
Rank Diminishing in Deep Neural Networks

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The rank of neural networks measures information flowing across layers. It is
2 an instance of a key structural condition that applies across broad domains of
3 machine learning. In particular, the assumption of low-rank feature representations
4 led to algorithmic developments in many architectures. For neural networks,
5 however, the intrinsic mechanism that yields low-rank structures remains vague and
6 unclear. To fill this gap, we perform a rigorous study on the behavior of network
7 rank, focusing particularly on the notion of rank deficiency. We theoretically
8 establish a universal monotone decreasing property of network ranks from the
9 basic rules of differential and algebraic composition, and uncover rank deficiency
10 of network blocks and deep function coupling. By virtue of our numerical tools,
11 we provide the first empirical analysis of the per-layer behavior of network ranks
12 in realistic settings, *i.e.*, ResNets, deep MLPs, and Transformers on ImageNet.
13 These empirical results are in direct accord with our theory. Furthermore, we reveal
14 a novel phenomenon of independence deficit caused by the rank deficiency of
15 deep networks, where classification confidence of a given category can be linearly
16 decided by the confidence of a handful of other categories. The theoretical results
17 of this work, together with the empirical findings, may advance understanding of
18 the inherent principles of deep neural networks.

19 1 Introduction

20 In mathematics, the rank of a smooth function measures the volume of independent information
21 captured by the function [20]. Deep neural networks are highly smooth functions, thus the rank
22 of a network has long been an essential concept in machine learning that underlies many tasks
23 such as information compression [46, 54, 35, 52, 47], network pruning [31, 53, 5, 24, 9], data
24 mining [6, 23, 10, 55, 17, 28], computer vision [57, 56, 30, 26, 28, 58], and natural language
25 processing [8, 27, 7, 11]. Numerous methods are either designed to utilize the mathematical property
26 of network ranks, or are derived from an assumption that low-rank structures are to be preferred.

27 Yet a rigorous investigation to the behavior of rank of general networks, combining both theoretical
28 and empirical arguments, is still absent in current research, weakening our confidence in the being able
29 to predict performance. To the best of our knowledge, there are only a few previous works discussing
30 the rank behavior of specific network architectures, like attention blocks [14] and BatchNorms [12, 4]
31 in pure MLP structures. The empirical validation of those methods are also limited to shallow
32 networks, specific architectures, or merely the final layers of deep networks, leaving the global
33 behavior of general deep neural networks mysterious due to prohibitive space-time complexity for
34 measuring them. Rigorous work on network rank that combines both strong theoretical and empirical
35 evidence would have significant implications.

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

52 2 Preliminaries

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

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

56 For simplicity, we further write the k -th sub-network¹ of \mathbf{F} as

$$\mathbf{F}_k = \mathbf{f}^k \circ \dots \circ \mathbf{f}^1, \quad (2)$$

57 and we use $\mathcal{F}_k = \mathbf{F}_k(\mathcal{X})$ to denote the feature space of the k -th sub-network on the data domain \mathcal{X} .
 58 We are more interested in the behavior of network rank in the feature spaces rather than scalar outputs
 59 (which trivially have rank 1). Thus for classification or regression networks that output a scalar value,
 60 we will consider $\mathbf{F} = \mathbf{F}_L$ as the transformation from the input space to the final feature space instead.
 61 Thus, we always have $n \gg 1$ and $d \gg 1$. For example, for ResNet-50 [18] architecture on ImageNet,
 62 we only consider the network slice from the inputs to the last feature layer of 2,048 units.

63 **Rank of Function** The rank of a function $\mathbf{f} = (\mathbf{f}_1, \dots, \mathbf{f}_d)^T : \mathbb{R}^n \rightarrow \mathbb{R}^d$ refers to the rank of its
 64 Jacobi matrix \mathbf{J}_f over its input domain \mathcal{X} , which is defined as

$$\text{Rank}(\mathbf{f}) = \text{Rank}(\mathbf{J}_f) = \text{Rank}((\partial \mathbf{f}_i(\mathbf{x}) / \partial \mathbf{x}_j)_{n \times d}). \quad (3)$$

65 The rank of a function represents the volume of information captured by it in the output [20]. That is
 66 why it is so important to investigate the behavior of neural networks and many practical applications.
 67 Theoretically, by the Rank Theorem and Sard’s Theorem of manifolds [20], we can know that rank of
 68 the function equals the intrinsic dimension of its output feature space, as captured by the following
 69 lemma.²

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

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

¹In 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.

²Due to space limitation, all the related proofs are attached in the supplementary material.

77 3 Numerical Tools

78 Validating the rank behavior of deep neural networks is a challenging task because it involves
 79 operations of high complexity on large-scale non-sparse matrices, which is infeasible both in time
 80 and space. Computing the full Jacobian representation of sub-networks of ResNet-50, for example,
 81 consumes over 150G GPU memory and several days at a single input point. In accuracy, this is
 82 even more challenging as rank is very sensitive to small perturbations. The digital accuracy of
 83 float32, $1.19e - 7$ [38], cannot be trivially neglected in computing matrix ranks. Thus, in this section
 84 we establish some numerical tools for validating our subsequent arguments, and provide rigorous
 85 theoretical support for them.

86 3.1 Numerical Rank: Stable Alternative to Rank

87 The rank of large matrices is known to be unstable: it varies significantly under even small noise
 88 perturbations [40]. Matrices perturbed by even small Gaussian noises are almost surely of full rank,
 89 regardless of the true rank of the original matrix. Thus in practice we have to use an alternative: we
 90 count the number of singular values larger than some given threshold ϵ as the numerical rank of the
 91 matrix. Let $\mathbf{W} \in \mathbb{R}^{n \times d}$ be a given matrix. Its numerical rank with tolerance ϵ is

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

92 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
 93 values, and $\#$ is the counting measure for finite sets. We can prove that the numerical rank is stable
 94 under small perturbations. Based on Weyl inequalities [48], we have the following theorem.

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

99 This property of the numerical rank metric makes it a suitable tool for investigating the rank behavior
 100 of neural networks. Possible small noises can be filtered out in Jacobian matrices of networks by
 101 using numerical rank. It is worth mentioning that random matrices no longer have full rank almost
 102 surely under the numerical rank. Instead their rank distribution can be inferred from the well-known
 103 Marcenko–Pastur distribution [33] of random matrices. So under numerical rank, low-rank matrices
 104 will be commonly seen. In this paper, we always use the numerical rank when measuring ranks.

105 3.2 Partial Rank of the Jacobian: Estimating Lower Bound of Lost Rank in Deep Networks

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

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

111 where $1 \leq j_1 < \dots < j_K \leq n$. We can efficiently compute sub-matrix of the Jacobian by zero
 112 padding to small patches of input images. For any data point $\mathbf{x} \in \mathbb{R}^n$, let $\text{Sub}(\mathbf{x}, j_1, \dots, j_K) =$
 113 $(\mathbf{x}_{j_1}, \dots, \mathbf{x}_{j_K})^T \in \mathbb{R}^K$, and ψ pad $\text{Sub}(\mathbf{x}, j_1, \dots, j_K)$ to the spatial size of \mathbf{x} with zeros:
 114 $\psi(\text{Sub}(\mathbf{x}, j_1, \dots, j_K)) = (0, \dots, 0, \mathbf{x}_{j_1}, 0, \dots, \mathbf{x}_{j_K}, 0, \dots, 0)^T \in \mathbb{R}^n$ with $\psi(\text{Sub}(\mathbf{x}, j_1, \dots, j_K))_{j_k} =$
 115 $\mathbf{x}_{j_k}, k = 1, \dots, K$. We then have $\mathbf{J}_{f \circ \psi} = \text{Sub}(\mathbf{J}_f, j_1, \dots, j_K)$. As K can be very small compared
 116 with n , computing $\mathbf{J}_{f \circ \psi}$ can be very cheap in time and space. The partial rank of Jacobian matrices
 117 of the network layers measures information captured among the spatial footprint j_1, \dots, j_K of the
 118 original input. They inherit the order relation of the rank of full Jacobian matrices. Thus we can
 119 validate the rank diminishing of network Jacobian matrices through the partial rank.

120 **Lemma 2.** *For differentiable $\mathbf{f}_1, \mathbf{f}_2$, $|\text{Rank}(\mathbf{f}_1) - \text{Rank}(\mathbf{f}_2 \circ \mathbf{f}_1)| \geq |\text{Rank}(\text{Sub}(\mathbf{f}_1, j_1, \dots, j_K)) -$
 121 $\text{Rank}(\text{Sub}(\mathbf{f}_2 \circ \mathbf{f}_1, j_1, \dots, j_K))|, \forall 1 \leq K \leq n, 1 \leq j_1, \dots, j_K \leq n$. Thus variance of partial
 122 ranks of adjacent sub-networks gives a lower bound on the variance of their ranks.*

123 3.3 Classification Dimension: Estimating Final Feature Dimension

124 Measuring the intrinsic dimension of feature manifolds is known to be intractable. So we turn to
 125 an approximation procedure. For most classification networks, a linear regression over the final
 126 feature manifold decides the final network prediction and accuracy. So we can estimate the intrinsic
 127 dimension as the minimum number of principal components in the final feature space to preserve a
 128 high classification accuracy. Given network slice $F : \mathbb{R}^n \rightarrow \mathbb{R}^d$ from input $\mathcal{X} \subset \mathbb{R}^n$ to final feature
 129 space $F(\mathcal{X}) \subset \mathbb{R}^d$, we independently sample N points from random variable $F(\mathbf{x})$, $\mathbf{x} \sim \mathbb{P}_{\mathcal{X}}$, where
 130 $\mathbb{P}_{\mathcal{X}}$ is the distribution of validation set of data \mathcal{X} . We then compute the covariance matrix Σ of those
 131 N samples, and eigenvectors $\mathbf{q}^1, \dots, \mathbf{q}^d$ of Σ , sorted by their eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$.
 132 Let $\text{cls} : \mathbb{R}^d \rightarrow \mathbb{R}^c$ be the classification predictions based on the final feature representation $F(\mathbf{x})$,
 133 Pro_k be projection operator in Euclidean space that projects to the linear subspace spanned by top- K
 134 eigenvectors $\mathbf{q}^1, \dots, \mathbf{q}^k$, $k \leq d$, $\mathbb{P}_{\mathbf{x}, \mathbf{y}}$ be the joint distribution of sample \mathbf{x} and its label \mathbf{y} , and $\mathbf{1}_{\text{cond}}$
 135 the indicator for condition cond . The classification dimension is then defined as

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

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

138 4 Principle of Rank Diminishing

139 We turn to the *principle of rank diminishing*. We first give a universal justification with minimum
 140 limitation on the network, so that we can safely apply this principle to many practical scenarios.

141 The principle of rank diminishing describes the behavior of general neural networks with almost
 142 everywhere smooth components, which exhibits the monotone decrease of network ranks and intrinsic
 143 dimensionality of feature manifolds as follows.

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

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

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

148 **Short Argument that the Principle Should Hold Universally.** Theorem 2 is ultra intrinsic for
 149 deep neural networks. It comes directly from the chain rules of differential and basic rules of matrix
 150 multiplications. The basic rule of matrix multiplication tells that, for any two matrices \mathbf{A} and \mathbf{B} , we
 151 have $\text{Rank}(\mathbf{AB}) \leq \min\{\text{Rank}(\mathbf{A}), \text{Rank}(\mathbf{B})\}$ [21]. Taking this into the chain rule of differential
 152 of $\mathbf{J}_F = \mathbf{J}_{\mathbf{f}^L} \mathbf{J}_{\mathbf{f}^{L-1}} \dots \mathbf{J}_{\mathbf{f}^1}$, we then have $\text{Rank}(\mathbf{J}_{F_k}) = \text{Rank}(\mathbf{J}_{\mathbf{f}^k \circ F_{k-1}}) = \text{Rank}(\mathbf{J}_{\mathbf{f}^k} \mathbf{J}_{F_{k-1}}) \leq$
 153 $\text{Rank}(\mathbf{J}_{F_{k-1}})$, $k = 2, \dots, L$, which is Eq. (7). Applying Lemma 1 to Eq. (7) then yields Eq. (8).

154 **Chance of Equal Signs Holding is Small.** A hypothetical but not practical concern would be that,
 155 is it possible that most of the equal signs of Eqs. (7) and (8) hold, so that the rank of network remains
 156 no significant dropping throughout the network? This concern can be mitigated by empirical and
 157 theoretical arguments. In what follows we will find that, 1) in practice, the rank of sub-networks
 158 decreases significantly after applying subsequent layers as shown in Fig. 1, and 2) in theory, there are
 159 two strong impetuses in deep neural networks to enforce strict decreasing of ranks which we will
 160 discuss in Secs. 4.1 and 4.2.

161 4.1 Structural Impetus of Strict Decreasing

162 Numerous explicit structures of the network layers can lead to a strict decrease in network ranks.
 163 Specifically, the following theorem gives a condition for the strictly greater signs to hold in the
 164 principle of rank diminishing.

Arch.	Network	Activ.	#Param.	Main Block	#Layer	Top-1 Acc.
ResNets	ResNet-18 [18]	ReLU [36]	11.7M	Bottleneck	11	69.8%
	ResNet-50 [18]	ReLU [36]	25.6M	Bottleneck	19	76.1%
MLP-like	GluMixer-24 [41]	SiLU [19]	25.0M	Mixer-Block	24	78.1%
	ResMLP-S24 [44]	GELU [19]	30.0M	Mixer-Block	24	79.4%
Transformer	ViT-T [15]	GELU [19]	5.7M	ViT-Block	13	75.5%
	Swin-T [32]	GELU [19]	29.0M	Swin-Block	18	81.3%

Table 1: Information of networks used in empirical validations. All pretrained on ImageNet.

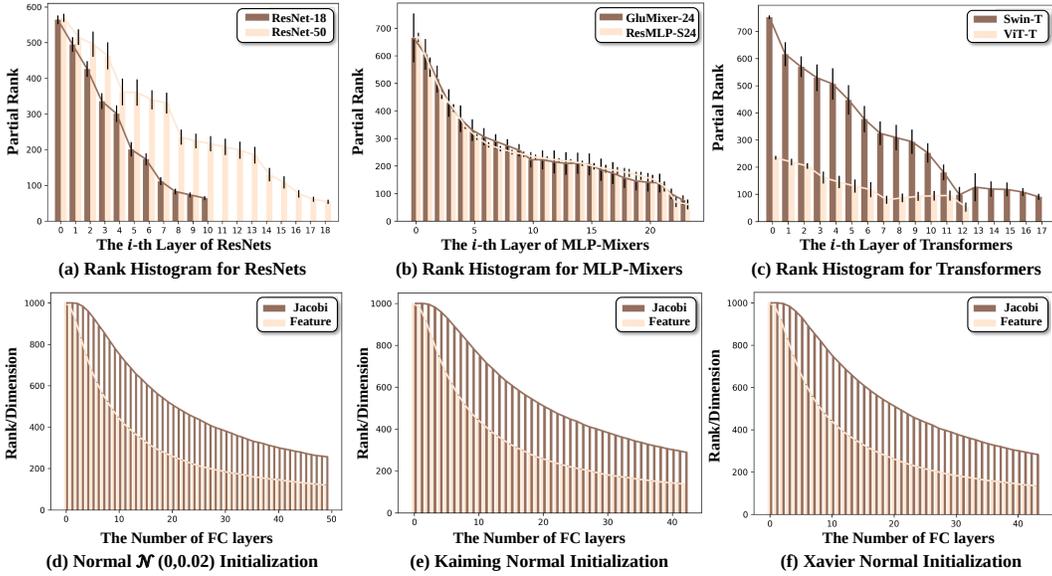


Figure 1: Partial rank of Jacobian matrices of CNN, MLP, and Transformer architecture networks for different layers on ImageNet (top row); rank of Jacobian matrices and feature dimensions of linear MLP network following conditions of Theorem 5 (bottom row). All the models show a similar trend of exponential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

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

169 By this theorem, one can immediately find that most frequently used layer designs have high
170 risk in inducing strict decreasing of network ranks. Normalization layers like LayerNorm [2],
171 InstanceNorm [45], and BatchNorm [25] may lose dimensions modestly, as the output feature
172 remains invariant along the normalized direction at each point. Linear layers like convolutions, linear
173 transformations (e.g. dense layers), and attentions, can lose rank considerably according to the rank
174 of their weight matrices. They constitute the explicit structural impetus to decrease network ranks
175 and intrinsic dimensions of feature manifolds.

176 4.2 Implicit Impetus of Strict Decreasing

177 Apart from the structural impetus we propose in Sec. 4.1, there is a more intrinsic strength to pull
178 down network ranks, which we call the implicit impetus. Deep neural networks repeatedly apply
179 layer networks from a fixed function pool (ReLU, MLP, CNN, attention, ResNet block, etc.) to
180 the input data and intermediate features to get outputs. Such paradigm accords with the cocycle
181 dynamic systems studied by Lyapunov *et al.* [42, 50], where the Furstenberg–Kesten theorem [16]
182 and multiplicative ergodic theorem [39] prove that logarithms of singular values divided by evolution

³The rigorous version is given in the supplementary material.

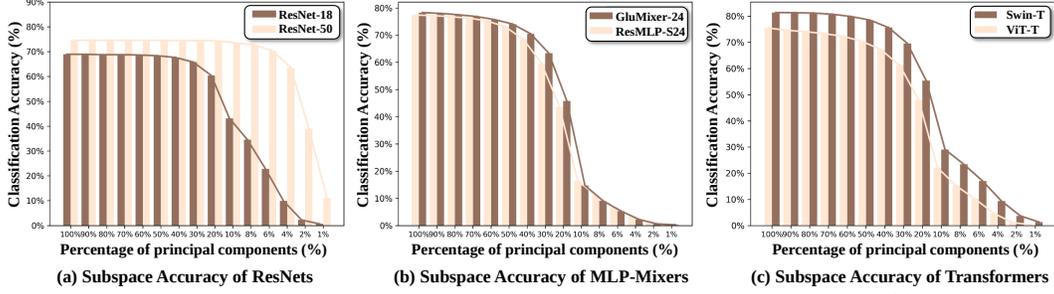


Figure 2: Classification Accuracy (top-1) of using subspaces spanned by top- k % eigenvectors (principal components) of the final feature manifolds. For all networks a small percentage (see Tab. 2) of eigenvectors are enough to reproduce the classification accuracy of the whole network, indicating a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

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 2: 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.

183 time of such chaos system converge to stable constants when time goes to infinity. While products of
 184 long chains of matrices are the simplest form of cocycle dynamic systems [29], we can get an intrinsic
 185 impetus of rank collapse tendency of Jacobian matrices independent of network architectures.

186 **Theorem 4** (Spectrum Centralization of Function Coupling). *Let the network be $\mathbf{F} = \mathbf{f}^L \circ \dots \circ \mathbf{f}^1$,
 187 and all the ambient dimensions of feature manifolds be the same as the ambient dimension of inputs,
 188 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
 189 follows some distribution μ , and $\mathbb{E}_\mu[\max\{\log \|\mathbf{J}_{\mathbf{f}^k}\|_2, 0\}] < \infty$. Let σ_k denotes the k -th largest
 190 singular value of $\mathbf{J}_{\mathbf{F}}$. Then there is an integer $r < n$ and positive constants μ_r, \dots, μ_n that only
 191 depend on μ such that we have for μ -almost everywhere,*

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

192 *That means for any tolerance $\epsilon > 0$, $\text{Rank}_\epsilon(\mathbf{F})$ drops below $r + 1$ with an exponential speed as
 193 $L \rightarrow \infty$.*

194 If further assuming Gaussian distributions of Jacobian matrices, we can prove $r = 1$ and give a more
 195 accurate estimation of the constant μ_r, \dots, μ_n . As a consequence, we can find that rank of networks
 196 collapses to 1 almost surely, which is formalized in the following theorem.

197 **Theorem 5.** *Let the network be $\mathbf{F} = \mathbf{f}^L \circ \dots \circ \mathbf{f}^1$, and all the ambient dimensions of feature
 198 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
 199 that $\mathbf{J}_{\mathbf{f}^i}$ independently follows the standard Gaussian distribution. Let σ_k denotes the k -th largest
 200 singular value of $\mathbf{J}_{\mathbf{F}}$. Then almost surely*

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

201 *where $\psi = \Gamma/\Gamma'$ and Γ is the Gamma function. That means for a large L and any tolerance ϵ ,
 202 $\text{Rank}_\epsilon(\mathbf{F})$ drops to 1 exponentially with speed nC^L , where $C < 1$ is a positive constant that only
 203 depends on n .*

204 **Connection with Gradient Explosion** Bengio *et. al.* [3, 37] discuss the gradient explosion issue of
 205 deep neural networks, where the largest singular value of the Jacobian matrix tends to infinity when
 206 the layer gets deeper. This problem could be viewed as a special case of Theorem 5 that investigates
 207 the behavior of all singular values of deep neural networks. The behavior of network ranks in fact
 208 manipulates the well-known gradient explosion issue. Rigorously, we have the following conclusion.

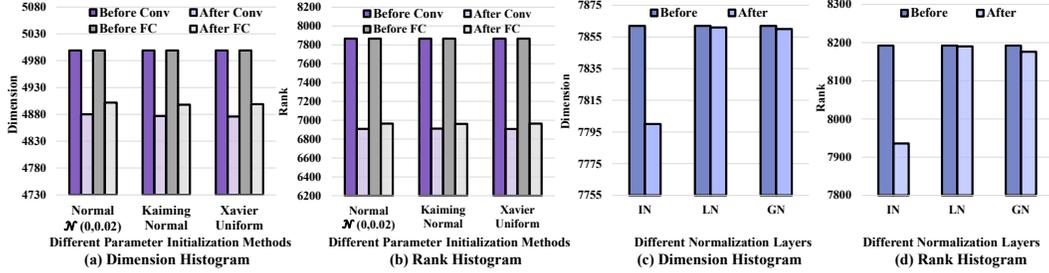


Figure 3: PCA dimension of feature spaces and rank of Jacobian matrix for commonly seen network components under standard Gaussian inputs and randomized weights. Convolution and FC layers tend to lose rank considerably; normalization layers, like InstanceNorm (IN) [45], LayerNorm (LN) [2], and GroupNorm (GN) [51], lose rank modestly. But none can preserve rank.

209 **Corollary 1.** *Under the condition of Theorem 5, then almost surely gradient explosion happens at an*
 210 *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.*

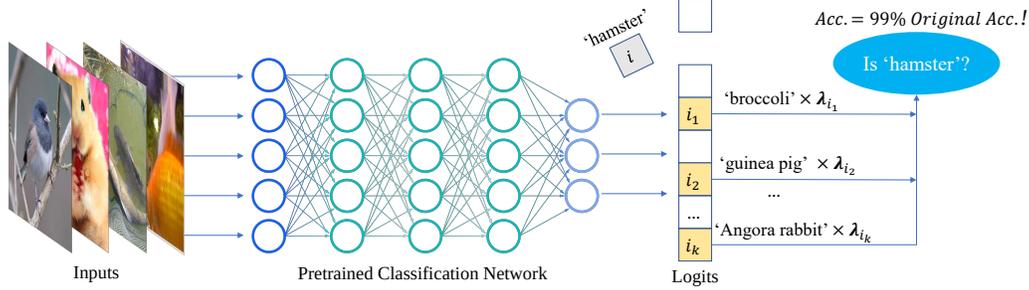
211 4.3 Validations

212 **Setup** In this section, we validate our theory in three types of architectures of benchmark deep
 213 neural networks, CNNs, MLPs, and Transformers, in the ImageNet [13] data domain. Information of
 214 those networks is listed in Tab. 1. For validating the tendency of network rank of Jacobian matrices,
 215 we use the numerical rank of sub-matrices of Jacobian on the central $16 \times 16 \times 3$ image patch of input
 216 images. We report the results of other choices of patches in the Appendix. When measuring rank, we
 217 set $\epsilon = \text{eps} \times N$, where eps is the digital accuracy of float32 (i.e., $1.19e - 7$) and N is the number
 218 of singular values of the matrix to measure. This threshold represents the minimum digital accuracy
 219 of numerical rank we can capture in data stored as float32. All the experiments are conducted on the
 220 validation set of ImageNet and NVIDIA A100-SXM-80G GPUs.

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

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

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



(a) Independence deficit: Classification confidence of some categories are linearly decided by a few other categories with fixed coefficients.

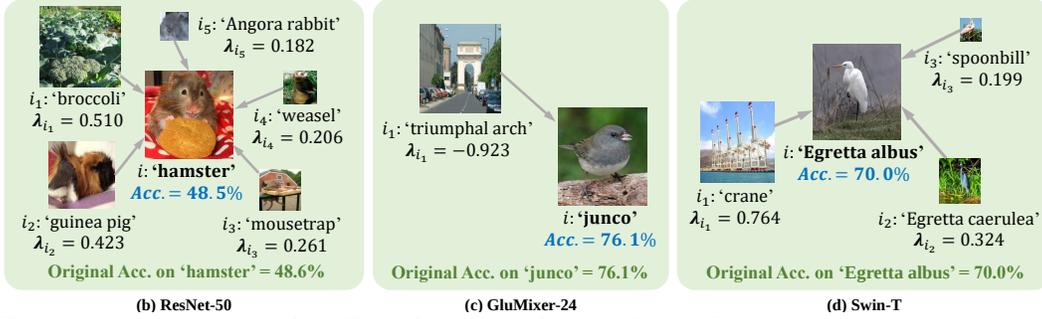


Figure 4: Independence deficit. Classification confidence of some ImageNet categories are linearly decided by a few other categories with fixed coefficients in the whole data domain. We illustrate this phenomenon in (a). Here we present some results from ResNet-50, GluMixer-24, and Swin-T. In the (b,c,d) 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.

245 a striking resemblance to the exponential law of our toy setting, which confirms the proposed implicit
 246 impetus in those models.

247 **Structural Impetus** We validate the structural impetus in Fig. 3. To give an estimation for general
 248 cases, here we use Gaussian noises with the size of $128 \times 8 \times 8$ as inputs, and randomize weights of the
 249 network components to be validated. We plug those components into a simple fully-connected (FC)
 250 layer of 8,192 hidden units. As the structure is simple, we directly measure the intrinsic dimension of
 251 feature spaces and the full rank of Jacobian matrices before and after the features pass the network
 252 components to be measured. The dimension is determined by the number of PCA eigenvalues [22, 49]
 253 larger than $1.19e - 7 \times N \times \sigma_{\max}$, where N is the number of PCA eigenvalues, and σ_{\max} is the
 254 largest PCA eigenvalue. The batch size is set to 5,000. We find convolution (the kernel size is 3×3)
 255 and FC layers (the weight size is 8,192) tend to lose rank considerably, while different normalization
 256 layers also lose rank modestly. But none of them can preserve rank invariant.

257 **Possible Remission Approaches to Rank Diminishing** There are quite some techniques, at least
 258 in theory, can remiss the network rank diminishing. Typical examples are skip connection [14] and
 259 BatchNorm [12], which we will discuss in the Appendix due to page limitation.

260 5 Independence Deficit of Final Feature Manifolds

261 In this section, we provide a further perspective to study the low-rank structure of the final feature
 262 manifold, which induces an interesting finding of independence deficit in deep neural networks. We

263 have already known that the final feature representations of deep neural networks admit a very low
 264 intrinsic dimension. Thus there are only a few independent representations to decide the classification
 265 scores for all the 1,000 categories of ImageNet. It is then curious whether we can predict the outputs
 266 of the network for some categories based on the outputs for a few other categories, as illustrated in
 267 Fig. 4 (a). And if we can, will those categories be strongly connected to each other? A surprising fact
 268 is that, we can find many counter examples of irrelevant categories dominating the network outputs
 269 for given categories regarding various network architectures. This interesting phenomenon indicates
 270 a rather drastic competing in the final feature layer for the tight rank budgets of all categories, which
 271 yields non-realistic dependencies of different categories.

272 To find the dependencies of categories in final features, we can solve the following Lasso problem [43],

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

273 where $\mathbf{F}(\mathbf{x}) \in \mathbb{R}^{1000}$ is the slice of network from inputs to the final feature representation, \mathbf{x} is the
 274 sample from ImageNet \mathcal{X} , and \mathbf{W} is the final dense layer. The solution λ^* will be a sparse vector,
 275 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

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

276 where $\text{logits}(\mathbf{x}, i_j), j = 1, \dots, k$ is the logits of network for category i_j , *i.e.*, $\text{logits}(\mathbf{x}, i_j) =$
 277 $\mathbf{W}_{i_j} \mathbf{F}(\mathbf{x})$. It is easy to see that outputs for category i is linearly decided by outputs for i_1, \dots, i_k and
 278 is dominated by outputs for i_1 .

279 In Fig. 4 we demonstrate the solutions of Eq. (12) for three different categories in ImageNet with
 280 $\eta = 20$, and network architectures ResNet-50, GluMixer-24, and Swin-T. The results are surprising.
 281 It shows that many categories of the network predictions are in fact ‘redundant’, as they are purely
 282 decided by the predictions of the other categories with simple linear coefficients. In this case, the
 283 entanglement of different categories cannot be avoided, thus the network may perform poorly under
 284 domain shift. An even more surprising finding is that, some very irrelevant categories hold the largest
 285 weights when deciding the predictions of the redundant categories. This means that the networks
 286 just neglect the unique representations of those categories in training and yield over-fitting when
 287 predicting them.

288 6 Related Work

289 Previous studies of rank deficiency in deep neural networks follow two parallel clues. One is the
 290 study of rank behavior in specific neural network architectures. [14] studies deep networks consisting
 291 of pure self-attention networks, and proves that they converge exponentially to a rank-1 matrix
 292 under the assumption of globally bounded weight matrices. [12] studies the effect of BatchNorm on
 293 MLPs and shows that BatchNorm can prevent drastic diminishing of network ranks in some small
 294 networks and datasets. Both of those works avoid directly validating the behavior of network ranks in
 295 intermediate layers due to the lacking of efficient numerical tools. An independent clue is the study
 296 of implicit self-regularization, which finds that weight matrices tend to lose ranks after training. [34]
 297 studies this phenomenon in infinitely-wide, over-parametric neural networks with tools from random
 298 matrix theory. [1] studies this phenomenon in deep matrix decomposition. Those works focus on the
 299 theoretical behavior of network ranks induced by the training process instead of network depth.

300 7 Conclusion

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

References

- 310
- 311 [1] S. Arora, N. Cohen, W. Hu, and Y. Luo. Implicit regularization in deep matrix factorization. *Adv. Neural*
312 *Inform. Process. Syst.*, 32, 2019.
- 313 [2] J. L. Ba, J. R. Kiros, and G. E. Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- 314 [3] Y. Bengio, P. Simard, and P. Frasconi. Learning long-term dependencies with gradient descent is difficult.
315 *IEEE Trans. Neural Netw.*, 5(2):157–166, 1994.
- 316 [4] N. Bjorck, C. P. Gomes, B. Selman, and K. Q. Weinberger. Understanding batch normalization. *Adv.*
317 *Neural Inform. Process. Syst.*, 31, 2018.
- 318 [5] C. Blakenev, Y. Yan, and Z. Zong. Is pruning compression?: Investigating pruning via network layer
319 similarity. In *IEEE Winter Conf. Appl. Comput. Vis.*, pages 914–922, 2020.
- 320 [6] X. Cai, C. Ding, F. Nie, and H. Huang. On the equivalent of low-rank linear regressions and linear
321 discriminant analysis based regressions. In *Proc. ACM SIGKDD Int. Conf. Knowl. Discov. Data Min.*,
322 pages 1124–1132, 2013.
- 323 [7] B. Chen, T. Dao, E. Winsor, Z. Song, A. Rudra, and C. Ré. Scatterbrain: Unifying sparse and low-rank
324 attention. *Adv. Neural Inform. Process. Syst.*, 34, 2021.
- 325 [8] P. Chen, S. Si, Y. Li, C. Chelba, and C.-J. Hsieh. Groupreduce: Block-wise low-rank approximation for
326 neural language model shrinking. *Adv. Neural Inform. Process. Syst.*, 31, 2018.
- 327 [9] P. Chen, H.-F. Yu, I. Dhillon, and C.-J. Hsieh. Drone: Data-aware low-rank compression for large NLP
328 models. *Adv. Neural Inform. Process. Syst.*, 34:29321–29334, 2021.
- 329 [10] Y. Cheng, L. Yin, and Y. Yu. Lorslim: Low rank sparse linear methods for top-n recommendations. In
330 *IEEE Int. Conf. on Data Min.*, pages 90–99. IEEE, 2014.
- 331 [11] J. Chiu, Y. Deng, and A. Rush. Low-rank constraints for fast inference in structured models. *Adv. Neural*
332 *Inform. Process. Syst.*, 34, 2021.
- 333 [12] H. Daneshmand, J. Kohler, F. Bach, T. Hofmann, and A. Lucchi. Batch normalization provably avoids
334 ranks collapse for randomly initialised deep networks. *Adv. Neural Inform. Process. Syst.*, 33:18387–18398,
335 2020.
- 336 [13] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A large-scale hierarchical image
337 database. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 248–255. Ieee, 2009.
- 338 [14] Y. Dong, J.-B. Cordonnier, and A. Loukas. Attention is not all you need: Pure attention loses rank doubly
339 exponentially with depth. In *Int. Conf. Mach. Learn.*, pages 2793–2803. PMLR, 2021.
- 340 [15] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani,
341 M. Minderer, G. Heigold, S. Gelly, et al. An image is worth 16x16 words: Transformers for image
342 recognition at scale. In *Int. Conf. Learn. Represent.*, 2020.
- 343 [16] H. Furstenberg and H. Kesten. Products of random matrices. *Ann. Math. Stat.*, 31(2):457–469, 1960.
- 344 [17] D. Goldfarb and Z. Qin. Robust low-rank tensor recovery: Models and algorithms. *SIAM J. Matrix Anal.*
345 *Appl.*, 35(1):225–253, 2014.
- 346 [18] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *IEEE Conf. Comput.*
347 *Vis. Pattern Recog.*, pages 770–778, 2016.
- 348 [19] D. Hendrycks and K. Gimpel. Gaussian error linear units (GELUs). *arXiv preprint arXiv:1606.08415*,
349 2016.
- 350 [20] M. W. Hirsch. *Differential topology*, volume 33. Springer Science & Business Media, 2012.
- 351 [21] K. Hoffman. *Linear algebra*. Englewood Cliffs, NJ, Prentice-Hall, 1971.
- 352 [22] H. Hotelling. Relations between two sets of variates. *Biometrika*, 28(3/4):321–377, 1936.
- 353 [23] C.-J. Hsieh, K.-Y. Chiang, and I. S. Dhillon. Low rank modeling of signed networks. In *Proc. ACM*
354 *SIGKDD Int. Conf. Knowl. Discov. Data Min.*, pages 507–515, 2012.
- 355 [24] Y.-C. Hsu, T. Hua, S. Chang, Q. Lou, Y. Shen, and H. Jin. Language model compression with weighted
356 low-rank factorization. In *Int. Conf. Learn. Represent.*, 2021.

- 357 [25] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal
358 covariate shift. In *Int. Conf. Mach. Learn.*, pages 448–456. PMLR, 2015.
- 359 [26] L. Jing, L. Yang, J. Yu, and M. K. Ng. Semi-supervised low-rank mapping learning for multi-label
360 classification. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1483–1491, 2015.
- 361 [27] R. Karimi Mahabadi, J. Henderson, and S. Ruder. Compacter: Efficient low-rank hypercomplex adapter
362 layers. *Adv. Neural Inform. Process. Syst.*, 34, 2021.
- 363 [28] M. Kheirandishfard, F. Zohrizadeh, and F. Kamangar. Deep low-rank subspace clustering. In *IEEE Conf.
364 Comput. Vis. Pattern Recog. Worksh.*, pages 864–865, 2020.
- 365 [29] J. F. C. Kingman. Subadditive ergodic theory. *Ann. Probab.*, pages 883–899, 1973.
- 366 [30] S. Kong and C. Fowlkes. Low-rank bilinear pooling for fine-grained classification. In *IEEE Conf. Comput.
367 Vis. Pattern Recog.*, pages 365–374, 2017.
- 368 [31] M. Lin, R. Ji, Y. Wang, Y. Zhang, B. Zhang, Y. Tian, and L. Shao. Hrank: Filter pruning using high-rank
369 feature map. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1529–1538, 2020.
- 370 [32] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo. Swin transformer: Hierarchical vision
371 transformer using shifted windows. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 10012–10022, 2021.
- 372 [33] V. A. Marčenko and L. A. Pastur. Distribution of eigenvalues for some sets of random matrices. *Math.
373 USSR Sb.*, 1(4):457, 1967.
- 374 [34] C. H. Martin and M. W. Mahoney. Implicit self-regularization in deep neural networks: Evidence from
375 random matrix theory and implications for learning. *J. Mach. Learn. Res.*, 22(165):1–73, 2021.
- 376 [35] J. Mu, R. Xiong, X. Fan, D. Liu, F. Wu, and W. Gao. Graph-based non-convex low-rank regularization for
377 image compression artifact reduction. *IEEE Trans. Image Process.*, 29:5374–5385, 2020.
- 378 [36] V. Nair and G. E. Hinton. Rectified linear units improve restricted boltzmann machines. In *Int. Conf. Mach.
379 Learn.*, pages 807–814. PMLR, 2010.
- 380 [37] R. Pascanu, T. Mikolov, and Y. Bengio. On the difficulty of training recurrent neural networks. In *Int. Conf.
381 Mach. Learn.*, pages 1310–1318. PMLR, 2013.
- 382 [38] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein,
383 L. Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Adv. Neural Inform.
384 Process. Syst.*, 32, 2019.
- 385 [39] Y. B. Pesin. Characteristic Lyapunov exponents and smooth ergodic theory. *Russ. Math. Surv.*, 32(4):55,
386 1977.
- 387 [40] A. Quarteroni, R. Sacco, and F. Saleri. *Numerical mathematics*, volume 37. Springer Science & Business
388 Media, 2010.
- 389 [41] N. Shazeer. Glu variants improve transformer. *arXiv preprint arXiv:2002.05202*, 2020.
- 390 [42] R. Temam. *Infinite-dimensional dynamical systems in mechanics and physics*, volume 68. Springer Science
391 & Business Media, 2012.
- 392 [43] R. Tibshirani. Regression shrinkage and selection via the lasso. *J. R. Stat. Soc. Series B Stat. Methodol.*,
393 58(1):267–288, 1996.
- 394 [44] H. Touvron, P. Bojanowski, M. Caron, M. Cord, A. El-Nouby, E. Grave, G. Izacard, A. Joulin, G. Synnaeve,
395 J. Verbeek, et al. Resmlp: Feedforward networks for image classification with data-efficient training. *arXiv
396 preprint arXiv:2105.03404*, 2021.
- 397 [45] D. Ulyanov, A. Vedaldi, and V. Lempitsky. Instance normalization: The missing ingredient for fast
398 stylization. *arXiv preprint arXiv:1607.08022*, 2016.
- 399 [46] T. Vogels, S. P. Karimireddy, and M. Jaggi. Practical low-rank communication compression in decentralized
400 deep learning. *Adv. Neural Inform. Process. Syst.*, 33:14171–14181, 2020.
- 401 [47] H. Wang and N. Ahuja. Rank-r approximation of tensors using image-as-matrix representation. In *IEEE
402 Conf. Comput. Vis. Pattern Recog.*, volume 2, pages 346–353. IEEE, 2005.

- 403 [48] H. Weyl. Das asymptotische verteilungsgesetz der eigenwerte linearer partieller differentialgleichungen.
404 *Math. Ann.*, 71:441–479, 1912.
- 405 [49] S. Wold, K. Esbensen, and P. Geladi. Principal component analysis. *Chemom. Intell. Lab. Syst*, 2(1-3):
406 37–52, 1987.
- 407 [50] A. Wolf, J. B. Swift, H. L. Swinney, and J. A. Vastano. Determining Lyapunov exponents from a time
408 series. *Physica D*, 16(3):285–317, 1985.
- 409 [51] Y. Wu and K. He. Group normalization. In *Eur. Conf. Comput. Vis.*, pages 3–19, 2018.
- 410 [52] J. Ye. Generalized low rank approximations of matrices. *Mach. Learn.*, 61(1):167–191, 2005.
- 411 [53] X. Yu, T. Liu, X. Wang, and D. Tao. On compressing deep models by low rank and sparse decomposition.
412 In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 7370–7379, 2017.
- 413 [54] Z. Zha, X. Yuan, B. Wen, J. Zhou, J. Zhang, and C. Zhu. From rank estimation to rank approximation:
414 Rank residual constraint for image restoration. *IEEE Trans. Image Process.*, 29:3254–3269, 2019.
- 415 [55] M. Zhan, S. Cao, B. Qian, S. Chang, and J. Wei. Low-rank sparse feature selection for patient similarity
416 learning. In *IEEE Int. Conf. on Data Min.*, pages 1335–1340. IEEE, 2016.
- 417 [56] T. Zhang, B. Ghanem, S. Liu, C. Xu, and N. Ahuja. Low-rank sparse coding for image classification. In
418 *Int. Conf. Comput. Vis.*, pages 281–288, 2013.
- 419 [57] Y. Zhang, Z. Jiang, and L. S. Davis. Learning structured low-rank representations for image classification.
420 In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 676–683, 2013.
- 421 [58] H. Zou, T. Hastie, and R. Tibshirani. Sparse principal component analysis. *J. Comput. Graph. Stat.*, 15(2):
422 265–286, 2006.

423 Checklist

- 424 1. For all authors...
- 425 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
426 contributions and scope? [Yes]
- 427 (b) Did you describe the limitations of your work? [No]
- 428 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 429 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
430 them? [Yes]
- 431 2. If you are including theoretical results...
- 432 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 433 (b) Did you include complete proofs of all theoretical results? [Yes] Proofs are included in
434 the appendix.
- 435 3. If you ran experiments...
- 436 (a) Did you include the code, data, and instructions needed to reproduce the main
437 experimental results (either in the supplemental material or as a URL)? [No]
- 438 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
439 were chosen)? [Yes]
- 440 (c) Did you report error bars (e.g., with respect to the random seed after running
441 experiments multiple times)? [Yes]
- 442 (d) Did you include the total amount of compute and the type of resources used (e.g., type
443 of GPUs, internal cluster, or cloud provider)? [Yes]
- 444 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 445 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 446 (b) Did you mention the license of the assets? [No]
- 447 (c) Did you include any new assets either in the supplemental material or as a URL? [No]

- 448 (d) Did you discuss whether and how consent was obtained from people whose data you're
449 using/curating? [N/A]
- 450 (e) Did you discuss whether the data you are using/curating contains personally identifiable
451 information or offensive content? [N/A]
- 452 5. If you used crowdsourcing or conducted research with human subjects...
- 453 (a) Did you include the full text of instructions given to participants and screenshots, if
454 applicable? [N/A]
- 455 (b) Did you describe any potential participant risks, with links to Institutional Review
456 Board (IRB) approvals, if applicable? [N/A]
- 457 (c) Did you include the estimated hourly wage paid to participants and the total amount
458 spent on participant compensation? [N/A]