

PAPER • OPEN ACCESS

Learning sparse features can lead to overfitting in neural networks^{*}

To cite this article: Leonardo Petrini *et al* *J. Stat. Mech.* (2023) 114003

View the [article online](#) for updates and enhancements.

You may also like

- [Phase diagram of stochastic gradient descent in high-dimensional two-layer neural networks](#)
Rodrigo Veiga, Ludovic Stephan, Bruno Loureiro et al.
- [Self-consistent dynamical field theory of kernel evolution in wide neural networks](#)
Blake Bordelon and Cengiz Pehlevan
- [Short-time large deviations of the spatially averaged height of a Kardar–Parisi–Zhang interface on a ring](#)
Timo Schorlepp, Pavel Sasorov and Baruch Meerson

PAPER: ML 2023

Learning sparse features can lead to overfitting in neural networks*

Leonardo Petrini^{1,3,**}, Francesco Cagnetta^{1,3},
Eric Vanden-Eijnden² and Matthieu Wyart¹

¹ Institute of Physics École Polytechnique Fédérale de Lausanne, Lausanne, Switzerland

² Courant Institute of Mathematical Sciences New York University, New York, United States of America

E-mail: leonardo.petrini@epfl.ch

Received 22 May 2023

Accepted for publication 13 September 2023

Published 15 November 2023



CrossMark

Online at stacks.iop.org/JSTAT/2023/114003
<https://doi.org/10.1088/1742-5468/ad01b9>

Abstract. It is widely believed that the success of deep networks lies in their ability to learn a meaningful representation of the features of the data. Yet, understanding when and how this feature learning improves performance remains a challenge. For example, it is beneficial for modern architectures to be trained to classify images, whereas it is detrimental for fully-connected networks to be trained on the same data. Here, we propose an explanation for this puzzle, by showing that feature learning can perform worse than lazy training (via the random feature kernel or the neural tangent kernel) as the former can lead to a sparser neural representation. Although sparsity is known to be essential for learning anisotropic data, it is detrimental when the target function is constant or smooth along certain directions of the input space. We illustrate this phenomenon in two settings: (i) regression of Gaussian random functions on the d -dimensional unit sphere and (ii) classification of benchmark data sets

³Equal contribution (a coin was flipped).

*This article is an updated version of: Petrini L, Cagnetta F, Vanden-Eijnden E, Wyart M 2022 Learning sparse features can lead to overfitting in neural networks *Advances in Neural Information Processing Systems* vol 35, ed S Koyejo, S Mohamed, A Agarwal, D Belgrave, K Cho and A Oh Curran Associates, Inc pp 9403–16

**Author to whom any correspondence should be addressed.



Original Content from this work may be used under the terms of the [Creative Commons Attribution 4.0 licence](https://creativecommons.org/licenses/by/4.0/). Any further distribution of this work must maintain attribution to the author(s) and the title of the work, journal citation and DOI.

of images. For (i), we compute the scaling of the generalization error with the number of training points and show that methods that do not learn features generalize better, even when the dimension of the input space is large. For (ii), we show empirically that learning features can indeed lead to sparse and thereby less smooth representations of the image predictors. This fact is plausibly responsible for deteriorating the performance, which is known to be correlated with smoothness along diffeomorphisms.

Keywords: deep learning, neuronal networks, machine learning

Contents

1. Introduction 3
 1.1. Our contribution3
 1.2. Related work5
 2. Problem and notation 5
 3. Asymptotic analysis of generalization 8
 4. Numerical tests of the theory 12
 5. Evidence for overfitting along diffeomorphisms in image data sets 12
 6. Conclusion 15
 Acknowledgments 15
 Appendix A. Quick recap of spherical harmonics 16
 A.1. Expansion of ReLU and combinations thereof 17
 A.2. Dot-product kernels on the sphere 18
 Appendix B. Uniqueness and sparsity of the L1 minimizer 20
 Appendix C. Proof of Proposition 1 21
 Appendix D. Asymptotics of generalization in $d=2$ 24
 Appendix E. Asymptotic of generalization via the spectral bias ansatz 26
 Appendix F. Spectral bias via the replica calculation 27
 Appendix G. Training wide neural networks: does GD find the
 minimal-norm solution? 29
 Appendix H. Sensitivity of the predictor to transformations other
 than diffeomorphisms..... 31
 Appendix I. Maximum-entropy model of diffeomorphisms 32
 References 33

1. Introduction

Neural networks are responsible for a technological revolution in a variety of machine learning tasks. Many such tasks require learning functions of high-dimensional inputs from a finite set of examples, and thus should be generically hard due to the *curse of dimensionality* [1, 2]: the exponent that controls the scaling of the generalization error with the number of training examples is inversely proportional to the input dimension d . For instance, for standard image classification tasks with d in the range of $10^3 - 10^5$, the exponent should be practically vanishing, in contrast to what is observed in practice [3]. In this respect, understanding the success of neural networks is still an open question. A popular explanation is that, during training, neurons adapt to features in the data that are relevant for the task [4], effectively reducing the input dimension and making the problem tractable [5–7]. However, understanding quantitatively if this intuition is true and how it depends on the structure of the task remains a challenge.

Recently, much progress has been made in characterizing the conditions that lead to feature learning, in the overparameterized setting where networks generally perform best. When the initialization scale of the network parameters is large [8] one encounters the *lazy training regime*, where neural networks behave as kernel methods [9, 10] (coined neural tangent kernel or NTK) and features are not learned. By contrast, when the initialization scale is small, a *feature learning regime* is found [11–13] where the network parameters evolve significantly during training. This limit is much less well understood apart from very simple architectures, where it can be shown to lead to sparse representations where a limited number of neurons are active after training [14]. These sparse representations can also be obtained by regularizing the weights during training [2, 15].

In terms of performance, most theoretical works have focussed on fully-connected networks. For these architectures, feature learning was shown to significantly outperform lazy training [11, 16–19] for certain tasks, including approximating a function that depends only on a subset or a linear combination of the input variables. However, when these primitive networks are trained on image data sets, learning features are detrimental [20, 21], as illustrated in figure 1 (see [19, figure 3] for the analogous plot in the case of a target function depending on just one of the input variables, where learning features are beneficial). A similar result was observed in simple models of data [22]. These facts are unexplained, yet central to understanding the implicit bias of the feature learning regime.

1.1. Our contribution

Our main contribution is to provide an account of the drawbacks of learning sparse representations based on the following set of ideas. Consider, for concreteness, an image classification problem: *(i)* image class varies little along smooth deformations of the image; *(ii)* because tasks like image classification require a continuous distribution of neurons to be represented; *(iii)* thus, requiring sparsity can be detrimental to performance. We build our argument as follows.

Learning sparse features can lead to overfitting in neural networks

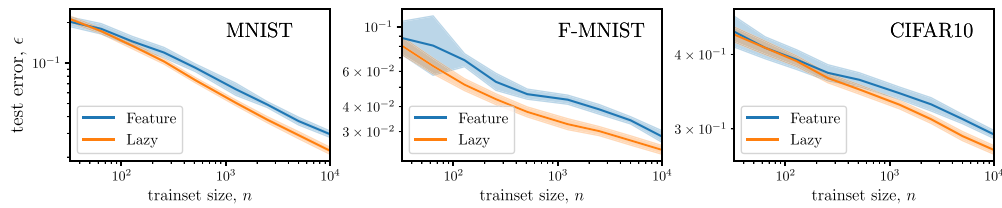


Figure 1. Feature versus lazy in image classification. Generalization error as a function of the training-set size n for infinite-width fully connected networks (FCNs) trained in the feature (blue) and lazy regime (orange). In the latter case the limit is arrived at exactly by training an SVC algorithm with the analytical NTK [23]. In the former case, the infinite-width limit can be accurately approximated for these data sets by considering very wide nets ($H = 10^3$), and performing ensemble averaging on different initial conditions of the parameters shown in [24, 25]. Panels correspond to different benchmark image data sets [26–28]. Results are averaged over ten different initializations of the networks and data sets.

- In order to find a quantitative description of the phenomenon, we start from the problem of regression of a random target function of controlled smoothness on the d -dimensional unit sphere and study the properties of the minimizers of the empirical loss with n observations, both in the lazy and the feature learning regimes. More specifically, we consider two extreme limits—the NTK limit and mean-field limit—as representatives of lazy and feature regimes, respectively (section 2). Both these limits admit a simple formulation that allows us to predict generalization performance. In particular, our results on feature learning rely on solutions that have atomic support. This property can be justified for one-hidden-layer neural networks with ReLU activations and weight decay. Yet, we also find such sparsity empirically using gradient descent (GD) in the absence of regularization, if the weights are initialized to be small enough.
- We find that lazy training leads to smoother predictors than feature learning. As a result, lazy training outperforms feature learning when the target function is also sufficiently smooth. Otherwise, the performance of the two methods is comparable, in the sense that they display the same asymptotic decay of generalization error with the number of training examples. Our predictions are obtained from asymptotic arguments that we systematically back up with numerical studies.
- For image data sets, it is believed that diffeomorphisms of images are key transformations along which the predictor function should only mildly vary to obtain good performance [29]. Based on the results above, a natural explanation as to why lazy outperforms feature for fully connected networks is that it leads to predictors with smaller variations along diffeomorphisms. We confirm that this is indeed the case empirically on benchmark data sets.

Numerical experiments are performed in PyTorch [30], and the code for reproducing the experiments is available online at github.com/pctl-epfl/regressionsphere.

1.2. Related work

The property training ReLU networks in the feature regime leads to a sparse representation that was observed empirically [31]. This property can be justified for one-hidden-layer networks by casting training as an L1 minimization problem [2, 32], then using a representer theorem [15, 33, 34]. This is analogous to what is commonly done in predictive sparse coding [35–38].

Many works have investigated the benefits of learning sparse representations in neural networks [2, 16–19, 39, 40] and study cases in which the true function only depends on a linear subspace of the input space, and show that feature learning profitably captures such property. Even for more general problems, sparse representations of the data might emerge naturally during deep network training—a phenomenon coined *neural collapse* [41]. Similar sparsification phenomena, for instance, have been found to allow for learning convolutional layers from scratch [42, 43]. Our study builds on this body of the literature by pointing out that learning sparse features can be detrimental, if the task does not allow for it.

There is currently no general framework to rigorously predict the learning curve exponent β defined as $\epsilon(n) = \mathcal{O}(n^{-\beta})$ for kernels. Some of our asymptotic arguments can be obtained by other approximations, such as assuming that data points lie on a lattice in \mathbb{R}^d [44] or by using the non-rigorous replica method of statistical physics [45–47]. In the case $d=2$, we provide a more explicit mathematical formulation of our results, which leads to analytical results for certain kernels. We systematically back up our predictions with numerical tests as d varies.

Finally, in the context of image classification, the connection between performance and ‘stability’ or smoothness towards small diffeomorphisms of the inputs has been conjectured by [29, 48]. Empirically, a strong correlation between these two quantities was shown to hold across various architectures for real data sets [49]. In this reference, it was found that fully connected networks lose their stability over training. Here, we show that this effect is much less pronounced in the lazy regime.

2. Problem and notation

Task. We consider a supervised learning scenario with n training points $\{\mathbf{x}_i\}_{i=1}^n$ uniformly drawn on the d -dimensional unit sphere \mathbb{S}^{d-1} . We assume that the target function f^* is an isotropic Gaussian random process on \mathbb{S}^{d-1} and control its statistics via the spectrum, by introducing the decomposition of f^* into spherical harmonics (see appendix A for definitions),

$$f^*(\mathbf{x}) = \sum_{k \geq 0} \sum_{\ell=1}^{\mathcal{N}_{k,d}} f_{k,\ell}^* Y_{k,\ell}(\mathbf{x}) \quad \text{with} \quad \mathbb{E}[f_{k,\ell}^*] = 0, \quad \mathbb{E}[f_{k,\ell}^* f_{k',\ell'}^*] = c_k \delta_{k,k'} \delta_{\ell,\ell'}. \quad (2.1)$$

We assume that all the c_k with k odd vanish apart from c_1 . This is required to guarantee that f^* can be approximated as well as desired with a one-hidden-layer ReLU network

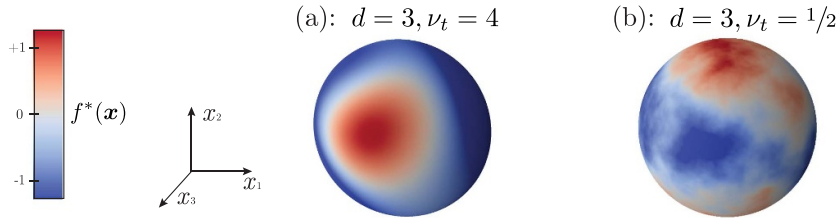


Figure 2. Gaussian random process on the sphere. We show two samples of the task introduced in section 2 when the target function $f^*(\mathbf{x})$ is defined on the 3D unit sphere. (a) and (b) show samples of large and small smoothness coefficient ν_t , respectively.

with no biases, as discussed in appendix A. We also assume that the non-zero c_k decay as a power of k for $k \gg 1$, $c_k \sim k^{-2\nu_t - (d-1)}$. The exponent $\nu_t > 0$ controls the (weak) differentiability of f^* on the sphere (see appendix A) and also the statistics of f^* in real space:

$$\mathbb{E} [|f^*(\mathbf{x}) - f^*(\mathbf{y})|^2] = O(|\mathbf{x} - \mathbf{y}|^{2\nu_t}) = O((1 - \mathbf{x} \cdot \mathbf{y})^{\nu_t}) \quad \text{as } \mathbf{x} \rightarrow \mathbf{y}. \quad (2.2)$$

Examples of such a target function for $d = 3$ and different values of ν_t are reported in figure 2.

Neural network representation in the feature regime. In this regime, we aim to approximate the target function $f^*(x)$ via a *one-hidden-layer neural network* of width H ,

$$f_H(\mathbf{x}) = \frac{1}{H} \sum_{h=1}^H w_h \sigma(\boldsymbol{\theta}_h \cdot \mathbf{x}), \quad (2.3)$$

where $\{\boldsymbol{\theta}_h\}_{h=1}^H$ (the features) and $\{w_h\}_{h=1}^H$ (the weights) are the network parameters to be optimized, and $\sigma(x)$ denotes the ReLU function, $\sigma(x) = \max\{0, x\}$. If we assume that $\{\boldsymbol{\theta}_h, w_h\}_{h=1}^H$ are independently drawn from a probability measure μ on $\mathbb{S}^{d-1} \times \mathbb{R}$ so that the Radon measure $\gamma = \int_{\mathbb{R}} w \mu(\cdot, dw)$ exists, then as $H \rightarrow \infty$,

$$\lim_{H \rightarrow \infty} f_H(\mathbf{x}) = \int_{\mathbb{S}^{d-1}} \sigma(\boldsymbol{\theta} \cdot \mathbf{x}) d\gamma(\boldsymbol{\theta}) \quad \text{a.e. on } \mathbb{S}^{d-1}. \quad (2.4)$$

This is the so-called mean-field limit [11, 12], and it is then natural to determine the optimal γ via,

$$\gamma^* = \arg \min_{\gamma} \int_{\mathbb{S}^{d-1}} |d\gamma(\boldsymbol{\theta})| \quad \text{subject to: } \int_{\mathbb{S}^{d-1}} \sigma(\boldsymbol{\theta} \cdot \mathbf{x}_i) d\gamma(\boldsymbol{\theta}) = f^*(\mathbf{x}_i) \quad \forall i = 1, \dots, n. \quad (2.5)$$

In practice, we can approximate this minimization problem using a network with large but finite width, constraining the feature to be on the sphere $|\boldsymbol{\theta}_h| = 1$, and minimizing the following empirical loss with L1 regularization on the weights,

$$\min_{\substack{\{w_h, \boldsymbol{\theta}_h\}_{h=1}^H \\ |\boldsymbol{\theta}_h|=1}} \frac{1}{2n} \sum_{i=1}^n \left(f^*(\mathbf{x}_i) - \frac{1}{H} \sum_{h=1}^H w_h \sigma(\boldsymbol{\theta}_h \cdot \mathbf{x}_i) \right)^2 + \frac{\lambda}{H} \sum_{h=1}^H |w_h|. \quad (2.6)$$

This minimization problem leads to (2.5) when $H \rightarrow \infty$ and $\lambda \rightarrow 0$. Note that, by homogeneity of ReLU (2.6) can be shown to be equivalent to imposing a regularization on the L2 norm of all parameters [32, theorem 10], i.e. the usual weight decay.

To proceed, we make the following assumption about the minimizer γ^* :

Assumption 1. The minimizer γ^* of (2.5) is unique and atomic, with $n_A \leq n$ atoms, i.e. $\{w_i^*, \boldsymbol{\theta}_i^*\}_{i=1}^{n_A}$ exists so that,

$$\gamma^* = \sum_{i=1}^{n_A} w_i^* \delta_{\boldsymbol{\theta}_i^*}. \quad (2.7)$$

The main component of the assumption is the uniqueness of γ^* ; if it holds, the sparsity of γ^* follows from the representer theorem, see e.g. [33]. Both the uniqueness and sparsity of the minimizer can be justified to hold generically using asymptotic arguments involving recasting the L1 minimization problem (2.5) as a linear programming one. These arguments are standard (see e.g. [50]) and are presented in appendix B for the reader's convenience. From our arguments below to deduce the scaling of the generalization error, we will mainly use that $n_A = O(n)$ —we shall confirm this fact numerically even in the absence of regularization if the weights are initialized to be small enough. Note that from assumption 1 it follows that the predictor in the feature regime corresponding to the minimizer γ^* takes the following form:

$$f^{\text{FEATURE}}(\mathbf{x}) = \sum_{i=1}^{n_A} w_i^* \sigma(\boldsymbol{\theta}_i^* \cdot \mathbf{x}). \quad (2.8)$$

Neural network representation in the lazy regime. In this regime, we approximate the target function $f^*(x)$ via,

$$f^{\text{NTK}}(\mathbf{x}) = \sum_{i=1}^n g_i K^{\text{NTK}}(\mathbf{x}_i \cdot \mathbf{x}), \quad (2.9)$$

where the weights $\{g_i\}_{i=1}^n$ solve,

$$f^*(\mathbf{x}_j) = \sum_{i=1}^n g_i K^{\text{NTK}}(\mathbf{x}_i \cdot \mathbf{x}_j), \quad j = 1, \dots, n. \quad (2.10)$$

and $K^{\text{NTK}}(\mathbf{x} \cdot \mathbf{y})$ is the NTK [9]:

$$K^{\text{NTK}}(\mathbf{x} \cdot \mathbf{y}) = \int_{\mathbb{S}^{d-1} \times \mathbb{R}} (\sigma(\boldsymbol{\theta} \cdot \mathbf{x}) \sigma(\boldsymbol{\theta} \cdot \mathbf{y}) + w^2 \mathbf{x} \cdot \mathbf{y} \sigma'(\boldsymbol{\theta} \cdot \mathbf{x}) \sigma'(\boldsymbol{\theta} \cdot \mathbf{y})) d\mu_0(\boldsymbol{\theta}, w). \quad (2.11)$$

Here, μ_0 is a fixed probability distribution which, in the NTK training regime [9], is the distribution of the features and weights at initialization. It is well-known [51] that the solution to the kernel ridge regression problem 2.10 can also be expressed via the kernel trick as,

$$f^{\text{NTK}}(\mathbf{x}) = \int_{\mathbb{S}^{d-1} \times \mathbb{R}} (g_w(\boldsymbol{\theta}, w) \sigma(\boldsymbol{\theta} \cdot \mathbf{x}) + w \mathbf{x} \cdot \mathbf{g}_\theta(\boldsymbol{\theta}, w) \sigma'(\boldsymbol{\theta} \cdot \mathbf{x})) d\mu_0(\boldsymbol{\theta}, w), \quad (2.12)$$

where \mathbf{g}_θ and g_w are the solutions of,

$$\begin{aligned} & \min_{g_w, \mathbf{g}_\theta} \int_{\mathbb{S}^{d-1} \times \mathbb{R}} (g_w^2(w, \boldsymbol{\theta}) + |\mathbf{g}_\theta(w, \boldsymbol{\theta})|^2) d\mu_0(\boldsymbol{\theta}, w) \\ \text{subject to: } & \int_{\mathbb{S}^{d-1} \times \mathbb{R}} (g_w(w, \boldsymbol{\theta}) \sigma(\boldsymbol{\theta} \cdot \mathbf{x}_i) + w \mathbf{x}_i \cdot \mathbf{g}_\theta(w, \boldsymbol{\theta}) \sigma'(\boldsymbol{\theta} \cdot \mathbf{x}_i)) d\mu_0(\boldsymbol{\theta}, w) = f^*(\mathbf{x}_i) \\ & \forall i = 1, \dots, n. \end{aligned} \quad (2.13)$$

Another lazy limit can be obtained equivalently by training only the weights while keeping the features at their initialization value. This is equivalent to forcing $\mathbf{g}_\theta(\boldsymbol{\theta}, w)$ to vanish in equation (2.13), again resulting in a kernel method. The kernel, in this case, is called the *random feature kernel* (K^{RFK}) and can be obtained from equation (2.11) by setting $d\mu_0(\boldsymbol{\theta}, w) = \delta_{w=0} d\tilde{\mu}_0(\boldsymbol{\theta})$. The minimizer can then be written as in equation (2.9) with K^{NTK} replaced by K^{RFK} .

3. Asymptotic analysis of generalization

In this section, we characterize the asymptotic decay of the generalization error $\bar{\epsilon}(n)$ averaged over several realizations of the target function f^* . Denoting with $d\tau^{d-1}(\mathbf{x})$ the uniform measure on \mathbb{S}^{d-1} ,

$$\bar{\epsilon}(n) = \mathbb{E}_{f^*} \left[\int d\tau^{d-1}(\mathbf{x}) (f^n(\mathbf{x}) - f^*(\mathbf{x}))^2 \right] = \mathcal{A}_d n^{-\beta} + o(n^{-\beta}), \quad (3.1)$$

for some constant \mathcal{A}_d , which might depend on d but not on n . For both the lazy (see equation (2.9)) and feature regimes (see equation (2.8)) the predictor can be written as the sum of $\mathcal{O}(n)$ terms:

$$f^n(\mathbf{x}) = \sum_{j=1}^{\mathcal{O}(n)} g_j \varphi(\mathbf{x} \cdot \mathbf{y}_j) := \int_{\mathbb{S}^{d-1}} g^n(\mathbf{y}) \varphi(\mathbf{x} \cdot \mathbf{y}) d\tau(\mathbf{y}). \quad (3.2)$$

In the feature regime, the g_j 's (\mathbf{y}_j) coincide with the optimal weights w_j^* (features $\boldsymbol{\theta}_j^*$), and φ with the activation function σ . In the lazy regime, \mathbf{y}_j are the training points \mathbf{x}_j , φ is the NTK or RFK and the g_j 's are the weights solving equation (2.10). We have defined the density $g^n(\mathbf{x}) = \sum_j |\mathbb{S}^{d-1}| g_j \delta(\mathbf{x} - \mathbf{y}_j)$ in order to cast the predictor as a convolution on the sphere. As a result, the projections of f^n onto spherical harmonics $Y_{k,\ell}$ read $f_{k,\ell}^n = g_{k,\ell}^n \varphi_k$, where $g_{k,\ell}^n$ is the projection of $g^n(\mathbf{x})$ and φ_k that of $\varphi(\mathbf{x} \cdot \mathbf{y})$. For ReLU neurons one has (as shown in appendix A):

$$\varphi_k^{\text{LAZY}} \sim k^{-(d-1)-2\nu} \quad \text{with } \nu = 1/2(\text{NTK}), 3/2(\text{RFK}), \quad \varphi_k^{\text{FEATURE}} \sim k^{-\frac{d-1}{2}-3/2}. \quad (3.3)$$

Main result. Consider a target function f^* with smoothness exponent ν_t as defined above, with data lying on \mathbb{S}^{d-1} . If f^* is learnt with a one-hidden-layer network with ReLU neurons in the regimes specified above, then the generalization error follows $\bar{\epsilon}(n) \sim n^{-\beta}$ with:

$$\beta^{\text{LAZY}} = \frac{\min\{2(d-1) + 4\nu, 2\nu_t\}}{d-1} \quad \text{with } \nu = \begin{cases} 1/2 \text{ for NTK,} \\ 3/2 \text{ for RFK,} \end{cases} \quad (3.4a)$$

$$\beta^{\text{FEATURE}} = \frac{\min\{(d-1) + 3, 2\nu_t\}}{d-1}. \quad (3.4b)$$

This is our central result. This implies that if the target function is a smooth isotropic Gaussian field (realized for large ν_t), then lazy outperforms feature, in the sense that training the network in the lazy regime leads to a better scaling of the generalization performance with the number of training points.

Strategy. There is no general framework for a rigorous derivation of the generalization error in the ridgeless limit $\lambda \rightarrow 0$. Predictions such as those of equation (3.4) can be obtained by either assuming that training points (for equation (3.4a)) and neurons (for equation (3.4b)) lie on a periodic lattice [44] or (for equation (3.4a)) using the replica method from physics [45], as shown in appendix F. Here, we follow a different route, by first characterizing the form of the predictor for $d=2$ (proof in appendix C). This property alone allows us to determine the asymptotic scaling of the generalization error. We use it to analytically obtain the generalization error in the NTK case with a slightly simplified function φ (details in appendix D). This calculation motivates a simple ansatz for the form of $g^n(\mathbf{x})$ entering equation (3.2) and its projections onto spherical harmonics, which extends naturally to arbitrary dimensions. We systematically confirm in numerical experiments the predictions resulting from this ansatz.

Properties of the predictor in $d = 2$. On the unit circle \mathbb{S}^1 all points are identified by a polar angle $x \in [0, 2\pi)$. Hence, both the target function and the estimated predictor are functions of the angle, and all functions of the scalar product are in fact functions of the difference in angle. In particular, introducing $\tilde{\varphi}(x) = \varphi(\cos(x))$,

$$f^n(x) = \sum_j g_j \tilde{\varphi}(x - x_j) \equiv \int_0^{2\pi} \frac{dy}{2\pi} g^n(y) \tilde{\varphi}(x - y), \quad (3.5)$$

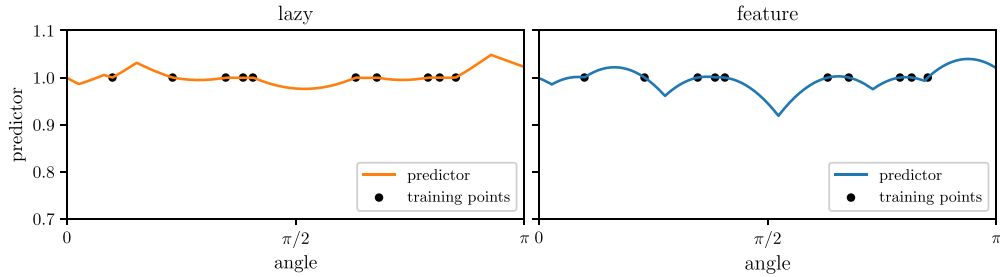


Figure 3. Feature versus lazy predictor. Predictor of the lazy (left) and feature (right) regime when learning the constant function on the ring with eight uniformly sampled training points.

where we defined,

$$g^n(x) = \sum_{j=1}^n (2\pi g_j) \delta(y - x_j). \tag{3.6}$$

For both the feature regime and the NTK limit, the first derivative of $\tilde{\varphi}(x)$ is continuous except for two values of x (0 and π for lazy, $-\pi/2$ and $\pi/2$ for feature) so that $\tilde{\varphi}(x)''$ has a singular part consisting of two Dirac delta functions.

As a result, the second derivative of the predictor $(f^n)''$ has a singular part consisting of many Dirac deltas. If we denote with $(f^n)''_r$ the regular part, obtained by subtracting all the delta functions, we show that (see appendix C):

Proposition 1 (informal). *As $n \rightarrow \infty$, $(f^n)''_r$ converges to a function having a finite second moment, i.e.*

$$\lim_{n \rightarrow \infty} \mathbb{E}_{f^*} [(f^n)''_r(x)]^2 = \text{const.} < \infty. \tag{3.7}$$

In the large n limit, the predictor displays a singular second derivative at $O(n)$ points. Proposition 1 implies that outside of these singular points the second derivative is well defined. Thus, as n becomes large and the singular points approach each other, the predictor can be approximated by a chain of parabolas, as highlighted in figure 3 and noted in [47] for a Laplace kernel. This property alone allows us to determine the asymptotic scaling of the error in $d=2$. In simple terms, Proposition 1 follows from the convergence of g^n to the function satisfying $f^*(x) = \int \frac{dy}{2\pi} g(y) \tilde{\varphi}_r(x-y)$, which is guaranteed under our assumptions on the target function—a detailed proof is given in appendix C.

Decay of the error in $d=2$ (sketch). The full calculation is in appendix D. Consider a slightly simplified problem where $\tilde{\varphi}$ has a single discontinuity in its derivative, located at $x=0$. In this case, $f^n(x)$ is singular if and only if x is a data point. Consider then the interval $x \in [x_i, x_{i+1}]$ and set $\delta_i = x_{i+1} - x_i$, $x_{i+1/2} = (x_{i+1} + x_i)/2$. If the target function is smooth enough ($\nu_t > 2$), then a Taylor expansion implies $|f^*(x_{i+1/2}) - f^n(x_{i+1/2})| \sim \delta_i^2$. Since the distances δ_i between adjacent singular points are random variables with

mean of order $1/n$ and finite moments, it is straightforward to obtain that $\bar{\epsilon}(n) \sim \sum_i (f^*(x_{i+1/2}) - f^n(x_{i+1/2}))^2 \sim \sum_i \delta_i^4 \sim n^{-4}$. By contrast, if f^* is not sufficiently smooth ($\nu_t \leq 2$), then $|f^*(x_{i+1/2}) - f^n(x_{i+1/2})| \sim \delta_i^{2\nu_t}$, leading to $\bar{\epsilon}(n) \sim n^{-2\nu_t}$. Note that for this asymptotic argument to apply to the feature learning regime, one must ensure that the distribution of the rescaled distance between adjacent singularities $n\delta_i$ has a finite fourth moment. This is obvious in the lazy regime, where the δ_i 's are controlled by the position of the training points, but not in the feature regime, where the distribution of singular points is determined by that of the neuron's features. Nevertheless, we show that this must be the case in our setup in appendix D.

Interpretation in terms of spectral bias. From the discussion above, it is evident that there is a length scale δ of order $1/n$ so that $f^n(x)$ is a good approximation of $f^*(x)$ over scales larger than δ . In terms of Fourier modes³, one has (i) $\widehat{f^n}(k)$, which matches $\widehat{f^n}(k)$ at long wavelengths, i.e. for $k \ll k_c \sim 1/n$. (ii) In addition, since the phases $\exp(ikx_j)$ become effectively random phases for $k \gg k_c$, $\widehat{g^n}(k) = \sum_j g_j \exp(ikx_j)$ becomes a Gaussian random variable with zero mean and fixed variance and thus (iii) $\widehat{f^n}(k) = \widehat{g^n}(k)\widehat{\varphi}(k)$ decorrelates from f^* for $k \gg k_c$. Therefore,

$$\bar{\epsilon}(n) \sim \sum_{|k| > k_c} \mathbb{E}_{f^*} \left[\left(\widehat{g^n}(k)\widehat{\varphi}(k) - \widehat{f^n}(k) \right)^2 \right] \sim \sum_{|k| \geq k_c} \mathbb{E}_{f^*} \left[\left(\widehat{g^n}(k) \right)^2 \right] \widehat{\varphi}(k)^2 + \mathbb{E}_{f^*} \left[\left(\widehat{f^n}(k) \right)^2 \right]. \quad (3.8)$$

For $\nu_t > 2$, one has $\sum_j g_j^2 \sim n^{-1} \lim_{n \rightarrow \infty} \int g^n(x)^2 dx \sim n^{-1}$. It follows (see appendix E for details) that the sum is dominated by the first term, hence entirely controlled by the Fourier coefficients of $\widehat{f^n}(k)$ at large k . A smoother predictor corresponds to a faster decay of $\widehat{f^n}(k)$ with k , thus a faster decay of the error with n . Plugging the relevant decays yields $\bar{\epsilon} \sim n^{-4}$ for feature regime and lazy regime with the NTK, and n^{-6} for lazy regime with the RFK (which is smoother than the NTK). For $\nu_t \leq 2$, the two terms have comparable magnitude (see appendix E), thus $\bar{\epsilon} \sim n^{-2\nu_t}$.

Generalization to higher dimensions. The argument above can be generalized for any d by replacing Fourier modes with projections onto spherical harmonics. Thus, the characteristic distance between training points scales as $n^{-1/(d-1)}$, thus $k_c \sim n^{-1/(d-1)}$. Our ansatz is that, as in $d=2$: (i) for $k \ll k_c$, the predictor modes coincide with those of the target function $f^n_{k,l} \approx f^*_{k,l}$ (this corresponds to the spectral bias result of kernel methods, stating that the predictor reproduces the first $O(n)$ projections of the target in the kernel eigenbasis [45]); (ii) for $k \gg k_c$, $g^n_{k,l}$ is a sum of uncorrelated terms, thus a Gaussian variable with zero mean and fixed variance; (iii) $f^n_{k,l} = g^n_{k,l}\varphi_k$ decorrelates from $f^*_{k,l}$ for $k \gg k_c$. (i), (ii) and (iii) imply that,

$$\bar{\epsilon}(n) \sim \sum_{k \geq k_c} \sum_{l=1}^{N_{k,d}} \mathbb{E}_{f^*} \left[\left(f^n_{k,l} - f^*_{k,l} \right)^2 \right] \sim \sum_{k \geq k_c} \sum_{l=1}^{N_{k,d}} \mathbb{E}_{f^*} \left[\left(g^n_{k,l} \right)^2 \right] \varphi_k^2 + k^{-2\nu_t - (d-1)}. \quad (3.9)$$

³ The Fourier transform of a function $f(x)$ is indicated by the hat, $\widehat{f}(k)$.

As shown in appendix E, from this expression it is straightforward to obtain equation (3.4). Note again that when the target is sufficiently smooth so that the predictor-dependent term dominates, the error is determined by the smoothness of the predictor. In particular, since $d > 2$, the predictor of feature learning is less smooth than both the NTK and RfK ones, due to the slower decay of the corresponding φ_k .

4. Numerical tests of the theory

We successfully test our predictions by computing the learning curves of both lazy and feature regimes when (i) the target function is constant on the sphere for varying d , see figure 4, and (ii) the target is a Gaussian random field with varying smoothness ν_i , as shown in figure G1 of appendix G. For the lazy regime, we perform kernel regression using the analytical expression of the NTK [52] (see also equation (A.19)). For the feature regime, we find that our predictions hold when having a small regularization, although it takes unreachable times for GD to exactly recover the minimal-norm solution—a more in-depth discussion can be found in appendix G. An example of the atomic distribution of neurons found after training, which contrasts with the initial distribution, is displayed in figure 5(a), left panel.

Another way to obtain sparse features is to initialize the network with very small weights [14], as proposed in [8]. As in the presence of an infinitesimal weight decay, this scheme also leads to sparse solutions with $n_A = \mathcal{O}(n)$ —an asymptotic dependence confirmed in figure G3 of appendix G. This observation implies that our predictions must apply in that case too, which we confirm in figure G3.

5. Evidence for overfitting along diffeomorphisms in image data sets

For fully-connected networks, the feature regime is well adapted to learn anisotropic tasks [16]. If the target function does not depend on a certain linear subspace of the input space, e.g. the pixels at the corner of an image, then neurons align perpendicularly to these directions [19]. In contrast, our results highlight a drawback of this regime when the target function is constant or smooth along directions in input space that require a continuous distribution of neurons to be represented. In this case, the adaptation of the weights to the training points leads to a predictor with a sparse representation. This predictor would be less smooth than in the lazy regime and thus underperform.

Does this view hold for images and explain why learning their features is detrimental for fully connected networks? The first positive empirical evidence is that the neurons' distribution of networks trained on image data indeed becomes sparse in the feature regime, as illustrated in figure 5(a), right, for CIFAR10 [28]. This observation raises the question of which are the directions in input space (i) along which the target should vary smoothly, and (ii) that are not easily represented by a discrete set of neurons. An example of these directions is global translations, which conserve the norm of the input and do not change the image class. The lazy regime predictor is indeed smoother than the feature one with respect to translations of the input (see appendix H). Yet, these

Learning sparse features can lead to overfitting in neural networks

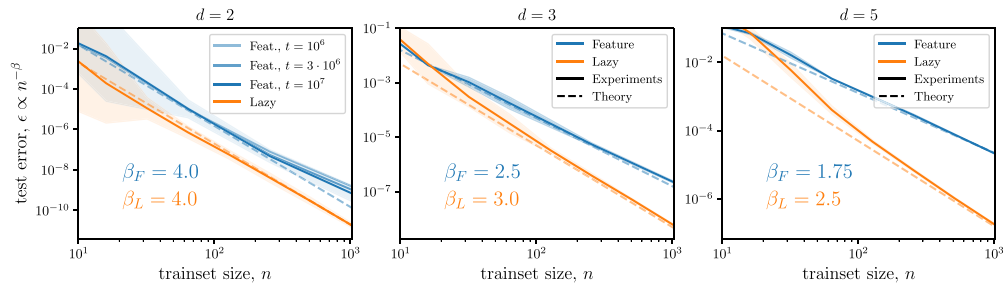


Figure 4. Generalization error for a constant function $f^*(\mathbf{x}) = 1$. Generalization error as a function of the training set size n for a network trained in the feature regime with L1 regularization (blue) and kernel regression corresponding to the infinite-width lazy regime (orange). Numerical results (full lines) and the exponents predicted by the theory (dashed) are plotted. Panels correspond to different input-space dimensions ($d = 2, 3, 5$). Results are averaged over ten different initializations of the networks and data sets. For $d = 2$ and large n , the gap between experiments and prediction for the feature regime is due to the finite training time t . Indeed our predictions become more accurate as t increases, as illustrated on the left.

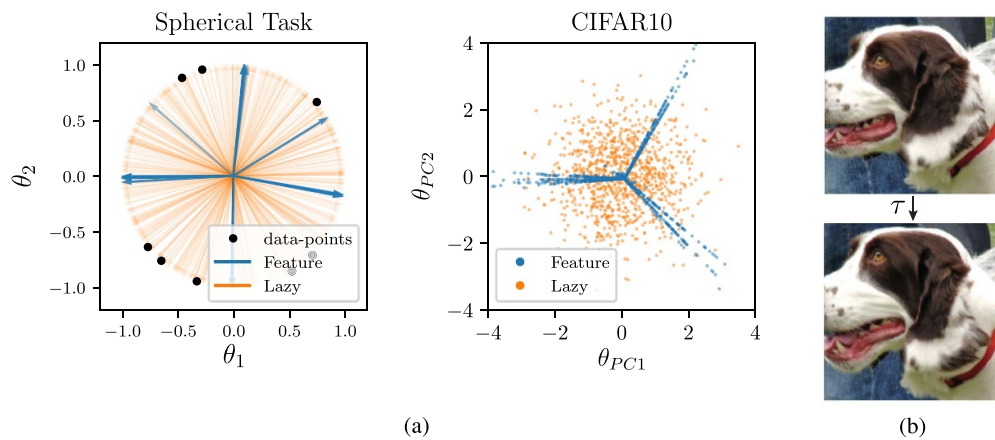


Figure 5. Features sparsification and example of a diffeomorphism. (a) Features sparsification. 1st Panel: Distribution of neuron’s feature for the task of learning a constant function on the sphere in 2D. Arrows represent a subset of the network features $\{\theta_h\}_{h=1}^H$ after training in the lazy and feature regimes. Training is performed on $n = 8$ data points (black dots). 2nd Panel: FCN trained on CIFAR10. On the axes the first two principal components of the features $\{\theta_h\}_{h=1}^H$ after training on $n = 32$ points in the feature (blue) and lazy (orange) regimes. Similar to what is observed when learning a constant function, the θ_h angular distribution becomes sparse with training in the feature regime. (b) Example of diffeomorphism. Sample of a max-entropy deformation τ [49] when applied to a natural image, illustrating that it does not change the image class for the human brain.

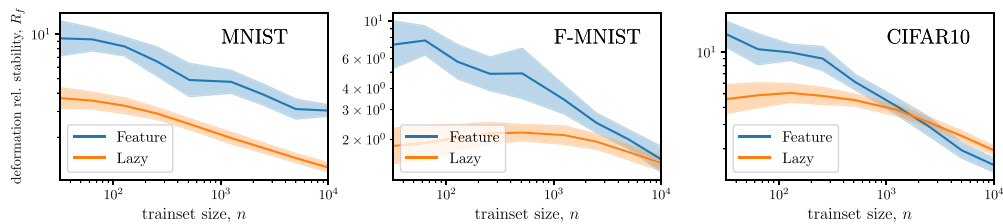


Figure 6. Sensitivity to diffeomorphisms versus number of training points. Relative sensitivity of the predictor to small diffeomorphisms of the input images, in the two regimes, for varying number of training points n and different image data sets. Smaller values correspond to a smoother predictor, on average. Results are computed using the same predictors as in figure 1.

transformations live in a space of dimension 2, which is small in comparison with the full dimensionality d of the data and thus may play a negligible role.

A much larger class of transformations believed to have little effect on the target are small diffeomorphisms [29]. A diffeomorphism τ acting on an image is illustrated in figure 5(b), which highlights that our brain still perceives the content of the transformed image as in the original one. Near-invariance of the task to these transformations is believed to play a key role in the success of deep learning, and in explaining how neural networks beat the curse of dimensionality [48]. Indeed, if modern architectures can become insensitive to these transformations, then the dimensionality of the problem is considerably reduced. In fact, it was found that the architectures displaying the best performance are precisely those that learn to vary smoothly along these transformations [49].

Small diffeomorphisms are likely the directions we are looking for. In order to test this hypothesis, following [49], we characterize the smoothness of a function along these diffeomorphisms, relative to that of random directions in the input space. Specifically, we use the *relative sensitivity*:

$$R_f = \frac{\mathbb{E}_{x,\tau} \|f(\tau x) - f(x)\|^2}{\mathbb{E}_{x,\eta} \|f(x + \eta) - f(x)\|^2}. \quad (5.1)$$

In the numerator, the average is given over the test set and over an ensemble of diffeomorphisms, reviewed in appendix I. The magnitude of the diffeomorphisms is chosen so that each pixel is shifted by one on average. In the denominator, the average runs over the test set and the vectors η sampled uniformly on the sphere of radius $\|\eta\| = \mathbb{E}_{x,\tau} \|\tau x - x\|$, and this fixes the transformation magnitude.

We measure R_f as a function of n for three benchmark data sets of images, as shown in figure 6. We indeed find that R_f is consistently smaller in the lazy training regime, where features are not learned. Overall, this observation supports the view that learning sparse features is detrimental when data present (near) invariance to transformations that cannot be represented sparsely by the architecture considered. Figure 1 supports the idea that—for benchmark image data sets—this negative effect

overcomes the well-known positive effects of learning features, e.g. becoming insensitive to pixels on the edges of images (see appendix H for evidence of this effect).

6. Conclusion

Our central result is that learning sparse features can be detrimental if the task presents invariance or smooth variations along transformations that are not adequately captured by the neural network architecture. For fully connected networks, these transformations can be rotations of the input, but also continuous translations and diffeomorphisms.

Our analysis relies on the sparsity of the features learned by a shallow fully connected architecture. Even in the infinite width limit, when trained in the feature learning regime, these networks behave as $\mathcal{O}(n)$ neurons. The asymptotic analysis that we perform for random Gaussian fields on the sphere leads to predictions for the learning curve exponent β in different training regimes, which we verify. These kinds of results are scarce in the literature.

Note that our analysis focuses on ReLU neurons because (i) these are very often used in practice and (ii) in that case, β will depend on the training regime, allowing for stringent numerical tests. If smooth activations (e.g. softplus) are considered, we expect that learning features will still be detrimental for generalization. Yet, the difference will not appear in the exponent β , but in other aspects of the learning curves (including numerical coefficients and pre-asymptotic effects) that are harder to predict.

Most fundamentally, our results underline that the success of feature learning for modern architectures still lacks a sufficient explanation. Indeed, most of the theoretical studies that previously emphasized the benefits of learning features have been considering fully connected networks, for which learning features can, in practice, be a drawback. It is tempting to argue that, in modern architectures, learning features are not at a disadvantage because smoothness along diffeomorphisms can be enforced from the start—due to the locally connected, convolutional and pooling layers [29, 53]. Yet, the best architectures often do not perform pooling and are not stable towards diffeomorphisms at initialization. *During training*, learning features lead to more stable and smoother solutions along diffeomorphisms [49, 54]. Understanding why building sparse features enhances stability in these architectures may ultimately explain the magical feat of deep CNNs, learning tasks in high dimensions.

Acknowledgments

We gratefully acknowledge Lénaïc Chizat, Antonio Sclocchi and Umberto M Tomasini for helpful discussions. The work of MW is supported by a grant from the Simons Foundation (Grant No. 454953) and from the NSF under Grant No. 200021-165509. The work of EVE is supported by the National Science Foundation under Awards DMR-1420073, DMS-2012510 and DMS-2134216, by the Simons Collaboration on Wave Turbulence, Grant No. 617006, and by a Vannevar Bush Faculty Fellowship.

Appendix A. Quick recap of spherical harmonics

Spherical harmonics. This appendix collects some introductory background on spherical harmonics and dot-product kernels on the sphere [55]. See [56, 57] for an expanded treatment. Spherical harmonics are homogeneous polynomials on the sphere $\mathbb{S}^{d-1} = \{\mathbf{x} \in \mathbb{R}^d \mid \|\mathbf{x}\| = 1\}$, with $\|\cdot\|$ denoting the L2 norm. Given the polynomial degree $k \in \mathbb{N}$, there are $\mathcal{N}_{k,d}$ linearly independent spherical harmonics of degree k on \mathbb{S}^{d-1} , with,

$$\mathcal{N}_{k,d} = \frac{2k+d-2}{k} \binom{d+k-3}{k-1}, \quad \begin{cases} \mathcal{N}_{0,d} = 1 & \forall d, \\ \mathcal{N}_{k,d} \asymp A_d k^{d-2} & \text{for } k \gg 1, \end{cases} \quad (\text{A.1})$$

where \asymp means logarithmic equivalence for $k \rightarrow \infty$ and $A_d = \sqrt{2/\pi}(d-2)^{\frac{3}{2}-d}e^{d-2}$. Thus, we can introduce a set of $\mathcal{N}_{k,d}$ spherical harmonics $Y_{k,\ell}$ for each k , with ℓ ranging in $1, \dots, \mathcal{N}_{k,d}$, which are orthonormal with respect to the uniform measure on the sphere $d\tau(\mathbf{x})$,

$$\{Y_{k,\ell}\}_{k \geq 0, \ell = 1, \dots, \mathcal{N}_{k,d}}, \quad \langle Y_{k,\ell}, Y_{k,\ell'} \rangle_{\mathbb{S}^{d-1}} := \int_{\mathbb{S}^{d-1}} Y_{k,\ell}(\mathbf{x}) Y_{k,\ell'}(\mathbf{x}) d\tau(\mathbf{x}) = \delta_{\ell,\ell'}. \quad (\text{A.2})$$

Because of the orthogonality of homogeneous polynomials with different degree, the set is a complete orthonormal basis for the space of square-integrable functions on \mathbb{S}^{d-1} . For any function $f: \mathbb{S}^{d-1} \rightarrow \mathbb{R}$, then,

$$f(\mathbf{x}) = \sum_{k \geq 0} \sum_{\ell=1}^{\mathcal{N}_{k,d}} f_{k,\ell} Y_{k,\ell}(\mathbf{x}), \quad f_{k,\ell} = \int_{\mathbb{S}^{d-1}} f(\mathbf{x}) Y_{k,\ell}(\mathbf{x}) d\tau(\mathbf{x}). \quad (\text{A.3})$$

Furthermore, spherical harmonics are eigenfunctions of the Laplace–Beltrami operator Δ , which is nothing but the restriction of the standard Laplace operator to \mathbb{S}^{d-1} ,

$$\Delta Y_{k,\ell} = -k(k+d-2)Y_{k,\ell}. \quad (\text{A.4})$$

Legendre polynomials. By fixing a direction \mathbf{y} in \mathbb{S}^{d-1} , one can select, for each k , the only spherical harmonic of degree k , which is invariant for rotations that leave \mathbf{y} unchanged. This particular spherical harmonic is, in fact, a function of $\mathbf{x} \cdot \mathbf{y}$ and is called the Legendre polynomial of degree k , $P_{k,d}(\mathbf{x} \cdot \mathbf{y})$ (also referred to as Gegenbauer polynomial). Legendre polynomials can be written as a combination of the orthonormal spherical harmonics $Y_{k,\ell}$ via the addition theorem [56, theorem 2.9],

$$P_{k,d}(\mathbf{x} \cdot \mathbf{y}) = \frac{1}{\mathcal{N}_{k,d}} \sum_{\ell=1}^{\mathcal{N}_{k,d}} Y_{k,\ell}(\mathbf{x}) Y_{k,\ell}(\mathbf{y}). \quad (\text{A.5})$$

Alternatively, $P_{k,d}$ is given explicitly as a function of $t = \mathbf{x} \cdot \mathbf{y} \in [-1, 1]$ via the Rodrigues formula [56, theorem 2.23],

$$P_{k,d}(t) = \left(-\frac{1}{2}\right)^k \frac{\Gamma\left(\frac{d-1}{2}\right)}{\Gamma\left(k + \frac{d-1}{2}\right)} (1-t^2)^{\frac{3-d}{2}} \frac{d^k}{dt^k} (1-t^2)^{k+\frac{d-3}{2}}. \tag{A.6}$$

Here, Γ denotes the Gamma function, and $\Gamma(z) = \int_0^\infty x^{z-1} e^{-x} dx$. Legendre polynomials are orthogonal on $[-1, 1]$ with respect to the measure with density $(1-t^2)^{(d-3)/2}$, which is the probability density function of the scalar product between two points on \mathbb{S}^{d-1} :

$$\int_{-1}^{+1} P_{k,d}(t) P_{k',d}(t) (1-t^2)^{\frac{d-3}{2}} dt = \frac{|\mathbb{S}^{d-1}| \delta_{k,k'}}{|\mathbb{S}^{d-2}| \mathcal{N}_{k,s}}. \tag{A.7}$$

Here, $|\mathbb{S}^{d-1}| = 2\pi^{\frac{d}{2}}/\Gamma(\frac{d}{2})$ denotes the surface area of the d -dimensional unit sphere ($|\mathbb{S}^0|=2$ by definition).

To sum up, given $\mathbf{x}, \mathbf{y} \in \mathbb{S}^{d-1}$, functions of \mathbf{x} or \mathbf{y} can be expressed as a sum of projections on the orthonormal spherical harmonics, whereas functions of $\mathbf{x} \cdot \mathbf{y}$ can be expressed as a sum of projections on the Legendre polynomials. The relationship between the two expansions is elucidated in the Funk–Hecke formula [56, theorem 2.22]:

$$\int_{\mathbb{S}^{d-1}} f(\mathbf{x} \cdot \mathbf{y}) Y_{k,\ell}(\mathbf{y}) d\tau(\mathbf{y}) = Y_{k,\ell}(\mathbf{x}) \frac{|\mathbb{S}^{d-2}|}{|\mathbb{S}^{d-1}|} \int_{-1}^{+1} f(t) P_{k,d}(t) (1-t^2)^{\frac{d-3}{2}} dt := f_k Y_{k,\ell}(\mathbf{x}). \tag{A.8}$$

A.1. Expansion of ReLU and combinations thereof

We can apply equation (A.8) to have an expansion of neurons $\sigma(\boldsymbol{\theta} \cdot \mathbf{x})$ in terms of spherical harmonics [2, appendix D]. After defining,

$$\varphi_k := \frac{|\mathbb{S}^{d-2}|}{|\mathbb{S}^{d-1}|} \int_{-1}^{+1} \sigma(t) P_{k,d}(t) (1-t^2)^{\frac{d-3}{2}} dt, \tag{A.9}$$

one has,

$$\sigma(\boldsymbol{\theta} \cdot \mathbf{x}) = \sum_{k \geq 0} \mathcal{N}_{k,d} \varphi_k P_{k,d}(\boldsymbol{\theta} \cdot \mathbf{x}) = \sum_{k \geq 0} \varphi_k \sum_{\ell=1}^{\mathcal{N}_{k,d}} Y_{k,\ell}(\boldsymbol{\theta}) Y_{k,\ell}(\mathbf{x}). \tag{A.10}$$

For ReLU activations, in particular, $\sigma(t) = \max(0, t)$, thus,

$$\varphi_k^{\text{ReLU}} = \frac{|\mathbb{S}^{d-2}|}{|\mathbb{S}^{d-1}|} \int_0^{+1} t P_{k,d}(t) (1-t^2)^{\frac{d-3}{2}} dt. \tag{A.11}$$

Note that when k is odd, $P_{k,d}$ is an odd function of t , thus the integrand $t P_{k,d}(t) (1-t^2)^{\frac{d-3}{2}}$ is an even function of t . As a result, the integral on the right-hand side of equation (A.11) coincides with half the integral over the full domain $[-1, 1]$:

$$\int_0^{+1} t P_{k,d}(t) (1-t^2)^{\frac{d-3}{2}} dt = \frac{1}{2} \int_{-1}^{+1} t P_{k,d}(t) (1-t^2)^{\frac{d-3}{2}} dt = 0 \text{ for } k > 1, \quad (\text{A.12})$$

because, due to equation (A.7), $P_{k,d}$ is orthogonal to all polynomials with degree strictly lower than k . For even k we can use equation (A.6) and obtain [2] (see equation (3.3), main text):

$$\begin{aligned} \int_0^{+1} t P_{k,d}(t) (1-t^2)^{\frac{d-3}{2}} dt &= \left(-\frac{1}{2}\right)^k \frac{\Gamma\left(\frac{d-1}{2}\right)}{\Gamma\left(k + \frac{d-1}{2}\right)} \int_0^1 t \frac{d^k}{dt^k} (1-t^2)^{k+\frac{d-3}{2}} dt \\ &= -\left(-\frac{1}{2}\right)^k \frac{\Gamma\left(\frac{d-1}{2}\right)}{\Gamma\left(k + \frac{d-1}{2}\right)} \frac{d^{k-2}}{dt^{k-2}} (1-t^2)^{k+\frac{d-3}{2}} \Big|_{t=0}^{t=1} \\ &\Rightarrow \varphi_k^{\text{ReLU}} \sim k^{-\frac{d-1}{2}-\frac{3}{2}} \text{ for } k \gg 1 \text{ and even.} \end{aligned} \quad (\text{A.13})$$

Because all φ_k^{ReLU} with $k > 1$ and odd vanish, even summing an infinite number of neurons $\sigma(\boldsymbol{\theta} \cdot \mathbf{x})$ with varying $\boldsymbol{\theta}$ does not allow us to approximate any function on \mathbb{S}^{d-1} , but only those that have vanishing projections on all the spherical harmonics $Y_{k,\ell}$ with $k > 1$ and odd. This is why we set the odd coefficients of the target function spectrum to zero in equation (2.1).

A.2. Dot-product kernels on the sphere

In addition, general dot-product kernels on the sphere admit an expansion such as equation (A.10),

$$\mathcal{C}(\mathbf{x} \cdot \mathbf{y}) = \sum_{k \geq 0} \mathcal{N}_{k,d} c_k P_{k,d}(\boldsymbol{\theta} \cdot \mathbf{x}) = \sum_{k \geq 0} c_k \sum_{\ell=1}^{\mathcal{N}_{k,d}} Y_{k,\ell}(\boldsymbol{\theta}) Y_{k,\ell}(\mathbf{x}), \quad (\text{A.14})$$

with,

$$c_k = \frac{|\mathbb{S}^{d-2}|}{|\mathbb{S}^d|} \int_{-1}^1 \mathcal{C}(t) P_{k,d}(t) (1-t^2)^{\frac{d-3}{2}} dt. \quad (\text{A.15})$$

The asymptotic decay of c_k for large k is controlled by the behavior of $\mathcal{C}(t)$ near $t = \pm 1$, [58]. More precisely [58, theorem 1], if \mathcal{C} is infinitely differentiable in $(-1, 1)$ and has the following expansion around ± 1 ,

$$\begin{cases} \mathcal{C}(t) = p_1(1-t) + c_1(1-t)^\nu + o((1-t)^\nu) \text{ near } t = +1; \\ \mathcal{C}(t) = p_{-1}(-1+t) + c_{-1}(-1+t)^\nu + o((-1+t)^\nu) \text{ near } t = -1, \end{cases} \quad (\text{A.16})$$

where $p_{\pm 1}$ are polynomials and ν is not an integer, then,

$$\begin{aligned} k \text{ even: } c_k &\sim (c_1 + c_{-1}) k^{-2\nu-(d-1)}; \\ k \text{ odd: } c_k &\sim (c_1 - c_{-1}) k^{-2\nu-(d-1)}. \end{aligned} \quad (\text{A.17})$$

The result above implies that that if $c_1 = c_{-1}$ ($c_1 = -c_{-1}$), then the eigenvalues with k odd (even) decay faster than $k^{-2\nu-(d-2)}$. Moreover, if \mathcal{C} is infinitely differentiable in $[-1, 1]$, then c_k decays faster than any polynomial.

NTK and RFK of one-hidden-layer ReLU networks. Let \mathbb{E}_θ denote expectation over a multivariate normal distribution with zero mean and unitary covariance matrix. For any $\mathbf{x}, \mathbf{y} \in \mathbb{S}^{d-1}$, the RFK of a one-hidden-layer ReLU network equation (2.3) with all parameters initialized as independent Gaussian random numbers with zero mean and unit variance reads,

$$\begin{aligned} K^{\text{RFK}}(\mathbf{x} \cdot \mathbf{y}) &= \mathbb{E}_\theta[\sigma(\boldsymbol{\theta} \cdot \mathbf{x})\sigma(\boldsymbol{\theta} \cdot \mathbf{y})] \\ &= \frac{(\pi - \arccos(t))t + \sqrt{1-t^2}}{2\pi}, \text{ with } t = \mathbf{x} \cdot \mathbf{y}. \end{aligned} \tag{A.18}$$

The NTK of the same network reads, with σ' denoting the derivative of ReLU or Heaviside function,

$$\begin{aligned} K^{\text{NTK}}(\mathbf{x} \cdot \mathbf{y}) &= \mathbb{E}_\theta[\sigma(\boldsymbol{\theta} \cdot \mathbf{x})\sigma(\boldsymbol{\theta} \cdot \mathbf{y})] + (\mathbf{x} \cdot \mathbf{y})\mathbb{E}_\theta[\sigma'(\boldsymbol{\theta} \cdot \mathbf{x})\sigma'(\boldsymbol{\theta} \cdot \mathbf{y})] \\ &= \frac{2(\pi - \arccos(t))t + \sqrt{1-t^2}}{2\pi}, \text{ with } t = \mathbf{x} \cdot \mathbf{y}. \end{aligned} \tag{A.19}$$

As functions of a dot-product on the sphere, both NTK and RFK admit a decomposition in terms of spherical harmonics as equation (A.15). For dot-product kernels, this expansion coincides with the Mercer’s decomposition of the kernel [55], that is, the coefficients of the expansion are the eigenvalues of the kernel. The asymptotic decay of the eigenvalues of these kernels φ_k^{NTK} and φ_k^{RFK} can be obtained by applying equation (A.16) [58, theorem 1]. Equivalently, one can see that K^{RFK} is proportional to the convolution on the sphere of ReLU with itself, therefore $\varphi_k^{\text{RFK}} = (\varphi_k^{\text{ReLU}})^2$. Similarly, the asymptotic decay of φ_k^{NTK} can be related to that of the coefficients of σ' , derivative of ReLU, $\varphi_k(\sigma') \sim k\varphi_k(\sigma)$, thus $\varphi_k^{\text{NTK}} \sim k^2(\varphi_k^{\text{ReLU}})^2$. Both methods lead to equation (3.3) of the main text.

Gaussian random fields and equation (2.2). Consider a Gaussian random field f^* on the sphere with covariance kernel $\mathcal{C}(\mathbf{x} \cdot \mathbf{y})$,

$$\mathbb{E}[f^*(\mathbf{x})] = 0, \quad \mathbb{E}[f^*(\mathbf{x})f^*(\mathbf{y})] = \mathcal{C}(\mathbf{x} \cdot \mathbf{y}), \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{S}^{d-1}. \tag{A.20}$$

f^* can be equivalently specified via the statistics of the coefficients $f_{k,\ell}^*$,

$$\mathbb{E}[f_{k,\ell}^*] = 0, \quad \mathbb{E}[f_{k,\ell}^*f_{k',\ell'}^*] = c_k\delta_{k,k'}\delta_{\ell,\ell'}, \tag{A.21}$$

with c_k denoting the eigenvalues of \mathcal{C} in equation (A.15). Note that the eigenvalues are degenerate with respect to ℓ because the covariance kernel is a function $\mathbf{x} \cdot \mathbf{y}$. As a result, the random function f^* is isotropic in law.

If c_k decays as a power of k , then this power controls the weak differentiability (in the mean-squared sense) of the random field f^* . In fact, from equation (A.4),

$$\left\| \Delta^{m/2} f^* \right\|^2 = \sum_{k \geq 0} \sum_{\ell} (-k(k+d-2))^m (f_{k,\ell}^*)^2. \tag{A.22}$$

Upon averaging over f^* one gets,

$$\mathbb{E} \left[\left\| \Delta^{m/2} f^* \right\|^2 \right] = \sum_{k \geq 0} (-k(k+d-2))^m \sum_{\ell} \mathbb{E} \left[(f_{k,\ell}^*)^2 \right] = \sum_{k \geq 0} (-k(k+d-2))^m \mathcal{N}_{k,d} c_k. \tag{A.23}$$

From equation (A.16) [58, theorem 1], if $\mathcal{C}(t) \sim (1-t)^{\nu_t}$ for $t \rightarrow 1$ and/or $\mathcal{C}(t) \sim (-1+t)^{\nu_t}$ for $t \rightarrow -1$, then $c_k \sim k^{-2\nu_t-(d-1)}$ for $k \gg 1$. In addition, for finite but arbitrary d , $(-k(k+d-2))^m \sim k^{2m}$ and $\mathcal{N}_{k,s} \sim k^{d-2}$ (see equation (A.1)). Hence, the summand in the right-hand side of equation (A.23) is $\sim k^{2(m-\nu_t)-1}$, thus,

$$\mathbb{E} \left[\left\| \Delta^{m/2} f^* \right\|^2 \right] < \infty \quad \forall m < \nu_t. \tag{A.24}$$

Alternatively, one can think of ν_t as controlling the scaling of the difference δf^* over inputs separated by a distance δ . From equation (A.20),

$$\begin{aligned} \mathbb{E} [|f^*(\mathbf{x}) - f^*(\mathbf{y})|^2] &= 2\mathcal{C}(1) - 2\mathcal{C}(\mathbf{x} \cdot \mathbf{y}) = 2\mathcal{C}(1) + O((1 - \mathbf{x} \cdot \mathbf{y})^{\nu_t}) \\ &= 2\mathcal{C}(1) + O(|\mathbf{x} - \mathbf{y}|^{2\nu_t}) \end{aligned} \tag{A.25}$$

Appendix B. Uniqueness and sparsity of the L1 minimizer

Recall that we want to find the γ^* that solves,

$$\gamma^* = \arg \min_{\gamma} \int_{\mathbb{S}^{d-1}} |d\gamma(\boldsymbol{\theta})| \quad \text{subject to} \quad \int_{\mathbb{S}^{d-1}} \sigma(\boldsymbol{\theta} \cdot \mathbf{x}_i) d\gamma(\boldsymbol{\theta}) = f^*(\mathbf{x}_i) \quad \forall i = 1, \dots, n. \tag{B.1}$$

In this appendix, we argue that the uniqueness of γ^* , which implies that it is atomic with at most n atoms, is a natural assumption. We start by discretizing the measure γ into H atoms, with H arbitrarily large. Then, the problem equation (B.1) can be rewritten as,

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \|\mathbf{w}\|_1, \quad \text{subject to} \quad \boldsymbol{\Phi} \mathbf{w} = \mathbf{y}, \tag{B.2}$$

with $\boldsymbol{\Phi} \in \mathbb{R}^{H \times n}$, $\Phi_{h,i} = \sigma(\boldsymbol{\theta}_h \cdot \mathbf{x}_i)$ and $y_i = f^*(\mathbf{x}_i)$.

Given $\mathbf{w} \in \mathbb{R}^H$, let $\mathbf{u} = \max(\mathbf{w}, 0) \geq \mathbf{0}$ and $\mathbf{v} = -\max(-\mathbf{w}, 0) \geq \mathbf{0}$ so that $\mathbf{w} = \mathbf{u} - \mathbf{v}$. It is well known (see e.g. [50]) that the minimization problem in (B.2) can be recast in terms of \mathbf{u} and \mathbf{v} into a linear programming problem. That is, $\mathbf{w}^* = \mathbf{u}^* - \mathbf{v}^*$ with,

$$(\mathbf{u}^*, \mathbf{v}^*) = \operatorname{argmin}_{\mathbf{u}, \mathbf{v}} \mathbf{e}^T (\mathbf{u} + \mathbf{v}), \quad \text{subject to } \Phi \mathbf{u} - \Phi \mathbf{v} = \mathbf{y}, \quad \mathbf{u} \geq \mathbf{0}, \quad \mathbf{v} \geq \mathbf{0}, \quad (\text{B.3})$$

where $\mathbf{e} = [1, 1, \dots, 1]^T$. Assuming that this problem is feasible (i.e. there is at least one solution to $\Phi \mathbf{u} - \Phi \mathbf{v} = \mathbf{y}$ so that $\mathbf{u} \geq \mathbf{0}, \mathbf{v} \geq \mathbf{0}$), it is known that it admits an extremal solution, i.e. solutions so that at most n entries of $(\mathbf{u}^*, \mathbf{v}^*)$ (and hence \mathbf{w}^*) are non-zero. The issue is whether such an extremal solution is unique. Assume that there are two, say $(\mathbf{u}_1^*, \mathbf{v}_1^*)$ and $(\mathbf{u}_2^*, \mathbf{v}_2^*)$. Then, by convexity,

$$(\mathbf{u}_t^*, \mathbf{v}_t^*) = (\mathbf{u}_1^*, \mathbf{v}_1^*) t + (\mathbf{u}_2^*, \mathbf{v}_2^*) (1 - t), \quad (\text{B.4})$$

is also a minimizer of (B.3) for all $t \in [0, 1]$, with the same minimum value $\mathbf{u}_t^* + \mathbf{v}_t^* = \mathbf{u}_1^* + \mathbf{v}_1^* = \mathbf{u}_2^* + \mathbf{v}_2^*$. Generalizing this argument to the case of more than two extremal solutions, we conclude that all minimizers are global, with the same minimum value, and they live on the simplex where $\mathbf{e}^T (\mathbf{u} + \mathbf{v}) = \mathbf{e}^T (\mathbf{u}_1 + \mathbf{v}_1)$. Therefore, nonuniqueness requires that this simplex has a nontrivial intersection with the feasible set where $\Phi \mathbf{u} - \Phi \mathbf{v} = \mathbf{y}$ with $\mathbf{u} \geq \mathbf{0}, \mathbf{v} \geq \mathbf{0}$. We argue that, generically, this will not be the case, i.e. the intersection will be trivial, and the extremal solution unique. In particular, since in our case we are in fact interested in the problem (B.1), we can always perturb slightly the discretization into H atoms of γ to guarantee that the extremal solution is unique. Since this is true no matter how large H is, and any Radon measure can be approached to arbitrary precision using this discretization, we conclude that the minimizer of (B.1) should be unique as well, with at most n atoms.

Appendix C. Proof of Proposition 1

In this section, we provide the formal statement and proof of Proposition 1. Let us recall the general form of the predictor for both lazy and feature regimes in $d=2$. From equation (3.6),

$$f^n(x) = \sum_{j=1}^n g_j \tilde{\varphi}(x - x_j) = \int \frac{dy}{2\pi} g^n(y) \tilde{\varphi}(x - y). \quad (\text{C.1})$$

where n is the number of training points for the lazy regime and the number of atoms for the feature regime and, for $x \in (-\pi, \pi]$,

$$\tilde{\varphi}(x) = \begin{cases} \max\{0, \cos(x)\} & \text{(feature regime),} \\ \frac{2(\pi - |x|)\cos(x) + \sin(|x|)}{2\pi} & \text{(lazy regime, NTK),} \\ \frac{(\pi - |x|)\cos(x) + \sin(|x|)}{2\pi} & \text{(lazy regime, RFK).} \end{cases} \quad (\text{C.2})$$

All these functions $\tilde{\varphi}$ have jump discontinuities on some derivative. The first for feature and NTK, the third for RFK. If the l th derivative has jump discontinuities, the $l+1$ th only exists in a distributional sense and it can be generically written as a sum of a regular function and a sequence of Dirac masses located at the discontinuities. With m denoting the number of these discontinuities and $\{x_j\}_j$ their locations, $f^{(l)}$ denoting the l th derivative of f , for some $c_j \in \mathbb{R}$,

$$f^{(l+1)}(x) = f_r^{(l+1)}(x) + \sum_{j=1}^m c_j \delta(x - x_j), \tag{C.3}$$

where f_r denotes the *regular* part of f .

Proposition 2. Consider a random target function f^* satisfying equation (2.1) and the predictor f^n obtained by training a one-hidden-layer ReLU network on n samples $(x_i, f^*(x_i))$ in the feature or in the lazy regime (equation (C.1)). Then, with $\widehat{f}(k)$ denoting the Fourier transform of $f(x)$, one has,

$$\lim_{|k| \rightarrow \infty} \lim_{n \rightarrow \infty} \frac{\widehat{(f^n)_r''}(k)}{\widehat{f^*}(k)} = c, \tag{C.4}$$

where c is a constant (different for every regime). This result implies that as $n \rightarrow \infty$, $(f^n)_r''(x)$ converges to a function having a finite second moment, i.e.

$$\begin{aligned} \lim_{n \rightarrow \infty} \mathbb{E}_{f^*} [(f^n)_r''(x)]^2 &= \lim_{n \rightarrow \infty} \mathbb{E}_{f^*} \left[\int dx ((f^n)_r''(x))^2 \right] \\ &= \lim_{n \rightarrow \infty} \mathbb{E}_{f^*} \left[\sum_k \widehat{(f^n)_r''}(k)^2 \right] = \text{const.} < \infty, \end{aligned} \tag{C.5}$$

using the fact that $\mathbb{E}_{f^*} [(f^n)_r''(x)]^2$ does not depend on x and $\mathbb{E}_{f^*} [\sum_k \widehat{(f^*)}^2(k)] = \text{const.}$

Proof. Because our target functions are random fields that are in L_2 with probability one, and the Reproducing Kernel Hilbert Space of our kernels are dense in that space, we know that the test error vanishes as $n \rightarrow \infty$ [59]. As a result,

$$f^*(x) = \lim_{n \rightarrow \infty} f^n(x) = \lim_{n \rightarrow \infty} \int \frac{dy}{2\pi} g^n(y) \tilde{\varphi}(x - y). \tag{C.6}$$

Consider first the feature regime and the NTK lazy regime. In both cases $\tilde{\varphi}$ has two jump discontinuities in the first derivative, located at $x=0, \pi$ for the NTK and at $x = \pm \pi/2$. Therefore, we can write the second derivative as the sum of a regular function and two Dirac masses,

$$\begin{aligned} (\tilde{\varphi}^{\text{FEATURE}})'' &= -\max\{0, \cos(x)\} + \delta(x - \pi/2) + \delta(x + \pi/2), \\ (\tilde{\varphi}^{\text{NTK}})'' &= \frac{-2(\pi - |x|)\cos(x) + 3\sin(|x|)}{2\pi} - \frac{1}{2\pi}\delta(x) + \frac{1}{2\pi}\delta(x - \pi). \end{aligned} \tag{C.7}$$

As a result, the second derivative of the predictor can be written as the sum of a regular part $(f^n)_r''$ and a sequence of $2n$ Dirac masses. After subtracting the Dirac masses, both sides of equation (C.1) can be differentiated twice and yield,

$$(f^n)_r''(x) = \int \frac{dy}{2\pi} g^n(y) \tilde{\varphi}_r''(x-y). \tag{C.8}$$

Hence, in the Fourier representation we have,

$$\widehat{(f^n)_r''}(k) = \widehat{g}^n(k) \left(-k^2 \widehat{\tilde{\varphi}}_r(k) \right), \tag{C.9}$$

where we defined,

$$\widehat{\tilde{\varphi}}(k) = \int_{-\pi}^{\pi} \frac{dx}{\sqrt{2\pi}} e^{ikx} \tilde{\varphi}(x), \quad \widehat{\tilde{\varphi}}_r(k) = \int_{-\pi}^{\pi} \frac{dx}{\sqrt{2\pi}} e^{ikx} \tilde{\varphi}_r(x), \tag{C.10}$$

and used $\widehat{\tilde{\varphi}_r''}(k) = -k^2 \widehat{\tilde{\varphi}}_r(k)$. By universal approximation we have,

$$\widehat{f}^*(k) = \int_{-\pi}^{\pi} \frac{dx}{\sqrt{2\pi}} e^{ikx} f^*(x) = \lim_{n \rightarrow \infty} \widehat{g}^n(k) \widehat{\tilde{\varphi}}(k) \quad \Rightarrow \quad \lim_{n \rightarrow \infty} \widehat{g}^n(k) = \frac{\widehat{f}^*(k)}{\widehat{\tilde{\varphi}}(k)}. \tag{C.11}$$

As a result, by combining equations (C.9) and (C.11) we deduce,

$$\lim_{n \rightarrow \infty} \widehat{(f^n)_r''}(k) = -\frac{k^2 \widehat{\tilde{\varphi}}_r(k)}{\widehat{\tilde{\varphi}}(k)} \widehat{f}^*(k). \tag{C.12}$$

To complete the proof using this result it remains to estimate the scaling of $\widehat{\tilde{\varphi}}_r(k)$ and $\widehat{\tilde{\varphi}}(k)$ in the large $|k|$ limit.

For the feature regime, a direct calculation shows that $\tilde{\varphi}_r'' = -\tilde{\varphi}$, implying that $\widehat{\tilde{\varphi}}_r(k) = -\widehat{\tilde{\varphi}}(k)$. This proves that equation (C.4) is satisfied with $c = -1$.

For the NTK lazy regime $\tilde{\varphi}_r''$ and $-\tilde{\varphi}$ are different but they have similar singular expansions near $x=0$ and π . Therefore, their Fourier coefficients display the same asymptotic decay. More specifically, with $t = \cos(x)$ (or $x = \arccos(t)$), so that $\tilde{\varphi}(x) = \varphi(t)$, one has,

$$\begin{cases} \varphi^{\text{NTK}}(t) = t - \frac{1}{\sqrt{2\pi}} (1-t)^{1/2} + O\left((1-t)^{3/2}\right) \text{ near } t = +1; \\ \varphi^{\text{NTK}}(t) = -\frac{1}{\sqrt{2\pi}} (-1+t)^{1/2} + O\left((-1+t)^{3/2}\right) \text{ near } t = -1, \end{cases} \tag{C.13}$$

and

$$\begin{cases} (\varphi^{\text{NTK}})_r''(t) = -t + \frac{5}{\sqrt{2\pi}} (1-t)^{1/2} + O\left((1-t)^{3/2}\right) \text{ near } t = +1; \\ (\varphi^{\text{NTK}})_r''(t) = +\frac{5}{\sqrt{2\pi}} (-1+t)^{1/2} + O\left((-1+t)^{3/2}\right) \text{ near } t = -1. \end{cases} \tag{C.14}$$

Therefore, due to equations (A.17) and (C.4) is satisfied with $c = -5$. The same procedure can be applied to the RfK lazy regime, with the exception that it is the fourth derivative of $\tilde{\varphi}^{\text{RfK}}$, which can be written as a regular part plus Dirac masses, but one

can still obtain the Fourier coefficients of the second derivative's regular part by dividing those of the fourth derivative's regular part by k^2 . \square

Appendix D. Asymptotics of generalization in $d=2$

In this section, we compute the decay of generalization error $\bar{\epsilon}$ with the number of samples n in the following 2D setting:

$$f^n(x) = \sum_{j=1}^n g_j \tilde{\varphi}(x - x_j), \quad (\text{D.1})$$

where the x_j 's are the training points (as in the NTK case) and φ has a single discontinuity on the first derivative in 0.

Let us order the training points clockwise on the ring so that $x_1=0$ and $x_{i+1} > x_i$ for all $i=1, \dots, n$, with $x_{n+1} := 2\pi$. On each of the x_i the predictor coincides with the target,

$$f^n(x_i) = f^*(x_i) \quad \forall i = 1, \dots, n. \quad (\text{D.2})$$

For large enough n , the difference $x_{i+1} - x_i$ is small enough so that, within (x_i, x_{i+1}) , $f^n(x)$ can be replaced with its Taylor series expansion up to the second order. In practice, the predictors appear like the cable of a suspension bridge with the pillars located on the training points. In particular, we can consider an expansion around $x_i^+ := x_i + \epsilon$ for any $\epsilon > 0$ and then let $\epsilon \rightarrow 0$ from above:

$$f^n(x) = f^n(x_i^+) + (x - x_i^+) f^{n'}(x_i^+) + \frac{(x - x_i^+)^2}{2} (f^n)''(x_i^+) + \mathcal{O}\left((x - x_i^+)^3\right). \quad (\text{D.3})$$

By differentiability of f^n in (x_i, x_{i+1}) , the second derivative can be computed at any point inside (x_i, x_{i+1}) without changing the order of approximation in equation (D.3). In particular, we can replace $(f^n)''(x_i^+)$ with c_i , the mean curvature of f^n in (x_i, x_{i+1}) . Moreover, since $\epsilon \rightarrow 0$, $f^n(x_i^+) \rightarrow f^*(x_i)$ and $f^n(x_{i+1}^-) \rightarrow f^*(x_{i+1})$. By introducing the limiting slope $m_i^+ := \lim_{x \rightarrow 0^+} f^{n'}(x_i + x)$, we can write,

$$f^n(x) = f^*(x_i) + (x - x_i) m_i^+ + \frac{(x - x_i)^2}{2} c_i + \mathcal{O}\left((x - x_i^+)^3\right). \quad (\text{D.4})$$

Computing equation (D.4) at $x=x_{i+1}$ yields a closed form for the limiting slope m_i^+ as a function of the mean curvature c_i , the interval length $\delta_i := (x_{i+1} - x_i)$ and $\Delta f_i := f^*(x_{i+1}) - f^*(x_i)$. Specifically,

$$m_i^+ = \frac{\Delta f_i}{\delta_i} - \frac{\delta_i}{2} c_i. \quad (\text{D.5})$$

The generalization error can then be split into contributions from all the intervals. If $\nu_t > 2$, a Taylor expansion leads to the following:

$$\begin{aligned}
 \epsilon(n) &= \int_0^{2\pi} \frac{dx}{2\pi} (f^n(x) - f^*(x))^2 \\
 &= \sum_{i=1}^n \int_{x_i}^{x_{i+1}} \frac{dx}{2\pi} \left[(x - x_i) (m_i^+ - (f^*)'(x_i)) \right. \\
 &\quad \left. + \frac{(x - x_i)^2}{2} (c_i - (f^*)''(x_i)) + o\left((x - x_i^+)^2\right) \right]^2 \\
 &= \sum_{i=1}^n \int_0^{\delta_i} \frac{d\delta}{2\pi} \left[\delta (m_i^+ - (f^*)'(x_i)) + \frac{\delta^2}{2} (c_i - (f^*)''(x_i)) + o(\delta^2) \right]^2 \\
 &= \sum_{i=1}^n \frac{1}{2\pi} \left[\frac{\delta_i^3}{3} (m_i^+ - (f^*)'(x_i))^2 \right. \\
 &\quad \left. + \frac{\delta_i^5}{20} (c_i - (f^*)''(x_i))^2 + \frac{\delta_i^4}{4} (m_i^+ - (f^*)'(x_i)) (c_i - (f^*)''(x_i)) + o(\delta_i^5) \right]. \quad (D.6)
 \end{aligned}$$

In addition, since $\Delta f_i = (f^*)'(x_i)\delta_i + (f^*)''(x_i)\delta_i^2/2 + O(\delta_i^3)$,

$$m_i^+ - (f^*)'(x_i) = \frac{\delta_i}{2} ((f^*)''(x_i) - c_i) + o(\delta_i)^2, \quad (D.7)$$

thus,

$$\epsilon(n) = \frac{1}{2\pi} \sum_{i=1}^n \left[\frac{\delta_i^5}{120} (c_i - (f^*)''(x_i))^2 + o(\delta_i^5) \right]. \quad (D.8)$$

This implies that,

$$\bar{\epsilon}(n) = \frac{n^{-4} \left(n^{-1} \sum_{i=1}^n (n\delta_i)^5 \right)}{240\pi} \lim_{n \rightarrow \infty} \int \mathbb{E}_{f^*} \left[((f^n)''(x) - (f^*)''(x))^2 \right] dx + o(n^{-4}) \sim \frac{1}{n^4}, \quad (D.9)$$

where we used that (i) the integral converges to some finite value, due to Proposition 2. From appendix C, this integral can be estimated as $\sum_k \mathbb{E}_{f^*} [(cf^*(k) - k^2 f^*(k))^2]$, which indeed converges for $\nu_t > 2$. (ii) $(n^{-1} \sum_{i=1}^n (n\delta_i)^5)$ has a deterministic limit for large n . It is clear for the lazy regime since the distance between adjacent singularities δ_i follows an exponential distribution of mean $\sim \frac{1}{n}$. We expect this result to also be true for the feature regime in our set-up. Indeed, in the limit $n \rightarrow \infty$, the predictor approaches a parabola between singular points, which generically cannot fit more than three random points. There must thus be a singularity at least every two data points with a probability approaching unity as $n \rightarrow \infty$, which implies that $(n^{-1} \sum_{i=1}^n (n\delta_i)^5)$ converges to a constant for large n .

Finally, for $\nu_t < 2$, the same decomposition in intervals applies, but a Taylor expansion to second order does not hold. The error is then dominated by the fluctuations of f^* on the scale of the intervals, as indicated in the main text.

Appendix E. Asymptotic of generalization via the spectral bias ansatz

According to the spectral bias ansatz, the first n modes of the predictor $f_{k,\ell}^n$ coincide with the modes of the target function $f_{k,\ell}^*$. Therefore, the asymptotic scaling of the error with n is entirely controlled by the remaining modes,

$$\epsilon(n) \sim \sum_{k \geq k_c} \sum_{\ell=1}^{\mathcal{N}_{k,d}} (f_{k,\ell}^n - f_{k,\ell}^*)^2 \text{ with } \sum_{k \leq k_c} \mathcal{N}_{k,d} \sim n. \quad (\text{E.1})$$

Since $\mathcal{N}_{k,d} \sim k^{d-2}$ for $k \gg 1$, one has that, for large n , $k_c \sim n^{\frac{1}{d-1}}$. After averaging the error over target functions we obtain,

$$\bar{\epsilon}(n) \sim \sum_{k \geq k_c} \sum_{\ell=1}^{\mathcal{N}_{k,d}} \left\{ \mathbb{E}_{f^*} \left[(f_{k,\ell}^n)^2 \right] + \mathbb{E}_{f^*} \left[(f_{k,\ell}^*)^2 \right] - 2 \mathbb{E}_{f^*} \left[(f_{k,\ell}^n f_{k,\ell}^*) \right] \right\}. \quad (\text{E.2})$$

Let us recall that, with the predictor having the general form in equation (3.2), then,

$$f_{k,\ell}^n = g_{k,\ell}^n \varphi_k \quad \text{with} \quad g_{k,\ell}^n = \sum_{j=1}^n g_j Y_{k,\ell}(\mathbf{y}_j), \quad (\text{E.3})$$

where the \mathbf{y}_j 's denote the training points for the lazy regime and the neuron features for the feature regime. For $k \ll k_c$, where $f_{k,\ell}^n = f_{k,\ell}^*$, $g_{k,\ell}^n = f_{k,\ell}^*/\varphi_k$. For $k \gg k_c$, due to the highly oscillating nature of $Y_{k,\ell}$, the factors $Y_{k,\ell}(\mathbf{y}_j)$ are essentially decorrelated random numbers with zero mean and finite variance since the values of $(Y_{k,\ell}(\mathbf{y}_j))^2$ are limited by the addition theorem equation (A.5). Let us denote the variance with σ_Y . By the central limit theorem, $g_{k,\ell}^n$ converges to a Gaussian random variable with zero mean and finite variance $\sigma_Y^2 \sum_{j=1}^n g_j^2$. As a result,

$$\begin{aligned} \bar{\epsilon}(n) &\sim \sum_{k \geq k_c} \sum_{\ell=1}^{\mathcal{N}_{k,d}} \left\{ \left(\sum_{j=1}^n g_j^2 \right) \varphi_k^2 + \mathbb{E}_{f^*} \left[(f_{k,\ell}^*)^2 \right] \right\} \\ &= \left(\sum_{j=1}^n g_j^2 \right) \sum_{k \geq k_c} \mathcal{N}_{k,d} \varphi_k^2 + \sum_{k \geq k_c} \mathcal{N}_{k,d} c_k, \end{aligned} \quad (\text{E.4})$$

where we have used the definition of f^* (equation (2.1)) to set the expectation of $(f_{k,\ell}^*)^2$ to c_k .

Large ν_t case. When f^* is smooth enough, the error is controlled by the predictor term proportional to $\sum_{j=1}^n g_j^2$. More specifically, if,

$$\sum_{k \geq 0} \sum_{\ell=1}^{\mathcal{N}_{k,d}} \frac{c_k}{\varphi_k^2} < +\infty, \quad (\text{E.5})$$

then the function $g^n(\mathbf{x})$ converges to the square-summable function $g^*(\mathbf{x})$ so that $f^*(\mathbf{x}) = \int g^*(\mathbf{y})\varphi(\mathbf{x} \cdot \mathbf{y}) d\tau(\mathbf{y})$. With $c_k \sim k^{-2\nu_t - (d-1)}$ and $\mathcal{N}_{k,d} \sim k^{d-2}$, in the lazy regime $\varphi_k \sim k^{-(d-1)-2\nu}$ equation (E.5) is satisfied when $2\nu_t > 2(d-1) + 4\nu$ ($\nu = 1/2$ for the NTK and $3/2$ for the RFK). In the feature regime $\varphi_k \sim k^{-(d-1)/2-3/2}$, equation (E.5) is satisfied when $2\nu_t > (d-1) + 3$. If $g^n(\mathbf{x})$ converges to a square-summable function, then,

$$\sum_{j=1}^n g_j^2 = \frac{1}{n} \int g^n(\mathbf{x})^2 d\tau(\mathbf{x}) + o(n^{-1}) = \frac{1}{n} \sum_{k \geq 0} \mathcal{N}_{k,d} \frac{c_k}{\varphi_k^2} + o(n^{-1}), \quad (\text{E.6})$$

which is proportional to n^{-1} . In addition, since $\mathcal{N}_{k,d} \sim k^{d-2}$ and $k_c \sim n^{\frac{1}{d-1}}$, one has,

$$n^{-1} \sum_{k \geq k_c} \mathcal{N}_{k,d} \varphi_k \sim \begin{cases} n^{-1} k^{d-1} k^{-2(d-1)-4\nu} \Big|_{k=n^{\frac{1}{d-1}}} \sim n^{-2-\frac{4\nu}{d-1}} \text{ (Lazy)}, \\ n^{-1} k^{d-1} k^{-(d-1)-3} \Big|_{k=n^{\frac{1}{d-1}}} \sim n^{-1-\frac{3}{d-1}} \text{ (Feature)}, \end{cases} \quad (\text{E.7})$$

and

$$\sum_{k \geq k_c} \mathcal{N}_{k,d} c_k \sim k^{d-1} k^{-2\nu_t - (d-1)} \Big|_{k=n^{\frac{1}{d-1}}} \sim n^{-\frac{2\nu_t}{d-1}}. \quad (\text{E.8})$$

Hence, if ν_t is large enough so that equation (E.5) is satisfied, the asymptotic decay of the error is given by equation (E.7).

Small ν_t case. If equation (E.7) does not hold, then $g^n(\mathbf{x})$ is not square-summable in the limit $n \rightarrow \infty$. However, for large but finite n only the modes up to the k_c th are correctly reconstructed. Therefore,

$$\sum_{j=1}^n g_j^2 \sim \frac{1}{n} \sum_{k \leq k_c} \mathcal{N}_{k,d} \frac{c_k}{\varphi_k^2} \sim \begin{cases} n^{-1} k^{-2\nu_t} k^{2(d-1)+4\nu} \Big|_{k=n^{\frac{1}{d-1}}} \sim n^{-\frac{2\nu_t}{d-1}} n^{1+\frac{4\nu}{d-1}} \text{ (Lazy)}, \\ n^{-1} k^{-2\nu_t} k^{(d-1)+3} \Big|_{k=n^{\frac{1}{d-1}}} \sim n^{-\frac{2\nu_t}{d-1}} n^{\frac{3}{d-1}} \text{ (Feature)}. \end{cases} \quad (\text{E.9})$$

For both feature and lazy, multiplying the term above by $\sum_{k \geq k_c} \mathcal{N}_{k,d} \varphi_k$ from equation (E.7) yields $\sim n^{-2\nu_t/(d-1)}$. This is also the scaling of the target function term equation (E.8), implying that for small ν_t one has,

$$\bar{\epsilon}(n) \sim n^{-\frac{2\nu_t}{d-1}}, \quad (\text{E.10})$$

in both the feature and lazy regimes.

Appendix F. Spectral bias via the replica calculation

Due to the equivalence with kernel methods, the asymptotic decay of the test error in the lazy regime can be computed with the formalism of [45], which also provides a non-rigorous justification for the spectral bias ansatz. By ranking the eigenvalues from

the largest to the smallest, so that φ_ρ denotes the ρ th eigenvalue and denoting with c_ρ the variance of the projections of the target onto the ρ th eigenfunction, one has,

$$\epsilon(n) = \sum_{\rho} \epsilon_{\rho}(n), \quad \epsilon_{\rho}(n) = \frac{\kappa(n)^2}{(\varphi_{\rho} + \kappa(n))^2} c_{\rho}, \quad \kappa(n) = \frac{1}{n} \sum_{\rho} \frac{\varphi_{\rho} \kappa(n)}{\varphi_{\rho} + \kappa(n)}. \quad (\text{F.1})$$

It is convenient to introduce the eigenvalue density,

$$\mathcal{D}(\varphi) := \sum_{k \geq 0} \sum_{l=1}^{\mathcal{N}_{k,d}} \delta(\varphi - \varphi_k) = \sum_{k \geq 0} \mathcal{N}_{k,d} \delta(\varphi - \varphi_k) \sim \int_0^{\infty} k^{d-2} \delta\left(\varphi - k^{-(d-1)-2\nu}\right) \text{ for } k \gg 1. \quad (\text{F.2})$$

After changing variables in the delta function, one finds,

$$\mathcal{D}(\varphi) \sim \varphi^{-\frac{2(d-1)+2\nu}{(d-1)+2\nu}} \text{ for } \varphi \ll 1. \quad (\text{F.3})$$

This can be used for inferring the asymptotics of $\kappa(n)$,

$$\begin{aligned} \kappa(n) &= \frac{1}{n} \sum_{\rho} \frac{\varphi_{\rho} \kappa(n)}{\varphi_{\rho} + \kappa(n)} \sim \frac{1}{n} \int d\varphi \mathcal{D}(\varphi) \frac{\varphi \kappa(n)}{\varphi + \kappa(n)} \\ &\sim \frac{1}{n} \int_0^{\kappa(n)} d\varphi \mathcal{D}(\varphi) \varphi + \frac{\kappa(n)}{n} \int_{\kappa(n)}^{\varphi_0} d\varphi \mathcal{D}(\varphi) \\ &\sim \frac{1}{n} \kappa(n)^{1-\frac{(d-1)}{(d-1)+2\nu}} \Rightarrow \kappa(n) \sim n^{-1-\frac{2\nu}{d-1}}. \end{aligned} \quad (\text{F.4})$$

Once the scaling of $\kappa(n)$ has been determined, the modal contributions to the error can be split according to whether $\varphi_{\rho} \ll \kappa(n)$ or $\varphi_{\rho} \gg \kappa(n)$. The scaling of φ_{ρ} with the rank ρ is determined self-consistently,

$$\rho \sim \int_{\varphi_{\rho}}^{\varphi_1} d\varphi \mathcal{D}(\varphi) \sim \varphi_{\rho}^{-\frac{d-1}{(d-1)+2\nu}} \Rightarrow \varphi_{\rho} \sim \rho^{-1-\frac{2\nu}{d-1}} \Rightarrow \varphi_{\rho} \gg (\ll) \kappa(n) \Leftrightarrow \rho \ll (\gg) n. \quad (\text{F.5})$$

Therefore,

$$\epsilon(n) \sim \kappa(n)^2 \sum_{\rho \ll n} \frac{c_{\rho}}{\varphi_{\rho}^2} + \sum_{\rho \gg n} c_{\rho}. \quad (\text{F.6})$$

Note that $\kappa(n)^2$ scales as $n^{-1} \sum_{k \geq k_c} \mathcal{N}_{k,s} \varphi_k$ in equation (F.6), whereas $\sum_{\rho \ll n} c_{\rho} / \varphi_{\rho}^2$ corresponds to $n \sum_j g_j^2$ in equation (E.9) so that the first term on the right-hand side of equation (F.6) matches that of equation (E.4). The same matching is found for the second term on the right-hand side of equation (F.6) so that the replica calculation justifies the spectral bias ansatz.

Appendix G. Training wide neural networks: does GD find the minimal-norm solution?

In the main text, we provided predictions for the asymptotics of the test error of the minimal norm solution that fits all the training data. Does the prediction hold when the solution of equations (2.5) and (2.13) is approximately found by GD? More specifically, is the solution found by GD the minimal-norm one?

Feature learning. We answer these questions by performing full-batch GD in two settings (further details about the trainings are provided in the code repository, `experiments.md` file),

1. **Min-L1.** Here, we update weights and features of equation (2.3), with $\xi = 0$, by following the negative gradient of,

$$\mathcal{L}_{\text{Min-L1}} = \frac{1}{2n} \sum_{i=1}^n (f^*(\mathbf{x}_i) - f(\mathbf{x}_i))^2 + \frac{\lambda}{H} \sum_{h=1}^H |w_h|, \quad (\text{G.1})$$

with $\lambda \rightarrow 0^+$. The weights w_h are initialized to zero and the features are initialized uniformly and constrained to be on the unit sphere.

2. **α -trick.** Following [8], here we minimize,

$$\mathcal{L}_{\alpha\text{-trick}} = \frac{1}{2n\alpha} \sum_{i=1}^n (f^*(\mathbf{x}_i) - \alpha f(\mathbf{x}_i))^2, \quad (\text{G.2})$$

with $\alpha \rightarrow 0$. This trick allows us to be far from the lazy regime by forcing the weights to evolve to $\mathcal{O}(1/\alpha)$, when fitting a target of order 1.

In both cases, the solution found by GD is sparse, in the sense that it is supported on a finite number of neurons—in other words, the measure $\gamma(\boldsymbol{\theta})$ becomes atomic, satisfying Assumption 1. Furthermore, we find that

1. for **Min-L1**, the generalization error prediction holds (figures 4 and G1) as the minimal norm solution if effectively recovered, see figure G2. Such clean results in terms of features position are difficult to achieve for large n because the training dynamics becomes very slow and reaching convergence becomes computationally infeasible. Furthermore, we observe the test error to plateau and reach its infinite-time limit much earlier than the parameters, which allows for the scaling predictions to hold.
2. **α -trick**, however, does not recover the minimal-norm solution, figure G2. Moreover, the solution found is of the type (2.7) since it is sparse and supported on a number of atoms that scales linearly with n , figure G3, left. Therefore, we find that our predictions for the generalization error also hold in this case, see figure G3, right.

Lazy learning. In this case, the correspondence between the solution found by GD and the minimal-norm one is well established [9]. Therefore, numerical experiments are performed here via kernel regression and the analytical NTK equation (A.19). Given a

Learning sparse features can lead to overfitting in neural networks

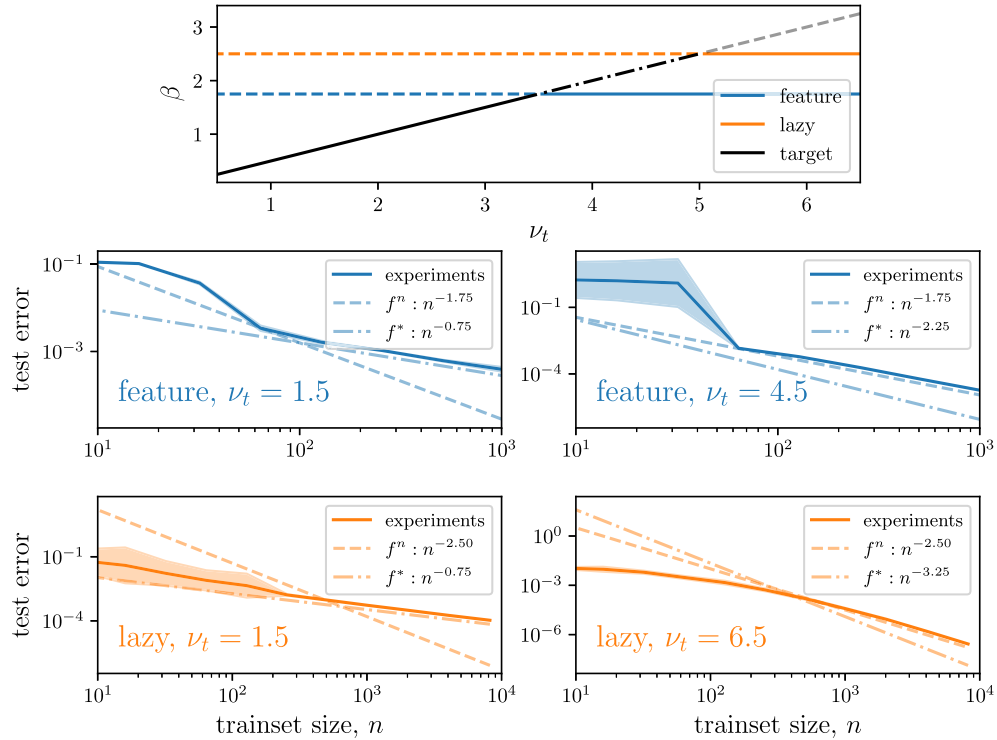


Figure G1. Generalization error decay versus target smoothness and training regime. Here, data points are sampled uniformly from the spherical surface in $d = 5$ and the target function is an infinite-width FCN with activation function $\sigma(\cdot) = |\cdot|^{\nu_t - \frac{1}{2}}$, corresponding to a Gaussian random process of smoothness ν_t . 1st row: generalization error decay exponent as a function of the target smoothness ν_t . Three curves correspond to the target contribution to the generalization error (black) and the predictor contribution in either the feature (blue) or lazy (orange) regime. Full lines highlight the dominating contributions to the generalization error. 2nd row: agreement between predictions and experiments in the feature regime for a non-smooth (left) and smooth (right) target. In the first case, the error is dominated by the target f^* , in the second by the predictor f^n —predicted exponents β are indicated in the legends. 3rd row: analogous of the previous row for the lazy regime.

data set $\{\mathbf{x}_i, y_i = f^*(\mathbf{x}_i)\}_{i=1}^n$, we define the gram matrix $\mathbf{K} \in \mathbb{R}^{n \times n}$ with elements $\mathbf{K}_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$ and the vector of target labels $\mathbf{y} = [y_1, y_2, \dots, y_n]$. The q_i 's in equation (2.9) can easily be recovered by solving the linear system:

$$\mathbf{y} = \frac{1}{n} \mathbf{K} \mathbf{q}. \tag{G.3}$$

Experiments. Numerical experiments are run with PyTorch on GPUs NVIDIA V100 (university internal cluster). Details for reproducing experiments are provided in the [code repository](#), [experiments.md](#) file. Individual trainings are run in 1 min to 1 h of wall time. We estimate a total of a thousand hours of computing time for running the preliminary and actual experiments present in this study.

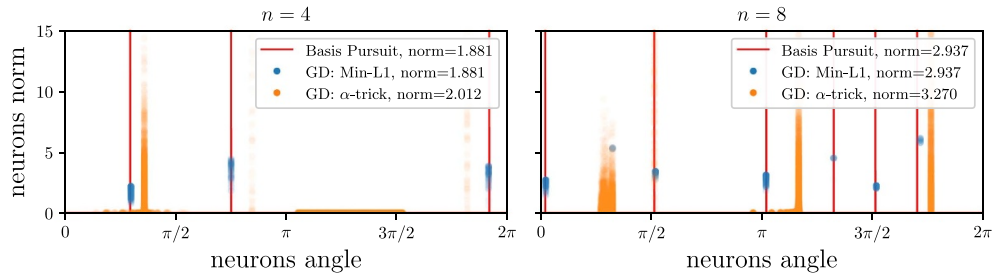


Figure G2. Comparing solutions. Solutions to the spherically symmetric task in $d = 2$ for $n = 4$ (left) and $n = 8$ (right) training points. Minimal norm solution in red (equation (2.5)), as found by basis pursuit [50]. Solutions found by GD in the Min-L1 and α -trick setting, respectively, are shown in blue and orange. Dots correspond to single neurons in the network. X -axis reports their angular position while the y -axis reports their norm, $|w_h| \|\theta_h\|_2$. Total norm of the solutions, $\frac{\alpha}{H} \sum_{h=1}^H |w_h| \|\theta_h\|_2$, is indicated in the legend.

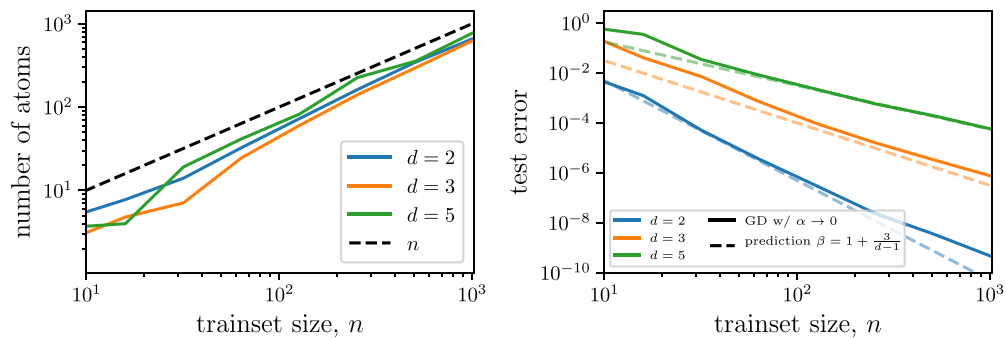


Figure G3. Solution found by the α -trick. We consider the case of approximating the constant target function on \mathbb{S}^{d-1} with an FCN. Training is performed starting from small initialization through the α -trick. Left: number of atoms n_A as a function of the number of training points n . Neurons that are active on the same subset of the training set are grouped together and we consider each group to be a distinct atom for the counting. Right: generalization error in the same setting (full), together with the theoretical predictions (dashed). Different colors correspond to different input dimensions. The case of $d = 2$ and large n suffers from the same finite time effects discussed in figure 4. Results are averaged over ten different initializations of the networks and data sets.

Appendix H. Sensitivity of the predictor to transformations other than diffeomorphisms

This section reports experiments to integrate the discussion of section 5. In particular, we (i) show that the lazy regime predictor is less sensitive to image translations than the feature regime one (as is the case for deformations, from figure 6) and (ii) provide evidence of the positive effects of learning features in image classifications, namely becoming invariant to pixels at the border of images that are unrelated to the task.

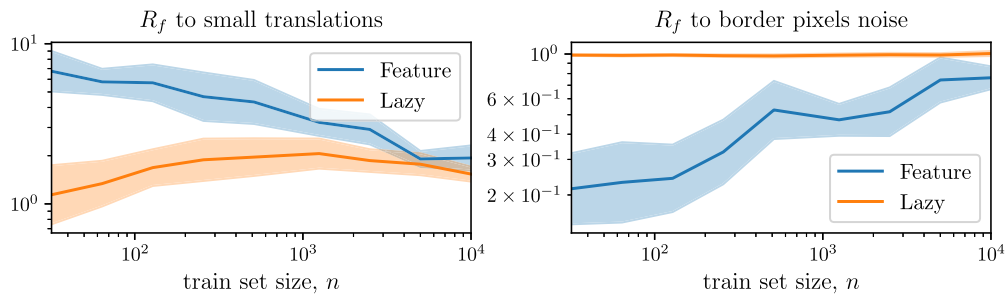


Figure H1. Sensitivity to input transformations versus number of training points. Relative sensitivity of the predictor to (left) random one-pixel translations and (right) white noise added at the boundary of the input images, in the two regimes, for varying number of training points n and when training on FashionMNIST. Smaller values correspond to a smoother predictor, on average. Results are computed using the same predictors as in figure 1. Left: for small translations, the behavior is the same compared to applying diffeomorphisms. Right: the lazy regime does not distinguish between noise added at the boundary or on the whole image ($R_f = 1$), while the feature regime becomes more insensitive to the former.

To prove the above points we consider, as in figure 6, the relative sensitivity of the predictors of lazy and feature regime with respect to global translations for point (i) and corruption of the boundary pixels for point (ii) . The relative sensitivity to translations is obtained from equation (5.1) after replacing the transformation τ with a one-pixel translation of the image in a random direction. For the relative sensitivity to boundary corruption, the transformation consists of adding zero-mean and unit-variance Gaussian numbers to the boundary pixels. Both relative sensitivities are plotted in figure H1, with translations on the left and boundary pixel corruption on the right.

In section 5, we then argue that differences in performance between the two training regimes can be explained by gaps in sensitivities with respect to input transformations that do not change the label. For (i) , the gap is similar to the one observed for diffeomorphisms (figure 6). Furthermore, the space of translations has negligible size with respect to input space, hence we expect the diffeomorphisms to have a more prominent effect. In case (ii) , the feature regime is less sensitive with respect to irrelevant pixel corruption and this would give it an advantage over the lazy regime. The fact that the performance difference is in favor of the lazy regime instead, means that these transformations only play a minor role.

Appendix I. Maximum-entropy model of diffeomorphisms

Here, we briefly review the maximum-entropy model of diffeomorphisms as introduced in [49].

An image can be thought of as a function $x(s)$ describing intensity in position $s = (u, v) \in [0, 1]^2$, where u and v are the horizontal and vertical (pixel) coordinates.

We denote τx the image deformed by τ , i.e. $[\tau x](s) = x(s - \tau(s))$ [49] and propose an ensemble of diffeomorphisms $\tau(s) = (\tau_u, \tau_v)$ with i.i.d. τ_u and τ_v defined as,

$$\tau_u = \sum_{i,j \in \mathbb{N}^+} C_{ij} \sin(i\pi u) \sin(j\pi v), \quad (\text{I.1})$$

where the C_{ij} 's are Gaussian variables of zero mean and variance $T/(i^2 + j^2)$ and T is a parameter controlling the deformation magnitude. Once τ is generated, pixels are displaced to random positions. See figure 5(b) for an example of such a transformation.

References

- [1] von Luxburg U and Bousquet O 2004 Distance-based classification with lipschitz functions *J. Mach. Learn. Res.* **5** 669–95
- [2] Bach F 2017 Breaking the curse of dimensionality with convex neural networks *J. Mach. Learn. Res.* **18** 629–81
- [3] Hestness J, Narang S, Ardalani N, Diamos G, Jun H, Kianinejad H, Patwary Md, Ali M, Yang Y and Zhou Y 2017 Deep learning scaling is predictable, empirically (arXiv:1712.00409)
- [4] Le Q V 2013 Building high-level features using large scale unsupervised learning *2013 IEEE Int. Conf. on Acoustics, Speech and Signal Processing* (IEEE) pp 8595–8
- [5] Shwartz-Ziv R and Tishby N 2017 Opening the black box of deep neural networks via information (arXiv:1703.00810)
- [6] Ansuini A, Laio A, Macke J H and Zoccolan D 2019 Intrinsic dimension of data representations in deep neural networks *Advances in Neural Information Processing Systems* pp 6111–22
- [7] Recanatani S, Farrell M, Advani M, Moore T, Lajoie G and Shea-Brown E 2019 Dimensionality compression and expansion in deep neural networks. (arXiv:1906.00443)
- [8] Chizat L, Oyallon E and Bach F 2019 On lazy training in differentiable programming *Advances in Neural Information Processing Systems* pp 2937–47
- [9] Jacot A, Gabriel F and Hongler C 2018 Neural tangent kernel: Convergence and generalization in neural networks *Proc. 32Nd Int. Conf. on Neural Information Processing Systems, (NIPS')* vol 18 (Curran Associates Inc.) pp 8580–9
- [10] Du S S, Zhai X, Póczos B and Singh A 2019 Gradient descent provably optimizes over-parameterized neural networks *Int. Conf. on Learning Representations*
- [11] Rotskoff G M and Vanden-Eijnden E 2018 Neural networks as interacting particle systems: Asymptotic convexity of the loss landscape and universal scaling of the approximation error (arXiv:1805.00915)
- [12] Mei S, Montanari A and Nguyen P-M 2018 A mean field view of the landscape of two-layer neural networks *Proc. Natl Acad. Sci.* **115** E7665–71
- [13] Sirignano J and Spiliopoulos K 2020 Mean field analysis of neural networks: a law of large numbers *SIAM J. Appl. Math.* **80** 725–52
- [14] Woodworth B, Gunasekar S, Lee J D, Moroshko E, Savarese P, Golan I, Soudry D and Srebro N 2020 Kernel and rich regimes in overparametrized models *Conf. on Learning Theory* (PMLR) pp 3635–73
- [15] de Dios J and Bruna J 2020 On sparsity in overparametrised shallow ReLU networks (arXiv:2006.10225)
- [16] Chizat L and Bach F 2020 Implicit bias of gradient descent for wide two-layer neural networks trained with the logistic loss *Conf. on Learning Theory* (PMLR) pp 1305–38
- [17] Ghorbani B, Mei S, Misiakiewicz T and Montanari A 2020 When do neural networks outperform Kernel methods? *Advances in Neural Information Processing Systems* p 33
- [18] Refinetti M, Goldt S, Krzakala F and Zdeborov L 2021 Classifying high-dimensional Gaussian mixtures: where kernel methods fail and neural networks succeed (arXiv:2102.11742)
- [19] Paccolat J, Petrini L, Geiger M, Tyloo K and Wyart M 2021 Geometric compression of invariant manifolds in neural networks *J. Stat. Mech.* **044001**
- [20] Geiger M, Spigler S, Jacot A and Wyart M 2020 Disentangling feature and lazy training in deep neural networks *J. Stat. Mech.* **113301**
- [21] Lee J, Schoenholz S S, Pennington J, Adlam B, Xiao L, Novak R and Sohl-Dickstein J 2020 Finite versus infinite neural networks: an empirical study (arXiv:2007.15801)

- [22] Ortiz-Jiménez G, Moosavi-Dezfooli S-M and Frossard P 2021 What can linearized neural networks actually say about generalization? *Advances in Neural Information Processing Systems* p 34
- [23] Chen Y, Huang W, Nguyen L M and Weng T-W 2021 On the equivalence between neural network and support vector machine (arXiv:2111.06063)
- [24] Geiger M, Jacot A, Spigler S, Gabriel F, Sagun L, d'Ascoli S, Biroli G, Hongler C and Wyart M 2020 Scaling description of generalization with number of parameters in deep learning *J. Stat. Mech.* **2020** 023401
- [25] Geiger M, Petrini L and Wyart M 2021 Landscape and training regimes in deep learning *Phys. Rep.* **924** 1–18
- [26] LeCun Y, Bottou L, Bengio Y and Haffner P 1998 Gradient-based learning applied to document recognition *Proc. IEEE* **86** 2278–324
- [27] Xiao H, Rasul K and Vollgraf R 2017 Fashion-MNIST: a novel image dataset for benchmarking machine learning algorithms (arXiv:1708.07747)
- [28] Krizhevsky A 2009 *Learning multiple layers of features from tiny images* (University of Toronto)
- [29] Bruna J and Mallat S 2013 Invariant scattering convolution networks *IEEE Trans. Pattern Anal. Mach. Intell.* **35** 1872–86
- [30] Paszke A *et al* 2019 PyTorch: an imperative style, high-performance deep learning library *Advances in Neural Information Processing Systems* vol 32 (Curran Associates, Inc.)
- [31] Maennel H, Bousquet O and Gelly S 2018 Gradient descent quantizes relu network features (arXiv:1803.08367)
- [32] Neyshabur B, Tomioka R and Srebro N 2015 Norm-based capacity control in neural networks *Conf. on Learning Theory* (PMLR) pp 1376–401
- [33] Boyer C, Chambolle A, De Castro Y, Duval V, De Gournay F'eric and Weiss P 2019 On representer theorems and convex regularization *SIAM J. Optim.* **29** 1260–81
- [34] Chizat L 2022 Sparse optimization on measures with over-parameterized gradient descent *Math. Program.* **194** 487–532
- [35] Olshausen B A and Field D J 1996 Emergence of simple-cell receptive field properties by learning a sparse code for natural images *Nature* **381** 607–9
- [36] Mairal J, Ponce J, Sapiro G, Zisserman A and Bach F 2008 Supervised dictionary learning *Advances in Neural Information Processing Systems* vol 21 (Curran Associates, Inc.)
- [37] Mehta N A and Gray A G 2013. Sparsity-based generalization bounds for predictive sparse coding *Proc. 30th Int. Conf. on Machine Learning, PMLR* vol 28 pp 36–44
- [38] Sulam J, Muthukumar R and Arora R 2021 Adversarial robustness of supervised sparse coding (arXiv:2010.12088)
- [39] Yehudai G and Shamir O 2019 On the power and limitations of random features for understanding neural networks *Advances in Neural Information Processing Systems* pp 6598–608
- [40] Ghorbani B, Mei S, Misiakiewicz T and Montanari A 2019 Limitations of lazy training of two-layers neural network *Advances in Neural Information Processing Systems* pp 9111–21
- [41] Vardan Papyan X Y H and Donoho D L 2020 Prevalence of neural collapse during the terminal phase of deep learning training *Proc. Natl Acad. Sci.* **117** 24652–63
- [42] Neyshabur B Towards learning convolutions from scratch 2020 (arXiv:2007.13657)
- [43] Ingrosso A and Goldt S 2022 Data-driven emergence of convolutional structure in neural networks *Proc. Natl Acad. Sci.* **119** e2201854119
- [44] Spigler S, Geiger M and Wyart M 2020 Asymptotic learning curves of kernel methods: empirical data versus teacher-student paradigm *J. Stat. Mech.* **124001**
- [45] Bordelon B, Canatar A and Pehlevan C 2020 Spectrum dependent learning curves in kernel regression and wide neural networks *Int. Conf. on Machine Learning* (PMLR) pp 1024–34
- [46] Cui H, Loureiro B, Krzakala F and Zdeborová L 2021 Generalization error rates in kernel regression: the crossover from the noiseless to noisy regime *Advances in Neural Information Processing Systems* p 34
- [47] Tomasini U M, Sclocchi A and Wyart M 2022 Failure and success of the spectral bias prediction for kernel ridge regression: the case of low-dimensional data (arXiv:2202.03348)
- [48] Mallat S 2016 Understanding deep convolutional networks *Phil. Trans. R. Soc. A* **374** 20150203
- [49] Petrini L, Favero A, Geiger M and Wyart M 2021 Relative stability toward diffeomorphisms indicates performance in deep nets *Advances in Neural Information Processing Systems* vol 34 (Curran Associates, Inc.) pp 8727–39
- [50] Shaobing Chen S, Donoho D L and Saunders M A 1998 Atomic decomposition by basis pursuit *SIAM J. Sci. Comput.* **20** 33–61
- [51] Scholkopf B and Smola A J 2001 *Learning With Kernels: Support Vector Machines, Regularization, Optimization and Beyond* (MIT press)

- [52] Cho Y and Lawrence K S 2009 Kernel methods for deep learning *Advances in Neural Information Processing Systems* vol 22 (Curran Associates, Inc.) pp 342–50
- [53] Bietti A and Mairal J 2019 Group invariance, stability to deformations and complexity of deep convolutional representations *J. Mach. Learn. Res.* **20** 876–924
- [54] Ruderman A, Rabinowitz N C, Morcos A S and Zoran D 2018 Pooling is neither necessary nor sufficient for appropriate deformation stability in CNNs (arXiv:[1804.04438](https://arxiv.org/abs/1804.04438))
- [55] Smola A, Ovári Z and Williamson R C 2000 Regularization with dot-product kernels *Advances in Neural Information Processing Systems* p 13
- [56] Atkinson K and Han W 2012 *Spherical Harmonics and Approximations on the Unit Sphere: an Introduction* vol 2044 (Springer Science & Business Media)
- [57] Efthimiou C and Frye C 2014 *Spherical Harmonics in p Dimensions* (Singapore: World Scientific)
- [58] Bietti A and Bach F 2021 Deep equals shallow for ReLU networks in kernel regimes (arXiv:[2009.14397](https://arxiv.org/abs/2009.14397))
- [59] Bach F 2022 Learning theory from first principles (in preparation) (MIT Press)