

---

# Principal Components Bias in Deep Neural Networks

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Recent work suggests that convolutional neural networks of different architectures  
2 learn to classify images in the same order. To understand this phenomenon, we  
3 revisit the over-parametrized deep linear network model. Our asymptotic analysis,  
4 assuming that the hidden layers are wide enough, reveals that the convergence rate  
5 of this model’s parameters is exponentially faster along directions corresponding  
6 to the larger principal components of the data, at a rate governed by the singular  
7 values. We term this convergence pattern the *Principal Components bias (PC-bias)*.  
8 We show how the *PC-bias* streamlines the order of learning of both linear and non-  
9 linear networks, more prominently at earlier stages of learning. We then compare  
10 our results to the spectral bias, showing that both biases can be seen independently,  
11 and affect the order of learning in different ways. Finally, we discuss how the  
12 *PC-bias* may explain some benefits of early stopping and its connection to PCA,  
13 and why deep networks converge more slowly when given random labels.

## 14 1 Introduction

15 The dynamics of learning in deep neural networks is an intriguing subject, not yet sufficiently  
16 understood. Diverse empirical data seems to support the hypothesis that neural networks start by  
17 learning a simple model, which then gains complexity as learning proceeds (Gunasekar et al., 2018;  
18 Soudry et al., 2018; Hu et al., 2020; Nakkiran et al., 2019; Gissin et al., 2019; Heckel & Soltanolkotabi,  
19 2019; Ulyanov et al., 2018; Valle-Perez et al., 2018). This phenomenon is sometimes called *simplicity*  
20 *bias* (Dingle et al., 2018; Shah et al., 2020).

21 Recent work additionally shows that neural networks learn the training examples of natural datasets  
22 in a consistent order, and further impose a consistent order on the test set (Hacohen et al., 2020;  
23 Pliushch et al., 2021). Below we call this effect *Learning Order Constancy (LOC)*. Currently, the  
24 characteristics of visual data, which may explain this consistently imposed order, remain unclear.  
25 Surprisingly, this universal order persists despite the variability introduced into the training of different  
26 models and architectures.

27 To understand this phenomenon, we start by analyzing the deep linear network model (Saxe et al.,  
28 2013, 2019), defined by the concatenation of linear operators. While not a universal approximator, it  
29 is nevertheless trained by minimizing a non-convex objective function with a multitude of minima.  
30 The investigation of such networks is often employed to shed light on the learning dynamics when  
31 complex geometric landscapes are explored by GD (Fukumizu, 1998; Arora et al., 2018).

32 In Section 2, we prove that the convergence of the weights of deep linear networks is governed  
33 by the eigendecomposition of the raw data in a phenomenon we term *PC-bias*. These asymptotic  
34 results, valid when the hidden layers are wide enough, can be seen as an extension of the known  
35 behavior of the single-layer convex linear model (Le Cun et al., 1991). Our work is closely related to  
36 (Saxe et al., 2013, 2019), where the deep linear model’s dynamics is analyzed as a function of the  
37 input and input-output statistics. Importantly, the analysis in (Saxe et al., 2013, 2019; Arora et al.,

38 2018) incorporates the simplifying assumption that the data’s singular values are identical (whitened  
39 data), an assumption which unfortunately **obscures the main result of our analysis** – the direct  
40 dependence of convergence rate on the singular values of the data.

41 In Section 3, we empirically show that this pattern of convergence is indeed observed in deep linear  
42 networks, validating the plausibility of our assumptions. We continue by showing that the *LOC-effect*  
43 in deep linear network is determined solely by their *PC-bias*. We prove a similar (weaker) result for  
44 the non-linear two-layer ReLU model introduced by Allen-Zhu et al. (2018), where this model is  
45 presented as a certain extension of NTK (Jacot et al., 2020). In this framework, convergence is fastest  
46 along the largest **kernel’s** principal components, a result related to the *Spectral bias* discussed below.

47 In Section 4, we extend the study empirically to non-linear networks, and investigate the relation  
48 between the *PC-bias* and the *LOC-effect* in general deep networks. We first show that the order  
49 by which examples are learned by linear networks is highly correlated with the order induced by  
50 prevalent deep CNN models. We then show directly that the learning order of non-linear CNN models  
51 is affected by the principal decomposition of the data. Moreover, the *LOC-effect* diminishes when  
52 data is whitened, indicating a tight connection between the *PC-bias* and the *LOC-effect*.

53 Our results are reminiscent of another phenomenon, termed *Spectral bias* (Rahaman et al., 2019;  
54 Cao et al., 2019), which associates the learning dynamics of neural networks with the Fourier  
55 decomposition of functions in the hypothesis space. Rahaman et al. (2019) empirically demonstrated  
56 that the complexity of classifiers learned by ReLU networks increases with time. Basri et al. (2019,  
57 2020) showed theoretically, by way of analyzing elementary neural network models, that these models  
58 first fit the data with low-frequency functions, and gradually add higher frequencies to improve the fit.

59 Nevertheless, the *spectral bias* and *PC-bias* are inherently different. Indeed, the eigendecomposition  
60 of raw images is closely related to the Fourier analysis of images as long as the statistical properties  
61 of images are (approximately) translation-invariant (Simoncelli & Olshausen, 2001; Torralba & Oliva,  
62 2003). Still, the *PC-bias* is guided by spectral properties of the raw data and is additionally blind to  
63 class labels. On the other hand, the *spectral bias*, as well as the related *frequency bias* that has been  
64 shown to characterize NTK models (Basri et al., 2020), are all guided by spectral properties of the  
65 learned hypothesis, which strongly depends on label assignment.

66 In Section 4.3 we investigate the relation between the *PC-bias*, *spectral bias*, and the *LOC-effect*.  
67 We find that the *LOC-effect* is very robust: (i) when we neutralize the *spectral bias* by using low  
68 complexity models such as deep linear networks, the effect is still observed; (ii) when we neutralize  
69 the *PC-bias* by using whitened data, the *LOC-effect* persists. We hypothesize that at the beginning of  
70 learning, the learning dynamics of neural models is controlled by the eigendecomposition of the raw  
71 data. As learning proceeds, control of the dynamics slowly shifts to other factors.

72 The *PC-bias* has implications beyond the *LOC-effect*, as expanded in Section 5 and Suppl. §A:

73 **1. Early stopping.** It is often observed that when training deep networks with real data, the highest  
74 generalization accuracy is obtained before convergence. Consequently, early stopping is often  
75 prescribed to improve generalization. Following the commonly used assumption that in natural  
76 images the lowest principal components correspond to noise (Torralba & Oliva, 2003), our results  
77 predict the benefits of early stopping, and relate it to PCA. In Section 5 we investigate the relevance  
78 of this conclusion to real non-linear networks (see, e.g., Basri et al. (2019); Li et al. (2020) for  
79 complementary accounts).

80 **2. Slower convergence with random labels.** Zhang et al. (2016) showed that neural networks  
81 can learn any label assignment. However, training with random label assignments is known to  
82 converge slower as compared to training with the original labels (Krueger et al., 2017). We report a  
83 similar phenomenon when training deep linear networks. Our analysis shows that when the principal  
84 eigenvectors are correlated with class identity, as is often the case in natural images, the loss decreases  
85 faster when given true label assignments as against random label assignments. In Section 5 we  
86 investigate this hypothesis empirically in linear and non-linear networks.

87 **3. Weight initialization.** Different weight initialization schemes have been proposed to stabilize the  
88 learning and minimize the hazard of "exploding gradients" (e.g., Glorot & Bengio, 2010; He et al.,  
89 2015). Our analysis (see Suppl. §A) identifies a related variant, which eliminates the hazard when  
90 all the hidden layers are roughly of equal width. In the deep linear model, it can be proven that the  
91 proposed normalization variant in a sense minimizes repeated gradient amplification.

92 **2 Theoretical analysis**

93 **Notations.** Let  $\mathbb{X} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$  denote the training data, where  $\mathbf{x} \in \mathbb{R}^q$  denotes the  $i$ -th data  
 94 point and  $\mathbf{y} \in \{0, 1\}^K$  its corresponding label. Let  $\frac{1}{n_i} \mathbf{m}_i$  denote the centroid (mean) of class  $i$  with  
 95  $n_i$  points, and  $M = [\mathbf{m}_1 \dots \mathbf{m}_K]^\top$ . Finally, let  $X$  and  $Y$  denote the matrices whose  $i^{\text{th}}$  column  
 96 is  $\mathbf{x}_i$  and  $\mathbf{y}_i$  respectively.  $\Sigma_{XX} = XX^\top$  and  $\Sigma_{YX} = YX^\top$  denote the covariance matrix of  $X$   
 97 and cross-covariance of  $X$  and  $Y$  respectively. We note that  $\Sigma_{XX}$  captures the structure of the data  
 98 irrespective of class identity.

99 **Definition 1** (Principal coordinate system). *The coordinate system obtained by rotating the data in  $\mathbb{R}^q$   
 100 by an orthonormal matrix  $U^\top$ , where  $SVD(\Sigma_{XX}) = UDU^\top$ . Now  $\Sigma_{XX} = D$ , a diagonal matrix whose  
 101 elements are the singular values of  $XX^\top$ , arranged in decreasing order  $d_1 \geq d_2 \geq \dots \geq d_q \geq 0$ .*

102 **Definition 2** (Compact representation). *Let  $f(\mathbf{x})$  denote a deep linear network. Then  $f(\mathbf{x}) =$   
 103  $\left(\prod_{l=L}^1 W_l\right)\mathbf{x} = \mathbf{W}\mathbf{x}$ , where  $\mathbf{W} \in \mathbb{R}^{K \times q}$  is called the compact representation of the network.*

104 **Definition 3** (Error matrix). *For a deep linear network whose compact representation is  $\mathbf{W}$ , the  
 105 error matrix is  $Er = \mathbf{W}\Sigma_{XX} - \Sigma_{YX}$ . In the principal coordinate system,  $Er = \mathbf{W}D - M$ .*

106 **Assumptions.** Our analysis assumes that the learning rate  $\mu$  is infinitesimal, and therefore terms  
 107 of size  $O(\mu^2)$  can be neglected. We further assume that the width of the hidden layers lies in  
 108  $[\mathfrak{m}, \mathfrak{m} + M_b]$ , where  $\mathfrak{m} \rightarrow \infty$  denotes a very large number and  $M_b$  is fixed. Thus terms of size  $O(\frac{1}{\mathfrak{m}})$   
 109 can also be neglected. In Fig. 1 we show the plausibility of these assumptions, where the predicted  
 110 dynamics is seen throughout the training of deep linear networks, even for small values of  $\mathfrak{m}$ .

111 **2.1 The dynamics of deep over-parametrized linear networks**

112 Consider a deep linear network with  $L$  layers, and let

$$L(\mathbb{X}) = \frac{1}{2} \|\mathbf{W}X - Y\|_F^2 \quad \mathbf{W} := \prod_{l=L}^1 W_l, \quad W_l \in \mathbb{R}^{m_l \times m_{l-1}} \quad (1)$$

113 Above  $m_l$  denotes the number of neurons in layer  $l$ , where  $m_0 = q$  and  $m_L = K$ .

114 **Theorem 1.** *In each time point  $s$ , the compact matrix representation  $\mathbf{W}$  obeys the following dynamics,  
 115 when using the notation  $Er^s$  defined in Def. 3:*

$$\mathbf{W}^{s+1} = \mathbf{W}^s - \mu \sum_{l=1}^L A_l^s \cdot Er^s \cdot B_l^s + O(\mu^2) \quad (2)$$

116 Above  $\mu$  denotes the learning rate.  $A_l^s$  and  $B_l^s$  are called gradient scale matrices, and are defined as

$$A_l^s := \left( \prod_{j=L}^{l+1} W_j^s \right) \left( \prod_{j=L}^{l+1} W_j^s \right)^\top \in \mathbb{R}^{K \times K} \quad B_l^s := \left( \prod_{j=l-1}^1 W_j^s \right)^\top \left( \prod_{j=l-1}^1 W_j^s \right) \in \mathbb{R}^{q \times q} \quad (3)$$

117 The proof can be found in Suppl. §B.

118 **Gradient scale matrices.** Some statistical properties of such matrices are established in Suppl. §A.  
 119 Note that when the number of hidden layers is 0 ( $L = 1$ ), both gradient scale matrices reduce to the  
 120 identity matrix and the dynamics in (2) is reduced to the following known result (e.g., Le Cun et al.,  
 121 1991):  $\mathbf{W}^{s+1} = \mathbf{W}^s - \mu Er^s$ . Recall, however, that the focus of this paper is the over-parameterized  
 122 linear model with  $L > 1$ , in which the loss is not convex. Since the difference between the convex  
 123 linear model and the over-parametrized deep model boils down to these matrices, our convergence  
 124 analysis henceforth focuses on the dynamics of the *gradient scale matrices*.

125 In accordance, we analyze the evolution of the *gradient scale matrices* as learning proceeds. Let  
 126  $\mathfrak{m} = \min(m_1, \dots, m_{L-1})$  denote the size of the smallest hidden layer. Initially for  $s = 0$ , all weight  
 127 matrices  $W_l^0$  are assumed to be initialized by sampling from a distribution with mean 0 and variance  
 128  $\sigma_l^2 = O(\frac{1}{\mathfrak{m}})$ . The specific normalization factor, alluded to in  $O(\frac{1}{\mathfrak{m}})$ , is a variant of the Glorot  
 129 initialization. Details and justification can be found in Suppl. §A.1.

130 At time  $s$ , let  $A_l^s(\mathfrak{m})$  and  $B_l^s(\mathfrak{m})$  denote a sequence of random *gradient scale matrices*, corresponding  
 131 to networks whose smallest hidden layer has  $\mathfrak{m}$  neurons. From Suppl. §A we deduce that:

132 **Theorem 2.** Using  $\xrightarrow{p}$  to denote convergence in probability as  $m \rightarrow \infty$ , and  $\forall s, l$ :

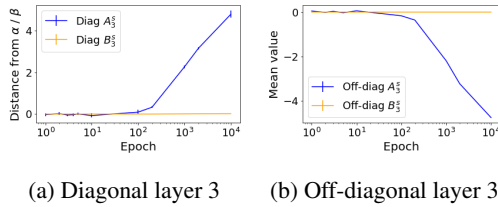
$$B_l^s(m) \xrightarrow{p} I, \quad \text{var}[B^l(m)] = O\left(\frac{1}{m}\right) \quad A_l^s(m) \xrightarrow{p} I, \quad \text{var}[A^l(m)] = O\left(\frac{1}{m}\right)$$

133 *Proof.* Proof by induction on  $s$ . Initially when  $s = 0$ , the claim follows from Thm 4 and Corr 5.1.  
 134 The induction step validity follows from Thm 6 and Thm 7 (see Suppl. §A.2).  $\square$

135 The detailed proof shows that the relevant constants are amplified with  $s$ . While they remain moderate  
 136 and  $m$  is sufficiently large,  $B_l^s(m) \approx I$  and  $A_l^s(m) \approx I \forall l$ . In this case, the dynamics of the over-  
 137 parameterized model is identical to the dynamics of the convex linear model,  $\mathbf{W}^{s+1} = \mathbf{W}^s - \mu E r^s$ .

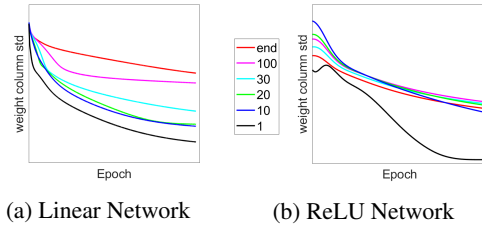
138 **Convergence rate.** In §A.2 we show that the convergence of  $B_l^s(m)$  to  $I$  is governed to some extent  
 139 by  $O\left(\frac{K}{m}\right)$ , while the convergence of  $A_l^s(m)$  is governed by  $O\left(\frac{q}{m}\right)$ . Recall that while  $m \rightarrow \infty$ ,  
 140  $q$  is the dimension of the data space which is fixed in advance and can be fairly large, while  $K$  is  
 141 the number of classes which is fixed and quite small. Typically,  $K \ll q$ . Thus we expect the right  
 142 *gradient scale matrices*  $B_l^s(m)$  to remain approximately  $I$  much longer than the left matrices  $A_l^s(m)$ .

143 **Empirical validation.** Since the results above are asymptotic, and to envision the difference between  
 144 convergence governed by  $O\left(\frac{K}{m}\right)$  vs.  $O\left(\frac{q}{m}\right)$ , we resort to simulations whose results are shown in  
 145 Fig. 1. These empirical results, recounting linear networks with 4 hidden layers of width 1024, clearly  
 146 show that during a significant part of the training both *gradient scale matrices* remain approximately  
 147  $I$ . The difference between the convergence rate of  $B_l^s$  and  $A_l^s$  is seen later on, when  $\Delta A_l^s$  starts to  
 148 increase shortly before convergence, while  $\Delta B_l^s$  remains essentially 0 throughout.



(a) Diagonal layer 3 (b) Off-diagonal layer 3

Figure 1: The dynamics of  $A_l^s$  and  $B_l^s$  when training 10 5-layered linear networks on the small-mammals dataset. (a) Mean distance of the diagonal elements of  $A_l^s$  and  $B_l^s$  from  $\alpha_l^s$  and  $\beta_l^s$  (as defined in Thm 3, §A.1). (b) Mean value of the off-diagonal elements of  $A_l^s$  and  $B_l^s$ . The networks reach maximal test accuracy at epoch  $s = 100$ , before the divergence of  $A_l^s$ . All layers behave similarly.



(a) Linear Network (b) ReLU Network

Figure 2: Empirical confirmation of the theoretical results reported below, showing the std of  $w_j$  over 10 independently trained networks as a function of the epoch, for 6 specific principal components (identified in the legend). Left: two-layer linear network. Right: two-layer non-linear network with ReLU activation.

## 149 2.2 Weight evolution

150  $K \ll q$  entails that  $B_l^s(m)$  remains approximately equal to  $I$  much longer than  $A_l^s(m)$ . This is  
 151 substantiated by the simulation results in Fig. 1. Consequently, while earlier on it is safe to assume  
 152 that both  $A_l^s \approx I$  and  $B_l^s \approx I$ , as learning proceeds only  $B_l^s \approx I$  is safe to assume.

153 With this in mind, we obtain expressions for the evolution of  $\mathbf{W}^s$  separately for earlier and later in  
 154 learning. We first shift to the principal coordinate system defined in Def 1. In this system we can  
 155 analyze each column of  $\mathbf{W}^s$  separately, where  $w_j^s$  and  $m_j$  denote the respective columns of  $\mathbf{W}^s$  and  
 156  $M$ . At the beginning of learning when both  $A_l^s \approx I$  and  $B_l^s \approx I$  (see §B.3 for a detailed derivation):

$$w_j^{s+1} = (\lambda_j)^s w_j^0 + [1 - (\lambda_j)^s] \frac{m_j}{d_j} \quad \lambda_j = 1 - \mu d_j L \quad (4)$$

157 Eq. 4 is reminiscent of the well understood dynamics of training the convex one layer linear model. It  
 158 is composed of two additive terms, revealing two parallel and independent processes:

- 160 1. The dependence on random initialization tends to 0 exponentially with decline rate  $\lambda_j$ .
- 161 2. The final value is the sum of a geometrical series with a common ratio  $\lambda_j$ .

162 In either case, convergence is fastest for the largest singular eigenvalue, or the first column of  $\mathbf{W}$ ,  
 163 and slowest for the smallest singular value. This behavior is visualized in Fig. 2a. Importantly,  
 164 the rate of convergence depends on the singular value  $d_j$ , the number of layers  $L$ , and the learning rate  $\mu$ .

165 In later stages of learning, when we can only assume that  $B_i^s \approx I$ , the dynamic becomes:

$$w_j^{s+1} = \prod_{\nu=1}^s (I - \mu d_j A^\nu) w_j^0 + \mu \left[ \sum_{\nu=1}^s \prod_{\rho=\nu+1}^s (I - \mu d_j A^\rho) A^\nu \right] m_j \quad (5)$$

166 where  $A^s = \sum_{l=1}^L A_l^s$ . The proof is provided in §B.3. Although the dynamics now depends on  
 167 matrices  $A^s$  as well, it is still the case that the convergence of each column is governed by its singular  
 168 value  $d_j$ . This suggests that while the *PC-bias* is more pronounced in earlier stages of learning, its  
 169 effect persists throughout.

170 The analysis above is extended to a simple non-linear ReLU model (cf. Arora et al., 2019) as detailed  
 171 in §B.2, with qualitatively similar results (albeit under unrealistic assumptions). Empirical results,  
 172 shown in Fig. 2b, indicate that the results are indicative beyond the assumed circumstances.

### 173 3 PC-bias: empirical study

174 In this section, we first analyze deep linear networks, showing that the convergence rate is indeed  
 175 governed by the principal singular values of the data, which demonstrates the plausibility of the  
 176 assumptions made in Section 2. We continue by extending the scope of the investigation to non-linear  
 177 neural networks, finding there evidence for the *PC-bias* mostly in the earlier stages of learning.

#### 178 3.1 Methodology

179 We say that a linear network is  $L$ -layered when it has  $L - 1$  hidden fully connected (FC) layers  
 180 (without convolutional layers). In our empirical study we relaxed some assumptions of the theoretical  
 181 study, in order to increase the resemblance of the trained networks to networks in common use.  
 182 Specifically, we changed the initialization to the commonly used Glorot initialization, replaced the  
 183  $L_2$  loss with the cross-entropy loss, and employed SGD instead of the deterministic GD. Notably,  
 184 the original assumptions yielded similar results. The results presented summarize experiments with  
 185 networks of equal width across all hidden layers, specifically the moderate value of  $m = 1024$ ,  
 186 keeping in mind that we test the relevance of asymptotic results for  $m \rightarrow \infty$ . Using a different width  
 187 for each layer yielded similar qualitative results. Details regarding the hyper-parameters, architectures,  
 188 and datasets can be found in §D.1, §D.3 and §D.4 respectively.

#### 189 3.2 PC-bias in deep linear networks

190 In this section, we train  $L$ -layered linear networks, then compute their compact representations  
 191  $W$  rotated to align with the canonical coordinate system (Def. 1). Note that each row  $w_r$  in  $W$   
 192 essentially defines the one-vs-all separating hyper-plane corresponding to class  $r$ .

193 To examine both the variability between models and their convergence rate, we inspect  $w_r$  at different  
 194 time points during learning. The rate of convergence can be measured directly, by observing the  
 195 changes in the weights of each element in  $w_r$ . These weight values<sup>1</sup> should be compared with  
 196 the optimal values in each row  $w_r$  of  $W_{opt} = YX^T(XX^T)$ . The variability between models is  
 197 measured by calculating the standard deviation (std) of each  $w_r$  across  $N$  models.

198 We begin with linear networks. We trained 10 5-layered FC linear networks, and 10 linear st-VGG  
 199 convolutional networks. When analyzing the compact representation of such networks we observe  
 200 similar behavior – weights corresponding to larger principal components converge faster to the  
 201 optimal value, and their variability across models converges faster to 0 (Figs. 3a,3b). Thus, while the  
 202 theoretical results are asymptotic, *PC-bias* is empirically seen throughout the entire learning process  
 203 of deep linear networks.

204 **Whitened data.** The *PC-bias* is neutralized when the data is whitened, at which point  $\Sigma_{XX}$  is the  
 205 scaled identity matrix. In Fig. 3c, we plot the results of the same experimental protocol while using a  
 206 ZCA-whitened dataset. As predicted, the networks no longer show any bias towards any principal  
 207 direction. Weights in all directions are scaled similarly, and the std over all models is the same in  
 208 each epoch, irrespective of the principal direction. (Additional experiments show that this is *not* an  
 209 artifact of the lack of uniqueness when deriving the principal components of a white signal).

<sup>1</sup>We note that the weights tend to start larger for smaller principal components, as can be seen in Fig. 3a left.

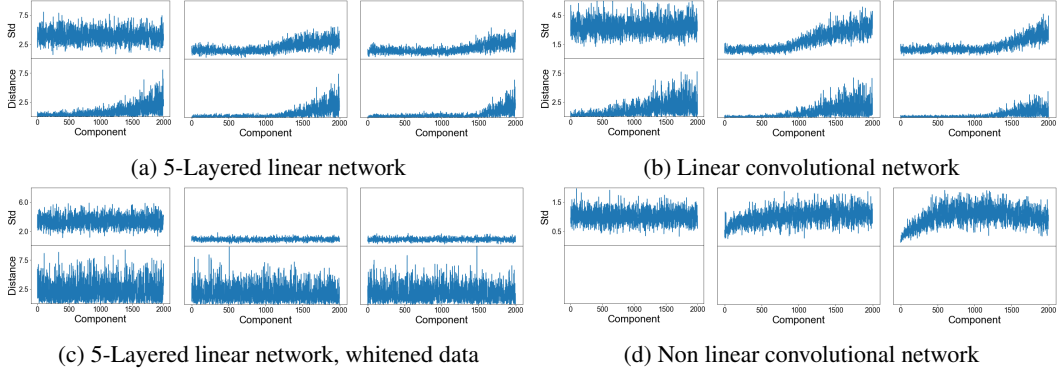


Figure 3: Convergence of the compact representation along the principal directions in different epochs. The value of the  $X$ -axis corresponds to the index of a principal eigenvalue, from the most significant to the least significant. (a) 10 5-layered linear networks trained on the cats and dogs dataset. 3 plots are provided, corresponding to snapshots taken at different stages of learning: the beginning (epoch 0, left), intermediate stage (middle), and close to convergence (right). Bottom panel: average distance of the weights in  $w_1$  from the optimal linear classifier; top panel: respective std. (b) Similarly, for 10 linear st-VGG convolutional networks, trained on CIFAR-10. (c) Similarly, for 10 5-layered linear networks, trained on the cats and dogs dataset, with ZCA-whitening. (d) Similarly, for 10 **non-linear** st-VGG networks trained on the cats and dogs dataset. Here the distance to the optimal solution is not well defined and we therefore only show the std.

### 210 3.3 PC-bias in general CNNs

211 In this section, we investigate the manifestation of the *PC-bias* in non-linear deep convolutional  
 212 networks. As we cannot directly track the learning dynamics separately in each principal direction of  
 213 non-linear networks, we adopt two different evaluation mechanisms:

214 **Linear approximation.** We considered several linear approximations, but since all of them showed  
 215 the same qualitative behavior, we report results with the simplest one. Specifically, to obtain a linear  
 216 approximation of a non-linear network, without max-pooling or batch-normalization layers, we  
 217 follow the definition of the compact representation from Section 2 while ignoring any non-linear  
 218 activation. We then align this matrix with the canonical coordinate system (Def. 1), and observe the  
 219 evolution of the weights and their std across models along the principal directions during learning.  
 220 Note that now the networks do not converge to the same compact representation, which is not unique.  
 221 Nevertheless, we see that the *PC-bias* governs the weight dynamics to a noticeable extent.

222 More specifically, in these networks a large fraction of the lowest principal components hardly changes  
 223 during learning, as good as being ignored. Nevertheless, the *PC-bias* affects the higher principal  
 224 components, most notably at the beginning of training (see Fig. 3d). Thus weights corresponding to  
 225 higher principal components converge faster, and the std across models of such weights decreases  
 226 faster for higher principal components.

227 **Projection to higher PC's.** We created a modified *test-set*, by project-  
 228 ing each test example on the span of the first  $P$  principal components.  
 229 This is equivalent to reducing the dimensionality of the test set to  $P$  using  
 230 PCA. We trained an ensemble of  $N=100$  st-VGG networks on the  
 231 original small mammals training set, then evaluated these networks dur-  
 232 ing training on 4 versions of the test-set, reduced to  $P=1,10,100,1000$   
 233 dimensions respectively. Mean accuracy is plotted in Fig. 4. Similar  
 234 results are obtained when training VGG-19 networks on CIFAR-10,  
 235 see §C.3.

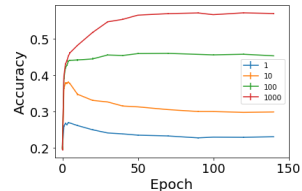


Figure 4: Mean accuracy of 10 st-VGG networks evaluated on test data projected to dimensionality  $\{1, 10, 100, 1000\}$ .

236 Taking a closer look at Fig. 4, we see that when evaluated on lower  
 237 dimensionality test-data ( $P=1,10$ ), the networks' accuracy peaks after  
 238 a few epochs, at which point performance starts to decrease. This result suggests that the networks  
 239 rely more heavily on these dimensions in the earlier phases of learning, and then continue to learn  
 240 other things. In contrast, when evaluated on higher dimensionality test-data ( $P=100,1000$ ), accuracy  
 241 continues to rise, longer so for larger  $P$ . This suggests that significant learning of the additional  
 242 dimensions continues in later stages of the learning.

243 **4 PC-bias: Learning Order Constancy**

244 In this section, we show that the *PC-bias* is significantly correlated with the learning order of deep  
 245 neural networks, and can therefore partially account for the *LOC-effect* described in Section 1.  
 246 Following Hacoen et al. (2020), we measure the "speed of learning" of each example by computing  
 247 its *accessibility* score. This score is given per example, and characterizes how fast an ensemble of  
 248  $N$  networks learns it. Formally,  $accessibility(\mathbf{x}) = \mathbb{E} [\mathbb{1}(f_i^e(\mathbf{x}) = y(\mathbf{x}))]$ , where  $f_i^e(\mathbf{x})$  denotes  
 249 the outcome of the  $i$ -th network trained over  $e$  epochs, and the mean is taken over networks and  
 250 epochs. For the set of datapoints  $\{(\mathbf{x}_j, \mathbf{y}_j)\}_{j=1}^n$ , *Learning Order Constancy* is manifested by the high  
 251 correlation between 2 instances of  $accessibility(\mathbf{x})$ , each computed from a different ensemble.

252 *PC-bias* is shown to pertain to *LOC* in two ways: First, in Section 4.1 we show high correlation  
 253 between the learning order in deep linear and non-linear networks. Since the *PC-bias* fully accounts  
 254 for *LOC* in deep linear networks, this suggests it also accounts (at least partially) for the observed  
 255 *LOC* in non-linear networks. Comparison with the *critical principal component* verifies this assertion.  
 256 Second, we show in Section 4.2 that when the *PC-bias* is neutralized, *LOC* diminishes as well. In  
 257 Section 4.3 we discuss the relationship between the *spectral bias*, *PC-bias* and the *LOC-effect*.

258 **4.1 PC-Bias is correlated with LOC**

259 We first compare the order of learning of non-linear models and deep linear networks by computing  
 260 the correlation between the *accessibility* scores of both models. This comparison reveals high  
 261 correlation ( $r = 0.85, p < 10^{-45}$ ), as seen in Fig. 5a. To investigate directly the connection between  
 262 the *PC-bias* and *LOC*, we define the *critical principal component* of an example to be the first  
 263 principal component  $P$ , such that a linear classifier trained on the original data can classify the  
 264 example correctly when projected to  $P$  principal components. We trained  $N=100$  st-VGG networks  
 265 on the cats and dogs dataset, and computed for each example its *accessibility* score and *critical*  
 266 *principal component*. In Fig. 5b we see strong negative correlation between the two scores ( $p=-0.93,$   
 267  $r < 10^{-4}$ ), suggesting that the *PC-bias* affects the order of learning as measured by *accessibility*.

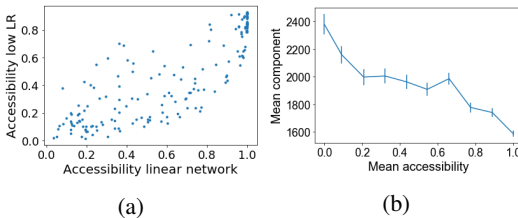


Figure 5: (a) Correlation between the *accessibility* score of  $N=100$  st-VGG networks trained with a low learning rate<sup>2</sup>, and  $N=100$  linear st-VGG networks, trained on small mammals. (b) Correlation between the accessibility score of  $N=100$  st-VGG networks trained on cats and dogs, and the *critical principal component* score. The *accessibility* plot is smoothed by moving average of width 10. Error bars indicate standard error.

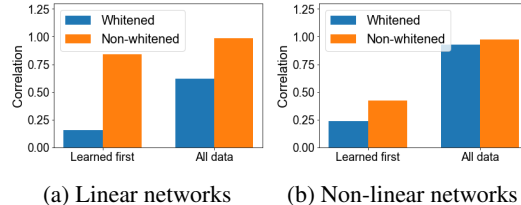


Figure 6: *LOC* measured with and without *PC-bias*. Each bar represents the correlation between the learning order of 2 collections of 10 networks trained on CIFAR-10. Orange bars represent natural images, in which the *PC-bias* is present, while blue bars represent whitened data, in which the *PC-bias* is eliminated. As *PC-bias* is more prominent earlier on, we compare these correlations for the entire data (right 2 bars), and for the subset of 20% "fastest learned" examples (left 2 bars).

268 **4.2 Neutralizing the PC-bias leads to diminishing LOC**

269 Whitening the data eliminates the *PC-bias* as shown in Fig. 3c, since all the singular values are now  
 270 identical. Here we use this observation to further probe into the dependency of the *Learning Order*  
 271 *Constancy* on the *PC-bias*. Starting with the linear case, we train 4 ensembles of  $N=10$  2-layered  
 272 linear networks on the cats and dogs dataset, 2 with and 2 without ZCA-whitening. We compute the  
 273 *accessibility* score for each ensemble separately, and correlate the scores of the 2 ensembles in each  
 274 test case. Each correlation captures the consistency of the *LOC-effect* for the respective condition.  
 275 This correlation is expected to be very high for natural images. Low correlation implies that the  
 276 *LOC-effect* is weak, as training the same network multiple times yields a different learning order.

<sup>2</sup>As non-linear models achieve the accuracy of linear models within an epoch or 2, low learning rate is used.



277 Fig. 6a shows the results for deep linear networks. As expected, the correlation when using natural  
 278 images is very high. However, when using whitened images, correlation plummets, indicating that  
 279 the *LOC-effect* is highly dependent on the *PC-bias*. We note that the drop in the correlation is much  
 280 higher when considering only the 20% "fastest learned" examples, suggesting that the *PC-bias* affects  
 281 learning order more evidently at earlier stages of learning.

282 Fig. 6b shows the results when repeating this experiment with non-linear networks, training 2  
 283 collections of  $N=10$  VGG-19 networks on CIFAR-10. We find that the elimination of the *PC-bias*  
 284 in this case affects *LOC* much less, suggesting that the *PC-bias* can only partially account for the  
 285 *LOC-effect* in the non-linear case. However, we note that at the beginning of learning, when the  
 286 *PC-bias* is most pronounced, once again the drop is much larger and very significant (half).

### 287 4.3 Spectral bias, PC-bias and LOC

288 The *spectral bias* (Rahaman et al., 2019) characterizes the dynamics of learning in neural networks  
 289 differently, asserting that initially neural models can be described by low frequencies only. This may  
 290 provide an alternative explanation to LOC. Recall that LOC is manifested in the consistency of the  
 291 *accessibility* score across networks. To compare between the *spectral bias* and *accessibility* score,  
 292 we first need to estimate for each example whether it can be correctly classified by a low frequency  
 293 model. Accordingly, we define for each example a *discriminability* measure – the percentage out  
 294 of its  $k$  neighbors that share with it class identity. Intuitively, an example has a low *discriminability*  
 295 score when it is surrounded by examples from other classes, which forces the learned boundary to  
 296 incorporate high frequencies. In §C.2 we show that in the 2D case analyzed by Rahaman et al. (2019),  
 297 this measure strongly correlates ( $r=-0.8$ ,  $p < 10^{-2}$ ) with the spectral bias.

298 We trained several networks (VGG-19 and st-VGG) on several real datasets, including small-  
 299 mammals, STL-10, CIFAR-10/100 and a subset of ImageNet-20. For each network and dataset,  
 300 we computed the *accessibility* score as well as the *discriminability* of each example. The vector  
 301 space, in which discriminability is evaluated, is either the raw data or the network’s perceptual space  
 302 (penultimate layer activation). The correlation between these scores is shown in Table 1.

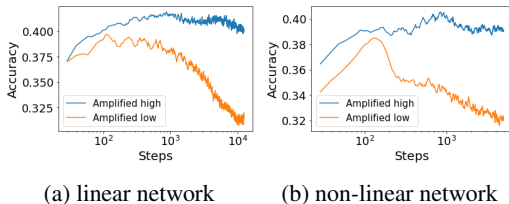


Figure 7: Effects of amplifying the highest (blue) and lowest (orange) principal components.

Table 1: Correlation between *accessibility* and *discriminability*.

Dataset	Raw data	Penultimate
Small mammals	0.46	0.85
ImageNet 20	0.01	0.54
CIFAR-100	0.51	0.85
STL10	0.44	0.7

303 Using raw data, low correlation is still seen between the *accessibility* and *discriminability* scores  
 304 when inspecting the smaller datasets (small mammals, CIFAR-100 and STL10). This correlation  
 305 vanishes when considering the larger ImageNet-20 dataset. It would appear that on its own, the  
 306 *spectral bias* cannot adequately explain the *LOC-effect*. On the other hand, in the perceptual space,  
 307 the correlation between *discriminability* and *accessibility* is quite significant for all datasets. Contrary  
 308 to our supposition, it seems that networks learn a representation where the *spectral bias* is evident,  
 309 but this bias does not necessarily govern its learning before the representation has been learned.

## 310 5 PC-bias: further implications

311 **Early Stopping and the Generalization Gap.** Considering natural images, it is often assumed that  
 312 the least significant principal components of the data represent noise (Torralba & Oliva, 2003). In  
 313 such cases, our analysis predicts that as noise dominates the components learned later in learning,  
 314 early stopping is likely to be beneficial. To test this hypothesis directly, we manipulated CIFAR-10  
 315 to amplify the signal in either the 1.5% most significant (higher) or 1.5% least significant (lower)  
 316 principal components (see examples in Fig. 16, Suppl. §D). Accuracy over the original test set,  
 317 after training 10 st-VGG and linear st-VGG networks on these manipulated images, can be seen  
 318 in Fig. 7. Both in linear and non-linear networks, early stopping is more beneficial when lower



319 principal components are amplified, and significantly less so when higher components are amplified,  
 320 as predicted by the *PC-bias*.

321 **Slower Convergence with Random Labels.** Deep neural models can learn any random label  
 322 assignment to a given training set (Zhang et al., 2016). However, when trained on randomly labeled  
 323 data, convergence appears to be much slower (Krueger et al., 2017). Assume, as before, that in natural  
 324 images the lower principal components are dominated by noise. We argue that the *PC-bias* now  
 325 predicts this empirical result, since learning randomly labeled examples requires signal present in  
 326 lower principal components. To test this hypothesis directly, we trained 10 2-layered linear networks  
 327 on datasets of natural images. Indeed, these networks converge slower with random labels (see  
 328 Fig. 8a). In Fig. 8b we repeat this experiment after having whitened the images, to neutralize the  
 329 *PC-bias*. Now convergence rate is identical, whether the labels are original or shuffled. Clearly, in  
 330 deep linear networks the *PC-bias* gives a full account of this phenomenon.

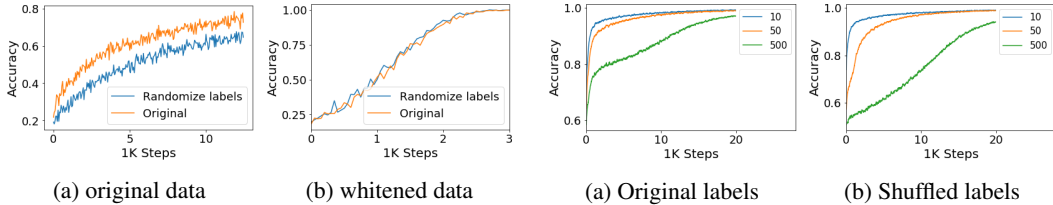


Figure 8: Learning curves of 10 2-layered linear networks, with real and shuffled labels, (a) before and (b) after whitening.

Figure 9: Learning curves of st-VGG networks trained on 3 datasets, which are linearly separable after projection to the highest  $P$  principal components (see legend).

331 To further check the relevance of this account to non-linear networks, we artificially generate datasets  
 332 where only the first  $P$  principal components are discriminative, while the remaining components  
 333 become noise by design. We constructed two such datasets: in one the labels are correlated with the  
 334 original labels, in the other they are not. Specifically, PCA is used to reduce the dimensionality of a  
 335 two-class dataset to  $P$ , and the optimal linear separator in the reduced representation is computed.  
 336 Next, all the labels of points that are incorrectly classified by the optimal linear separator are switched,  
 337 so that the train and test sets are linearly separable by this separator. Note that the modified labels  
 338 are still highly correlated with the original labels (for  $P = 500$ :  $p = 0.82$ ,  $r < 10^{-10}$ ). The  
 339 second dataset is generated by repeating the process while starting from randomly shuffled labels.  
 340 This dataset is likewise fully separable when projected to the first  $P$  components, but its labels are  
 341 uncorrelated with the original labels (for  $P = 500$ :  $p = 0.06$ ,  $r < 10^{-10}$ ).

342 The mean training accuracy of 10 non-linear networks with  $P=10,50,500$  is plotted in Fig. 9a (first  
 343 dataset) and Fig. 9b (second dataset). In both cases, the lower  $P$  is (namely, only the first few principal  
 344 components are discriminative), the faster the data is learned by the non-linear network. Whether the  
 345 labels are real or shuffled makes little qualitative difference, as predicted by the *PC-bias*.

## 346 6 Summary and discussion

347 When trained with gradient descent, the convergence rate of the over-parameterized deep linear  
 348 network model is provably governed by the eigendecomposition of the data, and specifically, pa-  
 349 rameters corresponding to the most significant principal components converge faster than the least  
 350 significant components. Empirical evidence is provided for the relevance of these results to more  
 351 realistic non-linear networks. We term this effect *PC-bias*. This result provides a complementary  
 352 account for some prevalent empirical observations, including the benefit of early stopping and the  
 353 slower convergence rate with shuffled labels.

354 We use the *PC-bias* to explicate the *Learning Order Constancy (LOC)*, showing that examples  
 355 learned at earlier stages are more distinguishable by the higher principal components, demonstrating  
 356 that networks' training relies more heavily on higher principal components early on. A causal link  
 357 between the *PC-bias* and the *LOC-effect* is demonstrated, as the *LOC-effect* diminishes when the  
 358 *PC-bias* is eliminated by whitening the images. We analyze these findings in view of a related  
 359 phenomenon termed *spectral bias*. While the *PC-bias* may be more prominent early on, the *spectral*  
 360 *bias* may be more important in later stages of learning.

## 361 **References**

- 362 Allen-Zhu, Z., Li, Y., and Liang, Y. Learning and generalization in overparameterized neural  
363 networks, going beyond two layers. *arXiv preprint arXiv:1811.04918*, 2018.
- 364 Arora, S., Cohen, N., and Hazan, E. On the optimization of deep networks: Implicit acceleration by  
365 overparameterization. In *International Conference on Machine Learning*, pp. 244–253, 2018.
- 366 Arora, S., Du, S., Hu, W., Li, Z., and Wang, R. Fine-grained analysis of optimization and generaliza-  
367 tion for overparameterized two-layer neural networks. In *International Conference on Machine*  
368 *Learning*, pp. 322–332, 2019.
- 369 Basri, R., Jacobs, D. W., Kasten, Y., and Kritchman, S. The convergence rate of neural networks for  
370 learned functions of different frequencies. In *Advances in Neural Information Processing Systems*,  
371 pp. 4761–4771, 2019.
- 372 Basri, R., Galun, M., Geifman, A., Jacobs, D., Kasten, Y., and Kritchman, S. Frequency bias in neural  
373 networks for input of non-uniform density. In *International Conference on Machine Learning*, pp.  
374 685–694. PMLR, 2020.
- 375 Cao, Y., Fang, Z., Wu, Y., Zhou, D.-X., and Gu, Q. Towards understanding the spectral bias of deep  
376 learning. *arXiv preprint arXiv:1912.01198*, 2019.
- 377 Dingle, K., Camargo, C. Q., and Louis, A. A. Input–output maps are strongly biased towards simple  
378 outputs. *Nature communications*, 9(1):1–7, 2018.
- 379 Fukumizu, K. Effect of batch learning in multilayer neural networks. *Gen*, 1(04):1E–03, 1998.
- 380 Gissin, D., Shalev-Shwartz, S., and Daniely, A. The implicit bias of depth: How incremental learning  
381 drives generalization. *arXiv preprint arXiv:1909.12051*, 2019.
- 382 Glorot, X. and Bengio, Y. Understanding the difficulty of training deep feedforward neural networks.  
383 In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*,  
384 pp. 249–256, 2010.
- 385 Gunasekar, S., Lee, J., Soudry, D., and Srebro, N. Implicit bias of gradient descent on linear  
386 convolutional networks. *arXiv preprint arXiv:1806.00468*, 2018.
- 387 Hacoen, G., Choshen, L., and Weinshall, D. Let’s agree to agree: Neural networks share classification  
388 order on real datasets. In *International Conference on Machine Learning*, pp. 3950–3960. PMLR,  
389 2020.
- 390 He, K., Zhang, X., Ren, S., and Sun, J. Delving deep into rectifiers: Surpassing human-level  
391 performance on imagenet classification. In *Proceedings of the IEEE International Conference on*  
392 *Computer Vision (ICCV)*, December 2015.
- 393 Heckel, R. and Soltanolkotabi, M. Denoising and regularization via exploiting the structural bias of  
394 convolutional generators. *arXiv preprint arXiv:1910.14634*, 2019.
- 395 Hu, W., Xiao, L., Adlam, B., and Pennington, J. The surprising simplicity of the early-time learning  
396 dynamics of neural networks. *arXiv preprint arXiv:2006.14599*, 2020.
- 397 Jacot, A., Gabriel, F., and Hongler, C. Neural tangent kernel: Convergence and generalization in  
398 neural networks, 2020.
- 399 Krueger, D., Ballas, N., Jastrzebski, S., Arpit, D., Kanwal, M. S., Maharaj, T., Bengio, E., Fischer, A.,  
400 and Courville, A. Deep nets don’t learn via memorization. 2017.
- 401 Le Cun, Y., Kanter, I., and Solla, A. S. Second order properties of error surfaces learning time and  
402 generalization. *Advances in neural information processing systems*, 3:918–924, 1991.
- 403 Li, M., Soltanolkotabi, M., and Oymak, S. Gradient descent with early stopping is provably robust  
404 to label noise for overparameterized neural networks. In *International Conference on Artificial*  
405 *Intelligence and Statistics*, pp. 4313–4324. PMLR, 2020.

- 406 Nakkiran, P., Kaplun, G., Kalimeris, D., Yang, T., Edelman, B. L., Zhang, F., and Barak, B. Sgd  
407 on neural networks learns functions of increasing complexity. *arXiv preprint arXiv:1905.11604*,  
408 2019.
- 409 Pliushch, I., Mundt, M., Lupp, N., and Ramesh, V. When deep classifiers agree: Analyzing  
410 correlations between learning order and image statistics. *arXiv preprint arXiv:2105.08997*, 2021.
- 411 Rahaman, N., Baratin, A., Arpit, D., Draxler, F., Lin, M., Hamprecht, F., Bengio, Y., and Courville,  
412 A. On the spectral bias of neural networks. In *International Conference on Machine Learning*, pp.  
413 5301–5310. PMLR, 2019.
- 414 Saxe, A. M., McClelland, J. L., and Ganguli, S. Exact solutions to the nonlinear dynamics of learning  
415 in deep linear neural networks. *arXiv preprint arXiv:1312.6120*, 2013.
- 416 Saxe, A. M., McClelland, J. L., and Ganguli, S. A mathematical theory of semantic development in  
417 deep neural networks. *Proceedings of the National Academy of Sciences*, 116(23):11537–11546,  
418 2019.
- 419 Shah, H., Tamuly, K., Raghunathan, A., Jain, P., and Netrapalli, P. The pitfalls of simplicity bias in  
420 neural networks. *arXiv preprint arXiv:2006.07710*, 2020.
- 421 Simoncelli, E. P. and Olshausen, B. A. Natural image statistics and neural representation. *Annual*  
422 *review of neuroscience*, 24(1):1193–1216, 2001.
- 423 Soudry, D., Hoffer, E., Nacson, M. S., Gunasekar, S., and Srebro, N. The implicit bias of gradient  
424 descent on separable data. *The Journal of Machine Learning Research*, 19(1):2822–2878, 2018.
- 425 Torralba, A. and Oliva, A. Statistics of natural image categories. *Network: computation in neural*  
426 *systems*, 14(3):391–412, 2003.
- 427 Ulyanov, D., Vedaldi, A., and Lempitsky, V. Deep image prior. In *Proceedings of the IEEE conference*  
428 *on computer vision and pattern recognition*, pp. 9446–9454, 2018.
- 429 Valle-Perez, G., Camargo, C. Q., and Louis, A. A. Deep learning generalizes because the parameter-  
430 function map is biased towards simple functions. *arXiv preprint arXiv:1805.08522*, 2018.
- 431 Zhang, C., Bengio, S., Hardt, M., Recht, B., and Vinyals, O. Understanding deep learning requires  
432 rethinking generalization. *arXiv preprint arXiv:1611.03530*, 2016.

433 **Checklist**

- 434 1. For all authors...
- 435 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
436 contributions and scope? [Yes]
- 437 (b) Did you describe the limitations of your work? [Yes]
- 438 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 439 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
440 them? [Yes]
- 441 2. If you are including theoretical results...
- 442 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See the  
443 "Assumptions" paragraph as Section 2
- 444 (b) Did you include complete proofs of all theoretical results? [Yes] Each theorem reference  
445 to its proof. Proofs can be found in Suppl. §A,B
- 446 3. If you ran experiments...
- 447 (a) Did you include the code, data, and instructions needed to reproduce the main ex-  
448 perimental results (either in the supplemental material or as a URL)? [No] All data,  
449 instructions and hyper-parameters are explicitly written in the main paper and/or in the  
450 Suppl. (see §D.4,D.4). The code itself will be provided once the anonymity will be  
451 lifted.
- 452 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
453 were chosen)? [Yes]
- 454 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
455 ments multiple times)? [Yes]
- 456 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
457 of GPUs, internal cluster, or cloud provider)? [N/A]
- 458 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 459 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 460 (b) Did you mention the license of the assets? [N/A]
- 461 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 462
- 463 (d) Did you discuss whether and how consent was obtained from people whose data you're  
464 using/curating? [N/A]
- 465 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
466 information or offensive content? [N/A]
- 467 5. If you used crowdsourcing or conducted research with human subjects...
- 468 (a) Did you include the full text of instructions given to participants and screenshots, if  
469 applicable? [N/A]
- 470 (b) Did you describe any potential participant risks, with links to Institutional Review  
471 Board (IRB) approvals, if applicable? [N/A]
- 472 (c) Did you include the estimated hourly wage paid to participants and the total amount  
473 spent on participant compensation? [N/A]