
Implicit Regularization via Neural Feature Alignment

Aristide Baratin^{1*}

Thomas George^{1*}

César Laurent¹

R Devon Hjelm^{2,1}

Guillaume Lajoie¹

Pascal Vincent^{1,3}

Simon Lacoste-Julien^{1,3}

¹ Mila, Université de Montréal ² Microsoft Research ³ Canada CIFAR AI chair

Abstract

We approach the problem of implicit regularization in deep learning from a geometrical viewpoint. We highlight a regularization effect induced by a dynamical alignment of the neural tangent features introduced by Jacot et al. (2018), along a small number of task-relevant directions. This can be interpreted as a combined mechanism of feature selection and compression. By extrapolating a new analysis of Rademacher complexity bounds for linear models, we motivate and study a heuristic complexity measure that captures this phenomenon, in terms of sequences of tangent kernel classes along optimization paths. The code for our experiments is available at https://github.com/tfjgeorge/ntk_alignment.

1 Introduction

One important property of deep neural networks is their ability to generalize well on real data. Surprisingly, this is even true with very high-capacity networks *without explicit regularization* (Neyshabur et al., 2015; Zhang et al., 2017; Hoffer et al., 2017). This seems at odds with the usual understanding of the bias-variance trade-off (Geman et al., 1992; Neal et al., 2018; Belkin et al., 2019): highly complex models are expected to overfit the training data and perform poorly on test data (Hastie et al., 2009). Solving this apparent paradox requires understanding the various learning biases induced by the training procedure, which can act as implicit regularizers (Neyshabur et al., 2015, 2017b).

In this paper, we help clarify one such implicit regu-

larization mechanism, by examining the evolution of the *neural tangent features* (Jacot et al., 2018) learned by the network along the optimization paths. Our results can be understood from two complementary perspectives: a *geometric* perspective – the (uncentered) covariance of the tangent features defines a metric on the function class, akin to the Fisher information metric (e.g., Amari, 2016); and a *functional* perspective – through the tangent kernel and its RKHS. In standard supervised classification settings, our main observation is a dynamical alignment of the tangent features along a small number of task-relevant directions during training. We interpret this phenomenon as a combined mechanism of *feature selection* and *compression*. The intuition motivating this work is that such a mechanism allows large models to adapt their capacity to the task, which in turn underpins their generalization abilities.

Specifically, our main contributions are as follows:

1. Through experiments with various architectures on MNIST and CIFAR10, we give empirical insights on how the tangent features and their kernel adapt to the task during training (Section 3). We observe in particular a sharp increase of the anisotropy of their spectrum early in training, as well as an increasing similarity with the class labels, as measured by *centered kernel alignment* (Cortes et al., 2012).
2. Drawing upon intuitions from linear models (Section 4.1), we argue that such a dynamical alignment acts as *implicit regularizer*. We motivate a new heuristic complexity measure which captures this phenomenon, and empirically show better correlation with generalization compared to various measures proposed in the recent literature (Section 4).

2 Preliminaries

Let \mathcal{F} be a class of functions (e.g a neural network) parametrized by $\mathbf{w} \in \mathbb{R}^P$. We restrict here to *scalar*

*Equal contribution. Proceedings of the 24th International Conference on Artificial Intelligence and Statistics (AISTATS) 2021, San Diego, California, USA. PMLR: Volume 130. Copyright 2021 by the author(s).

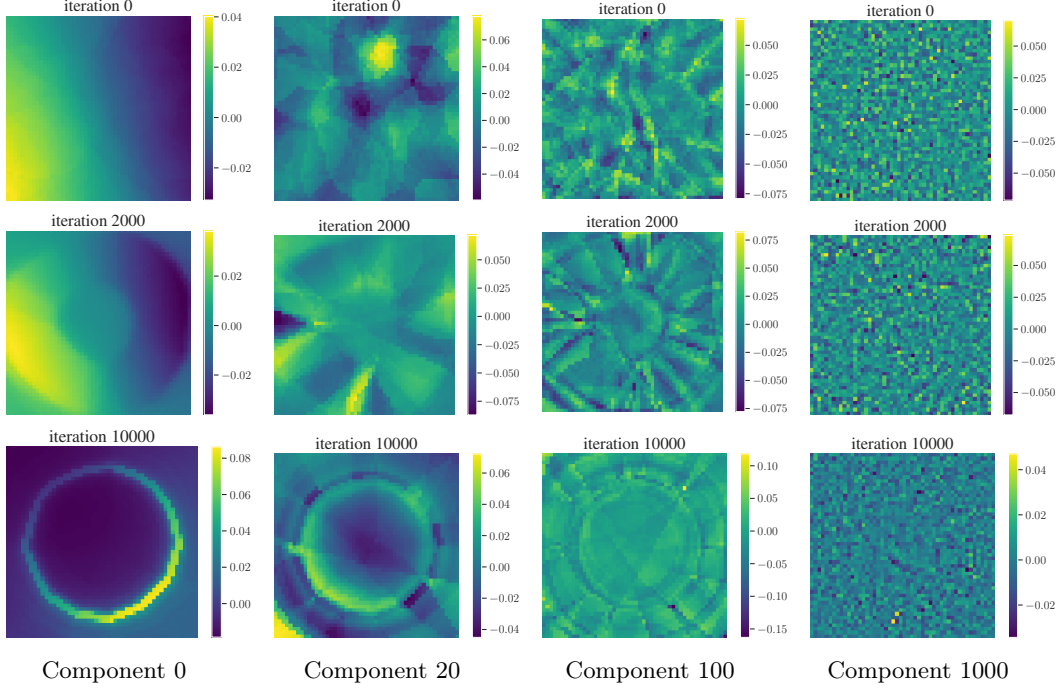


Figure 1: Evolution of eigenfunctions of the tangent kernel, ranked in nonincreasing order of the eigenvalues (**in columns**), at various iterations during training (**in rows**), for the $2d$ Disk dataset. After a number of iterations, we observe modes corresponding to the class structure (e.g. boundary circle) in the top eigenfunctions. Combined with an increasing anisotropy of the spectrum (e.g. $\lambda_{20}/\lambda_1 = 1.5\%$ at iteration 0, 0.2% at iteration 2000), this illustrates a stretch of the tangent kernel, hence a (soft) compression of the model, along a small number of features that are highly correlated with the classes.

functions $f_{\mathbf{w}} : \mathcal{X} \rightarrow \mathbb{R}$ to keep notation light.¹

Tangent Features. We define the **tangent features** as the function gradients w.r.t the parameters,

$$\Phi_{\mathbf{w}}(\mathbf{x}) := \nabla_{\mathbf{w}} f_{\mathbf{w}}(\mathbf{x}) \in \mathbb{R}^P. \quad (1)$$

The corresponding kernel $k_{\mathbf{w}}(\mathbf{x}, \tilde{\mathbf{x}}) = \langle \Phi_{\mathbf{w}}(\mathbf{x}), \Phi_{\mathbf{w}}(\tilde{\mathbf{x}}) \rangle$ is the **tangent kernel** (Jacot et al., 2018). Intuitively, the tangent features govern how small changes in parameter affect the function’s outputs,

$$\delta f_{\mathbf{w}}(\mathbf{x}) = \langle \delta \mathbf{w}, \Phi_{\mathbf{w}}(\mathbf{x}) \rangle + O(\|\delta \mathbf{w}\|^2). \quad (2)$$

More formally, the (uncentered) covariance matrix $g_{\mathbf{w}} = \mathbb{E}_{\mathbf{x} \sim \rho} [\Phi_{\mathbf{w}}(\mathbf{x}) \Phi_{\mathbf{w}}(\mathbf{x})^\top]$ w.r.t the input distribution ρ acts as a **metric tensor** on \mathcal{F} : assuming $\mathcal{F} \subset L^2(\rho)$, this is the metric induced on \mathcal{F} by pullback of the L^2 scalar product. It characterizes the geometry of the function class \mathcal{F} . Metric (as symmetric $P \times P$ matrices) and tangent kernels (as rank P integral operators) share the same spectrum (see Prop 4 in Appendix A.3).

Spectral Bias. The structure of the tangent features impacts the evolution of the function during training. To formalize this, we introduce the covariance eigenvalue decomposition $g_{\mathbf{w}} = \sum_{j=1}^P \lambda_{\mathbf{w}j} \mathbf{u}_{\mathbf{w}j} \mathbf{u}_{\mathbf{w}j}^\top$, which

¹The extension to vector-valued functions, relevant for the multiclass classification setting, is presented in Appendix A, along with more mathematical details.

summarizes the predominant directions in parameter space. Given n input samples (\mathbf{x}_i) and $\mathbf{f}_{\mathbf{w}} \in \mathbb{R}^n$ the vector of outputs $f_{\mathbf{w}}(\mathbf{x}_i)$, consider gradient descent updates $\delta \mathbf{w}_{\text{GD}} = -\eta \nabla_{\mathbf{w}} L$ for some cost function $L := L(\mathbf{f}_{\mathbf{w}})$. The following elementary result (see Appendix A.5) shows how the corresponding function updates in the linear approximation (2), $\delta f_{\text{GD}}(\mathbf{x}) := \langle \delta \mathbf{w}_{\text{GD}}, \Phi_{\mathbf{w}}(\mathbf{x}) \rangle$, decompose in the **eigenbasis**² of the tangent kernel:

$$u_{\mathbf{w}j}(\mathbf{x}) = \frac{1}{\sqrt{\lambda_{\mathbf{w}j}}} \langle \mathbf{v}_{\mathbf{w}j}, \Phi_{\mathbf{w}}(\mathbf{x}) \rangle \quad (3)$$

Lemma 1 (Local Spectral Bias). *The function updates decompose as $\delta f_{\text{GD}}(\mathbf{x}) = \sum_{j=1}^P \delta f_j u_{\mathbf{w}j}(\mathbf{x})$ with*

$$\delta f_j = -\eta \lambda_{\mathbf{w}j} (\mathbf{u}_{\mathbf{w}j}^\top \nabla_{\mathbf{f}_{\mathbf{w}}} L), \quad (4)$$

where $\mathbf{u}_{\mathbf{w}j} = [u_{\mathbf{w}j}(\mathbf{x}_1), \dots, u_{\mathbf{w}j}(\mathbf{x}_n)]^\top \in \mathbb{R}^n$ and $\nabla_{\mathbf{f}_{\mathbf{w}}}$ denotes the gradient w.r.t the sample outputs.

This illustrates how, from the point of view of function space, the metric/tangent kernel eigenvalues act as a mode-specific rescaling $\eta \lambda_{\mathbf{w}j}$ of the learning rate.³ This

²The functions $(u_{\mathbf{w}j})_{j=1}^P$ form an orthonormal family in $L^2(\rho)$, i.e. $\mathbb{E}_{\mathbf{x} \sim \rho} [u_{\mathbf{w}j} u_{\mathbf{w}j'}] = \delta_{jj'}$, and yield the spectral decomposition $k_{\mathbf{w}}(\mathbf{x}, \tilde{\mathbf{x}}) = \sum_{j=1}^P \lambda_{\mathbf{w}j} u_{\mathbf{w}j}(\mathbf{x}) u_{\mathbf{w}j}(\tilde{\mathbf{x}})$ of the tangent kernel as an integral operator (see Appendix A.3).

³Intuitively, the eigenvalue $\lambda_{\mathbf{w}j}$ can be thought of as defining a local ‘learning speed’ for the mode j .

is a local version of a well-known bias for linear models trained by gradient descent (e.g in linear regression, see Appendix A.5.2), which prioritizes learning functions within the top eigenspaces of the kernel. Several recent works (Bietti & Mairal, 2019; Basri et al., 2019; Yang & Salman, 2019) investigated such bias for neural networks, in *linearized* regimes where the tangent kernel remains constant during training (Jacot et al., 2018; Du et al., 2019; Allen-Zhu et al., 2019). As a simple example, for a randomly initialized MLP on 1D uniform data, Fig. 8 in Appendix A.5 shows an alignment of the tangent kernel eigenfunctions with Fourier modes of increasing frequency, in line with prior empirical observations (Rahaman et al., 2019; Xu et al., 2019) of a ‘spectral bias’ towards low-frequency functions.

Tangent Features Adapt to the Task. By contrast, our aim in this paper is to highlight and discuss *non-linear* effects, in the (standard) regime where the tangent features and their kernel evolve during training (e.g., Geiger et al., 2019; Woodworth et al., 2020).

As a first illustration of such effects, Fig. 1 shows visualizations of eigenfunctions of the tangent kernel (ranked in nonincreasing order of the eigenvalues), during training of a 6-layer deep 256-unit wide MLP by gradient descent of the binary cross entropy loss, on a simple classification task: $y(\mathbf{x}) = \pm 1$ depending on whether $\mathbf{x} \sim \text{Unif}[-1, 1]^2$ is in the centered disk of radius $\sqrt{2/\pi}$ (details in Appendix C.1). After a number of iterations, we observe (rotation invariant) modes corresponding to the class structure (e.g. boundary circle) showing up in the *top* eigenfunctions of the learned kernel. We also note an increasing spectrum anisotropy – for example, the ratio λ_{20}/λ_1 , which is 1.5% at iteration 0, has dropped to 0.2% at iteration 2000. The interpretation is that the tangent kernel (and the metric) *stretch* along a relatively small number of directions that are highly correlated with the classes during training. We quantify and investigate this effect in more detail below.

3 Neural Feature Alignment

In this section, we study in more detail the evolution of the tangent features during training. Our main results are to highlight (i) a sharp increase of the anisotropy of their spectrum early in training; (ii) an increasing similarity with the class labels, as measured by **centered kernel alignment** (CKA) (Cristianini et al., 2002; Cortes et al., 2012). We interpret this as a combined mechanism of feature selection and model compression.

3.1 Setup

We run experiments on MNIST (LeCun et al., 2010) and CIFAR10 (Krizhevsky & Hinton, 2009) with standard

MLPs, VGG (Simonyan & Zisserman, 2014) and Resnet (He et al., 2016) architectures, trained by stochastic gradient descent (SGD) with momentum, using cross-entropy loss. We use PyTorch (Paszke et al., 2019) and NNGeometry (George, 2021) for efficient evaluation of tangent kernels.

In multiclass settings, tangent kernels evaluated on n samples carry additional class indices $y \in \{1 \dots c\}$ and thus are $nc \times nc$ matrices, $(\mathbf{K}_{\mathbf{w}})_{ij}^{yy'} := k_{\mathbf{w}}(\mathbf{x}_i, y; \mathbf{x}_j, y')$ (details in Appendix A.4). In all our experiments, we evaluate tangent kernels on mini-batches of size $n = 100$ from both the training set and the test set; for $c = 10$ classes, this yields kernel matrices of size 1000×1000 . We report results obtained from *centered* tangent features $\Phi_{\mathbf{w}}(\mathbf{x}) \rightarrow \Phi_{\mathbf{w}}(\mathbf{x}) - \mathbb{E}_{\mathbf{x}} \Phi_{\mathbf{w}}(\mathbf{x})$, though we obtain qualitatively similar results for uncentered features (see plots in Appendix C.2).

3.2 Spectrum Evolution

We first investigate the evolution of the tangent kernel *spectrum* for a VGG19 on CIFAR 10, trained with and without label noise (Fig. 2). The take away is an anisotropic increase of the spectrum during training. We report results for kernels evaluated on training examples (solid line) and test examples (dashed line).⁴

The first observation is a significant *increase* of the spectrum, early in training (note the log scale for the x -axis). By the time the model reaches 100% training accuracy, the maximum and average eigenvalues (Fig. 2, 2nd row) have gained more than 2 orders of magnitude.

The second observation is that this evolution is highly *anisotropic*, i.e larger eigenvalues increase faster than lower ones. This results in a (sharp) increase of spectrum anisotropy, early in training. We quantify this using a notion of **effective rank** based on spectral entropy (Roy & Vetterli, 2007). Given a kernel matrix \mathbf{K} in $\mathbb{R}^{r \times r}$ with (strictly) positive eigenvalues $\lambda_1, \dots, \lambda_r$, let $\mu_j = \lambda_j / \sum_{i=1}^r \lambda_i$ be the trace-normalized eigenvalues. The effective rank is defined as $\text{erank} = \exp(H(\boldsymbol{\mu}))$ where $H(\boldsymbol{\mu})$ is the Shannon entropy,

$$\text{erank} = \exp(H(\boldsymbol{\mu})), \quad H(\boldsymbol{\mu}) = - \sum_{j=1}^r \mu_j \log(\mu_j). \quad (5)$$

This effective rank is a real number between 1 and r , upper bounded by $\text{rank}(\mathbf{K})$, which measures the ‘uniformity’ of the spectrum through the entropy. We also track the various **trace ratios**

$$T_k = \sum_{j < k} \lambda_j / \sum_j \lambda_j, \quad (6)$$

⁴The striking similarity of the plots for train and test kernels suggests that the spectrum of empirical tangent kernels is robust to sampling variations in our setting.

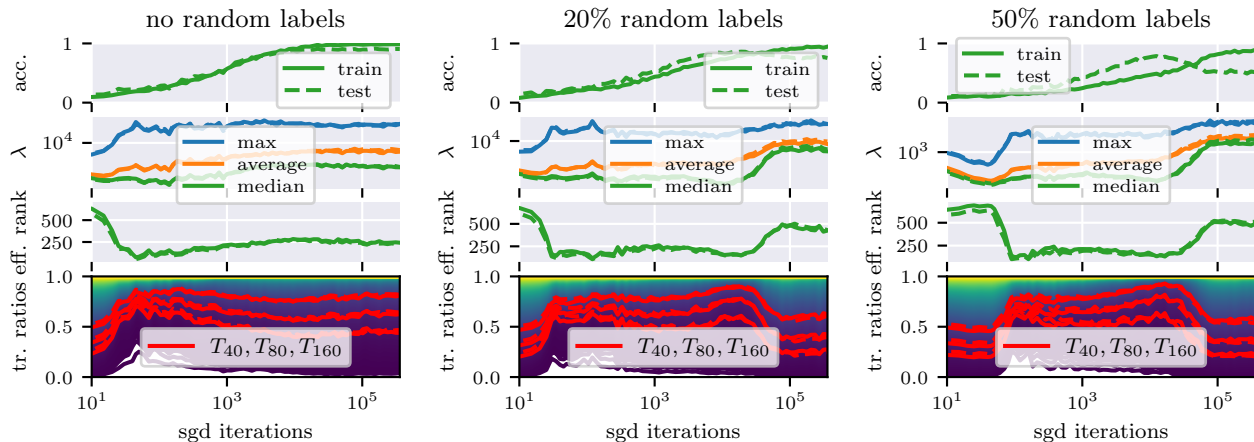


Figure 2: Evolution of the tangent kernel **spectrum** (max, average and median eigenvalues), **effective rank** (5) and **trace ratios** (6) during training of a VGG19 on CIFAR10 with various ratio of random labels, using cross-entropy and SGD with batch size 100, learning rate 0.01 and momentum 0.9. Tangent kernels are evaluated on batches of size 100 from both the training set (solid lines) and the test set (dashed lines). The plots in the top row show train/test accuracy.

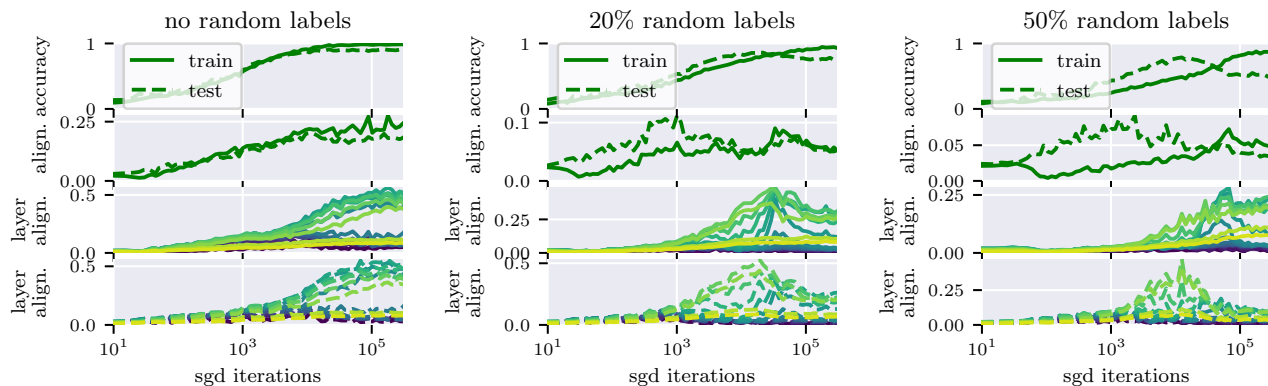


Figure 3: Evolution of the (tangent) **feature alignment with class labels** as measured by CKA (7), during training of a VGG19 on CIFAR10 (same setup as in Fig. 2). Tangent kernels and label vectors are evaluated on batches of size 100 from both the training set (solid lines) and the test set (dashed lines). The plots in the last two rows show the alignment of tangent features associated to *each layer*. Layers are mapped to colors sequentially from input layer (-), through intermediate layers (-), to output layer (-). See Fig. 11 and 13 in Appendix C for additional architectures and datasets.

which quantify the relative importance of the top k eigenvalues.

We note (Fig. 2, third row) a drop of the effective rank early in training (e.g. to less than 10% of its initial value in our experiments with no random labels; less than 20% when half of the labels are randomized). This can also be observed from the highlighted (in red) trace ratios T_{40} , T_{80} and T_{160} (Fig. 2, fourth row), e.g. the first top 40 eigenvalues (T_{40}), over 1000 in total, accounting for more than 70% of the total trace.

Remarkably, in the presence of high label noise, the effective rank of the tangent kernel (and hence that of the metric) evaluated on *training* examples (anti)-correlates nicely with the *test* accuracy: while decreasing and remaining relatively low during the learning

phase (increase of test accuracy), it begins to rise again when overfitting starts (decrease of test accuracy). This suggests that this effective rank already provides a good proxy for the effective capacity of the network.

3.3 Alignment to class labels

We now include the evolution of the eigenvectors in our study. We investigate the similarity of the learned tangent features with the class label through centered kernel alignment. Given two kernel matrices \mathbf{K} and \mathbf{K}' in $\mathbb{R}^{r \times r}$, it is defined as (Cortes et al., 2012)

$$\text{CKA}(\mathbf{K}, \mathbf{K}') = \frac{\text{Tr}[\mathbf{K}_c \mathbf{K}'_c]}{\|\mathbf{K}_c\|_F \|\mathbf{K}'_c\|_F} \in [0, 1] \quad (7)$$

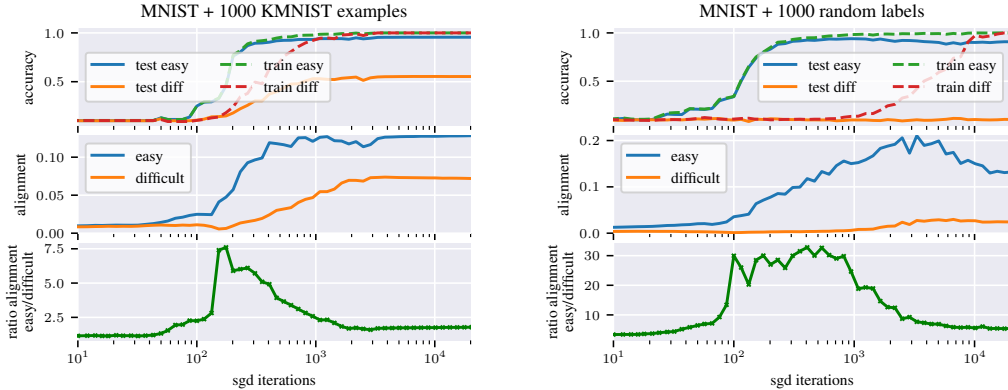


Figure 4: Alignment *easy* versus *difficult*: We augment a dataset composed of 10,000 *easy* MNIST examples with 1000 *difficult* examples from 2 different setups: **(left)** 1000 MNIST examples with random label **(right)** 1000 KMNIST examples. We train a MLP with 6 layers of 80 hidden units using SGD with learning rate=0.02, momentum=0.9 and batch size=100. We observe that the alignment to (train) labels increases faster and to a higher value for the easy examples.

where the subscript c denotes the feature centering operation, i.e. $\mathbf{K}_c = C\mathbf{K}C$ where $C = I_r - \frac{1}{r}\mathbf{1}\mathbf{1}^T$ is the centering matrix, and $\|\cdot\|_F$ is the Frobenius norm. CKA is a normalized version of the Hilbert-Schmidt Independence Criterion (Gretton et al., 2005) designed as a dependence measure for two sets of features. The normalization makes CKA invariant under isotropic rescaling.

Let $\mathbf{Y} \in \mathbb{R}^{nc}$ be the vector resulting from the concatenation of the one-hot label representations $\mathbf{Y}_i \in \mathbb{R}^c$ of the n samples. Similarity with the labels is measured through CKA with the rank-one kernel $\mathbf{K}_Y := \mathbf{Y}\mathbf{Y}^T$. Intuitively, $\text{CKA}(\mathbf{K}, \mathbf{K}_Y)$ is high when \mathbf{K} has low (effective) rank and such that the angle between \mathbf{Y} and its top eigenspaces is small.⁵ Maximizing such an index has been used as a criterion for kernel selection in the literature on learning kernels (Cortes et al., 2012).

With the same setup as in Section 3.2, we observe (Fig. 3, 2nd row) an increasingly high CKA between the tangent kernel and the labels as training progresses. The trend is similar for other architectures and datasets (e.g., Fig. 11 in Appendix C shows CKA plots for MLP on MNIST and Resnets 18 on CIFAR10).

Interestingly, in the presence of high level noise, the CKA reaches a much higher value during the learning phase (increase of test accuracy) for tangent kernels and labels evaluated for *test* than for *train* inputs (note test labels are not randomized). Together with Equ. 4, this suggests a stronger learning bias towards features predictive of the *clean* labels. This is line with empirical observations that, in the presence of noise, deep networks ‘learn patterns faster than noise’ (Arpit et al., 2017) (see Section 3.4 below for additional insights).

⁵In the limiting case $\text{CKA}(\mathbf{K}, \mathbf{K}_Y) = 1$, the features are all aligned with each other and parallel to \mathbf{Y} .

We also report the alignments of the *layer-wise* tangent kernels. By construction, the tangent kernel, obtained by pairing features $\Phi_{w_p}(\mathbf{x})\Phi_{w_p}(\tilde{\mathbf{x}})$ and summing over all parameters w_p of the network, can also be expressed as the sum of layer-wise tangent kernels, $\mathbf{K}_w = \sum_{\ell=1}^L \mathbf{K}_w^\ell$, where \mathbf{K}_w^ℓ results from summing only over parameters of the layer ℓ . We observe a high CKA, reaching more than 0.5 for a number of intermediate layers.⁶ In the presence of high label noise, we note that CKAs tend to peak when the test accuracy does.

3.4 Hierarchical Alignment

A key aspect of the generalization question concerns the articulation between learning and memorization, in the presence of noise (Zhang et al., 2017) or difficult examples (e.g., Sagawa et al., 2020). Motivated by this, we would like to probe the evolution of the tangent features *separately* in the directions of both types of examples in such settings. To do so, our strategy is to measure CKA for tangent kernels and label vectors evaluated on examples from two subsets of the same size in the training dataset – one with ‘easy’ examples, the other with ‘difficult’ ones. Our setup is to augment 10,000 MNIST training examples with 1000 difficult examples of 2 types: (i) examples with random labels and (ii) examples from the dataset KMNIST (Clanuwat et al., 2018). KMNIST images present features similar to MNIST digits (grayscale handwritten characters) but represent Japanese characters.

The results are shown in Fig. 4. As training progresses, we observe that the CKA on the easy examples increases faster (and to a higher value) than that on the

⁶We were expecting to see a gradually increasing CKA with ℓ ; we do not have any intuitive explanation for the relatively low alignment observed for the very top layers.

difficult ones; in the case of the (structured) difficult examples from KMNIST, we also note an increase of the CKA later in training. This demonstrates a hierarchy in the adaptation of the kernel, measured by the ratio between both alignments. From the intuition developed in the paper (see spectral bias in Equ.(4)), we interpret this aspect of the non-linear dynamics as favoring a sequentialization of learning across patterns of different complexity (‘easy patterns first’), a phenomenon analogous to one pointed out in the context of deep linear networks (Saxe et al., 2014; Lampinen et al., 2018; Gidel et al., 2019).

3.5 Ablation

Effect of depth. In order to study the influence of depth on alignment and test the robustness to the choice of seeds, we reproduce the experiment of the previous section for MLP with different depths, while varying parameter initialization and minibatch sampling. Our results, shown in Fig 13 (Appendix C), suggest that the alignment effect is magnified as depth increases. We also observe that the ratio of the maximum alignment between easy and difficult examples is increased with depth, but stays high for a smaller number of iterations.

Effect of the learning rate. We observed in our experiments that increasing the learning rate tend to enhance alignment effects.⁷ As an illustration, we reproduce in Fig. 14 the same plots as in Fig. 2, for a learning rate reduced to 0.003. We observe a similar drop of the effective rank as in Fig. 2 at the beginning of training, but to a much (about 3 times) higher value.

4 Measuring Complexity

In this section, drawing upon intuitions from linear models, we illustrate in a simple setting how the alignment of tangent features can act as implicit regularization. By extrapolating Rademacher complexity bounds for linear models, we also motivate a new complexity measure for neural networks and compare its correlation to generalization against various measures proposed in the literature. We refer to Appendix B for a review of classical results, further technical details, and proofs.

⁷Note that for wide enough networks and small enough learning rate, we expect to recover the linear regime where the tangent features are constant during training (Jacot et al., 2018; Du et al., 2019; Allen-Zhu et al., 2019).

4.1 Insights from Linear Models

4.1.1 Setup

We restrict here to scalar functions $f_{\mathbf{w}}(\mathbf{x}) = \langle \mathbf{w}, \Phi(\mathbf{x}) \rangle$ linearly parametrized by $\mathbf{w} \in \mathbb{R}^P$. Such a function class defines a *constant* (tangent) kernel and geometry, as defined in Section 2. Given n input samples, the n features $\Phi(\mathbf{x}_i) \in \mathbb{R}^P$ yield an $n \times P$ feature matrix Φ .

Our discussion will be based on the (empirical) **Rademacher complexity**, which shows up in generalization bounds (Bartlett & Mendelson, 2002); see Appendix B.2 for a review. It measures how well \mathcal{F} correlates with random noise on the sample set \mathcal{S} :

$$\widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}) = \mathbb{E}_{\sigma \in \{\pm 1\}^n} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(\mathbf{x}_i) \right]. \quad (8)$$

The Rademacher complexity depends on the size (or **capacity**) of the class \mathcal{F} . Constraints on the capacity, such as those induced by the implicit bias of the training algorithm, can reduce the Rademacher complexity and lead to sharper generalization bounds.

A standard approach for controlling capacity is in terms of the *norm* of the weight vector – usually the ℓ_2 -norm. In general, given any invertible matrix $A \in \mathbb{R}^{P \times P}$, we may consider the norm $\|\mathbf{w}\|_A := \sqrt{\mathbf{w}^\top g_A \mathbf{w}}$ induced by the metric $g_A = AA^\top$. Consider the (sub)classes of functions induced by balls of given radius:

$$\mathcal{F}_{M_A}^A = \{f_{\mathbf{w}} : \mathbf{x} \mapsto \langle \mathbf{w}, \Phi(\mathbf{x}) \rangle \mid \|\mathbf{w}\|_A \leq M_A\}. \quad (9)$$

A direct extension of standard bounds for the Rademacher complexity (see Appendix B.3) yields,

$$\widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_{M_A}^A) \leq (M_A/n) \|A^{-1} \Phi^\top\|_{\mathbb{F}} \quad (10)$$

where $\|A^{-1} \Phi^\top\|_{\mathbb{F}}$ is the Froebenius norm of the *rescaled* feature matrix.⁸

This freedom in the choice of rescaling matrix A raises the question of which of the norms $\|\cdot\|_A$ provide meaningful measures of the model’s capacity. Recent works (Belkin et al., 2018; Muthukumar et al., 2020) pointed out that using ℓ_2 norm is not coherently linked with generalization in practice. We discuss this issue in Appendix B.5, illustrating how meaningful norms critically depend on the geometry defined by the features.

4.1.2 Feature Alignment as Implicit Regularization

Here we describe a simple procedure making the geometry *adaptive* along optimization paths. The goal is to

⁸We also have $\|A^{-1} \Phi^\top\|_{\mathbb{F}} = \sqrt{\text{Tr} \mathbf{K}_A}$ in terms of the (rescaled) kernel matrix $\mathbf{K}_A = \Phi g_A^{-1} \Phi^\top$.

SuperNat update ($\tilde{A}_0 = I$, $\Phi_0 = \Phi$, $K_0 = K$):

1. Perform gradient step $\tilde{\mathbf{w}}_{t+1} \leftarrow \mathbf{w}_t + \delta \mathbf{w}_{\text{GD}}$
2. Find minimizer \tilde{A}_{t+1} of $\|\delta \mathbf{w}_{\text{GD}}\|_{\tilde{A}} \|\tilde{A}^{-1} \Phi_t^\top\|_{\text{F}}$
3. Reparametrize:

$$\mathbf{w}_{t+1} \leftarrow \tilde{A}_{t+1}^\top \tilde{\mathbf{w}}_{t+1}, \Phi_{t+1} \leftarrow \tilde{A}_{t+1}^{-1} \Phi_t$$

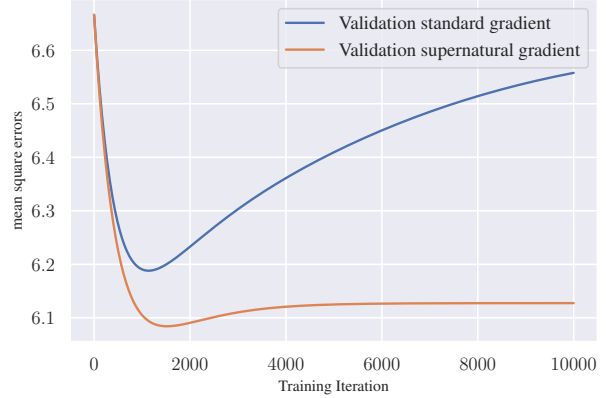


Figure 5: **(left)** **SuperNat** algorithm and **(right)** validation curves obtained with standard and **SuperNat** gradient descent, on the noisy linear regression problem. At each iteration, **SuperNat** identifies dominant features and stretches the kernel along them, thereby slowing down and eventually freezing the learning dynamics in the noise direction. This naturally yields better generalization than standard gradient descent on this problem.

illustrate in a simple setting how feature alignment can impact complexity and generalization, in a way that mimics the behaviour of a non-linear dynamics. The idea is to *learn* a rescaling metric at each iteration of our algorithm, using a local version of the bounds (10).

Complexity of Learning Flows. Since we are interested in functions $f_{\mathbf{w}}$ that result from an iterative algorithm, we consider functions $f_{\mathbf{w}} = \sum_t \delta f_{\mathbf{w}_t}$ written in terms of a sequence of updates⁹ $\delta f_{\mathbf{w}_t}(\mathbf{x}) = \langle \delta \mathbf{w}_t, \Phi(\mathbf{x}) \rangle$ (we set $f_0 = 0$ to keep the notation simple), with *local* constraints on the parameter updates:

$$\mathcal{F}_{\mathbf{m}}^{\mathbf{A}} = \{f_{\mathbf{w}} : \mathbf{x} \mapsto \sum_t \langle \delta \mathbf{w}_t, \Phi(\mathbf{x}) \rangle \mid \|\delta \mathbf{w}_t\|_{A_t} \leq m_t\} \quad (11)$$

The result (10) extends as follows.

Theorem 2 (Complexity of Learning Flows). *Given any sequences \mathbf{A} and \mathbf{m} of invertible matrices $A_t \in \mathbb{R}^{P \times P}$ and positive numbers $m_t > 0$, we have the bound*

$$\widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_{\mathbf{m}}^{\mathbf{A}}) \leq \sum_t (m_t/n) \|A_t^{-1} \Phi^\top\|_{\text{F}}. \quad (12)$$

Note that, by linear reparametrization invariance $\mathbf{w} \mapsto A^\top \mathbf{w}$, $\Phi \mapsto A^{-1} \Phi$, the *same* result can be formulated in terms of the sequence $\Phi = \{\Phi_t\}_t$ of feature maps $\Phi_t = A_t^{-1} \Phi$. The function class (11) can equivalently be written as

$$\mathcal{F}_{\mathbf{m}}^{\Phi} = \{f_{\mathbf{w}} : \mathbf{x} \mapsto \sum_t \langle \tilde{\delta \mathbf{w}}_t, \Phi_t(\mathbf{x}) \rangle \mid \|\tilde{\delta \mathbf{w}}_t\|_2 \leq m_t\} \quad (13)$$

⁹In order to not assume a specific upper bound on the number of iterations, we can think of the updates from an iterative algorithm as an infinite sequence $\{\delta \mathbf{w}_0, \dots, \delta \mathbf{w}_t, \dots\}$ such that for some T , $\delta \mathbf{w}_t = 0$ for all $t > T$.

In this formulation, the result (12) reads:

$$\widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_{\mathbf{m}}^{\Phi}) \leq \sum_t (m_t/n) \|\Phi_t\|_{\text{F}}. \quad (14)$$

Optimizing the Feature Scaling. To obtain learning flows with lower complexity, Thm. 2 suggests modification of the algorithm to include, at each iteration t , a reparametrization step with a suitable matrix \tilde{A}_t giving a low contribution to the bound (12). Applied to gradient descent (GD), this leads to a new update rule sketched in Fig. 5 (left), where the optimization in Step 2 is over a given class of reparametrization matrices. The successive reparametrizations yield a varying feature map $\Phi_t = A_t^{-1} \Phi$ where $A_t = \tilde{A}_0 \dots \tilde{A}_t$.¹⁰

In the original representation Φ , **SuperNat** amounts to natural gradient descent (Amari, 1998) with respect to the local metric $g_{A_t} = A_t A_t^\top$. By construction, we also have $\delta f_{\mathbf{w}_t}(\mathbf{x}) = \langle \delta \mathbf{w}_{\text{GD}}, \Phi_t(\mathbf{x}) \rangle$ where $\delta \mathbf{w}_{\text{GD}}$ are standard gradient descent updates in the linear model with feature map Φ_t .

As an example, let $\Phi = \sum_{j=1}^n \sqrt{\lambda_j} \mathbf{u}_j \mathbf{v}_j^\top$ be the SVD of the feature matrix. We restrict to the class of matrices

$$\tilde{A}_{\nu} = \sum_{j=1}^n \sqrt{\nu_j} \mathbf{v}_j \mathbf{v}_j^\top + \text{Id}_{\text{span}\{\mathbf{v}\}^\perp} \quad (15)$$

labelled by weights $\nu_j > 0, j = 1, \dots, n$. With such a class, the action $\Phi_t^\top \rightarrow A_{\nu}^{-1} \Phi_t^\top$ merely rescales the

¹⁰Note that upon training a non-linear model, the updates of the tangent features take the same form $\Phi_t = \tilde{A}_t^{-1} \Phi_{t-1}$ as in Step 3 of **SuperNat**, the difference being that \tilde{A}_t is now a differential operator, e.g. at first order $\tilde{A}_t = \text{Id} - \delta \mathbf{w}_t^\top \frac{\partial}{\partial \mathbf{w}_t}$.

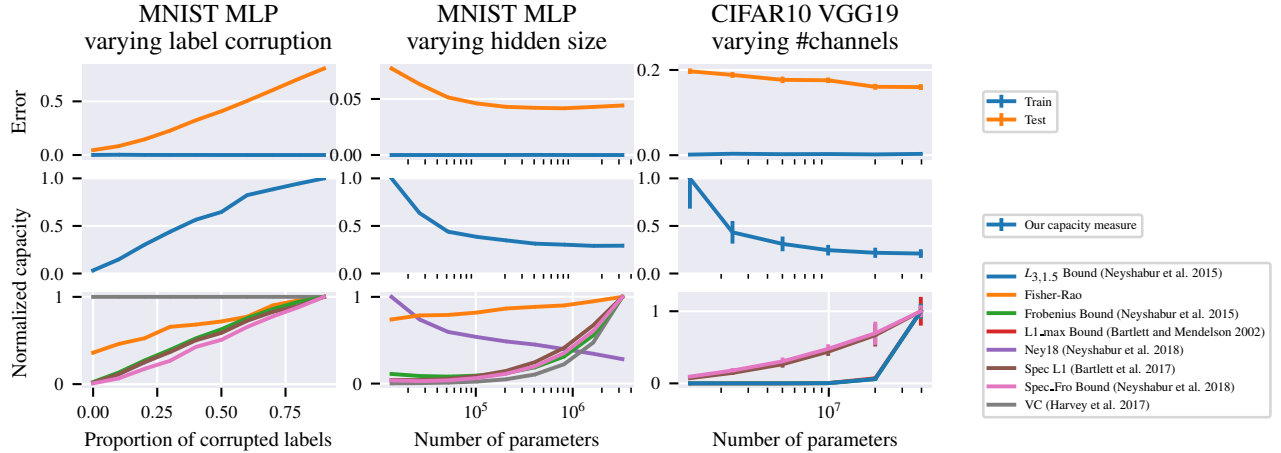


Figure 6: Complexity measures on MNIST with a 1 hidden layer MLP (**left**) as we increase the hidden layer size, (**center**) for a fixed hidden layer of 256 units as we increase label corruption and (**right**) for a VGG19 on CIFAR10 as we vary the number of channels. All networks are trained until cross-entropy reaches 0.01. Our proposed complexity measure and the one by Neyshabur et al. 2018 are the only ones to correctly reflect the shape of the generalization gap in these settings.

singular values $\lambda_{jt} \rightarrow \lambda_{jt}/\nu_j$, leaving the singular vectors unchanged. We work with gradient descent w.r.t a cost function L , so that $\delta \mathbf{w}_{\text{GD}} = -\eta \nabla_{\mathbf{w}} L$.

Proposition 3. *Any minimizer in Step 2 of SuperNat over matrices \mathbf{A}_ν in the class (15), takes the form*

$$\nu_{jt}^* = \kappa \frac{1}{|\mathbf{u}_j^\top \nabla_{\mathbf{f}_w} L|} \quad (16)$$

where $\nabla_{\mathbf{f}_w}$ denotes the gradient w.r.t the sample outputs $\mathbf{f}_w := [f_w(\mathbf{x}_1), \dots, f_w(\mathbf{x}_n)]^\top$, for some constant $\kappa > 0$.

In this context, this yields the following update rule, up to isotropic rescaling, for the singular values of Φ_t :

$$\lambda_{j(t+1)} = |\mathbf{u}_j^\top \nabla_{\mathbf{f}_w} L| \lambda_{jt}. \quad (17)$$

In this illustrative setting, we see how the feature map (or kernel) adapts to the task, by stretching (resp. contracting) its geometry in directions \mathbf{u}_j along which the residual $\nabla_{\mathbf{f}_w} L$ has large (resp. small) components. Intuitively, if a large component $|\mathbf{u}_j^\top \nabla_{\mathbf{f}_w} L|$ corresponds to signal and a small one $|\mathbf{u}_k^\top \nabla_{\mathbf{f}_w} L|$ corresponds to noise, then the ratio $\lambda_{jt}/\lambda_{kt}$ of singular values gets rescaled by the signal-to-noise ratio, thereby increasing the alignment of the learned features to the signal.

As a proof of concept, we consider the following regression setup. We consider a linear model with Gaussian features $\Phi = [\varphi, \varphi_{\text{noise}}] \in \mathbb{R}^{d+1}$ where $\varphi \sim \mathcal{N}(0, 1)$ and $\varphi_{\text{noise}} \sim \mathcal{N}(0, \frac{1}{d} I_d)$. Given n input samples, the n features $\Phi(\mathbf{x}_i)$ yield $\varphi \in \mathbb{R}^n$ and $\varphi_{\text{noise}} \in \mathbb{R}^{n \times d}$. We assume the label vector takes the form $\mathbf{y} = \varphi + P_{\text{noise}}(\boldsymbol{\epsilon})$, where Gaussian noise $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma^2 I_n)$ is projected onto the noise features through $P_{\text{noise}} = \varphi_{\text{noise}} \varphi_{\text{noise}}^\top$. The model is trained by gradient descent of the mean

squared loss and its **SuperNat** variant, where Step 2 uses the analytical solution of Proposition 3. We set $d = 10$, $\sigma^2 = 0.1$ and use $n = 50$ training points.

Fig 5 (right) shows test error obtained with standard and **SuperNat** gradient descent on this problem. At each iteration, **SuperNat** identifies dominant features (feature selection, here φ) and stretches the metric along them, thereby slowing down and eventually freezing the dynamics in the orthogonal (noise) directions (compression). The working hypothesis in this paper, supported by the observations of Section 3, is that for neural networks, such a (tangent) feature alignment is dynamically induced as an effect of non-linearity.

4.2 A New Complexity Measure for Neural Networks

Equ. (14) provides a bound of the Rademacher complexity for the function classes (11) specified by a *fixed* sequence of feature maps (see Appendix B.4 for a generalization to the multiclass setting). By extrapolation to the case of non-deterministic sequences of feature maps, we propose using

$$\mathcal{C}(f_w) = \sum_t \|\delta \mathbf{w}_t\|_2 \|\Phi_t\|_{\text{F}} \quad (18)$$

as a heuristic measure of complexity for neural networks, where Φ_t is the learned tangent feature matrix¹¹ at training iteration t , and $\|\delta \mathbf{w}_t\|_2$ is the norm of the SGD update. Following a standard protocol for studying complexity measures, (e.g., Neyshabur et al.,

¹¹In terms of tangent kernels, $\|\Phi_t\|_{\text{F}} = \sqrt{\text{Tr} \mathbf{K}_t}$ where \mathbf{K}_t is the tangent kernel (Gram) matrix.

2017a), Fig. 6 shows its behaviour for MLP on MNIST and VGG19 on CIFAR10 trained with cross entropy loss, with **(left)** fixed architecture and varying level of corruption in the labels and **(right)** varying hidden layer size/number of channels up to 4 millions parameters, against other capacity measures proposed in the recent literature. We observe that it correctly reflects the shape of the generalization gap.

5 Related Work

Role of Feature Geometry in Linear Models.

Analysis of the relation between capacity and feature geometry can be traced back to early work on kernel methods (Schölkopf et al., 1999a), which lead to data-dependent error bounds in terms of the eigenvalues of the kernel Gram matrix (Schölkopf et al., 1999b).

Recently, new analysis of minimum norm interpolators and max margin solutions for overparametrized linear models emphasize the key role of feature geometry, and specifically feature anisotropy, in the generalization performance (Bartlett et al., 2019; Muthukumar et al., 2019, 2020; Xie et al., 2020). Feature anisotropy combined to a high predictive power of the dominant features is the condition for a high centered alignment between kernel and class labels. In the context of neural networks, our results highlight the role of the non linear training dynamics in favouring such conditions.

Generalization Measures. There has been a large body of work on complexity/generalization measures for neural networks (see, Jiang et al., 2020, and references therein), some of which theoretically motivated by norm or margin based bounds (e.g., Neyshabur et al., 2019; Bartlett et al., 2017). Liang et al. (2019) proposed using the Fisher-Rao norm of the solution as a geometrically invariant complexity measure. By contrast, our approach to measuring complexity takes into account the geometry along the whole optimization trajectories. Since the geometry we consider is defined through the gradient second moments, our perspective is closely related to the notions of stiffness (Fort et al., 2019) and coherent gradients (Chatterjee, 2020).

Dynamics of Tangent Kernels. Several recent works investigated the ‘feature learning’ regime where neural tangent kernels evolve during training (Geiger et al., 2019; Woodworth et al., 2020). Independent concurrent works highlight alignment and compression phenomena similar to the one we study here (Kopitkov & Indelman, 2020; Paccolat et al., 2020). We offer various complementary empirical insights, and frame the alignment mechanism from the point of view of implicit regularization.

6 Conclusion

Through experiments with modern architectures, we highlighted an effect of dynamical alignment of the neural tangent features and their kernel along a small number of task-dependent directions during training, reflected by an early drop of the effective rank and an increasing similarity with the class labels, as measured by centered kernel alignment. We interpret this effect as a combined mechanism of feature selection and model compression of around dominant features.

Drawing upon intuitions from linear models, we argued that such a dynamical alignment acts as implicit regularizer. By extrapolating a new analysis of Rademacher complexity bounds for linear models, we also proposed a complexity measure that captures this phenomenon, and showed that it correlates with the generalization gap when varying the number of parameters, and when increasing the proportion of corrupted labels.

The results of this paper open several avenues for further investigation. The type of complexity measure we propose suggests new principled ways to design algorithms that learn the geometry in which to perform gradient descent (Srebro et al., 2011; Neyshabur et al., 2017b). Whether a procedure such as **SuperNat** can produce meaningful practical results for neural networks remains to be seen.

One of the consequences one can expect from the alignment effects highlighted here is to bias learning towards explaining most of the data with a small number of highly predictive features. While this feature selection ability might explain in part the performance of neural networks on a range of supervised tasks, it may also make them brittle under spurious correlation (e.g., Sagawa et al., 2020) and underpin their notorious weakness to generalize out-of-distribution (e.g., Geirhos et al., 2020). Resolving this tension is an important challenge towards building more robust models.

Acknowledgments

We thank X. Y Lu and V. Thomas for collaboration at an early stage of this project, A. Sordani for insightful discussions, G. Gidel, A. Mitra and M. Pezeshki for helpful feedback. This research was partially supported by the Canada CIFAR AI Chair Program (held at Mila); by NSERC through the Discovery Grants RGPIN-2017-06936 (S.LJ) and RGPIN-2018-04821 (G.L) and an Alexander Graham Bell Canada Graduate Scholarship (CGS D) award (A.B); by FRQNT Young Investigator Startup Program 2019- NC-253251 (G.L); and by a Google Focused Research award (S.LJ). S.LJ and P.V are CIFAR Associate Fellows in the Learning in Machines & Brains program.

References

- Madhu S Advani and Andrew M Saxe. High-dimensional dynamics of generalization error in neural networks. *arXiv preprint arXiv:1710.03667*, 2017.
- Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via overparameterization. volume 97 of *Proceedings of Machine Learning Research*, pp. 242–252, Long Beach, California, USA, 09–15 Jun 2019. PMLR.
- Shun-Ichi Amari. Natural gradient works efficiently in learning. *Neural computation*, 10(2):251–276, 1998.
- Shun-Ichi Amari. *Information Geometry and Its Applications*, volume 194. Springer, 2016.
- Devansh Arpit, Stanislaw Jastrzebski, Nicolas Ballas, David Krueger, Emmanuel Bengio, Maxinder S Kanwal, Tegan Maharaj, Asja Fischer, Aaron Courville, Yoshua Bengio, et al. A closer look at memorization in deep networks. *arXiv preprint arXiv:1706.05394*, 2017.
- Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *JMLR*, 2002.
- Peter L. Bartlett, Dylan J. Foster, and Matus Telgarsky. Spectrally-normalized margin bounds for neural networks. In *NIPS*, 2017.
- Peter L Bartlett, Philip M Long, Gábor Lugosi, and Alexander Tsigler. Benign overfitting in linear regression. *arXiv preprint arXiv:1906.11300[stat.ML]*, 2019.
- Ronen Basri, David Jacobs, Yoni Kasten, and Shira Kritchman. The convergence rate of neural networks for learned functions of different frequencies. In *Advances in Neural Information Processing Systems 32*, pp. 4761–4771. 2019.
- Mikhail Belkin, Siyuan Ma, and Soumik Mandal. To understand deep learning we need to understand kernel learning. In *ICML*, 2018.
- Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-learning practice and the classical bias–variance trade-off. *Proceedings of the National Academy of Sciences*, 116(32):15849–15854, 2019.
- Yoshua Bengio, Olivier Delalleau, Nicolas Le Roux, Jean-François Paiement, Pascal Vincent, and Marie Ouimet. Learning eigenfunctions links spectral embedding and kernel PCA. *Neural Computation*, 16(10):2197–2219, 2004.
- Alberto Bietti and Julien Mairal. On the inductive bias of neural tangent kernels. In *Advances in Neural Information Processing Systems 32*, pp. 12893–12904. 2019.
- Mikio L Braun. *Spectral properties of the kernel matrix and their relation to kernel methods in machine learning*. PhD thesis, Universitäts- und Landesbibliothek Bonn, 2005.
- Satrajit Chatterjee. Coherent gradients: An approach to understanding generalization in gradient descent-based optimization. In *International Conference on Learning Representations*, 2020.
- Tarin Clanuwat, Mikel Bober-Irizar, Asanobu Kitamoto, Alex Lamb, Kazuaki Yamamoto, and David Ha. Deep learning for classical japanese literature. 2018.
- Corinna Cortes, Mehryar Mohri, and Afshin Rostamizadeh. Algorithms for learning kernels based on centered alignment. *JMLR*, 13(1):795–828, 2012. ISSN 1532-4435.
- Nello Cristianini, John Shawe-Taylor, André Elisseeff, and Jaz S. Kandola. On kernel-target alignment. In *NIPS*. 2002.
- Simon S. Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes overparameterized neural networks. In *International Conference on Learning Representations*, 2019.
- Stanislav Fort, Pawel Krzyszttof Nowak, Stanislaw Jastrzebski, and Sridhar Narayanan. Stiffness: A new perspective on generalization in neural networks. *arXiv preprint arXiv:1901.09491*, 2019.
- Mario Geiger, Stephano Spigler, Arthur Jacot, and Matthieu Wyart. Disentangling feature and lazy training in deep neural networks. *arXiv:1906.08034 [cs.LG]*, 2019.
- Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *arXiv preprint arXiv:2004.07780 [cs.CV]*, 2020.
- Stuart Geman, Elie Bienenstock, and René Doursat. Neural networks and the bias/variance dilemma. *Neural Computation*, 4(1):1–58, 1992. doi: 10.1162/neco.1992.4.1.1.
- Thomas George. NNGeometry: Easy and Fast Fisher Information Matrices and Neural Tangent Kernels in PyTorch, February 2021.
- Gauthier Gidel, Francis Bach, and Simon Lacoste-Julien. Implicit regularization of discrete gradient dynamics in linear neural networks. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 32*, pp. 3202–3211. Curran Associates, Inc., 2019.

- Arthur Gretton, Olivier Bousquet, Alexander Smola, and Bernhard Schölkopf. Measuring statistical dependence with hilbert-schmidt norms, 2005.
- Suriya Gunasekar, Jason Lee, Daniel Soudry, and Nathan Srebro. Characterizing implicit bias in terms of optimization geometry. In Jennifer Dy and Andreas Krause (eds.), *ICML*, volume 80 of *Proceedings of Machine Learning Research*, pp. 1832–1841, Stockholm, Sweden, 10–15 Jul 2018. PMLR.
- Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The elements of statistical learning: data mining, inference and prediction*. Springer, 2009.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Elad Hoffer, Itay Hubara, and Daniel Soudry. Train longer, generalize better: Closing the generalization gap in large batch training of neural networks. In *NIPS*, 2017.
- Arthur Jacot, Franck Gabriel, and Clement Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *NIPS*, pp. 8571–8580, 2018.
- Yiding Jiang, Behnam Neyshabur, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantastic generalization measures and where to find them. In *ICLR*, 2020.
- Ryo Karakida, Shotaro Akaho, and Shun-ichi Amari. Pathological spectra of the fisher information metric and its variants in deep neural networks. *arXiv:1910.05992 [stat.ML]*, 2019a.
- Ryo Karakida, Shotaro Akaho, and Shun-ichi Amari. Universal statistics of fisher information in deep neural networks: Mean field approach. *AISTATS 2019*, 2019b.
- D. Kopitkov and V. Indelman. Neural spectrum alignment: Empirical study. In *International Conference on Artificial Neural Networks (ICANN)*, September 2020.
- Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009.
- Andrew K Lampinen, Andrew K Lampinen, and Surya Ganguli. An analytic theory of generalization dynamics and transfer learning in deep linear networks. *arXiv.org*, 2018.
- S. Lang. *Fundamentals of Differential Geometry*. Graduate Texts in Mathematics. Springer New York, 2012. ISBN 9781461205418.
- Yann LeCun, Corinna Cortes, and CJ Burges. Mnist handwritten digit database. *ATT Labs [Online]*. Available: <http://yann.lecun.com/exdb/mnist>, 2, 2010.
- M. Ledoux and M. Talagrand. *Probability in Banach Spaces: Isoperimetry and Processes*. Springer Science & Business, New York, 2013.
- Tengyuan Liang, Tomaso Poggio, Alexander Rakhlin, and James Stokes. Fisher-rao metric, geometry, and complexity of neural networks. In *Proceedings of Machine Learning Research*, volume 89, pp. 888–896, 2019.
- Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of Machine Learning*. The MIT Press, 2012. ISBN 026201825X, 9780262018258.
- Vidya Muthukumar, Kailas Vodrahalli, Vignesh Subramanian, and Anant Sahai. Harmless interpolation of noisy data in regression. *arXiv preprint arXiv:1903.09139[cs.LG]*, 2019.
- Vidya Muthukumar, Adhyayan Narang, Vignesh Subramanian, Mikhail Belkin, Daniel Sahai, Hsu, and Anant Sahai. Classification vs regression in overparameterized regimes: Does the loss function matter? *arXiv preprint arXiv:2005.08054 [cs.LG]*, 2020.
- Brady Neal, Sarthak Mittal, Aristide Baratin, Vinayak Tantia, Matthew Scicluna, Simon Lacoste-Julien, and Ioannis Mitliagkas. A modern take on the bias-variance tradeoff in neural networks. *arXiv:1810.08591 [cs.LG]*, 2018.
- Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. In search of the real inductive bias: On the role of implicit regularization in deep learning. *ICLR workshop track*, 2015.
- Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and Nati Srebro. Exploring generalization in deep learning. In *Advances in Neural Information Processing Systems*, pp. 5949–5958, 2017a.
- Behnam Neyshabur, Ryota Tomioka, Ruslan Salakhutdinov, and Nathan Srebro. Geometry of optimization and implicit regularization in deep learning. *arXiv:1705.03071 [cs.LG]*, 2017b.
- Behnam Neyshabur, Zhiyuan Li, Srinadh Bhojanapalli, Yann LeCun, and Nathan Srebro. Towards understanding the role of over-parametrization in generalization of neural networks. *International Conference on Learning Representations (ICLR)*, 2019.
- Jonas Paccolat, Leonardo Petrini, Mario Geiger, Kevin Tyloo, and Matthieu Wyart. Geometric compression of invariant manifolds in neural nets. *arXiv preprint arXiv:2007.11471*, 2020.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor

- Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 32*, pp. 8024–8035. Curran Associates, Inc., 2019.
- Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht, Yoshua Bengio, and Aaron Courville. On the spectral bias of neural networks. In *Proceedings of the 36th International Conference on Machine Learning*, 2019.
- Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *NIPS*, 2007.
- Olivier Roy and Martin Vetterli. The effective rank: A measure of effective dimensionality. In *2007 15th European Signal Processing Conference*, pp. 606–610. IEEE, 2007.
- Shiori Sagawa, Aditi Raghunathan, Pang Wei Koh, and Percy Liang. An investigation of why overparameterization exacerbates spurious correlations. *arXiv:2005.04345 [cs.LG]*, 2020.
- Andrew M. Saxe, James L. McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural network. In *International Conference on Learning Representations*, 2014.
- B. Schölkopf, S. Mika, C. J.C. Burges, P. Knirsch, K. R. Müller, G. Ratsch, and A. J. Smola. Input space versus feature space in kernel-based methods. *Trans. Neur. Netw.*, 10(5):1000–1017, September 1999a. ISSN 1045-9227.
- B. Schölkopf, J. Shawe-Taylor, A.J. Smola, and R.C. Williamson. Kernel-dependent support vector error bounds. In *Artificial Neural Networks, 1999. ICANN 99*, volume 470 of *Conference Publications*, pp. 103–108. Max-Planck-Gesellschaft, IEEE, 1999b.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- Nati Srebro, Karthik Sridharan, and Ambuj Tewari. On the universality of online mirror descent. In *Advances in Neural Information Processing Systems 24*. 2011.
- Blake Woodworth, Suriya Gunasekar, Jason D. Lee, Edward Moroshko, Pedro Savarese, Itay Golan, Daniel Soudry, and Nathan Srebro. Kernel and rich regimes in overparametrized models. *arXiv:2002.09277 [cs.LG]*, 2020.
- Yuege Xie, Rachel Ward, Holger Rauhut, and Chou Hung-Hsu. Weighted optimization: better generalization by smoother interpolation. *arXiv preprint arXiv:2006.08495*, 2020.
- Zhi-Qin John Xu, Yaoyu Zhang, and Yanyang Xiao. Training behavior of deep neural network in frequency domain. In Tom Gedeon, Kok Wai Wong, and Minhoo Lee (eds.), *Neural Information Processing*, pp. 264–274, Cham, 2019. Springer International Publishing. ISBN 978-3-030-36708-4.
- Greg Yang and Hadi Salman. A fine grained spectral perspective on neural networks. *arxiv preprint arXiv:1907.10599[cs.LG]*, 2019.
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. *ICLR*, 2017.

APPENDICES: Implicit Regularization via Neural Feature Alignment

A Tangent Features and Geometry

We describe in more formal detail some of the notions introduced in Section 2 of the paper. We will consider general classes of vector-valued predictors:

$$\mathcal{F} = \{f_{\mathbf{w}} : \mathcal{X} \rightarrow \mathbb{R}^c \mid \mathbf{w} \in \mathcal{W}\}, \quad (19)$$

where the parameter space \mathcal{W} is a finite dimensional manifold of dimension P (typically \mathbb{R}^P). For multiclass classification, $f_{\mathbf{w}}$ outputs a score $f_{\mathbf{w}}(\mathbf{x})[y]$ for each class $y \in \{1 \cdots c\}$. Each function can also be viewed as a scalar function on $\mathcal{X} \times \mathcal{Y}$ where $\mathcal{Y} = \{1 \cdots c\}$ is the set of classes.

A.1 Metric

We assume that $\mathbf{w} \rightarrow f_{\mathbf{w}}$ is a smooth mapping from \mathcal{W} to $L^2(\rho, \mathbb{R}^c)$, where ρ is some input data distribution. The inclusion $\mathcal{F} \subset L^2(\rho, \mathbb{R}^c)$ equips \mathcal{F} with the L^2 scalar product and corresponding norm:

$$\langle f, g \rangle_{\rho} := \mathbb{E}_{\mathbf{x} \sim \rho} [f(\mathbf{x})^{\top} g(\mathbf{x})], \quad \|f\|_{\rho} := \sqrt{\langle f, f \rangle_{\rho}} \quad (20)$$

The parameter space \mathcal{W} inherits a **metric tensor** $g_{\mathbf{w}}$ by pull-back of the scalar product $\langle f, g \rangle_{\rho}$ on \mathcal{F} . That is, given $\zeta, \xi \in \mathcal{T}_{\mathbf{w}}\mathcal{W} \cong \mathbb{R}^P$ on the tangent space at \mathbf{w} (Lang, 2012),

$$g_{\mathbf{w}}(\zeta, \xi) = \langle \partial_{\zeta} f_{\mathbf{w}}, \partial_{\xi} f_{\mathbf{w}} \rangle_{\rho} \quad (21)$$

where $\partial_{\zeta} f_{\mathbf{w}} = \langle df_{\mathbf{w}}, \zeta \rangle$ is the directional derivative in the direction of ζ . Concretely, in a given basis of \mathbb{R}^P , the metric is represented by the matrix of gradient second moments:

$$(g_{\mathbf{w}})_{pq} = \mathbb{E}_{\mathbf{x} \sim \rho} \left[\left(\frac{\partial f_{\mathbf{w}}(\mathbf{x})}{\partial w_p} \right)^{\top} \frac{\partial f_{\mathbf{w}}(\mathbf{x})}{\partial w_q} \right] \quad (22)$$

where $w_p, p = 1, \dots, P$ are the parameter coordinates. The metric shows up by spelling out the line element $ds^2 := \|df_{\mathbf{w}}\|_{\rho}^2$, since we have,

$$\|df_{\mathbf{w}}\|_{\rho}^2 = \sum_{p,q=1}^P \left\langle \frac{\partial f_{\mathbf{w}}}{\partial w_p} dw_p, \frac{\partial f_{\mathbf{w}}}{\partial w_q} dw_q \right\rangle_{\rho} = \sum_{p,q=1}^P (g_{\mathbf{w}})_{pq} dw_p dw_q \quad (23)$$

A.2 Tangent Kernels

This geometry has a dual description in function space in terms of *kernels*. The idea is to view the differential of the mapping $\mathbf{w} \rightarrow f_{\mathbf{w}}$ at each \mathbf{w} as a map $df_{\mathbf{w}} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{T}_{\mathbf{w}}^* \mathcal{W} \cong \mathbb{R}^P$ defining (joined) features in the (co)tangent space. In a given basis, this yields the **tangent features** given by the function derivatives w.r.t the parameters,

$$\Phi_{w_p}(\mathbf{x})[y] := \frac{\partial f_{\mathbf{w}}(\mathbf{x})[y]}{\partial w_p} \quad (24)$$

The tangent feature map $\Phi_{\mathbf{w}}$ can be viewed as a function mapping each pair (\mathbf{x}, y) to a vector in \mathbb{R}^P . It defines the so-called **tangent kernel** (Jacot et al., 2018) through the Euclidean dot product $\langle \cdot, \cdot \rangle$ in \mathbb{R}^P :

$$k_{\mathbf{w}}(\mathbf{x}, y; \tilde{\mathbf{x}}, y') = \langle \Phi_{\mathbf{w}}(\mathbf{x})[y], \Phi_{\mathbf{w}}(\tilde{\mathbf{x}})[y'] \rangle = \sum_{p=1}^P \Phi_{w_p}(\mathbf{x})[y] \Phi_{w_p}(\tilde{\mathbf{x}})[y'] \quad (25)$$

It induces an integral operator on $L^2(\rho, \mathbb{R}^c)$ acting as

$$(k_{\mathbf{w}} \triangleright f)(\mathbf{x})[y] = \langle k_{\mathbf{w}}(\mathbf{x}, y; \cdot), f \rangle \quad (26)$$

The metric tensor (22) is expressed in terms of the tangent features as $(g_{\mathbf{w}})_{pq} = \langle \Phi_{w_p}, \Phi_{w_q} \rangle_{\rho}$.

A.3 Spectral Decomposition

The local metric tensor (as symmetric $P \times P$ matrix) and tangent kernel (as rank P integral operator) share the same spectrum. More generally, let

$$g_{\mathbf{w}} = \sum_{j=1}^P \lambda_{\mathbf{w}j} \mathbf{v}_{\mathbf{w}j} \mathbf{v}_{\mathbf{w}j}^{\top} \quad (27)$$

be the eigenvalue decomposition of the positive (semi-)definite symmetric matrix (22), where $\mathbf{v}_{\mathbf{w}j}^{\top} \mathbf{v}_{\mathbf{w}j'} = \delta_{jj'}$. Assuming non-degeneracy, i.e $\lambda_{\mathbf{w}j} > 0$, let $u_{\mathbf{w}j}, j \in \{1 \cdots P\}$ be the functions in $L^2(\rho, \mathbb{R}^c)$ defined as:

$$u_{\mathbf{w}j}(\mathbf{x})[y] = \frac{1}{\sqrt{\lambda_{\mathbf{w}j}}} \mathbf{v}_{\mathbf{w}j}^{\top} \Phi_{\mathbf{w}}(\mathbf{x})[y] \quad (28)$$

The following result holds.

Proposition 4 (Spectral decomposition). *The functions $(u_{j\mathbf{w}})_{j=1}^P$ form an orthonormal family in $L^2(\rho, \mathbb{R}^c)$. They are eigenfunctions of the tangent kernel as an integral operator, which admits the spectral decomposition:*

$$k_{\mathbf{w}}(\mathbf{x}, y; \tilde{\mathbf{x}}, y') = \sum_{j=1}^P \lambda_{\mathbf{w}j} u_{\mathbf{w}j}(\mathbf{x})[y] u_{\mathbf{w}j}(\tilde{\mathbf{x}})[y'] \quad (29)$$

In particular metric tensor and tangent kernels share the same spectrum.

Proof. We first show orthonormality, i.e $\langle u_{\mathbf{w}j}, u_{\mathbf{w}j'} \rangle_{\rho} = \delta_{jj'}$. We have indeed,

$$\langle u_{\mathbf{w}j}, u_{\mathbf{w}j'} \rangle_{\rho} = \frac{1}{\sqrt{\lambda_{\mathbf{w}j} \lambda_{\mathbf{w}j'}}} \sum_{p,q=1}^P (\mathbf{v}_{\mathbf{w}j})_p (\mathbf{v}_{\mathbf{w}j'})_q \langle \Phi_{w_p}, \Phi_{w_q} \rangle_{\rho} \quad (30)$$

$$= \frac{1}{\sqrt{\lambda_{\mathbf{w}j} \lambda_{\mathbf{w}j'}}} \mathbf{v}_{\mathbf{w}j}^{\top} g_{\mathbf{w}} \mathbf{v}_{\mathbf{w}j'} \quad (31)$$

$$= \frac{1}{\lambda_{\mathbf{w}j}} \lambda_{\mathbf{w}j} \delta_{jj'} \quad (32)$$

$$= \delta_{jj'} \quad (33)$$

where we used the definition of the matrix $(g_{\mathbf{w}})_{pq}$ and its eigenvalue decomposition. Next, using the action (26) of the tangent kernel, we prove that the functions $u_{\mathbf{w}j}$ defined in (28) is an eigenfunction with eigenvalue $\lambda_{\mathbf{w}j}$:

$$(k_{\mathbf{w}} \triangleright u_{\mathbf{w}j})(\mathbf{x})[y] = \sum_{p=1}^P \Phi_{w_p}(\mathbf{x})[y] \langle \Phi_{w_p}, u_{\mathbf{w}j} \rangle_{\rho} \quad (34)$$

$$= \frac{1}{\sqrt{\lambda_{\mathbf{w}j}}} \sum_{p,q=1}^P (\mathbf{v}_{\mathbf{w}j})_q \Phi_{w_p}(\mathbf{x})[y] \langle \Phi_{w_p}, \Phi_{w_q} \rangle \quad (35)$$

$$= \frac{1}{\sqrt{\lambda_{\mathbf{w}j}}} \mathbf{v}_{\mathbf{w}j}^{\top} g_{\mathbf{w}} \Phi_{\mathbf{w}}(\mathbf{x})[y] \quad (36)$$

$$= \frac{1}{\sqrt{\lambda_{\mathbf{w}j}}} (\lambda_{\mathbf{w}j} \mathbf{v}_{\mathbf{w}j}^{\top}) \Phi_{\mathbf{w}}(\mathbf{x})[y] \quad (37)$$

$$= \lambda_{\mathbf{w}j} \frac{1}{\sqrt{\lambda_{\mathbf{w}j}}} \mathbf{v}_{\mathbf{w}j}^{\top} \Phi_{\mathbf{w}}(\mathbf{x})[y] \quad (38)$$

$$= \lambda_{\mathbf{w}j} u_{\mathbf{w}j} \quad (39)$$

Inserting the resolution of unity $\text{Id}_P = \sum_{j=1}^P \mathbf{v}_{\mathbf{w}j} \mathbf{v}_{\mathbf{w}j}^{\top}$ in the expression (25) of the tangent kernel directly yields the spectral decomposition (29). \square

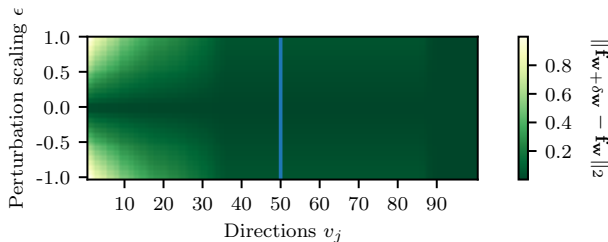


Figure 7: Variations of \mathbf{f}_w (evaluated on a test set) when perturbing the parameters in the directions given by the right singular vectors of the Jacobian (first 50 directions) or in randomly sampled directions (last 50 directions) on a VGG11 network trained for 10 epochs on CIFAR10. We observe that perturbations in most directions have almost no effect, except in those aligned with the top singular vectors.

A.4 Sampled Versions

Given n input samples $\mathbf{x}_1, \dots, \mathbf{x}_n$, any function $f: \mathcal{X} \rightarrow \mathbb{R}^c$ yields a vector $\mathbf{f} \in \mathbb{R}^{nc}$ obtained by concatenating the outputs $f(\mathbf{x}_i) \in \mathbb{R}^c$ of the n input samples \mathbf{x}_i . The sample output scores $f_w(\mathbf{x}_i)[y]$ thus yields $\mathbf{f}_w \in \mathbb{R}^{nc}$; and the tangent features $\Phi_{w_p}(\mathbf{x}_i)[y]$ are represented as a $nc \times P$ matrix Φ_w . Using this notation, (22) and (25) yield the sample covariance $P \times P$ matrix and kernel (Gram) $nc \times nc$ matrix:

$$\mathbf{G}_w = \Phi_w^\top \Phi_w, \quad \mathbf{K}_w = \Phi_w \Phi_w^\top \quad (40)$$

The eigenvalue decompositions of \mathbf{G}_w and \mathbf{K}_w follow from the (SVD) of Φ_w : assuming $P > nc$, we can write this SVD by indexing the singular values by a pair $J = (i, y)$ with $i = 1, \dots, n$ and $y = 1 \dots c$ as

$$\Phi_w = \sum_{J=1}^{nc} \sqrt{\hat{\lambda}_{wJ}} \hat{\mathbf{u}}_{wJ} \hat{\mathbf{v}}_{wJ}^\top \quad (41)$$

Such decompositions summarize the predominant directions both in parameter and feature space, in the neighborhood of \mathbf{w} : a small variation $\delta \mathbf{w}$ induces the first order variation $\delta \mathbf{f}_w$ of the function,

$$\delta \mathbf{f}_w := \Phi_w \delta \mathbf{w} = \sum_{J=1}^{nc} \sqrt{\hat{\lambda}_{wJ}} (\hat{\mathbf{v}}_{wJ}^\top \delta \mathbf{w}) \hat{\mathbf{u}}_{wJ} \quad (42)$$

Fig. 7 illustrates this ‘hierarchy’ for a VGG11 network (Simonyan & Zisserman, 2014) trained for 10 epoches on CIFAR10 (Krizhevsky & Hinton, 2009). We observe that perturbations in most directions have almost no effect, except in those aligned with the top singular vectors. This is reflected by a strong anisotropy of the tangent kernel spectrum. Recent analytical results for wide random neural networks also point to such a pathological structure of the spectrum (Karakida et al., 2019a,b).

A.5 Spectral Bias

A.5.1 Proof of Lemma 1

We consider parameter updates $\delta \mathbf{w}_{GD} := -\eta \nabla_{\mathbf{w}} L$ for gradient descent w.r.t a loss $L := L(\mathbf{f}_w)$, which is a function of the vector $\mathbf{f}_w \in \mathbb{R}^{nc}$ of sample output scores. We reformulate Lemma 1, extended to the multiclass setting.

Proposition 5 (Lemma 1 restated). *The gradient descent function updates in first order Taylor approximation, $\delta f_{GD}(\mathbf{x})[y] := \langle \delta \mathbf{w}_{GD}, \Phi_w(\mathbf{x})[y] \rangle$, decompose as,*

$$\delta f_{GD}(\mathbf{x})[y] = \sum_{j=1}^P \delta f_j u_{wj}(\mathbf{x})[y], \quad \delta f_j = -\eta \lambda_{wj} (\mathbf{u}_{wj}^\top \nabla_{\mathbf{f}_w} L) \quad (43)$$

where u_{wj} are the eigenfunctions (28) of the tangent kernel and $\mathbf{u}_{wj} \in \mathbb{R}^{nc}$ are their corresponding sample vector.

Proof. Inserting the resolution of unity $\text{Id}_P = \sum_{j=1}^P \mathbf{v}_{\mathbf{w}j} \mathbf{v}_{\mathbf{w}j}^\top$ in the expression for δf_{GD} yields

$$\delta f_{\text{GD}}(\mathbf{x})[y] = \sum_{j=1}^P (\mathbf{v}_{\mathbf{w}j}^\top \delta \mathbf{w}_{\text{GD}}) \mathbf{v}_{\mathbf{w}j}^\top \Phi_{\mathbf{w}}(\mathbf{x})[y] \quad (44)$$

$$= \sum_{j=1}^P \sqrt{\lambda_{\mathbf{w}j}} (\mathbf{u}_{\mathbf{w}j}^\top \delta \mathbf{w}_{\text{GD}}) u_{\mathbf{w}j}(\mathbf{x})[y] \quad (45)$$

Next, by the chain rule $\nabla_{\mathbf{w}} L = \Phi_{\mathbf{w}}^\top \nabla_{\mathbf{f}_{\mathbf{w}}} L$, so we can spell out:

$$\delta \mathbf{w}_{\text{GD}} = -\eta \sum_{j=1}^P \sqrt{\lambda_{\mathbf{w}j}} (\mathbf{u}_{\mathbf{w}j}^\top \nabla_{\mathbf{f}_{\mathbf{w}}} L) \mathbf{v}_{\mathbf{w}j}, \quad (46)$$

which implies that $(\mathbf{v}_{\mathbf{w}j}^\top \delta \mathbf{w}_{\text{GD}}) = \sqrt{\lambda_{\mathbf{w}j}} (\mathbf{u}_{\mathbf{w}j}^\top \nabla_{\mathbf{f}_{\mathbf{w}}} L)$. Substituting in (44) gives the desired result. \square

The decomposition (48) has a *sampled* version in terms of tangent feature and kernel matrices. Using the notation of SVD (41), let $\hat{\lambda}_{\mathbf{w}J}$, $\hat{\mathbf{u}}_{\mathbf{w}J}$ and $\hat{\mathbf{v}}_{\mathbf{w}J}$ be correspond to the (non-zero) eigenvalues and eigenvectors of the sample covariance and kernel (40). We consider the tangent kernel **principal components**, defined as the functions

$$\hat{u}_{\mathbf{w}J}(\mathbf{x})[y] = \frac{1}{\sqrt{\lambda_{\mathbf{w}J}}} \langle \hat{\mathbf{v}}_{\mathbf{w}J}, \Phi_{\mathbf{w}}(\mathbf{x})[y] \rangle, \quad (47)$$

which form an orthonormal family for the in-sample scalar product $\langle f, g \rangle_{\text{in}} = \sum_{i=1}^n f(\mathbf{x}_i) g(\mathbf{x}_i)$ and approximate the true kernel eigenfunctions (28) (e.g., Bengio et al., 2004; Braun, 2005). One can easily check from (41) that the vector $\hat{\mathbf{u}}_{\mathbf{w}J} \in \mathbb{R}^{nc}$ of sample outputs $\hat{u}(\mathbf{x}_i)[y]$ coincides with the J -th eigenvector of the tangent kernel matrix.

Proposition 6 (Sampled version of Prop 5). *The gradient descent function updates in first order Taylor approximation, $\delta f_{\text{GD}}(\mathbf{x})[y] := \langle \delta \mathbf{w}_{\text{GD}}, \Phi_{\mathbf{w}}(\mathbf{x})[y] \rangle$ decompose as,*

$$\delta f_{\text{GD}}(\mathbf{x})[y] = \sum_{j=1}^{nc} \delta f_J \hat{u}_{\mathbf{w}J}(\mathbf{x})[y], \quad \delta f_J = -\eta \hat{\lambda}_{\mathbf{w}J} (\hat{\mathbf{u}}_{\mathbf{w}J}^\top \nabla_{\mathbf{f}_{\mathbf{w}}} L) \quad (48)$$

in terms of the principal components (47) of the tangent kernel.

Proof. Same proof as for the previous Proposition, using the resolution of unity $\text{Id}_{nc} = \sum_{J=1}^{nc} \hat{\mathbf{v}}_{\mathbf{w}J} \hat{\mathbf{v}}_{\mathbf{w}J}^\top$. \square

A.5.2 The Case of Linear Regression

The previous Proposition gives a ‘local’ version of a classic decomposition of the training dynamics in linear regression (e.g., Advani & Saxe, 2017). In such a setting, $f_{\mathbf{w}} = \langle \mathbf{w}, \Phi(\mathbf{x}) \rangle$ are linearly parametrized scalar functions ($c = 1$) and $L = \frac{1}{2} \|\mathbf{f}_{\mathbf{w}} - \mathbf{y}\|^2$. We denote by $\Phi = \sum_{j=1}^n \hat{\lambda}_j \hat{\mathbf{u}}_j \hat{\mathbf{v}}_j^\top$ the $n \times P$ feature matrix and its SVD.

Proposition 7. *Gradient descent of the squared loss yields the function iterates,*

$$f_{\mathbf{w}_t} = f_{\mathbf{w}^*} + (\text{Id} - \eta K)^t (f_{\mathbf{w}_0} - f_{\mathbf{w}^*}) \quad (49)$$

where Id is the identity map and K is the operator acting on functions as $(K \triangleright f)(\mathbf{x}) = \sum_{i=1}^n k(\mathbf{x}, \mathbf{x}_i) f(\mathbf{x}_i)$ in terms of the kernel $k(\mathbf{x}, \tilde{\mathbf{x}}) = \langle \Phi(\mathbf{x}), \Phi(\tilde{\mathbf{x}}) \rangle$.

Proof. The updates $\delta \mathbf{w}_{\text{GD}} := -\eta \nabla_{\mathbf{w}} L$ induce the (exact) functional updates $\delta f_{\text{GD}} = f_{\mathbf{w}_{t+1}} - f_{\mathbf{w}_t}$ given by

$$\delta f_{\text{GD}}(\mathbf{x}) = -\eta \sum_{i=1}^n k(\mathbf{x}, \mathbf{x}_i) (f_{\mathbf{w}_t}(\mathbf{x}_i) - \mathbf{y}_i) \quad (50)$$

Substituting $\mathbf{y}_i = f_{\mathbf{w}^*}(\mathbf{x}_i)$ gives $f_{\mathbf{w}_{t+1}} - f_{\mathbf{w}^*} = (\text{id} - \eta K)(f_{\mathbf{w}_t} - f_{\mathbf{w}^*})$. Equ. 49 follows by induction. \square

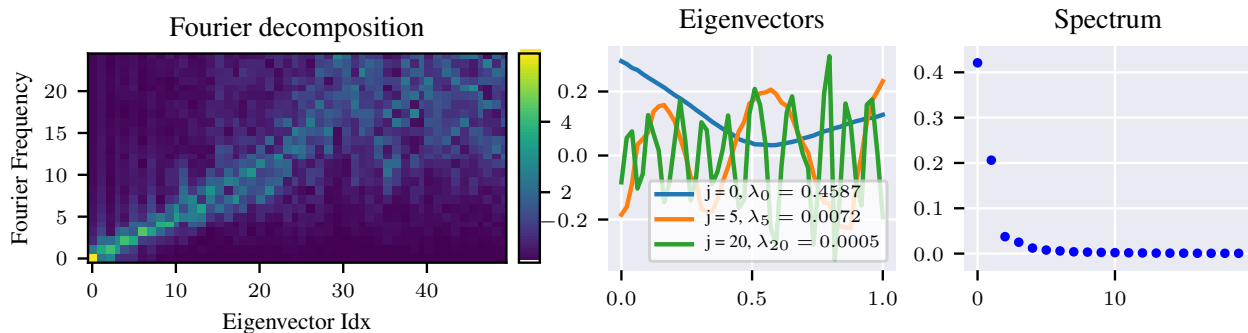


Figure 8: Eigendecomposition of the tangent kernel matrix of a random 6-layer deep 256-unit wide MLP on 1D uniform data (50 equally spaced points in $[0, 1]$). **(left)** Fourier decomposition (y -axis for frequency, colorbar for magnitude) of each eigenvector (x -axis), ranked in nonincreasing order of the eigenvalues. We observe that eigenvectors with increasing index j (hence decreasing eigenvalues) correspond to modes with increasing Fourier frequency. **(middle)** Plot of the j -th eigenvectors with $j \in \{0, 5, 20\}$ and **(right)** distribution of eigenvalues. We note the fast decay (e.g. $\lambda_{10}/\lambda_1 \approx 4\%$).

Lemma 8. *The kernel principal components $\hat{u}_j(\mathbf{x}) = \frac{1}{\sqrt{\hat{\lambda}_j}} \langle \hat{\mathbf{v}}_j, \Phi_{\mathbf{w}}(\mathbf{x}) \rangle$ are eigenfunctions of the operator K with corresponding eigenvalues $\hat{\lambda}_j$.*

Proof. By inserting $Id_n = \sum_j \hat{\mathbf{v}}_j \hat{\mathbf{v}}_j^\top$ in the expression of the kernel, one can write $k(\mathbf{x}, \mathbf{x}_i) = \sum_{j=1}^n \hat{u}_j(\mathbf{x}) \hat{u}_j(\mathbf{x}_i)$. Substituting in the definition of K and using the orthonormality of \hat{u}_j for the in-sample scalar product yield $K \triangleright \hat{u}_j = \hat{\lambda}_j \hat{u}_j$. \square

Together with (49), this directly leads to the decoupling of the training dynamics in the basis of kernel principal components.

Proposition 9 (Spectral Bias for Linear Regression). *By initializing $\mathbf{w}_0 = \Phi^\top \alpha_0$ in the span of the features, the function iterates in (49) uniquely decompose as,*

$$f_{\mathbf{w}_t}(\mathbf{x}) = \sum_{j=1}^n f_{jt} \hat{u}_j(\mathbf{x}), \quad f_{jt} = f_j^* + (1 - \eta \lambda_j)^t (f_{j0} - f_j^*) \quad (51)$$

where f_j^* are the coefficients of the (minimum ℓ_2 -norm) interpolating solution.

This standard result shows how each independent mode labelled by j has its own linear convergence rate. For example setting $\eta = 1/\lambda_1$, this gives $f_{jt} - f_j^* \propto e^{-t/\tau_j}$, where $\tau_j = -\log(1 - \frac{\lambda_j}{\lambda_1})$ is the time constant (number of iterations) for the mode j . Top modes f_j^* of the target function are learned faster than low modes.

In linearized regimes where deep learning reduces to kernel regression (Jacot et al., 2018; Du et al., 2019; Allen-Zhu et al., 2019), one can dwell further the nature of such a bias by analyzing the eigenfunctions of the neural tangent kernel (e.g., Yang & Salman, 2019). As a simple example, for a randomly initialized MLP on 1D uniform data, Fig. 8 shows the Fourier decomposition of such eigenfunctions, ranked in nonincreasing order of the eigenvalues. We observe that eigenfunctions with increasing index j (hence decreasing eigenvalues) correspond to modes with increasing Fourier frequency, with a remarkable alignment with Fourier modes for the first half of the spectrum. This is in line with observations (e.g., Rahaman et al., 2019) that deep networks tend to prioritize learning low frequency modes during training.

B Complexity Bounds

In this section, we spell out details and proofs for the content of Section 4.

B.1 Rademacher Complexity

Given a family $\mathcal{G} \subset \mathbb{R}^{\mathcal{Z}}$ of real-valued functions on a probability space (\mathcal{Z}, ρ) , the empirical Rademacher complexity of \mathcal{G} with respect to a sample $\mathcal{S} = \{\mathbf{z}_1, \dots, \mathbf{z}_n\} \sim \rho^n$ is defined as (Mohri et al., 2012):

$$\widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{G}) = \mathbb{E}_{\sigma \in \{\pm 1\}^n} \left[\sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i g(\mathbf{z}_i) \right], \quad (52)$$

where the expectation is over n i.i.d uniform random variables $\sigma_1, \dots, \sigma_n \in \{\pm 1\}$. For any $n \geq 1$, the Rademacher complexity with respect to samples of size n is then $\mathcal{R}_n(\mathcal{G}) = \mathbb{E}_{\mathcal{S} \sim \rho^n} \widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{G})$.

B.2 Generalization Bounds

Generalization bounds based on Rademacher complexity are standard (Bartlett et al., 2017; Mohri et al., 2012). We give here one instance of such a bound, relevant for classification task.

Setup. We consider a family \mathcal{F} of functions $f_{\mathbf{w}}: \mathcal{X} \rightarrow \mathbb{R}^c$ that output a score or probability $f_{\mathbf{w}}(\mathbf{x})[y]$ for each class $y \in \{1 \dots c\}$ (we take $c = 1$ for binary classification). The task is to find a predictor $f_{\mathbf{w}} \in \mathcal{F}$ with small expected classification error, which can be expressed e.g. as

$$L_0(f_{\mathbf{w}}) = \mathbb{P}_{(\mathbf{x}, y) \sim \rho} \{ \mu(f_{\mathbf{w}}(\mathbf{x}), y) < 0 \} \quad (53)$$

where $\mu(f(\mathbf{x}), y)$ denotes the **margin**,

$$\mu(f(\mathbf{x}), y) = \begin{cases} f(\mathbf{x})y & \text{binary case} \\ f(\mathbf{x})[y] - \max_{y' \neq y} f(\mathbf{x})[y'] & \text{multiclass case} \end{cases} \quad (54)$$

Margin Bound. We consider the **margin loss**,

$$\ell_{\gamma}(f_{\mathbf{w}}(\mathbf{x}), y) = \phi_{\gamma}(\mu(f_{\mathbf{w}}(\mathbf{x}), y)) \quad (55)$$

where $\gamma > 0$, and ϕ_{γ} is the **ramp** function: $\phi_{\gamma}(u) = 1$ if $u \leq 0$, $\phi(u) = 0$ if $u > \gamma$ and $\phi(u) = 1 - u/\gamma$ otherwise. We have the following bound for the expected error (53). With probability at least $1 - \delta$ over the draw $\mathcal{S} = \{\mathbf{z}_i = (\mathbf{x}_i, y_i)\}_{i=1}^n$ of size n , the following holds for all $f_{\mathbf{w}} \in \mathcal{F}$ (Mohri et al., 2012, Theorems 4.4. and 8.1):

$$L_0(f_{\mathbf{w}}) \leq \widehat{L}_{\gamma}(f_{\mathbf{w}}) + 2\widehat{\mathcal{R}}_{\mathcal{S}}(\ell_{\gamma}(\mathcal{F}, \cdot)) + 3\sqrt{\frac{\log \frac{2}{\delta}}{2n}} \quad (56)$$

where $\widehat{L}_{\gamma}(f_{\mathbf{w}}) = \frac{1}{n} \sum_{i=1}^n \ell_{\gamma}(f_{\mathbf{w}}(\mathbf{x}_i), y_i)$ is the empirical margin error and $\ell_{\gamma}(\mathcal{F}, \cdot)$ is the **loss class**,

$$\ell_{\gamma}(\mathcal{F}, \cdot) = \{(\mathbf{x}, y) \mapsto \ell_{\gamma}(f_{\mathbf{w}}(\mathbf{x}), y) \mid f_{\mathbf{w}} \in \mathcal{F}\} \quad (57)$$

For binary classifiers, because ϕ_{γ} is $1/\gamma$ -Lipschitz, we have in addition

$$\mathcal{R}_{\mathcal{S}}(\ell_{\gamma}(\mathcal{F}, \cdot)) \leq \frac{1}{\gamma} \mathcal{R}_{\mathcal{S}}(\mathcal{F}) \quad (58)$$

by Talagrand's contraction lemma (Ledoux & Talagrand, 2013) (see e.g. Mohri et al., 2012, lemma 4.2 for a detailed proof).

B.3 Complexity Bounds: Proofs

We first derive standard bounds for the linear classes of scalar functions,

$$\mathcal{F}_{M_A}^A = \{f_{\mathbf{w}}: \mathbf{x} \mapsto \langle \mathbf{w}, \Phi(\mathbf{x}) \rangle \mid \|\mathbf{w}\|_A \leq M_A\} \quad (59)$$

Proposition 10. *The empirical Rademacher complexity of $\mathcal{F}_{M_A}^A$ is bounded as,*

$$\widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_{M_A}^A) \leq (M_A/n) \sqrt{\text{Tr} \mathbf{K}_A} \quad (60)$$

where $(\mathbf{K}_A)_{ij} = k_A(\mathbf{x}_i, \mathbf{x}_j)$ is the kernel matrix associated to the rescaled features $A^{-1}\Phi$.

Proof. We use the notation of Section 4. For given Rademacher variables $\boldsymbol{\sigma} \in \{\pm 1\}^n$, we have,

$$\begin{aligned}
 \sup_{f \in \mathcal{F}_{M_A}^A} \sum_{i=1}^n \sigma_i f(\mathbf{x}_i) &= \sup_{\|\mathbf{w}\|_A \leq M_A} \sum_{i=1}^n \sigma_i \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle \\
 &= \sup_{\|A^\top \mathbf{w}\|_2 \leq M_A} \sum_{i=1}^n \sigma_i \langle A^\top \mathbf{w}, A^{-1} \Phi(\mathbf{x}_i) \rangle \\
 &= \sup_{\|\tilde{\mathbf{w}}\|_2 \leq M_A} \langle \tilde{\mathbf{w}}, \sum_{i=1}^n \sigma_i A^{-1} \Phi(\mathbf{x}_i) \rangle \\
 &= M_A \left\| \sum_{i=1}^n \sigma_i A^{-1} \Phi(\mathbf{x}_i) \right\|_2 \\
 &= M_A \sqrt{\boldsymbol{\sigma}^\top \mathbf{K}_A \boldsymbol{\sigma}}
 \end{aligned} \tag{61}$$

From (61) and the definition (52) we obtain:

$$\begin{aligned}
 \widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_{M_A}^A) &= \frac{M_A}{n} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sqrt{\boldsymbol{\sigma}^\top \mathbf{K}_A \boldsymbol{\sigma}} \right] \\
 &\leq \frac{M_A}{n} \sqrt{\mathbb{E}_{\boldsymbol{\sigma}} [\boldsymbol{\sigma}^\top \mathbf{K}_A \boldsymbol{\sigma}]} \\
 &\leq \frac{M_A}{n} \sqrt{\text{Tr} \mathbf{K}_A}
 \end{aligned} \tag{62}$$

where we used Jensen's inequality to pass $\mathbb{E}_{\boldsymbol{\sigma}}$ under the root, and that $\mathbb{E}[\sigma_i] = 0$ and $\sigma_i^2 = 1$ for all i . \square

We now extend the result to the families (11) of learning flows:

$$\mathcal{F}_{m_t}^A = \{f_{\mathbf{w}}: \mathbf{x} \mapsto \sum_t \langle \delta \mathbf{w}_t, \Phi(\mathbf{x}) \rangle \mid \|\delta \mathbf{w}_t\|_{A_t} \leq m_t\} \tag{63}$$

Theorem 11 (Theorem 2 restated). *The empirical Rademacher complexity of $\mathcal{F}_{m_t}^A$ is bounded as,*

$$\widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_{m_t}^A) \leq \sum_t (m_t/n) \sqrt{\text{Tr} \mathbf{K}_{A_t}} \tag{64}$$

where $(\mathbf{K}_{A_t})_{ij} = k_{A_t}(\mathbf{x}_i, \mathbf{x}_j)$ is the kernel matrix associated to the rescaled features $A_t^{-1} \Phi$.

Proof. This is simple extension of the previous proof:

$$\begin{aligned}
 \sup_{f \in \mathcal{F}_{m_t}^A} \sum_{i=1}^n \sigma_i f(\mathbf{x}_i) &= \sup_{\|\delta \mathbf{w}_t\|_{A_t} \leq m_t} \sum_{i=1}^n \sigma_i \sum_t \langle \delta \mathbf{w}_t, \Phi(\mathbf{x}_i) \rangle \\
 &= \sum_t \sup_{\|\tilde{\delta \mathbf{w}}_t\|_2 \leq m_t} \langle \tilde{\delta \mathbf{w}}_t, \sum_{i=1}^n \sigma_i A_t^{-1} \Phi(\mathbf{x}_i) \rangle \\
 &= \sum_t m_t \sqrt{\boldsymbol{\sigma}^\top \mathbf{K}_{A_t} \boldsymbol{\sigma}}
 \end{aligned} \tag{65}$$

and we conclude as in (62). \square

Finally, we note that the same result can be formulated in terms of an evolving feature map $\Phi_t = A_t^{-1} \Phi$ with kernel $k_t(\mathbf{x}, \tilde{\mathbf{x}}) = \langle \Phi_t(\mathbf{x}), \Phi_t(\tilde{\mathbf{x}}) \rangle$. In fact by reparametrization invariance, the function updates can also be written as $\delta f_{\mathbf{w}_t}(\mathbf{x}) = \langle \tilde{\delta \mathbf{w}}_t, \Phi_t(\mathbf{x}) \rangle$ where $\tilde{\delta \mathbf{w}}_t = A_t^{-1} \delta \mathbf{w}_t$. The function class (11) can equivalently be written as $\mathcal{F}_{m_t}^A = \mathcal{F}_{m_t}^{\Phi}$ where Φ denotes a fixed sequence of feature maps, $\Phi = \{\Phi_t\}_t$ and

$$\mathcal{F}_{m_t}^{\Phi} = \{f_{\mathbf{w}}: \mathbf{x} \mapsto \sum_t \langle \tilde{\delta \mathbf{w}}_t, \Phi_t(\mathbf{x}) \rangle \mid \|\tilde{\delta \mathbf{w}}_t\|_2 \leq m_t\} \tag{66}$$

In this formulation, the result (64) is expressed as,

$$\widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_{m_t}^{\Phi}) \leq \sum_t (m_t/n) \sqrt{\text{Tr} \mathbf{K}_t} \tag{67}$$

where $(\mathbf{K}_t)_{ij} = k_t(\mathbf{x}_i, \tilde{\mathbf{x}}_j)$ is the kernel matrix associated to the feature map Φ_t .

B.4 Bounds for Multiclass Classification

The generalization bound (56) is based on the **margin loss class** (57). In this section, we show how to bound $\widehat{\mathcal{R}}_{\mathcal{S}}(\ell_{\gamma}(\mathcal{F}, \cdot))$ in terms of tangent kernels for the original class \mathcal{F} of functions $f_{\mathbf{w}}: \mathcal{X} \rightarrow \mathbb{R}^c$ instead. Although the proof is adapted from standard techniques, to our knowledge Lemma 12 and Theorem 13 below are new results. In what follows, we denote by $\mu_{\mathcal{F}}$ the margin class,

$$\mu_{\mathcal{F}} = \{(\mathbf{x}, y) \rightarrow \mu(f_{\mathbf{w}}(\mathbf{x}), y) \mid f_{\mathbf{w}} \in \mathcal{F}\} \quad (68)$$

where $\mu(f_{\mathbf{w}}(\mathbf{x}), y)$ is the margin (54). We also define, for each $y \in \{1 \cdots c\}$,

$$\mathcal{F}_y = \{\mathbf{x} \mapsto f_{\mathbf{w}}(\mathbf{x})[y] \mid f_{\mathbf{w}} \in \mathcal{F}\}, \quad \mu_{\mathcal{F}, y} = \{\mathbf{x} \mapsto \mu(f_{\mathbf{w}}(\mathbf{x}), y) \mid f_{\mathbf{w}} \in \mathcal{F}\} \quad (69)$$

Lemma 12. *The following inequality holds:*

$$\widehat{\mathcal{R}}_{\mathcal{S}}(\ell_{\gamma}(\mathcal{F}, \cdot)) \leq \frac{c}{\gamma} \sum_{y=1}^c \widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_y) \quad (70)$$

Proof. We first follow the first steps of the proof of (Mohri et al., 2012, Theorem 8.1) to show that

$$\widehat{\mathcal{R}}_{\mathcal{S}}(\ell_{\gamma}(\mathcal{F}, \cdot)) \leq \frac{1}{\gamma} \sum_{y=1}^c \widehat{\mathcal{R}}_{\mathcal{S}}(\mu_{\mathcal{F}, y}) \quad (71)$$

We reproduce these steps here for completeness: first, it follows from the $1/\gamma$ -Lipschitzness of the ramp loss ϕ_{γ} in (55) and Talagrand's contraction lemma (Mohri et al., 2012, lemma 4.2) that

$$\widehat{\mathcal{R}}_{\mathcal{S}}(\ell_{\gamma}(\mathcal{F}, \cdot)) \leq \frac{1}{\gamma} \widehat{\mathcal{R}}_{\mathcal{S}}(\mu_{\mathcal{F}}) \quad (72)$$

Next, we write

$$\begin{aligned} \widehat{\mathcal{R}}_{\mathcal{S}}(\mu_{\mathcal{F}}) &:= \frac{1}{n} \mathbb{E}_{\sigma} \left[\sup_{f_{\mathbf{w}} \in \mathcal{F}} \sum_{i=1}^n \sigma_i \mu(f_{\mathbf{w}}(\mathbf{x}_i), y_i) \right] \\ &= \frac{1}{n} \mathbb{E}_{\sigma} \left[\sup_{f_{\mathbf{w}} \in \mathcal{F}} \sum_{i=1}^n \sigma_i \sum_{y=1}^c \mu(f_{\mathbf{w}}(\mathbf{x}_i), y) \delta_{y, y_i} \right] \\ &= \frac{1}{n} \sum_{y=1}^c \mathbb{E}_{\sigma} \left[\sup_{f_{\mathbf{w}} \in \mathcal{F}} \sum_{i=1}^n \sigma_i \mu(f_{\mathbf{w}}(\mathbf{x}_i), y) \delta_{y, y_i} \right] \end{aligned} \quad (73)$$

where $\delta_{y, y_i} = 1$ if $y = y_i$ and 0 otherwise; the second inequality follows from the sub-additivity of sup. Substituting $\delta_{y, y_i} = \frac{1}{2}(\epsilon_i + \frac{1}{2})$ where $\epsilon_i = 2\delta_{y, y_i} - 1 \in \{\pm 1\}$, we obtain

$$\begin{aligned} \widehat{\mathcal{R}}_{\mathcal{S}}(\mu_{\mathcal{F}}) &\leq \frac{1}{2n} \sum_{y=1}^c \mathbb{E}_{\sigma} \left[\sup_{f_{\mathbf{w}} \in \mathcal{F}} \sum_{i=1}^n (\epsilon_i \sigma_i) \mu(f_{\mathbf{w}}(\mathbf{x}_i), y) \right] + \frac{1}{2n} \sum_{y=1}^c \mathbb{E}_{\sigma} \left[\sup_{f_{\mathbf{w}} \in \mathcal{F}} \sum_{i=1}^n \sigma_i \mu(f_{\mathbf{w}}(\mathbf{x}_i), y) \right] \\ &= \sum_{y=1}^c \frac{1}{n} \mathbb{E}_{\sigma} \left[\sup_{f_{\mathbf{w}} \in \mathcal{F}} \sum_{i=1}^n \sigma_i \mu(f_{\mathbf{w}}(\mathbf{x}_i), y) \right] \\ &= \sum_{y=1}^c \widehat{\mathcal{R}}_{\mathcal{S}}(\mu_{\mathcal{F}, y}) \end{aligned} \quad (74)$$

Together with (72), this leads to (71).

Now, spelling out $\mu(f_{\mathbf{w}}(\mathbf{x}_i, y))$ gives

$$\begin{aligned}
 \widehat{\mathcal{R}}_{\mathcal{S}}(\mu_{\mathcal{F}, y}) &= \frac{1}{n} \mathbb{E}_{\sigma} \left[\sup_{f_{\mathbf{w}} \in \mathcal{F}} \sum_{i=1}^n \sigma_i (f_{\mathbf{w}}(\mathbf{x}_i)[y] - \max_{y' \neq y} f_{\mathbf{w}}(\mathbf{x}_i)[y']) \right] \\
 &= \widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_y) + \frac{1}{n} \mathbb{E}_{\sigma} \left[\sup_{f_{\mathbf{w}} \in \mathcal{F}} \sum_{i=1}^n (-\sigma_i) \max_{y' \neq y} f_{\mathbf{w}}(\mathbf{x}_i)[y'] \right] \\
 &= \widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_y) + \frac{1}{n} \mathbb{E}_{\sigma} \left[\sup_{f_{\mathbf{w}} \in \mathcal{F}} \sum_{i=1}^n \sigma_i \max_{y' \neq y} f_{\mathbf{w}}(\mathbf{x}_i)[y'] \right] \\
 &\leq \widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_y) + \widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{G}_y)
 \end{aligned} \tag{75}$$

where $\mathcal{G}_y = \{\max\{f_{y'} : y' \neq y\} \mid f_{y'} \in \mathcal{F}_{y'}\}$. Now (Mohri et al., 2012, lemma 8.1) show that $\widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{G}_y) \leq \sum_{y' \neq y} \widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_{y'})$. This leads to

$$\begin{aligned}
 \sum_{y=1}^c \widehat{\mathcal{R}}_{\mathcal{S}}(\mu_{\mathcal{F}, y}) &\leq \sum_{y=1}^c \widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_y) + \sum_{y=1}^c \sum_{\substack{y'=1 \\ y' \neq y}}^c \widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_{y'}) \\
 &= \sum_{y=1}^c \widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_y) + (c-1) \sum_{y=1}^c \widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_y) \\
 &= c \sum_{y=1}^c \widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_y)
 \end{aligned} \tag{76}$$

Substituting in (71) finishes the proof. \square

In the linear case, this results leads to analogous theorems as in B.3 in the multiclass setting. For example, considering the linear families of functions $\mathcal{X} \rightarrow \mathbb{R}^c$,

$$\mathcal{F}_{M_A}^A = \{\mathbf{x} \mapsto f_{\mathbf{w}}(\mathbf{x})[y] := \langle \mathbf{w}, \Phi(\mathbf{x})[y] \rangle \mid \|\mathbf{w}\|_A \leq M_A\} \tag{77}$$

where $(\mathbf{x}, y) \mapsto \Phi(\mathbf{x})[y]$ is some joint feature map, we have the following

Theorem 13. *The emp. Rademacher complexity of the margin loss class $\ell_{\gamma}(\mathcal{F}_{M_A}^A, \cdot)$ is bounded as,*

$$\widehat{\mathcal{R}}_{\mathcal{S}}(\ell_{\gamma}(\mathcal{F}_{M_A}^A, \cdot)) \leq (c^{3/2} M_A / \gamma n) \sqrt{\text{Tr} \mathbf{K}_A} \tag{78}$$

where $(\mathbf{K}_A)_{ij}^{yy'}$ is the kernel $nc \times nc$ matrix associated to the rescaled features $A^{-1} \Phi(\mathbf{x})[y]$.

Proof. Eq.70, and Theorem 13 applied to each linear family \mathcal{F}_y of (scalar) functions leads to

$$\widehat{\mathcal{R}}_{\mathcal{S}}(\ell_{\gamma}(\mathcal{F}_{M_A}^A, \cdot)) \leq \frac{c}{\gamma} \sum_{y=1}^c \frac{M_A}{n} \sqrt{\text{Tr} \mathbf{K}_A^{yy}} \tag{79}$$

where $\text{Tr} \mathbf{K}_A^{yy} := \sum_{i=1}^n (\mathbf{K}_A)_{ii}^{yy}$ is computed w.r.t to the indices $i = 1, \dots, n$ for fixed y . Passing the average $\frac{1}{c} \sum_{y=1}^c$ under the root using Jensen inequality, we conclude:

$$\begin{aligned}
 \widehat{\mathcal{R}}_{\mathcal{S}}(\ell_{\gamma}(\mathcal{F}_{M_A}^A, \cdot)) &\leq \frac{c^2 M_A}{\gamma n} \sqrt{\frac{1}{c} \sum_{y=1}^c \text{Tr} \mathbf{K}_A^{yy}} \\
 &= \frac{c^{3/2} M_A}{\gamma n} \sqrt{\text{Tr} \mathbf{K}_A}
 \end{aligned} \tag{80}$$

\square

The proof of the extension of these bounds to families learning flows follows the same line as in B.3.

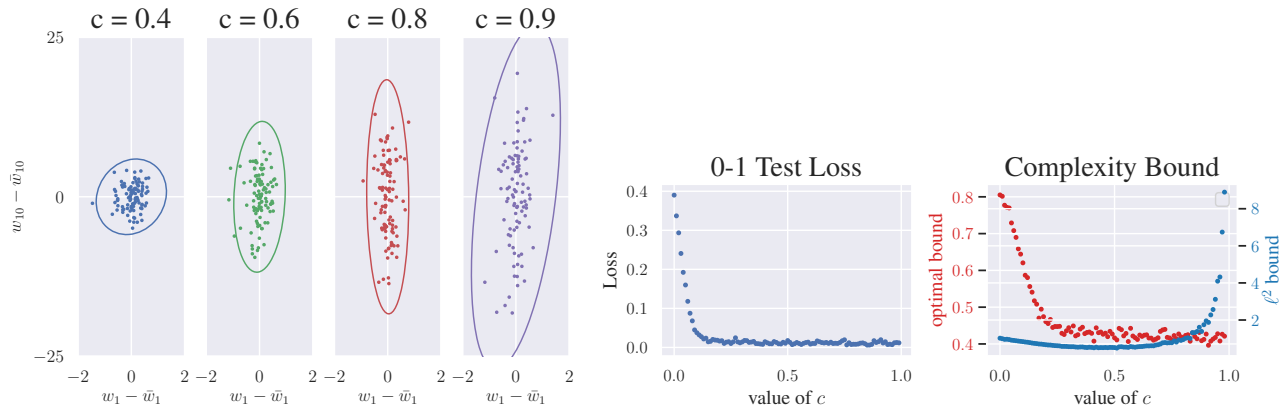


Figure 9: **Left:** 2D projection of the minimum ℓ^2 -norm interpolators $\mathbf{w}_{\mathcal{S}}^*$, $\mathcal{S} \sim \rho^n$, for linear models $f_{\mathbf{w}} = \langle \mathbf{w}, \Phi_c \rangle$, as the feature scaling factor varies from 0 (white features) to 1 (original, anisotropic features). For larger c , the solutions scatter in a very anisotropic way. **Right:** Average test classification loss and complexity bounds (78) with $A = \text{Id}$ (blue plot) for the solution vectors $\mathbf{w}_{\mathcal{S}}^*$, as we increase the scaling factor c . As feature anisotropy increases, the bound becomes increasingly loose and fails to reflect the shape of the test error. By contrast, the bound (10) with A optimized as in Proposition 14 (red plot) does not suffer from this problem.

B.5 Which Norm for Measuring Capacity?

Implicit biases of gradient descent are relatively well understood in linear models (e.g. Gunasekar et al. (2018)). For example when using square loss, it is well-known that gradient descent (initialized in the span of the data) converges to minimum ℓ^2 norm (resp. RKHS norm) solutions in parameter space (resp. function space). Yet, as pointed out by Belkin et al. (2018); Muthukumar et al. (2020), measuring capacity in terms of such norms is not coherently linked with generalization in practice. Here we discuss this issue by highlighting the critical dependence of meaningful norm-based capacity on the geometry defined by the features. We use the notation of Section 4.1: $\Phi = \sum_{j=1}^n \sqrt{\lambda_j} \mathbf{u}_j \mathbf{v}_j^\top$ denote the $n \times P$ feature matrix and its SVD decomposition.

A standard approach is to measure capacity in terms of the ℓ^2 norm the weight vector, e.g using bounds (10) with $A = \text{Id}$. If the distribution of solutions $\mathbf{w}_{\mathcal{S}}^*$, where $\mathcal{S} \sim \rho^n$ is sampled from the input distribution, is reasonably isotropic, taking the smallest ℓ^2 ball containing them (with high probability) gives an accurate description of the class of trained models. However for very anisotropic distributions, the solutions do not fill any such ball so describing trained models in terms of ℓ^2 balls is wasteful (Schölkopf et al., 1999a).

Now, for minimum ℓ^2 norm interpolators (Hastie et al., 2009),

$$\mathbf{w}^* = \Phi^\top \mathbf{K}^{-1} \mathbf{y} = \sum_{j=1}^n \frac{\mathbf{u}_j^\top \mathbf{y}}{\sqrt{\lambda_j}} \mathbf{v}_j, \quad (81)$$

where $\mathbf{K} = \Phi \Phi^\top$ is the kernel matrix, the solution distribution typically inherits the anisotropy of the features. For example, if $y_i = \bar{y}(\mathbf{x}_i) + \varepsilon_i$ where $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$, the covariance of the solutions with respect to noise is $\text{cov}_\varepsilon[\mathbf{w}^*, \mathbf{w}^*] = \sum_j \frac{\sigma^2}{\lambda_j} \mathbf{v}_j \mathbf{v}_j^\top$, which scales as $1/\lambda_j$ along \mathbf{v}_j .

To visualize this on a simple setting, we consider P random features of a RBF kernel¹², fit on 1D data \mathbf{x} modelled by N equally spaced points in $[-a, a]$. In this setting, the (true) feature map is represented by a $N \times P$ matrix with SVD $\Phi = \sum_j \sqrt{l_j} \psi_j \varphi_j^\top$. We assume the (true) labels are defined by the deterministic function $y(\mathbf{x}) = \text{sign}(\psi_1(\mathbf{x}))$. To highlight the effect of feature anisotropy, we further rescale the singular values as $l_j^c = 1 + c(l_j - 1)$ so as to interpolate between whitened features ($c=0$) and the original ones ($c=1$). We set $P=N=1000$. Fig 9 (left) shows 2D projections in the plane $(\varphi_1, \varphi_{10})$ of the (centered) minimum ℓ_2 norm solutions $\mathbf{w}_{\mathcal{S}}^* - \mathbb{E}_{\mathcal{S}} \mathbf{w}_{\mathcal{S}}^*$, for a pool of 100 training (sub)samples \mathcal{S} of size $n=50$, for increasing values of the scaling factor c . As c approaches 1, the solutions begin to scatter in a very anisotropic way in parameter space; as shown in Fig 9 (right), the complexity

¹²We used `RBFsampler` of `scikit-learn`, which implements a variant of Random Kitchen Sinks (Rahimi & Recht, 2007) to approximate the feature map of a RBF kernel with parameter $\gamma = 1$.

bound (60) based on the ℓ_2 norm, i.e $A = \text{Id}$ (blue plot), becomes increasingly loose and fails to reflect the shape of the test error.

To find a more meaningful capacity measure, Prop 13 suggests optimizing the bound (10) with $M_A = \|\mathbf{w}^*\|_A$, over a given class of rescaling matrices A . We give an example of this in the following Proposition.

Proposition 14. *Consider the class of matrices $A_\nu = \sum_{j=1}^n \sqrt{\nu_j} \mathbf{v}_j \mathbf{v}_j^\top + \text{Id}_{\text{span}\{\mathbf{v}\}^\perp}$, which act as mere rescaling of the singular values of the feature matrix. Any minimizer of the upper bound (60) for the minimum ℓ^2 -norm interpolator takes the form*

$$\nu_j^* = \kappa \frac{\sqrt{\lambda_j}}{|\mathbf{v}_j^\top \mathbf{w}^*|} = \kappa \frac{\lambda_j}{|\mathbf{u}_j^\top \mathbf{y}|} \quad (82)$$

where $\kappa > 0$ is a constant independent of j .

Proof. From (81) and the definition of A_ν , we first write

$$\|\mathbf{w}^*\|_{A_\nu}^2 = \sum_{j=1}^n \frac{\nu_j}{\lambda_j} (\mathbf{u}_j^\top \mathbf{y})^2, \quad \text{Tr} \mathbf{K}_{A_\nu} = \sum_{j=1}^n \frac{\lambda_j}{\nu_j} \quad (83)$$

The product of the above two terms has the critical points ν_j^* , $j = 1 \dots n$ which satisfy

$$\frac{(\mathbf{u}_j^\top \mathbf{y})^2}{\lambda_j} \text{Tr} \mathbf{K}_{A_\nu} - \frac{\lambda_j}{\nu_j^{*2}} \|\mathbf{w}^*\|_{A_\nu}^2 = 0 \quad (84)$$

giving the desired result $\nu_j^* \propto \lambda_j / |\mathbf{u}_j^\top \mathbf{y}|$. \square

In the context of Proposition 14, we see that the optimal norm $\|\cdot\|_{A_{\nu^*}}$ depends both on the feature geometry – through the singular values – and on the task – through the labels –. As shown in Fig 1 (right, red plot), in the above RBF feature setting, the resulting optimal bound on the Radecher complexity has a much nicer behaviour than the standard bound based on the ℓ^2 norm.¹³

B.6 SuperNat: Proof of Prop 3

Prop. 3 is a *local* version of Prop 14, where the feature rescaling factors are applied at each step of the training algorithm. The procedure is described in Fig 5 (left); the term to be optimized shows up in Step 2. With the chosen class of matrices described in Prop 3, the action $\Phi_t \rightarrow A_\nu^{-1} \Phi_t$ merely rescale its singular values $\lambda_{jt} \rightarrow \lambda_{jt} / \nu_j$, leaving its singular vectors $\mathbf{u}_j, \mathbf{v}_j$ unchanged.

Proposition 15 (Prop 3 restated). *For the class of rescaling matrices A_ν defined in Prop 14, any minimizer in Step 2 in Fig 5, where $\delta \mathbf{w}_{\text{GD}} = -\eta \nabla_{\mathbf{w}} L$, takes the form*

$$\nu_{jt}^* = \kappa \frac{1}{|\mathbf{u}_j^\top \nabla_{\mathbf{f}_w} L|} \quad (85)$$

where $\kappa > 0$ is a constant independent of j .

Proof. Using the chain rule and the SVD of the feature map Φ_t we write the gradient descent updates at iteration t of SuperNat as

$$\delta \mathbf{w}_{\text{GD}} = -\eta \Phi_t^\top \nabla_{\mathbf{f}_w} L \quad (86)$$

$$= -\eta \sum_{j=1}^n \sqrt{\lambda_{jt}} (\mathbf{u}_j^\top \nabla_{\mathbf{f}_w} L) \mathbf{v}_j, \quad (87)$$

¹³Note however that, since the optimal norm depends on the sample set \mathcal{S} , the resulting complexity bound does not directly yield a high probability bound on the generalization error as in (56). The more thorough analysis, which requires promoting (56) to uniform bounds over the choice of matrix A , is left for future work.

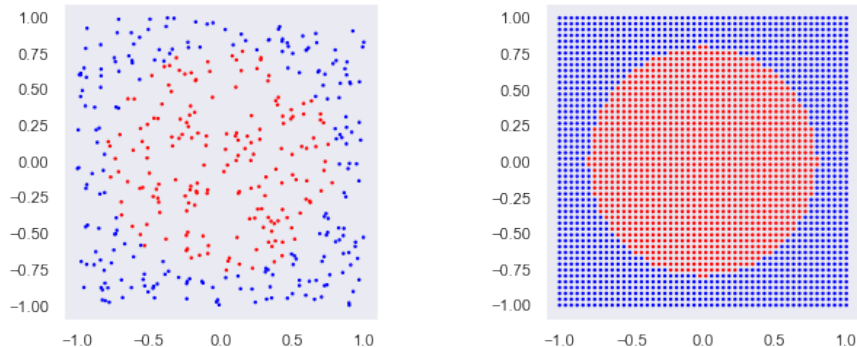


Figure 10: Disk dataset. **Left:** Training set of $n = 500$ points (\mathbf{x}_i, y_i) where $\mathbf{x} \sim \text{Unif}[-1, 1]^2$, $y_i = 1$ if $\|\mathbf{x}_i\|_2 \leq r = \sqrt{2/\pi}$ and -1 otherwise. **Right:** Large test sample (2500 points forming a 50×50 grid) used to evaluate the tangent kernel.

From the definition of A_ν , we then spell out

$$\|\delta \mathbf{w}_{\text{GD}}\|_{A_\nu}^2 = \eta^2 \sum_{j=1}^n (\nu_j \lambda_j) (\mathbf{u}_j^\top \nabla_{\mathbf{f}_w} L)^2, \quad \|A_\nu^{-1} \Phi_t\|_F := \text{Tr} \mathbf{K}_{t, A_\nu} = \sum_{j=1}^n \frac{\lambda_j}{\nu_j} \quad (88)$$

The product of the above two terms has the critical points ν_j^* , $j = 1 \cdots n$ which satisfy

$$\lambda_j (\mathbf{u}_j^\top \nabla_{\mathbf{f}_w} L)^2 \text{Tr} \mathbf{K}_{A_\nu} - \frac{\lambda_j}{\nu_j^{*2}} \|\delta \mathbf{w}_{\text{GD}}\|_{A_\nu}^2 = 0 \quad (89)$$

giving the desired result $\nu_j^* \propto 1/|\mathbf{u}_j^\top \nabla_{\mathbf{f}_w} L|$. \square

C Additional experiments

C.1 Synthetic Experiment: Fig. 1

To visualize the adaptation of the tangent kernel to the task during training, we perform the following synthetic experiment. We train a 6-layer deep 256-unit wide MLP on $n = 500$ points of the Disc dataset (\mathbf{x}, y) where $\mathbf{x} \sim \text{Unif}[-1, 1]^2$ and $y(\mathbf{x}) = \pm 1$ depending on whether is within the disk of center 0 and radius $\sqrt{2/\pi}$, see Fig 10. Fig. 1 in the main text shows visualizations of eigenfunctions sampled using a grid of $N = 2500$ points on the square, and ranked in non-increasing order of the spectrum $\lambda_1 \geq \cdots \geq \lambda_N$. After a number of iterations, we begin to see the class structure (e.g. boundary circle) emerge in the top eigenfunctions. We note also an increasingly fast spectrum decay (e.g. $\lambda_{20}/\lambda_1 = 1.5\%$ at iteration 0 and 0.2% at iteration 2000). The interpretation is that the kernel stretches in directions of high correlation with the labels.

C.2 More Alignment Plots

Varying datasets and architectures: Fig 11.

Uncentered kernel Experiments: Fig 12. The evolution of the alignment to the *uncentered* kernel, in order to assess whether this effect is consistent when removing centering. The experimental details are the same as in the main text; we also observe a similar increase of the alignment as training progresses.

C.3 Effect of depth on alignment

In order to study the influence of the architecture on the alignment effect, we measure the CKA for different networks and different initialization as we increase the depth. The results in Fig 13 suggest that the alignment effect is magnified as depth increases. We also observe that the ratio of the maximum alignment between easy and difficult examples is increased with depth, but stays high for a smaller number of iterations.

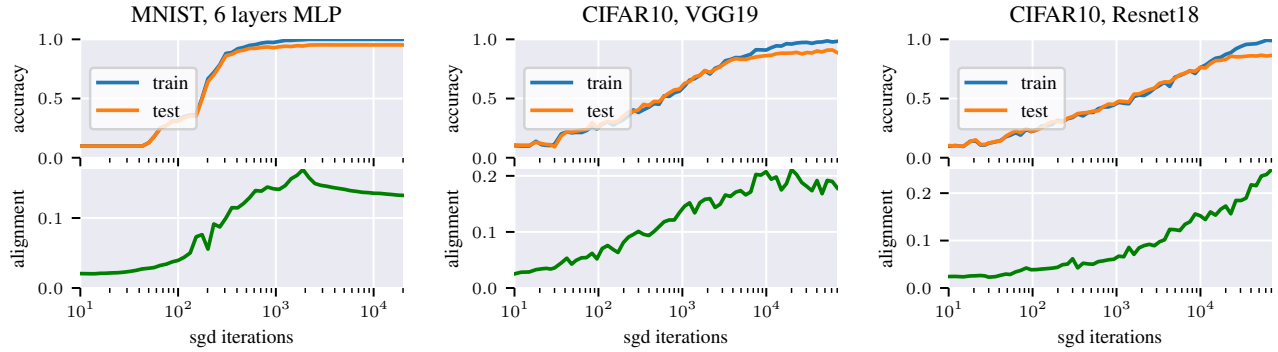


Figure 11: Evolution of the CKA between the tangent kernel and the class label kernel $K_Y = YY^T$ measured on a held-out test set for different architectures: **(left)** 6 layers of 80 hidden units MLP on MNIST **(middle)** VGG19 on CIFAR10 **(right)** Resnet18 on CIFAR10. We observe an increase of the alignment to the target function.

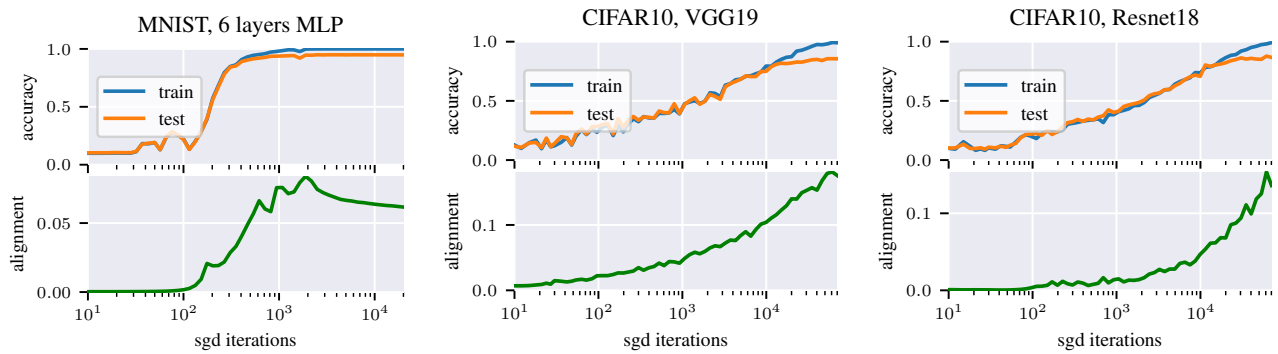


Figure 12: Same as figure 11 but without centering the kernel. Evolution of the uncentered kernel alignment between the tangent kernel and the class label kernel $K_Y = YY^T$ measured on a held-out test set for different architectures: **(left)** 6 layers of 80 hidden units MLP on MNIST **(middle)** VGG19 on CIFAR10 **(right)** Resnet18 on CIFAR10. We observe an increase of the alignment to the target function.

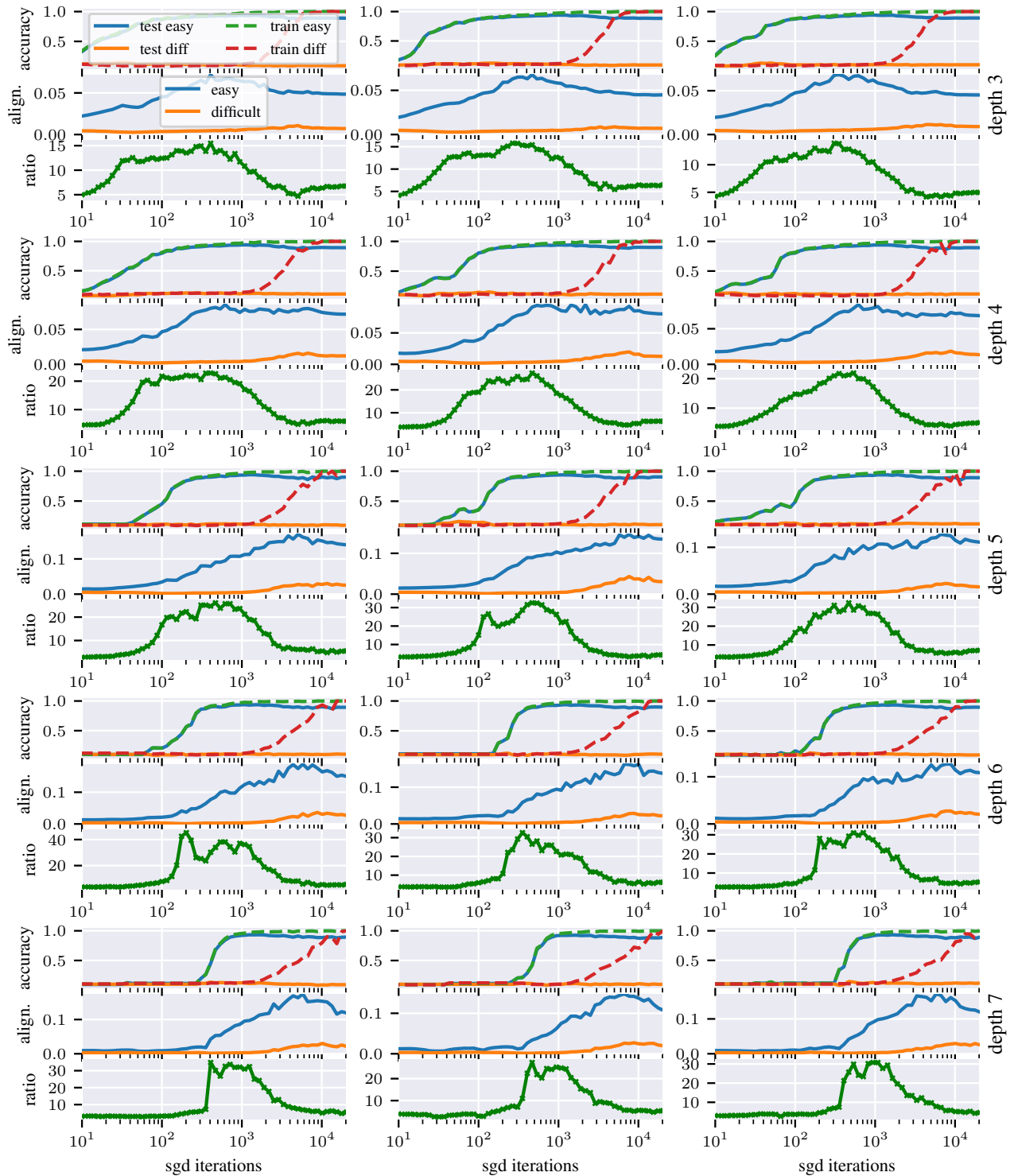


Figure 13: Effect of depth on alignment. 10,000 MNIST examples with 1000 random labels MNIST examples trained with learning rate=0.01, momentum=0.9 and batch size=100 for MLP with hidden layers size 60 and (in rows) varying depths (in columns) varying random initialization/minibatch sampling. As we increase the depth, the alignment starts increasing later in training and increases faster; and the ratio between easy and difficult alignments reaches a higher value.

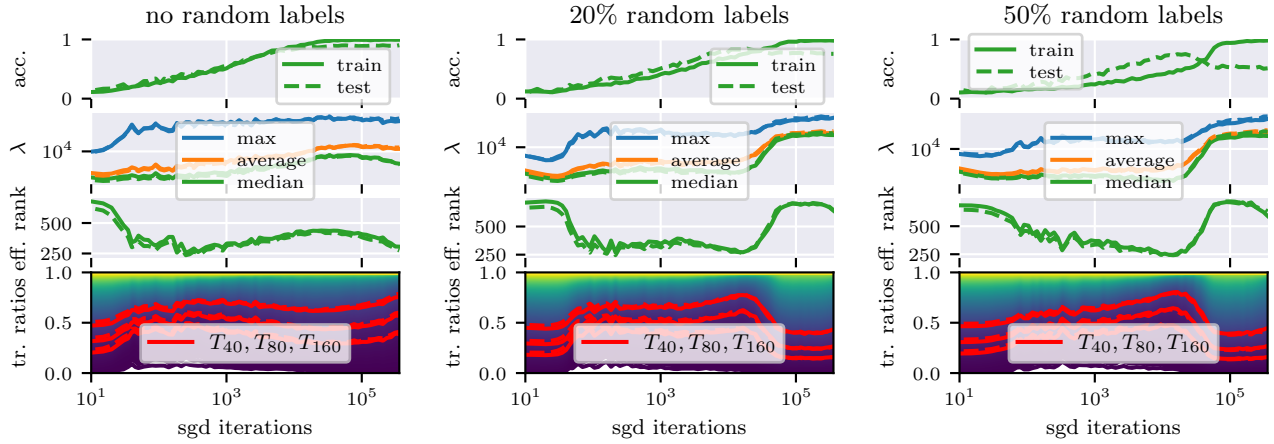


Figure 14: Evolution of tangent kernel spectrum, effective rank and trace ratios of a VGG19 trained by SGD with batch size 100, learning rate 0.003 and momentum 0.9 on dataset **(left)** CIFAR10 and **(right)** CIFAR10 with 50% random labels. We highlight the top 40, 80 and 160 trace ratios in red.

C.4 Spectrum Plots with lower learning rate : Fig. 14