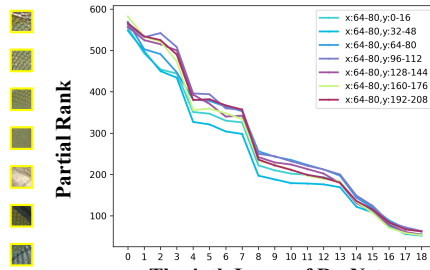
(a) Image



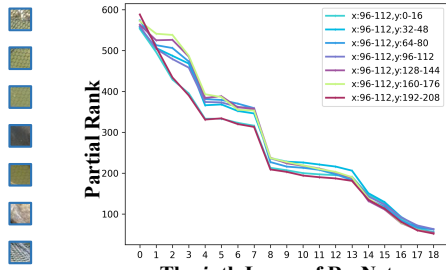
(b) Partial rank in the first column of the image



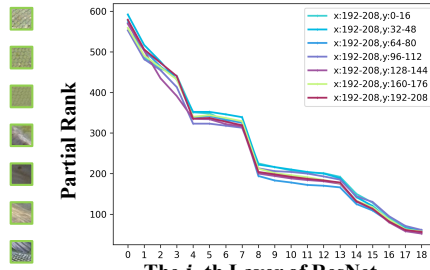
(c) Partial rank in the third column of the image



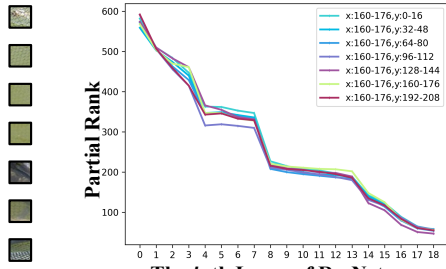
(d) Partial rank in the fifth column of the image



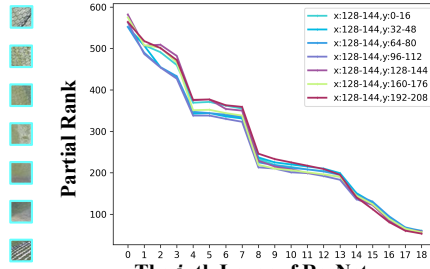
(e) Partial rank in the seventh column of the image



(f) Partial rank in the ninth column of the image



(g) Partial rank in the eleventh column of the image



(h) Partial rank in the thirteenth column of the image

Figure A7: The partial ranks of different input patches at the i -th layer of ResNet-50 on ImageNet.

Let F_k be the k -th sub-network of the whole network F . We want to measure the Perturbed PCA Dimension of $F_k(\mathcal{X})$, where \mathcal{X} is the input data domain. To this end, we compute

$$\text{PertDim} = \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\mathcal{X}}} [\text{PCADim}(\{F_k(\mathbf{x} + \epsilon) : \epsilon \sim \mathcal{N}(\mathcal{O}, \delta \mathcal{I})\})], \quad (\text{A104})$$

where PCADim for a set is the number of PCA eigenvalues larger than a threshold ξ . For each point \mathbf{x} , we sample 50,000 different perturbation ϵ to compute the PCA dimension of the neighborhood of $F_k(\mathbf{x})$. When computing the PCA dimension, we set $\delta = 1e-3$ and $\xi = 1.19e-7 \times 50000 \times \text{eig}_{\max}$, where eig_{\max} is the largest PCA eigenvalue. We then compute the mean value of PCA dimensions over the neighborhood of 100 random samples in the validation set of ImageNet as the final result.

We do not use PCA dimension of the feature manifolds directly as it is unable to cope with the highly non-linear structure of intermediate feature manifolds. However, the Perturbed PCA Dimension is able to estimate the dimensions of local neighborhoods of points in the feature manifolds. As local neighborhoods can be viewed as linear if the network is smooth, the Perturbed PCA Dimension could be more feasible than PCA dimension in our case. We provide the pseudo-code to compute the Perturbed PCA Dimension in Algorithm A2.

However, the perturbation is made in the ambient space of the input data manifold \mathcal{X} rather than the data manifold itself. Thus this estimation may considerably overestimate the intrinsic dimensions of feature manifolds. So we merely care about how many Perturbed PCA Dimensions are lost by a sub-network instead of its own Perturbed PCA Dimension. We call this quantity Δ Dimension, which is the difference between the Perturbed PCA Dimension of the current layer and that of the input layer for the given deep network. As shown in Fig. A8, we show the dropped dimensions of different feature layers of the CNN, MLP, and Transformer architectures on ImageNet. The results show that the Perturbed PCA Dimensions of feature manifolds of most networks decrease as the networks get deeper, thus confirming the rank diminishing principle we propose in Theorem 2.

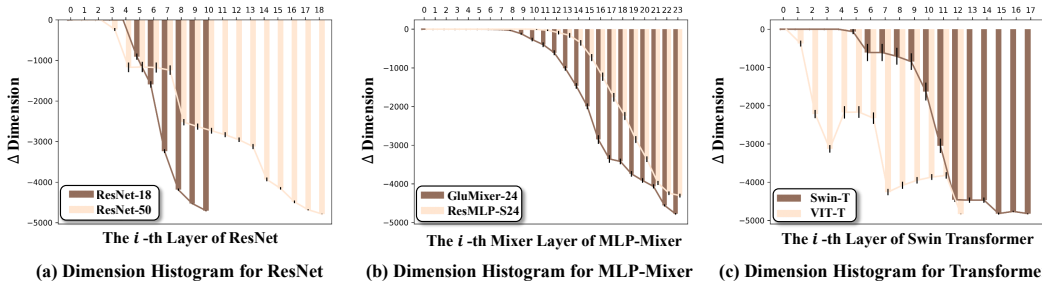


Figure A8: Dropped Perturbed PCA Dimension of different layers. Δ Dimension for the i -th layer is the difference between the Perturbed PCA Dimension of the i -th layer and that of the input layer of the CNN, MLP, and Transformer architectures on ImageNet.

Algorithm A1 Pseudocode of Partial Rank of the Jacobian.

```
# image: input images
# model: network
# row_idx, col_idx, patch_size: select a patch of the image to calculate Jacobian matrix

from functools import partial
import torch.nn.functional as functional

def Jacobian_rank(image, model):
    # select a patch of the image to calculate Jacobian matrix
    assert image.size(2) == image.size(3)
    image_size = image.size(2)
    image = image[:, :, row_idx:row_idx + patch_size, col_idx:col_idx + patch_size]
    zero_pad = partial(functional.pad, pad=[(image_size - patch_size) // 2 for _ in range(4)], value=0.)

    # calculate the jacobian matrix
    jacobian_matrix = jacobian(partial(net.forward, preprocess=zero_pad), image)

    # adopt trick to predict the singular values
    jacob = jacob.view(-1, image.size(1) * patch_size * patch_size)
    jacob = matmul(jacob.T, jacob)

    # calculate the partial rank of Jacobian matrix
    return matrix_rank(jacob, symmetric= True)
```

matmul: matrix multiplication; jacobian: calculate the jacobian matrix; matrix_rank: calculate the numerical rank of matrix.

Algorithm A2 Pseudocode of Perturbed PCA Dimension of Feature Spaces.

```
# image: input images
# model: network
# mag_perturb: magnitude of perturbations
# n_perturb: number of perturbations

def Perturbed_dimension(image, model, mag_perturb=1e-3, n_perturb=5000):
    # extract features with random perturbations
    features = []
    for _ in range(n_perturb):
        # sample random perturbation from Gaussian distribution
        perturb = randn_like(image) * mag_perturb
        # extract feature
        feature = model(image + perturb)
        features.append(feature)
    features = concatenation(features, dim=0)

    # calculate the covariance matrix
    x = input - mean(input, dim=0)
    x = x.view(x.size(0), -1)
    cov_matrix = matmul(x.T, x) # covariance matrix

    # calculate the perturbed PCA dimensions
    return matrix_rank(cov_matrix, symmetric= True)
```

matmul: matrix multiplication; randn_like: sample a random tensor from a Gaussian distribution; matrix_rank: calculate the numerical rank of matrix.

Algorithm A3 Pseudocode of the Classification Dimension of the Final Feature Manifold.

```
# image: input images
# model: network
# target: ground-truth labels
# acc_ratio: threshold for measuring intrinsic dimensions of final features

def PCA(X, n_components):
    n = X.shape[0]
    X_mean = mean(X, dim=0, keepdim=True)
    X = X - X_mean
    covariance_matrix = 1 / n * matmul(X.T, X)
    eigenvalues, eigenvectors = evd(covariance_matrix, eigenvectors=True)
    eigenvalues = norm(eigenvalues, dim=1) # modulus of complex numbers
    idx = argsort(-eigenvalues)
    eigenvectors = eigenvectors[:, idx]
    eigenvectors = eigenvectors[:, :n_components]
    return eigenvectors

def Feature_projection(X, V):
    X_proj = zeros_like(X)
    for component_idx in range(V.size(1)):
        eig_vec = V[:, component_idx].unsqueeze(-1)
        eig_vec_norm = eig_vec / norm(eig_vec, p=2, keepdim=True)
        w_proj = matmul(X, eig_vec_norm)
        X_proj_i = w_proj * eig_vec_norm.T
        X_proj += X_proj_i
    return X_proj

def Intrinsic_dimension(image, model, target, acc_ratio=0.95):
    # pre-extract features and calculate original classification accuracy
    feats = model(image) # [n_samples * n_channels]
    acc_ori = calc_acc(feats, target)

    for n_component in range(1, feats.size(1)):
        # compute the eigenvalues and eigenvectors of a real square matrix
        components = PCA(feats, n_component) # [n_channels * n_component]

        # reconstruct features with principal components
        feats_rec = Feature_projection(feats, components)

        # calculate classification accuracy
        acc = calc_acc(feats_rec, target)

        # return classification dimension
        if acc >= acc_ratio * acc_ori:
            return n_component
```

matmul: matrix multiplication; evd: eigen value decomposition; calc_acc: calculating classification accuracy.

Algorithm A4 Pseudocode of Independence Deficit.

```
# image: input images
# model: network
# target2index: dictionary mapping from category index to sample indices
# lr: learning rate for Lasso optimization
# n_iteration: number of iterations for Lasso optimization
# w_reg: weight of the L1 regularization term

def Feature_split(feats, class_i, target2index):
    sample_indices = target2index[class_i]
    start_idx, end_idx = sample_indices[0], sample_indices[-1]
    feats_i = feats[start_idx:end_idx+1, :]
    feats_i_n = concatenation((feats_i[:, :class_i], feats_i[:, class_i+1:]), dim=1)
    feats_i_p = feats_i[:, class_i:class_i+1]
    return feats_i_n, feats_i_p

def Independence_deficit(image, model, target2index, lr=1e-5, n_iteration=5000, w_reg=20.0):
    # pre-extract logits
    logits = model(image) # [n_samples * n_classes]

    # Lasso optimization
    for class_i in range(logits.size(1)):
        # split features by category index
        feats_n, feats_p = Feature_split(logits, class_i, target2index)

        # initialize the linear coefficients of category i
        param = Parameter(zeros(feats_n.size(1), 1))

        # start training
        for _ in range(n_iteration):
            loss = mse(matmul(feats_n, param), feats_p) + w_reg * l1_norm(param)
            loss.backward()
            param -= lr * param.grad

        # save trained coefficients
        save(param)
```

matmul: matrix multiplication; mse: mean squared error; l1_norm: sum of the magnitudes of the vectors in a space.