

1 We thank all the reviewers for their insightful comments and suggestions. We will add citations and discussions of all
2 the suggested related works in the full version. Reviewer-specific comments follow.

3 **Reviewer 1** We consider our mutual information framework to be a core contribution of our paper. In particular, our
4 formalization of “how well a linear classifier explains the performance of a model” has many advantages over prior
5 formalizations (e.g. see our response to Reviewer 3). Regarding our theory example (Section 5): We separate the
6 question of why SGD learns simple concepts first (Claim 1) from the question of why it does not forget them (Claim
7 2). Our current theory is only relevant to the second question, and it shows that in a simplified setting: SGD does not
8 “forget” the simple component even when trained to completion, provided it somehow learns the simple component
9 first. We argue that this simple example captures many properties of real settings (overparametrization, existence of
10 non-generalizable ERMs) and hence is valuable as a step towards more general theory. Regarding Claim 1, it is true
11 that we currently have no theoretical understanding of why linear learning occurs. We consider this one of the most
12 important open questions of our paper, and we are attempting to make progress on this in ongoing work.

13 **Reviewer 2** Thank you for pointing out the relevant papers. We also agree that it is scientifically valuable to describe
14 settings where these phenomenon fail to occur (bad initialization, bad architecture/parameterization, bad optimizer, or
15 pathological distributions). We plan on including a separate section with such examples in the final version.

16 **Reviewer 3** Regarding novelty of our claims: Although the idea that SGD learns functions of increasing complexity
17 has been informally floating in the community, our formalization has many advantages over prior formalizations, as
18 described below. Notably, our metrics respect the data distribution, are independent of network-parameterization, are
19 tractable to estimate in high dimensions, and are experimentally demonstrated for real-world distributions.

20 Concretely, regarding “On the spectral bias of neural networks” [1]: They consider measuring “simplicity” via the
21 Fourier spectrum of the learnt functions. However, the Fourier decomposition is taken with respect to a *uniform*
22 distribution on inputs. This notion is not as meaningful – if the data is not uniform (for example, natural images are
23 certainly not uniform in pixel-space), then a function which is highly correlated with a linear function when restricted
24 to the data distribution may appear highly non-linear with respect to the uniform distribution. And vice-versa – a
25 function which is nearly linear *under the uniform distribution* may in fact be highly non-linear when restricted to the
26 data distribution. Our metrics do not suffer from this issue – they are taken with respect to the true data distribution.

27 The synthetic experiments in [1] are all for 1-dimensional inputs, since to quote [1]: “explicitly evaluating the Fourier
28 coefficients ... becomes prohibitively expensive for larger d (e.g. on MNIST)”. Instead, their MNIST experiments
29 are only heuristically related to the metrics they propose. In contrast, the metrics in our paper are tractable even for
30 high-dimensional inputs, and we estimate them to high-precision on real datasets (MNIST and CIFAR).

31 Regarding “Understanding ... deep learning by Fourier analysis” [2]: This work also performs Fourier analysis with
32 respect to the uniform distribution on inputs, and so suffers from the same issues as [1]. Moreover [2] *requires* that
33 the input distribution is itself uniform to carry through the analysis¹. That is, the theorems of [2] do not hold for
34 non-uniform input distributions, such as images. The experiments in [2] are not relevant to our claims, since they
35 conflate the Fourier transform in the spatial domain (i.e. 2D Fourier transforming the input image, treated as a function
36 $\mathbb{R}^2 \rightarrow \mathbb{R}$) with the Fourier transform in function space. (i.e. Fourier transforming the classification function $\mathbb{R}^d \rightarrow \mathbb{R}$).
37 Finally, the work of [3] is largely unrelated to our work. The authors of [3] study how the internal layers of a network
38 vary with the *depth* of the layer, while we study how the end-to-end classification function evolves as a function of SGD
39 steps.

40 Reviewer 3 brings up two concerns with the performance correlation metric. First: why do we use $\mu_Y(F; L)$ instead of
41 simply $I(F; L)$? While it is true that $\mu_Y(F; L) \leq I(F; L)$, μ_Y captures the degree to which the information learned by
42 F about Y is explained by L – whereas $I(F; L)$ only captures the correlation of F and L , regardless of whether this
43 correlation is useful for predicting Y or not. For example, consider if $F(x) = L(x) \cdot \text{Bernoulli}(p)$. That is, F is a
44 linear classifier L with noisy outputs. Here, $I(F; L) \ll 1$, due to the noise in F . However, $\mu_Y(F; L) = I(F; Y)$, and
45 thus our metric recovers the fact that all the performance of F in predicting Y is explained by the linear L . Second,
46 Reviewer 3 describes a scenario where F first learns to classify the examples that are incorrectly classified by a linear
47 model L , and notes that μ treats this the same as learning the correct portion of L . This isn’t a problem, though, because
48 the incorrect portion of L is exactly the correct portion of the classifier $1 - L$, which is also linear. Contrast this with
49 the following scenario: the samples come from a *mixture* of L and an uncorrelated nonlinear model N , and F learns N
50 first. This is a true example of F not learning the linear part of the distribution, and accordingly $\mu_Y(F, L)$ will equal
51 zero. We will include formal examples to build intuition for our metric in the final version.

52 [1] On the Spectral Bias of Neural Networks (2018). By Rahaman, Baratin, Arpit, Draxler, Lin, Bengio, Courville.

53 [2] Understanding training and generalization in deep learning by Fourier analysis (2018). By Zhi-Qin John Xu.

54 [3] Understanding intermediate layers using linear classifier probes (2016). By Guillaume Alain and Yoshua Bengio.

¹For Equation (5) in <https://arxiv.org/pdf/1808.04295v4.pdf>, applying Parseval’s Theorem.