
Implicit Regularization via Feature Alignment

1 One important property of deep neural networks is their ability to generalize well on real data.
2 Surprisingly, this is even true with very high-capacity networks *without explicit regularization* [42,
3 61, 29]. This seems at odds with the usual understanding of the bias-variance trade-off [24, 41, 11].
4 Solving this apparent paradox requires understanding the various learning biases induced by the
5 training procedure, which can act as implicit regularizers [42, 44].

6 In this paper, we help clarify one such implicit regularization mechanism, by examining the evolution
7 of the *neural tangent features* [30] learned by the network along the optimization paths. Our results
8 can be understood from two complementary perspectives: a *geometric* perspective – the (uncentered)
9 covariance of the tangent features defines a metric on the model function class, akin to the Fisher
10 information metric [e.g., 2]; and a *functional* perspective – through the tangent kernel and its RKHS.

11 Our main observation is a dynamical alignment of the tangent features along a small number of task-
12 relevant directions during training (Section 3), which can be interpreted as a combined *feature selection*
13 and *compression* mechanism. The motivating intuition is that such a mechanism allows the model to
14 adapt its capacity to the task and underpins the generalization abilities of heavily overparametrized
15 models. Drawing upon intuitions from linear models, we motivate a new heuristic complexity measure
16 which captures this phenomenon, and empirically show correlation with generalization (Section 4).

17 **Preliminaries** Let \mathcal{F} be a class of functions (e.g a neural network) parametrized by $\mathbf{w} \in \mathbb{R}^P$. We
18 restrict here to *scalar* functions $f_{\mathbf{w}}: \mathcal{X} \rightarrow \mathbb{R}$ to keep notation light.¹ We define the **tangent features**
19 as the function gradients w.r.t the parameters, $\Phi_{\mathbf{w}}(\mathbf{x}) := \nabla_{\mathbf{w}} f_{\mathbf{w}}(\mathbf{x})$, which govern how small changes
20 in parameter affect the function’s outputs,

$$\delta f_{\mathbf{w}}(\mathbf{x}) = \langle \delta \mathbf{w}, \Phi_{\mathbf{w}}(\mathbf{x}) \rangle + O(\|\delta \mathbf{w}\|^2) \quad (1)$$

21 More formally, the (uncentered) covariance matrix $g_{\mathbf{w}} = \mathbb{E}_{\mathbf{x} \sim \rho} [\Phi_{\mathbf{w}}(\mathbf{x}) \Phi_{\mathbf{w}}(\mathbf{x})^{\top}]$ acts as a **metric**
22 **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
23 product (see Longer Version, Appendix A). It characterizes the geometry of the function class \mathcal{F} .
24 Metric (as symmetric matrices) and tangent kernels (as integral operators) share the same spectrum.

25 The structure of the tangent features impacts the evolution of the function during training. Given
26 n input samples, consider gradient descent updates $\delta \mathbf{w}_{\text{GD}} = -\eta \nabla_{\mathbf{w}} L$ for some cost function L . The
27 function updates $\delta f_{\text{GD}}(\mathbf{x}) := \langle \delta \mathbf{w}_{\text{GD}}, \Phi_{\mathbf{w}}(\mathbf{x}) \rangle$ in the linear approximation (9), decompose as

$$\delta f_{\text{GD}}(\mathbf{x}) = \sum_{j=1}^P \delta f_j u_{\mathbf{w}j}(\mathbf{x}), \quad \delta f_j = -\eta \lambda_{\mathbf{w}j} (\mathbf{u}_{\mathbf{w}j}^{\top} \nabla_{\mathbf{f}_{\mathbf{w}}} L), \quad (2)$$

28 where $(u_{\mathbf{w}j})_{j=1}^P$ is the **eigenbasis** of the tangent kernel and $\mathbf{u}_{\mathbf{w}j} = [u_{\mathbf{w}j}(\mathbf{x}_1), \dots, u_{\mathbf{w}j}(\mathbf{x}_n)]^{\top}$. From
29 the point of view of function space, the metric/tangent kernel eigenvalues act as a mode-specific
30 rescaling $\eta \lambda_{\mathbf{w}j}$ of the learning rate. This is a local version of a well-known bias for linear models
31 (see Longer Version, Appendix B.2), towards functions in the top eigenspaces of the kernel.

32 As a first illustration of *non-linear* effects, Fig. 3 (Longer Version) shows visualizations of eigenfunc-
33 tions of the tangent kernel of a MLP trained on a simple classification task: $y(\mathbf{x}) = \pm 1$ depending on
34 whether $\mathbf{x} \sim \text{Unif}[-1, 1]^2$ is in the centered disk of radius $\sqrt{2/\pi}$. After a number of iterations, we
35 observe (rotation invariant) modes corresponding to the class structure (e.g boundary circle) showing
36 up in the top eigenfunctions of the learned kernel. We also note an increasing spectrum anisotropy –
37 for example, the ratio λ_{20}/λ_1 , which is 1.5% at iteration 0, has dropped to 0.2% at iteration 2000.
38 The interpretation is that the tangent kernel *stretches* in the directions of the signal during training.

¹See Appendix A for the extension to vector-valued functions, along with further mathematical details.

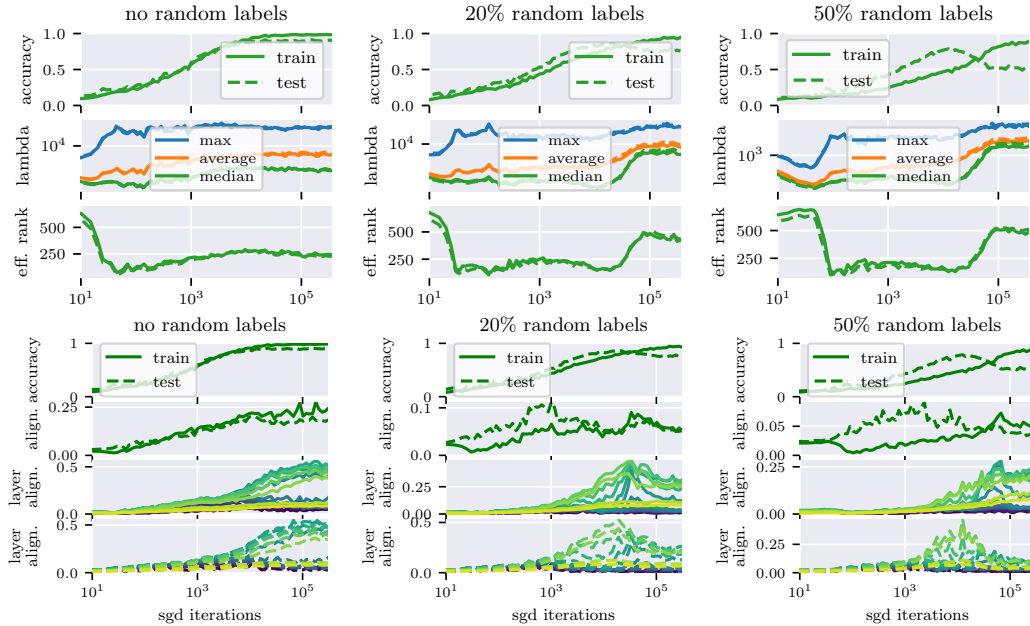


Figure 1: Evolution the *spectrum* and *effective rank* of the tangent kernel (**1st row**) and CKA and layer-wise CKA (**2nd row**) of a VGG19 on CIFAR10 with various ratios of random labels. For layer-wise alignment we map layers to colors sequentially from input layer (-), through intermediate layers (-), to output layer (-).

39 1 Neural Feature Alignment

40 We run experiments on MNIST [35] and CIFAR10 [33] with standard MLPs, VGG [55] and Resnet
 41 [28] architectures, using PyTorch [47] and NNGeometry [3] for efficient evaluation of tangent kernels.
 42 In multiclass settings, tangent kernels on n samples carry additional class indices $y \in \{1 \dots c\}$ and
 43 are treated as $nc \times nc$ matrices. We evaluate them on mini-batches (train or test) of size $n = 100$.

44 **Spectrum Evolution.** We report results (Fig. 1, 1st row) for tangent kernels evaluated on training
 45 examples (solid line) and test examples (dashed line). The main take away is an anisotropic increase
 46 of the kernel/metric spectrum during training. We quantify spectrum anisotropy through the various
 47 **trace ratios** $T_k = \sum_{j < k} \lambda_j / \sum_j \lambda_j$ as measures of the relative importance of the top k eigenvalues ;
 48 and using a notion of **effective rank** based on spectral entropy [50] (Longer Version, Appendix D).

49 We note an important decrease of the effective rank early in training, reaching a phase where only
 50 a few top eigenvalues account for most of the trace. This can be observed directly (Fig. 15) from
 51 the highlighted (in red) ratios T_{40} , T_{80} and T_{160} (Fig. 15), e.g. T_{80} accounting for 50% of the total
 52 trace (over 1000 eigenvalues). Remarkably, in the presence of high label noise, the effective rank
 53 of the tangent kernel evaluated on *training* examples (anti-)correlates nicely with the *test* accuracy,
 54 decreasing or remaining low during the learning phase and rising when overfitting starts. This
 55 suggests that the effective rank already provides a good proxy for the network’s effective capacity.

56 **Alignment to class labels.** We investigate the similarity of the tangent features with $\mathbf{Y} \in \mathbb{R}^{nc}$
 57 (concatenated one-hot vectors) through the **centered kernel alignment** (CKA) [19, 18] (Appendix
 58 D) $\text{CKA}(\mathbf{K}_w, \mathbf{K}_Y)$ with the rank-one kernel $\mathbf{K}_Y := \mathbf{Y}\mathbf{Y}^\top$. Intuitively, it is high when \mathbf{K}_w has low
 59 (effective) rank, and is such that the angle between \mathbf{Y} and its top eigenspaces is small.² Maximizing
 60 such an index has been used as a criterion for kernel selection in the literature on learning kernels [18].

61 We observe (Fig. 1, 2nd row) an increasingly high CKA as training progresses. The trend is similar
 62 for other architectures and datasets (Fig. 13 in Appendix E). Interestingly, in the presence of high
 63 level noise and during the learning phase, the CKA reaches a much higher value for *test* than for *train*
 64 kernels/labels (note that test labels are not randomized). Together with equation 11, this sheds lights
 65 on empirical observations that, in the presence of noise, deep networks ‘learn patterns first’ [5]

²In the limiting case $\text{CKA}(\mathbf{K}, \mathbf{K}_Y) = 1$, the features are all aligned with each other and parallel to \mathbf{Y} .

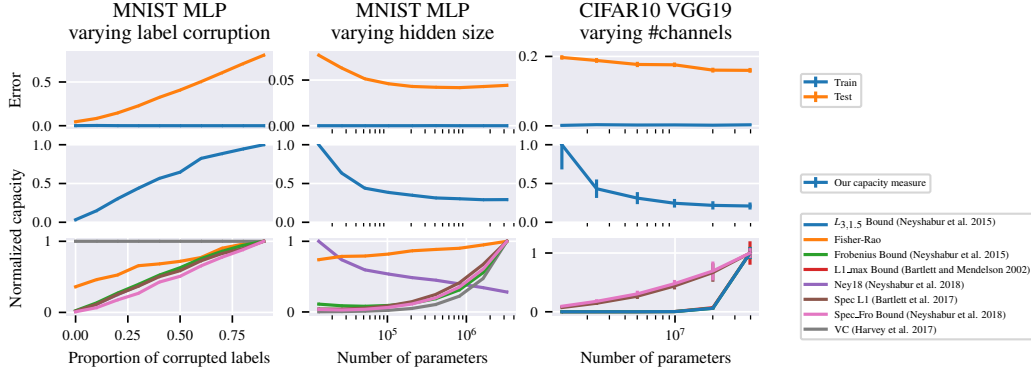


Figure 2: Normalized complexity measures on MNIST with a one hidden layer MLP (**left**) as we increase the hidden layer size, (**center**) for a fixed hidden layer of 256 units as we increase label corruption and (**right**) for a VGG19 on CIFAR10 as we vary the number of channels.

66 **Hierarchical Alignment.** A key aspect of the generalization question concerns the articulation
 67 between learning and memorization, in the presence of noise [61] or difficult examples [51]. In
 68 our next experiment, our setup is to augment 10,000 MNIST training examples with 1000 difficult
 69 examples of 2 types: (i) examples with random labels and (ii) examples from the dataset KMNIST
 70 [17]. Fig. 6 (Longer Version) shows that the (partial) CKA on the easy examples increases faster (and
 71 to a higher value) than that of the difficult examples. This suggests a hierarchy in the adaptation of the
 72 kernel, measured by the ratio between both alignments. This aspect of the non-linear dynamics favors
 73 a sequentialization of the learning (‘easy patterns first’) (see [52, 34, 25] for deep linear networks.)
 74 Fig 16 (Appendix E) shows that this effect is magnified as depth increases.

75 2 Measuring Complexity

76 2.1 Insights from Linear Models

77 **Setup.** We restrict here to functions $f_{\mathbf{w}}(\mathbf{x}) = \langle \mathbf{w}, \Phi(\mathbf{x}) \rangle$ linearly parametrized by $\mathbf{w} \in \mathbb{R}^P$. In this
 78 setting, (tangent) kernel and geometry are constant. Given n input samples, the n features $\Phi(\mathbf{x}_i) \in \mathbb{R}^P$
 79 yield a $n \times P$ feature matrix Φ . Our discussion is based on the **Rademacher complexity** showing
 80 up in generalization bounds [6]. It depends on the size (or **capacity**) of the function class.

81 A standard approach for controlling capacity is in terms of the *norm* of the weight vector – usually
 82 the ℓ_2 -norm. In general, given any invertible matrix $A \in \mathbb{R}^{P \times P}$, we may consider the norm
 83 $\|\mathbf{w}\|_A := \sqrt{\mathbf{w}^\top g_A \mathbf{w}}$ induced by the metric $g_A = AA^\top$. For $M_A > 0$, let $\mathcal{F}_{M_A}^A$ be the subclass of
 84 functions $f_{\mathbf{w}}$ such that $\|\mathbf{w}\|_A \leq M_A$. The Rademacher complexity can be bounded as,

$$\widehat{\mathcal{R}}(\mathcal{F}_{M_A}^A) \leq (M_A/n) \|A^{-1} \Phi^\top\|_F \quad (3)$$

85 in terms of the Frobenius norm of the *rescaled* feature matrix. This raises the question of which of
 86 the norms $\|\cdot\|_A$ provide meaningful capacity measures. Recent works [10, 40] pointed out that the
 87 ℓ_2 norm is not coherently linked with generalization in practice. We discuss this issue in Appendix
 88 C.5, illustrating how meaningful norms critically depend on the geometry defined by the features.

89 **Feature Alignment as Implicit Regularization.** The goal here is to illustrate in a simple setting how
 90 an *adaptive* geometry can act as implicit regularizer. In such setting, the idea is to *learn* a rescaling
 91 metric at each iteration of our algorithm, using a local version of the bounds (71). We consider
 92 functions $f_{\mathbf{w}} = \sum_t \delta f_{\mathbf{w}_t}$ written in terms of a sequence of updates³ $\delta f_{\mathbf{w}_t}(\mathbf{x}) = \langle \delta \mathbf{w}_t, \Phi(\mathbf{x}) \rangle$ (we set
 93 f_0 to keep the notation simple), with *local* constraints on the parameter updates:

$$\mathcal{F}_m^A = \{f_{\mathbf{w}}: \mathbf{x} \mapsto \sum_t \langle \delta \mathbf{w}_t, \Phi(\mathbf{x}) \rangle \mid \|\delta \mathbf{w}_t\|_{A_t} \leq m_t\} \quad (4)$$

³In 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 $t > T$.

94 **Theorem 1** (Complexity of Learning Flows). *Given any sequences \mathbf{A} and \mathbf{m} of invertible matrices*
 95 *$A_t \in \mathbb{R}^{P \times P}$ and positive numbers $m_t > 0$, we have the bound*

$$\widehat{\mathcal{R}}(\mathcal{F}_{\mathbf{m}}^{\mathbf{A}}) \leq \sum_t (m_t/n) \|A_t^{-1} \Phi^\top\|_{\text{F}} \quad (5)$$

96 The same result can be formulated in terms of the sequence of feature maps $\Phi_t = A_t^{-1} \Phi$. By
 97 reparametrization invariance, the function class (16) can equivalently be written as $\mathcal{F}_{\mathbf{m}}^{\mathbf{A}} = \mathcal{F}_{\mathbf{m}}^{\Phi}$ where
 98 $\Phi = \{\Phi_t\}_t$ and the norm constraints are $\|\tilde{\delta} \mathbf{w}_t\|_2 \leq m_t$; then (17) reads

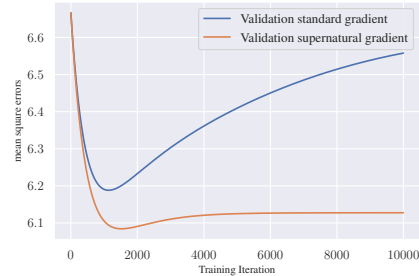
$$\widehat{\mathcal{R}}(\mathcal{F}_{\mathbf{m}}^{\Phi}) \leq \sum_t (m_t/n) \|\Phi_t\|_{\text{F}} \quad (6)$$

99 Thm. 3 suggests to include, at each iteration t , a reparametrization step with a suitable matrix
 100 \tilde{A}_t giving a low contribution to the bound (17). Applied to gradient descent, this leads to
 101 the new update rule below, where the optimization in Step 2 is over a given class of matrices.

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

1. Perform gradient step $\tilde{\mathbf{w}}_{t+1} \leftarrow \mathbf{w}_t + \delta \mathbf{w}_{\text{GD}}$
- 102 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:

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



103 The successive reparametrizations yield a varying feature map $\Phi_t = A_t^{-1} \Phi$ where $A_t = \tilde{A}_0 \cdots \tilde{A}_t$. In
 104 the original feature representation Φ , SuperNat amounts to performing natural gradient updates with
 105 respect to the local metric g_{A_t} ; and by construction, we also have $\delta f_{\mathbf{w}_t}(\mathbf{x}) = \langle \delta \mathbf{w}_{\text{GD}}, \Phi_t(\mathbf{x}) \rangle$ where
 106 $\delta \mathbf{w}_{\text{GD}}$ are standard gradient descent updates in the linear model with feature map Φ_t .

107 As an example, consider matrices \tilde{A}_ν acting diagonally in the right singular basis of the feature matrix,
 108 i.e by rescaling the singular vectors $\lambda_j \rightarrow \lambda_j/\nu_j$. Step 2 can be computed analytically (Longer
 109 Version, Prop. 4): up to isotropic rescaling, this yields the update rule $\lambda_{j(t+1)} = |\mathbf{u}_j^\top \nabla_{\mathbf{f}_w} L| \lambda_{jt}$
 110 for the singular values of Φ_t . This stretches (resp. contracts) the geometry in directions of large (resp.
 111 small) residual $\nabla_{\mathbf{f}_w} L$, thereby increases the alignment of the learned features to the signal. The
 112 working hypothesis in this paper, supported by the observations of Section 1, is that in the case of
 113 neural networks, such alignment of the features is dynamically induced as an effect of non-linearity.⁴

114 The plot shows the training curves for a simple model with Gaussian features $\Phi = [\varphi, \varphi_{\text{noise}}] \in \mathbb{R}^{d+1}$
 115 trained to regress $\mathbf{y} = \varphi + P_{\text{noise}}(\epsilon)$, with Gaussian noise is added in the direction of the noise
 116 features. SuperNat identifies dominant features (here φ) and stretches the metric along them, thereby
 117 slowing down and eventually freezing the dynamics in the orthogonal (noise) directions.

118 2.2 A New Complexity Measure for Neural Networks

119 Equ. (19) provides a bound of the Rademacher complexity for the function classes (16) specified by a
 120 fixed sequence of adaptive kernels (see Appendix C.4 for a generalization to the multiclass setting).
 121 By extrapolation to the case of non-deterministic sequences of kernels, we propose using

$$\mathcal{C}(f_{\mathbf{w}}) = \sum_t \|\delta \mathbf{w}_t\|_2 \|\Phi_t\|_{\text{F}} \quad (7)$$

122 where Φ_t is the tangent feature matrix⁵ at training iteration t , as a heuristic measure of complexity for
 123 neural networks. Following a standard protocol for studying complexity measures, [e.g., 43], Fig. 8
 124 shows its behaviour for MLP on MNIST and VGG19 on CIFAR10 trained with cross entropy loss,
 125 with (left) fixed architecture and varying level of corruption in the labels and (right) varying hidden
 126 layer size/number of channels up to 4 millions parameters, against other capacity measures proposed
 127 in the recent literature. We observe that it correctly reflects the shape of the generalization gap.

⁴For a non-linear model, the updates of the tangent feature take the same form $\Phi_t = \tilde{A}_t^{-1} \Phi_{t-1}$ as above, the
 difference being that \tilde{A}_t is no longer a matrix but a differential operator, e.g. at first order $A_t = \text{Id} - \delta \mathbf{w}_t^\top \frac{\partial}{\partial \mathbf{w}_t}$.

⁵In terms of tangent kernels, $\|\Phi_t\|_{\text{F}} = \sqrt{\text{Tr} \mathbf{K}_t}$ where \mathbf{K}_t is the tangent kernel matrix.

References

- [1] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. volume 97 of *Proceedings of Machine Learning Research*, pp. 242–252, Long Beach, California, USA, 09–15 Jun 2019. PMLR.
- [2] Shun-Ichi Amari. *Information Geometry and Its Applications*, volume 194. Springer, 2016.
- [3] Anonymous. {NNG}eometry: Easy and fast Fisher information matrices and neural tangent kernels in pytorch. In *Submitted to International Conference on Learning Representations*, 2021. under review.
- [4] Sanjeev Arora, Sanjeev Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In *ICML*, 2019.
- [5] 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.
- [6] Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *JMLR*, 2002.
- [7] Peter L. Bartlett, Dylan J. Foster, and Matus Telgarsky. Spectrally-normalized margin bounds for neural networks. In *NIPS*, 2017.
- [8] Peter L Bartlett, Philip M Long, Gábor Lugosi, and Alexander Tsigler. Benign overfitting in linear regression. *arXiv preprint arXiv:1906.11300[stat.ML]*, 2019.
- [9] 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.
- [10] Mikhail Belkin, Siyuan Ma, and Soumik Mandal. To understand deep learning we need to understand kernel learning. In *ICML*, 2018.
- [11] 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.
- [12] 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.
- [13] 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.
- [14] Yuan Cao, Zhiying Fang, Yue Wu, Ding-Xuan Zhou, and Quanquan Gu. Towards understanding the spectral bias of deep learning. *arXiv:1912.01198 [cs.LG]*, 2019.
- [15] Satrajit Chatterjee. Coherent gradients: An approach to understanding generalization in gradient descent-based optimization. In *International Conference on Learning Representations*, 2020.
- [16] L. Chizat and F Bach. A note on lazy training in supervised differentiable programming. *arXiv:1812.07956[math.OC]*, 2018.
- [17] Tarin Clanuwat, Mikel Bober-Irizar, Asanobu Kitamoto, Alex Lamb, Kazuaki Yamamoto, and David Ha. Deep learning for classical japanese literature. 2018.
- [18] Corinna Cortes, Mehryar Mohri, and Afshin Rostamizadeh. Algorithms for learning kernels based on centered alignment. *JMLR*, 13(1):795–828, 2012. ISSN 1532-4435.
- [19] Nello Cristianini, John Shawe-Taylor, André Elisseeff, and Jaz S. Kandola. On kernel-target alignment. In *NIPS*. 2002.

- 172 [20] Simon S. Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes
173 over-parameterized neural networks. In *International Conference on Learning Representations*,
174 2019.
- 175 [21] Stanislav Fort, Pawel Krzysztof Nowak, Stanislaw Jastrzebski, and Srini Narayanan. Stiffness: A
176 new perspective on generalization in neural networks. *arXiv preprint arXiv:1901.09491*, 2019.
- 177 [22] Mario Geiger, Stephano Spigler, Arthur Jacot, and Matthieu Wyart. Disentangling feature and
178 lazy training in deep neural networks. *arXiv:1906.08034 [cs.LG]*, 2019.
- 179 [23] Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel,
180 Matthias Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *arXiv*
181 *preprint arXiv:2004.07780 [cs.CV]*, 2020.
- 182 [24] Stuart Geman, Elie Bienenstock, and René Doursat. Neural networks and the bias/variance
183 dilemma. *Neural Computation*, 4(1):1–58, 1992. doi: 10.1162/neco.1992.4.1.1.
- 184 [25] Gauthier Gidel, Francis Bach, and Simon Lacoste-Julien. Implicit regularization of discrete
185 gradient dynamics in linear neural networks. In H. Wallach, H. Larochelle, A. Beygelzimer,
186 F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing*
187 *Systems 32*, pp. 3202–3211. Curran Associates, Inc., 2019.
- 188 [26] Arthur Gretton, Olivier Bousquet, Alexander Smola, and Bernhard Schölkopf. Measuring
189 statistical dependence with hilbert-schmidt norms, 2005.
- 190 [27] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The elements of statistical learning:*
191 *data mining, inference and prediction*. Springer, 2009.
- 192 [28] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
193 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
194 pp. 770–778, 2016.
- 195 [29] Elad Hoffer, Itay Hubara, and Daniel Soudry. Train longer, generalize better: Closing the
196 generalization gap in large batch training of neural networks. In *NIPS*, 2017.
- 197 [30] Arthur Jacot, Franck Gabriel, and Clement Hongler. Neural tangent kernel: Convergence and
198 generalization in neural networks. In *NIPS*, pp. 8571–8580. 2018.
- 199 [31] Yiding Jiang, Behnam Neyshabur, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantastic
200 generalization measures and where to find them. In *ICLR*, 2020.
- 201 [32] D. Kopitkov and V. Indelman. Neural spectrum alignment: Empirical study. In *International*
202 *Conference on Artificial Neural Networks (ICANN)*, September 2020.
- 203 [33] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images.
204 Technical report, Citeseer, 2009.
- 205 [34] Andrew K Lampinen, Andrew K Lampinen, and Surya Ganguli. An analytic theory of
206 generalization dynamics and transfer learning in deep linear networks. *arXiv.org*, 2018.
- 207 [35] Yann LeCun, Corinna Cortes, and CJ Burges. Mnist handwritten digit database. *ATT Labs*
208 *[Online]*. Available: <http://yann.lecun.com/exdb/mnist>, 2, 2010.
- 209 [36] M. Ledoux and M. Talagrand. *Probability in Banach Spaces: Isoperimetry and Processes*.
210 Springer Science & Business, New York, 2013.
- 211 [37] Tengyuan Liang, Tomaso Poggio, Alexander Rakhlin, and James Stokes. Fisher-rao metric,
212 geometry, and complexity of neural networks. In *Proceedings of Machine Learning Research*,
213 volume 89, pp. 888–896, 2019.
- 214 [38] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of Machine Learning*.
215 The MIT Press, 2012. ISBN 026201825X, 9780262018258.
- 216 [39] Vidya Muthukumar, Kailas Vodrahalli, Vignesh Subramanian, and Anant Sahai. Harmless
217 interpolation of noisy data in regression. *arXiv preprint arXiv:1903.09139[cs.LG]*, 2019.

- 218 [40] Vidya Muthukumar, Adhyyan Narang, Vignesh Subramanian, Mikhail Belkin, Daniel Sahai,
219 Hsu, and Anant Sahai. Classification vs regression in overparameterized regimes: Does the loss
220 function matter? *arXiv preprint arXiv:2005.08054 [cs.LG]*, 2020.
- 221 [41] Brady Neal, Sarthak Mittal, Aristide Baratin, Vinayak Tantia, Matthew Scicluna, Simon Lacoste-
222 Julien, and Ioannis Mitliagkas. A modern take on the bias-variance tradeoff in neural networks.
223 *arXiv:1810.08591 [cs.LG]*, 2018.
- 224 [42] Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. In search of the real inductive bias:
225 On the role of implicit regularization in deep learning. *ICLR workshop track*, 2015.
- 226 [43] Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and Nati Srebro. Exploring
227 generalization in deep learning. In *Advances in Neural Information Processing Systems*, pp.
228 5949–5958, 2017.
- 229 [44] Behnam Neyshabur, Ryota Tomioka, Ruslan Salakhutdinov, and Nathan Srebro. Geometry of
230 optimization and implicit regularization in deep learning. *arXiv:1705.03071 [cs.LG]*, 2017.
- 231 [45] Behnam Neyshabur, Zhiyuan Li, Srinadh Bhojanapalli, Yann LeCun, and Nathan Srebro.
232 Towards understanding the role of over-parameterization in generalization of neural networks.
233 *International Conference on Learning Representations (ICLR)*, 2019.
- 234 [46] Jonas Paccolat, Leonardo Petrini, Mario Geiger, Kevin Tyloo, and Matthieu Wyart. Geometric
235 compression of invariant manifolds in neural nets. *arXiv preprint arXiv:2007.11471*, 2020.
- 236 [47] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
237 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas
238 Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy,
239 Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-
240 performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc,
241 E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 32*, pp.
242 8024–8035. Curran Associates, Inc., 2019.
- 243 [48] Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht,
244 Yoshua Bengio, and Aaron Courville. On the spectral bias of neural networks. In *Proceedings*
245 *of the 36th International Conference on Machine Learning*, 2019.
- 246 [49] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *NIPS*,
247 2007.
- 248 [50] Olivier Roy and Martin Vetterli. The effective rank: A measure of effective dimensionality. In
249 *2007 15th European Signal Processing Conference*, pp. 606–610. IEEE, 2007.
- 250 [51] Shiori Sagawa, Aditi Raghunathan, Pang Wei Koh, and Percy Liang. An investigation of why
251 overparameterization exacerbates spurious correlations. *arXiv:2005.04345 [cs.LG]*, 2020.
- 252 [52] Andrew M. Saxe, James L. McClelland, and Surya Ganguli. Exact solutions to the nonlinear
253 dynamics of learning in deep linear neural network. In *International Conference on Learning*
254 *Representations*, 2014.
- 255 [53] B. Schölkopf, S. Mika, C. J.C. Burges, P. Knirsch, K. R. Muller, G. Ratsch, and A. J. Smola.
256 Input space versus feature space in kernel-based methods. *Trans. Neur. Netw.*, 10(5):1000–1017,
257 September 1999. ISSN 1045-9227.
- 258 [54] B. Schölkopf, J. Shawe-Taylor, AJ. Smola, and RC. Williamson. Kernel-dependent support
259 vector error bounds. In *Artificial Neural Networks, 1999. ICANN 99*, volume 470 of *Conference*
260 *Publications*, pp. 103–108. Max-Planck-Gesellschaft, IEEE, 1999.
- 261 [55] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale
262 image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- 263 [56] Nati Srebro, Karthik Sridharan, and Ambuj Tewari. On the universality of online mirror descent.
264 In *Advances in Neural Information Processing Systems 24*. 2011.

- 265 [57] Sharan Vaswani, Reza Babanezhad, Jose Gallego, Aaron Mishkin, Simon Lacoste-Julien, and
266 Nicolas Le Roux. To each optimizer a norm, to each norm its generalization. *arxiv preprint*
267 *arXiv:2006.06821[cs.LG]*, 2020.
- 268 [58] Blake Woodworth, Suriya Gunasekar, Jason D. Lee, Edward Moroshko, Pedro Savarese, Itay
269 Golan, Daniel Soudry, and Nathan Srebro. Kernel and rich regimes in overparametrized models.
270 *arXiv:2002.09277 [cs.LG]*, 2020.
- 271 [59] Zhi-Qin John Xu, Yaoyu Zhang, and Yanyang Xiao. Training behavior of deep neural network in
272 frequency domain. In Tom Gedeon, Kok Wai Wong, and Minhoo Lee (eds.), *Neural Information*
273 *Processing*, pp. 264–274, Cham, 2019. Springer International Publishing. ISBN 978-3-030-
274 36708-4.
- 275 [60] Greg Yang and Hadi Salman. A fine grained spectral perspective on neural networks. *arxiv*
276 *preprint arXiv:1907.10599[cs.LG]*, 2019.
- 277 [61] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding
278 deep learning requires rethinking generalization. *ICLR*, 2017.

Implicit Regularization via Feature Alignment (Longer Version)

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. [30], along a small number of task-relevant directions. This can be interpreted as a combined feature selection and compression mechanism. By extrapolating a new analysis of Rademacher complexity bounds for linear models, we propose and study a new heuristic measure of complexity which captures this phenomenon, in terms of sequences of tangent kernel classes along the learning trajectories.

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* [42, 61, 29]. This seems at odds with the usual understanding of the bias-variance trade-off [24, 41, 11]: highly complex models are expected to overfit the training data and perform poorly on test data [27]. Solving this apparent paradox requires understanding the various learning biases induced by the training procedure, which can act as implicit regularizers [42, 44].

In this paper, we help clarify one such implicit regularization mechanism, by examining the evolution of the *neural tangent features* [30] 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 model function class, akin to the Fisher information metric [e.g., 2]; and a *functional* perspective – through the tangent kernel and its RKHS.

Our main observation, in standard supervised classification settings, is a dynamical alignment of the tangent features along a small number of task-relevant directions during training. We interpret this phenomenon as combining a *feature selection* and a *compression* mechanisms. The intuition motivating this work is that such mechanisms are what allows the model to adapt its capacity to the task, which in turn underpins the generalization abilities of heavily overparametrized models.

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 an increasing similarity with the class labels, e.g. as measured by *centered kernel alignment* (CKA) [19, 18].
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* functions $f_{\mathbf{w}}: \mathcal{X} \rightarrow \mathbb{R}$ to keep notation light.⁶

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

$$\Phi_{\mathbf{w}}(\mathbf{x}) := \nabla_{\mathbf{w}} f_{\mathbf{w}}(\mathbf{x}) \quad (8)$$

The corresponding kernel $k_{\mathbf{w}}(\mathbf{x}, \bar{\mathbf{x}}) = \langle \Phi_{\mathbf{w}}(\mathbf{x}), \Phi_{\mathbf{w}}(\bar{\mathbf{x}}) \rangle$ is the **tangent kernel** [30]. Intuitively, the tangent features govern how small changes in parameter affect the function’s outputs,

$$\delta f_{\mathbf{w}}(\mathbf{x}) = \langle \delta \mathbf{w}, \Phi_{\mathbf{w}}(\mathbf{x}) \rangle + O(\|\delta \mathbf{w}\|^2) \quad (9)$$

⁶The extension to vector-valued functions, relevant for the multiclass classification setting, is presented in Appendix A, along with more mathematical details.

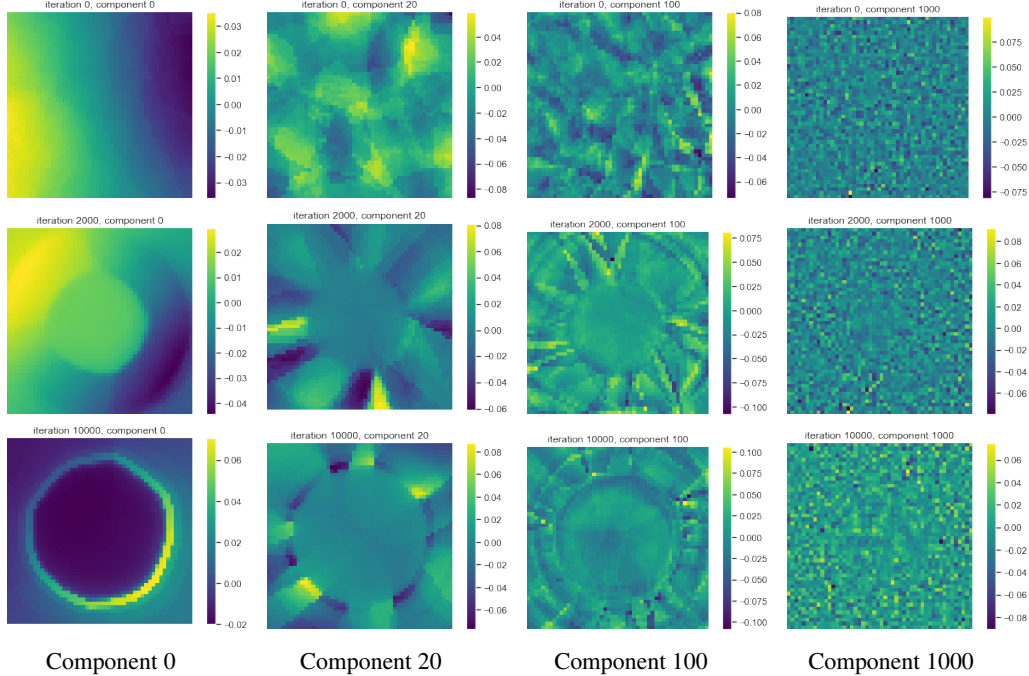


Figure 3: Evolution of eigenfunctions of the tangent kernel, ranked in nonincreasing order of the eigenvalues (in columns), at various iterations during training (in rows), for the $2d$ Disk dataset. After a number of iterations, we observe modes corresponding to the class structure (e.g. boundary circle) 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 in the directions of the signal.

321 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
 322 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}
 323 by pullback of the L^2 scalar product (see Appendix A). It characterizes the geometry of the function
 324 class \mathcal{F} .

325 **Spectral Bias.** The structure of the tangent features impacts the evolution of the function during train-
 326 ing. To formalize this, we introduce the covariance eigenvalue decomposition $g_{\mathbf{w}} = \sum_{j=1}^P \lambda_{\mathbf{w}j} \mathbf{v}_{\mathbf{w}j} \mathbf{v}_{\mathbf{w}j}^{\top}$,
 327 which summarizes the predominant directions in parameter space. Given n input samples (\mathbf{x}_i) and
 328 $\mathbf{f}_{\mathbf{w}} \in \mathbb{R}^n$ 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
 329 function $L := L(\mathbf{f}_{\mathbf{w}})$. The following elementary result (see Appendix B) shows how the corresponding
 330 function updates in the linear approximation (9), $\delta f_{\text{GD}}(\mathbf{x}) := \langle \delta \mathbf{w}_{\text{GD}}, \Phi_{\mathbf{w}}(\mathbf{x}) \rangle$, decompose in the
 331 **eigenbasis**⁷ of the tangent kernel:

$$u_{\mathbf{w}j}(\mathbf{x}) = \frac{1}{\sqrt{\lambda_{\mathbf{w}j}}} \langle \mathbf{v}_{\mathbf{w}j}, \Phi_{\mathbf{w}}(\mathbf{x}) \rangle \quad (10)$$

332 **Lemma 2** (Local Spectral Bias). *The function updates decompose as $\delta f_{\text{GD}}(\mathbf{x}) = \sum_{j=1}^P \delta f_j u_{\mathbf{w}j}(\mathbf{x})$*
 333 *with*

$$\delta f_j = -\eta \lambda_{\mathbf{w}j} (\mathbf{u}_{\mathbf{w}j}^{\top} \nabla_{\mathbf{f}_{\mathbf{w}}} L), \quad (11)$$

334 *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 $\nabla_{\mathbf{f}_{\mathbf{w}}}$ is the gradient w.r.t the sample outputs.*

335 This illustrates how, from the point of view of function space, the eigenvalues act as a mode-specific
 336 rescaling $\eta \lambda_{\mathbf{w}j}$ of the learning rate. This is a local version of a well-known bias for linear models
 337 trained by gradient descent (e.g in linear regression, see Appendix B.2), which prioritizes learning
 338 functions within the top eigenspaces of the kernel. Several recent works [12, 9, 60] investigated such

⁷The functions $u_{\mathbf{w}j}$, $j \in \{1 \dots P\}$ form an orthonormal family in $L^2(\rho)$, i.e. $\mathbb{E}_{\mathbf{x} \sim \rho} [u_{\mathbf{w}j} u_{\mathbf{w}j'}] = \delta_{jj'}$, yielding 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. Note that kernel and covariance share the same spectrum.

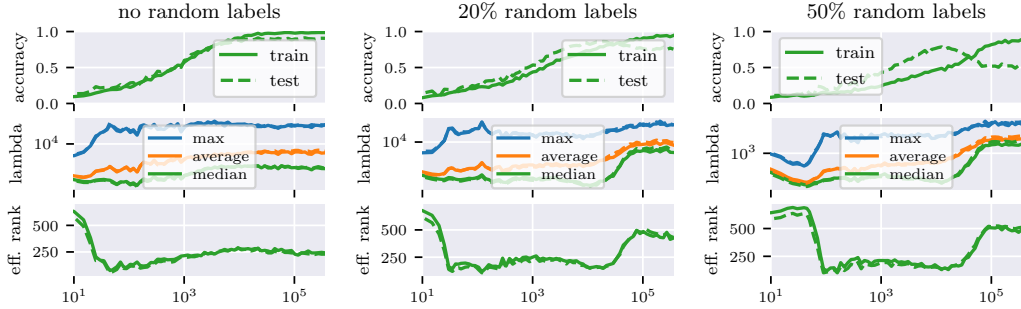


Figure 4: Evolution of tangent kernel spectrum and effective rank of a VGG19 trained by SGD with batch size 100, learning rate 0.01 and momentum 0.9 on CIFAR10 with various ratio of random labels. The small effective rank of the kernel biases the training procedure towards a few top eigenvectors.

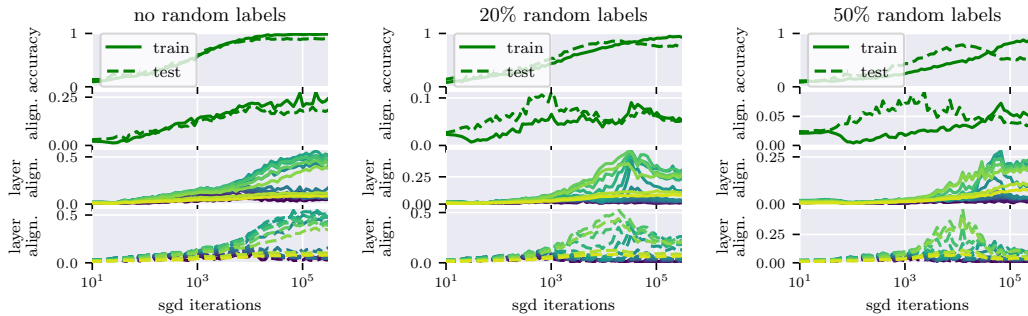


Figure 5: Evolution of kernel alignment and layer-wise kernel alignments of a VGG19 trained by SGD with batch size 100, learning rate 0.01 and momentum 0.9 on CIFAR10 with various ratios of random labels. For layer-wise alignment we map layers to colors sequentially from input layer (-), through intermediate layers (-), to output layer (-). See Figure 13 and 16 in appendix for additional architectures/datasets.

339 bias for neural networks, in *linearized* regimes where the tangent kernel remains constant during
 340 training [30, 20, 1]. As a simple example, for a randomly initialized MLP on 1D uniform data,
 341 Fig. 10 (Appendix B) shows an alignment of the tangent kernel eigenfunctions with Fourier modes of
 342 increasing frequency, explaining prior empirical observations [48, 59] of a ‘spectral bias’ towards
 343 low-frequency functions.

344 **Tangent Features Adapt to the Task.** By contrast, our aim in this paper is to highlight and discuss
 345 *non-linear* effects, in the (standard) regime where the tangent features and their kernel evolve during
 346 training [e.g., 22, 58].

347 As a first illustration of such effects, Fig. 3 shows visualizations of eigenfunctions of the tangent
 348 kernel (ranked in nonincreasing order of the eigenvalues), during training MLP by gradient descent
 349 of the binary cross entropy loss, on a simple classification task: $y(\mathbf{x}) = \pm 1$ depending on whether
 350 $\mathbf{x} \sim \text{Unif}[-1, 1]^2$ is in the centered disk of radius $\sqrt{2/\pi}$ (details in Appendix E). After a number of
 351 iterations, we observe (rotation invariant) modes corresponding to the class structure (e.g. boundary
 352 circle) showing up in the *top* eigenfunctions of the learned kernel. We also note an increasing
 353 spectrum anisotropy – for example, the ratio λ_{20}/λ_1 , which is 1.5% at iteration 0, has dropped to
 354 0.2% at iteration 2000. The interpretation is that the tangent kernel *stretches* along a small number of
 355 directions that are highly correlated with the signal during training. We quantify and investigate this
 356 alignment effect in more detail below.

357 3 Neural Feature Alignment

358 In this section, we perform experiments showing a dynamical alignment of the tangent features along
 359 a small number of task-relevant directions during training. We show in particular that networks learn

360 tangent features with increasing similarity with the class labels, as measured by **centered kernel**
 361 **alignment** (CKA) [19, 18]. We interpret this phenomenon as combining both a feature selection and
 362 a compression mechanism.

363 3.1 Setup

364 We run experiments on MNIST [35] and CIFAR10 [33] with standard MLPs, VGG [55] and Resnet
 365 [28] architectures, trained by stochastic gradient descent (SGD) with momentum, using cross-entropy
 366 loss. We use PyTorch [47] and NNGeometry [3] for efficient evaluation of tangent kernels.

367 In multiclass settings, tangent kernels evaluated on n samples carry additional class indices $y \in$
 368 $\{1 \dots c\}$ and thus are $nc \times nc$ matrices, $(\mathbf{K}_w)_{ij}^{yy'} := k_w(\mathbf{x}_i, y; \mathbf{x}_j, y')$. In all our experiments, we
 369 evaluate tangent kernels on mini-batches (either from the train or the test set) of size $n = 100$. For
 370 $c = 10$ classes, this yields kernel matrices of size 1000×1000 . We report results obtained from
 371 *centered* tangent features $\Phi_w(\mathbf{x}) \rightarrow \Phi_w(\mathbf{x}) - \mathbb{E}_x \Phi_w(\mathbf{x})$, though we obtain qualitatively similar
 372 results for uncentered features (see plots in Appendix E.2).

373 3.2 Spectrum Evolution

374 We first investigate the evolution of the tangent kernel *spectrum* for a VGG19 on CIFAR 10, trained
 375 with and without label noise (Fig. 4). The main take away is an anisotropic increase of the spectrum
 376 during training. We report results for kernels evaluated on training examples (solid line) and test
 377 examples (dashed line).⁸

378 The first observation is a significant *increase* of the spectrum, early in training (note the log scale for
 379 the number of iterations). By the time the model reaches 100% training accuracy, the maximum and
 380 average eigenvalues have gained more than 2 orders of magnitude.

381 The second observation is that this evolution is highly *anisotropic*. We quantify spectrum anisotropy
 382 using a notion of **effective rank** based on spectral entropy [50]. Given a kernel matrix \mathbf{K} in $\mathbb{R}^{r \times r}$ with
 383 (strictly) positive eigenvalues $\lambda_1, \dots, \lambda_r$, let $\mu_j = \lambda_j / \sum_{i=1}^r \lambda_i$ be the trace-normalized eigenvalues.
 384 The effective rank is defined as $\text{erank} = \exp(H(\boldsymbol{\mu}))$ where $H(\boldsymbol{\mu})$ is the Shannon entropy,

$$H(\boldsymbol{\mu}) = - \sum_{j=1}^r \mu_j \log(\mu_j) \quad (12)$$

385 This effective rank is a real number between 1 and r , upper bounded by $\text{rank}(\mathbf{K})$, which measures
 386 the ‘uniformity’ of the spectrum through the entropy. We also track the various **trace ratios**
 387 $T_k = \sum_{j < k} \lambda_j / \sum_j \lambda_j$ as measures of the relative importance of the top k eigenvalues (see Fig. 15
 388 in Appendix E.3).

389 We note an important decrease of the effective rank early in training (third row in Fig. 4), reaching a
 390 phase where only a few top eigenvalues account for most of the trace. This can be observed directly
 391 from the highlighted (in red) ratios T_{40}, T_{80} and T_{160} (Fig. 15), e.g. T_{80} accounting for 50% of the
 392 total trace (over 1000 eigenvalues). Remarkably, in the presence of high label noise, the effective rank
 393 of the tangent kernel evaluated on *training* examples (anti-)correlates nicely with the *test* accuracy,
 394 decreasing or remaining low during the learning phase (increase of test accuracy) and rising when
 395 overfitting starts (decrease of test accuracy). This suggests that the effective rank of the tangent kernel
 396 (and hence that of the metric) might already provide a good proxy for a measure of the effective
 397 capacity of the network.

398 3.3 Alignment to class labels

399 We now include the evolution of the eigenvectors in our analysis. We investigate the similarity of
 400 the learned tangent features with the class label through a similarity index called centered kernel
 401 alignment. Given two kernel matrices \mathbf{K} and \mathbf{K}' in $\mathbb{R}^{r \times r}$, it is defined as

$$\text{CKA}(\mathbf{K}, \mathbf{K}') = \frac{\text{Tr}[\mathbf{K}_c \mathbf{K}'_c]}{\|\mathbf{K}_c\|_F \|\mathbf{K}'_c\|_F} \in [0, 1] \quad (13)$$

⁸The 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.

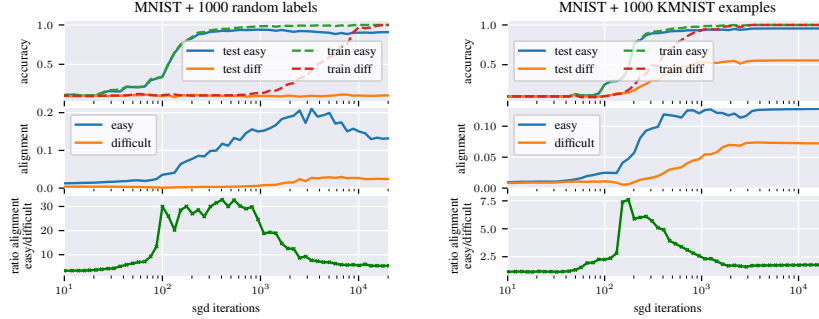


Figure 6: Alignment easy versus difficult: We augment a dataset composed of 10,000 *easy* MNIST examples with 1000 *difficult* examples from 2 different setups: **(left)** 1000 MNIST examples with random label **(right)** 1000 KMNIST examples. We train a MLP with 6 layers of 80 hidden units using SGD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

402 where the c subscript denotes the feature centering operation, i.e. $\mathbf{K}_c = C\mathbf{K}C$ where $C = I_r - \frac{1}{r}\mathbf{1}\mathbf{1}^T$
 403 is the centering matrix. CKA is a normalized version of the Hilbert-Schmidt Independence Criterion
 404 [26] designed as a dependence measure for two sets of features. The normalization by the Frobenius
 405 norms makes CKA invariant under isotropic rescaling.

406 Let $\mathbf{Y} \in \mathbb{R}^{nc}$ be the vector resulting from the concatenation of the one-hot label representations
 407 $\mathbf{Y}_i \in \mathbb{R}^c$ of the n samples. Similarity with the labels is measured through CKA with the rank-one
 408 kernel $\mathbf{K}_\mathbf{Y} := \mathbf{Y}\mathbf{Y}^T$. Intuitively, $\text{CKA}(\mathbf{K}, \mathbf{K}_\mathbf{Y})$ is high when \mathbf{K} has low (effective) rank and such
 409 that the angle between \mathbf{Y} and its top eigenspaces is small.⁹ Maximizing such index has been used as
 410 a criterion for kernel selection in the literature on learning kernels [18].

411 In the same setup as in Section 3.2, we observe (Fig. 5) an increasingly high CKA between tangent
 412 kernel and the labels as training progresses. The trend is similar for other architectures and datasets
 413 (Fig. 13 in Appendix E show CKA plots for MLP on MNIST and Resnets 18 on CIFAR10).

414 Interestingly, in the presence of high level noise and during the learning phase (increase of test
 415 accuracy), the CKA reaches a much higher value for kernels evaluated on test inputs than for kernels
 416 evaluated on training inputs (note that test labels are not randomized). Together with equation 11, the
 417 alignment of the tangent kernel along clean labels sheds lights on empirical observations that, in the
 418 presence of noise, deep networks ‘learn patterns first’ [5] (see Section 3.4 for additional insights).

419 We also report the alignments of the *layer-wise* tangent kernels \mathbf{K}_w^ℓ , obtained from the function
 420 gradients w.r.t parameters of layer ℓ . By construction, the tangent kernel is the sum of the layer-wise
 421 kernels over all layers of the network, $\mathbf{K}_w = \sum_{\ell=1}^L \mathbf{K}_w^\ell$. We observe a high CKA (reaching more
 422 than 0.5), especially for the *intermediate* layers¹⁰, suggesting the key role of depth in the overall
 423 alignment of the tangent kernel (see also Section 3.5).

424 3.4 Hierarchical Alignment

425 A key aspect of the generalization question for deep networks concerns the articulation between
 426 learning and memorization, in the presence of noise [61] or difficult examples [51]. Motivated by
 427 this, we would like to probe the evolution of the tangent features separately in the directions of both
 428 type of examples in such settings. To do so, our strategy is to measure partial CKA on examples from
 429 two subsets of the same size in the dataset – one with ‘easy’ examples, the other with ‘difficult’ ones.
 430 Our setup is to augment 10,000 MNIST training examples with 1000 difficult examples of 2 types: (i)
 431 examples with random labels and (ii) examples from the dataset KMNIST [17]. KMNIST images
 432 present similar features than MNIST digits (grayscale handwritten characters) but represent Japanese
 433 characters.

⁹In the limiting case $\text{CKA}(\mathbf{K}, \mathbf{K}_\mathbf{Y}) = 1$, the features are all aligned with each other and parallel to \mathbf{Y} .

¹⁰We 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.

434 The results are shown in Fig. 6. As training progresses, the CKA on the easy examples increases
 435 faster (and to a higher value); in the case of the (structured) difficult examples from KMNIST, we
 436 observe an increase of the CKA later in training. This demonstrates a hierarchy in the adaptation of
 437 the kernel, measured by the ratio between both alignments. From the intuition developed in the paper
 438 (see Section 2), this aspect of the non-linear dynamics favors a sequentialization of the learning ('easy
 439 patterns first'), a phenomenon analogous to one pointed out in the context of deep linear networks
 440 [52, 34, 25].

441 3.5 Ablation

442 In order to study the influence of depth on alignment and test the robustness to the choice of seeds,
 443 we reproduce the experiment of the previous section for MLP with different depths, while varying
 444 parameter initialization and minibatch sampling. Our results, shown in Fig 16 (Appendix E), suggest
 445 that the alignment effect is magnified as depth increases. We also observe that the ratio of the
 446 maximum alignment between easy and difficult examples is increased with depth, but stays high for a
 447 smaller number of iterations.

448 4 Measuring Complexity

449 In this section, drawing upon intuitions from linear models, we illustrate on a simple setting how
 450 the alignment effect highlighted in the previous section can act as implicit regularization. We also
 451 motivate a new complexity measure for neural networks and compare its correlation to generalization
 452 against various measures proposed in the recent literature.

453 4.1 Insights from Linear Models

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

457 Our discussion will be based on the **Rademacher complexity**, which shows up in generalization
 458 bounds [6]. It measures how well \mathcal{F} correlates with random noise on the sample set \mathcal{S} :

$$\widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}) = \mathbb{E}_{\sigma \in \{\pm 1\}^n} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(\mathbf{x}_i) \right] \quad (14)$$

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

462 A standard approach for controlling capacity is in terms of the *norm* of the weight vector – usually
 463 the ℓ_2 -norm. In general, given any invertible matrix $A \in \mathbb{R}^{P \times P}$, we may consider the norm
 464 $\|\mathbf{w}\|_A := \sqrt{\mathbf{w}^\top g_A \mathbf{w}}$ induced by the metric $g_A = AA^\top$. For $M_A > 0$, let $\mathcal{F}_{M_A}^A$ be the subclass of
 465 functions $f_{\mathbf{w}}$ such that $\|\mathbf{w}\|_A \leq M_A$. A direct extension of standard bounds for the Rademacher
 466 complexity (see Appendix C) yields,

$$\widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_{M_A}^A) \leq (M_A/n) \|A^{-1} \Phi^\top\|_{\text{F}} \quad (15)$$

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

468 This freedom in the choice of rescaling matrix A , due to linear reparametrization invariance, raises
 469 the question of which of the norms $\|\cdot\|_A$ provides meaningful measures of the model's capacity.
 470 Recent works [10, 40] pointed out that using ℓ_2 norm is not coherently linked with generalization in
 471 practice. We discuss this issue in Appendix C.5, illustrating how meaningful norms critically depend
 472 on the geometry defined by the features.¹¹

¹¹Analysis of the relation between capacity and feature geometry can be traced back to early work on kernel methods [53]

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

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:

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

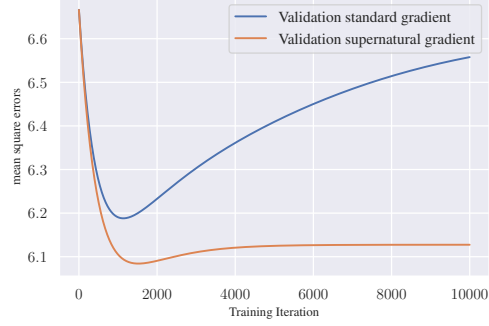


Figure 7: **(left)** SuperNat algorithm and **(right)** validation curves obtained with standard and SuperNat gradient descent, on the noisy linear regression problem. At each iteration, SuperNat identifies dominant features and stretches the kernel along them, thereby slowing down and eventually freezing the learning dynamics in the noise direction. This naturally yields better generalization than standard gradient descent on this problem.

473 4.1.1 Feature Alignment as Implicit Regularization

474 The goal here is to illustrate in a simple setting how an *adaptive* geometry along optimization
 475 trajectories can act as an implicit regularizer. In such setting, the idea is to *learn* a rescaling metric at
 476 each iteration of our algorithm, using a local version of the bounds (71).

477 **Complexity of Learning Flows.** Since we are interested in functions $f_{\mathbf{w}}$ that result from an iterative
 478 algorithm, we can assume they are written as $f_{\mathbf{w}} = f_0 + \sum_t \delta f_{\mathbf{w}_t}$ in terms of a sequence of updates
 479 $\delta f_{\mathbf{w}_t}(\mathbf{x}) = \langle \delta \mathbf{w}_t, \Phi(\mathbf{x}) \rangle$.¹² We set $f_0 = 0$ to keep the notation simple. Instead of considering classes
 480 of functions with direct constraints on the parameter, we consider functions resulting from a learning
 481 flow with *local* constraints on the parameter *updates*:

$$\mathcal{F}_m^{\mathbf{A}} = \{f_{\mathbf{w}} : \mathbf{x} \mapsto \sum_t \langle \delta \mathbf{w}_t, \Phi(\mathbf{x}) \rangle \mid \|\delta \mathbf{w}_t\|_{A_t} \leq m_t\} \quad (16)$$

482 The result (71) extends as follows.

483 **Theorem 3** (Complexity of Learning Flows). *Given any sequences \mathbf{A} and \mathbf{m} of invertible matrices*
 484 *$A_t \in \mathbb{R}^{P \times P}$ and positive numbers $m_t > 0$, we have the bound*

$$\widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_m^{\mathbf{A}}) \leq \sum_t (m_t/n) \|A_t^{-1} \Phi^\top\|_{\text{F}} \quad (17)$$

485 Equ. 17 provides us with bounds written in terms of local contributions at each iteration t . Note
 486 that the same result can be formulated in terms of the sequence of feature maps $\Phi_t = A_t^{-1} \Phi$. By
 487 reparametrization invariance, the function class (16) can equivalently be written as $\mathcal{F}_m^{\mathbf{A}} = \mathcal{F}_m^{\Phi}$ where
 488 $\Phi = \{\Phi_t\}_t$ and

$$\mathcal{F}_m^{\Phi} = \{f_{\mathbf{w}} : \mathbf{x} \mapsto \sum_t \langle \tilde{\delta} \mathbf{w}_t, \Phi_t(\mathbf{x}) \rangle \mid \|\tilde{\delta} \mathbf{w}_t\|_2 \leq m_t\} \quad (18)$$

489 In this formulation, the result (17) reads:

$$\widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_m^{\Phi}) \leq \sum_t (m_t/n) \|\Phi_t\|_{\text{F}} \quad (19)$$

490 **Optimizing the Feature Scaling.** To obtain learning flows with low complexity, Thm. 3 suggests to
 491 include, at each iteration t , a reparametrization step with a suitable matrix \tilde{A}_t giving a low contribution
 492 to the bound (17). Applied to gradient descent (GD), this leads to a new update rule sketched as in
 493 Fig 7 (left), where the optimization in Step 2 is over a given class of reparametrization matrices. As
 494 an example, we consider the class of matrices A_ν acting diagonally in the right singular basis of the
 495 feature matrix $\Phi = \sum_{j=1}^n \sqrt{\lambda_j} \mathbf{u}_j \mathbf{v}_j^\top$; which amounts to rescaling the singular vector $\lambda_j \rightarrow \lambda_j/\nu_j$.

¹²In 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 $t > T$.

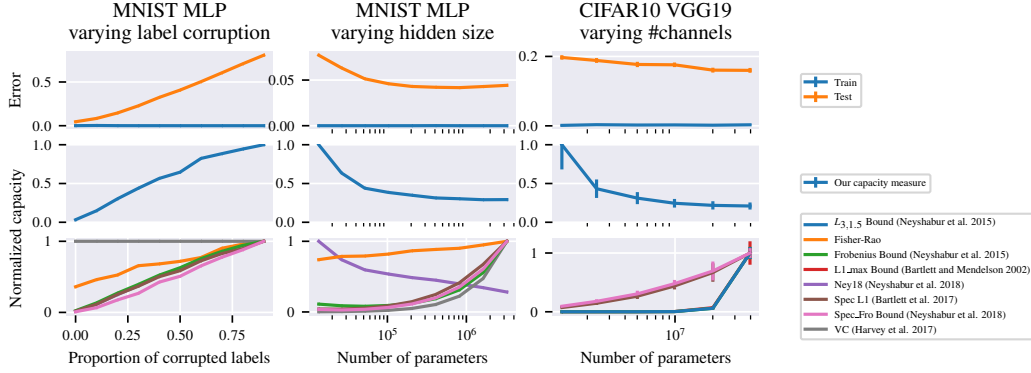


Figure 8: Complexity measures on MNIST with a one hidden layer MLP (left) as we increase the hidden layer size, (center) for a fixed hidden layer of 256 units as we increase label corruption and (right) for a VGG19 on CIFAR10 as we vary the number of channels. All networks are trained until cross-entropy loss reaches 0.01. Our proposed complexity measure and the one proposed by Neyshabur et al. 2018 are the only ones to correctly reflect the shape of the generalization gap.

496 **Proposition 4.** For the class of rescaling matrices A_ν defined above, any minimizer in Step 2 in Fig
 497 7, where $\delta \mathbf{w}_{GD} = -\eta \nabla_{\mathbf{w}} L$ is a GD updates w.r.t a loss L , takes the form

$$v_{jt}^* = \kappa \frac{1}{|\mathbf{u}_j^\top \nabla_{\mathbf{f}_w} L|} \quad (20)$$

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

499 The successive reparametrizations yield a varying feature map $\Phi_t = A_t^{-1} \Phi$ where $A_t = \tilde{A}_0 \cdots \tilde{A}_t$.
 500 In the original representation Φ , SuperNat amounts to natural gradient descent with respect to the
 501 local metric $g_{A_t} = A_t A_t^\top$. In the context of Proposition 4, this yields the following update rule, up to
 502 isotropic rescaling, for the singular values of Φ_t :

$$\lambda_{j(t+1)} = |\mathbf{u}_j^\top \nabla_{\mathbf{f}_w} L| \lambda_{jt} \quad (21)$$

503 In this illustrative setting, we see how the feature map (or kernel) adapts to the task, by stretching
 504 (resp. contracting) its geometry in directions \mathbf{u}_j along which the residual $\nabla_{\mathbf{f}_w} L$ has large (resp.
 505 small) components. Intuitively, if a large component $|\mathbf{u}_j^\top \nabla_{\mathbf{f}_w} L|$ corresponds to signal and a small
 506 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
 507 signal-to-noise ratio, thereby increasing the alignment of the learned feature matrix to the signal.

508 Fig 7 in Appendix ?? (right) shows results for the following regression setup. We consider Gaussian
 509 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 training
 510 features, we assume the label vector takes the form $\mathbf{y} = \varphi + P_{\text{noise}}(\epsilon)$, where Gaussian noise
 511 $\epsilon \sim \mathcal{N}(0, \sigma^2 I_n)$ is projected onto the noise features through $P_{\text{noise}} = \varphi_{\text{noise}} \varphi_{\text{noise}}^\top$. The model is
 512 trained by gradient descent of the mean square loss and its SuperNat variant, where Step 2 uses the
 513 analytical solution of Proposition 4. We set $d = 10, \sigma^2 = 0.1$ and use $n = 50$ training points. At
 514 each iteration, SuperNat identifies dominant features (feature selection, here φ) and stretches the
 515 metric along them, thereby slowing down and eventually freezing the dynamics in the orthogonal
 516 (noise) directions (compression).

517 4.2 A New Complexity Measure for Neural Networks

518 Equ. (19) provides a bound of the Rademacher complexity for the function classes (16) specified by a
 519 fixed sequence of adaptive kernels (see Appendix C.4 for a generalization to the multiclass setting).
 520 By extrapolation to the case of non-deterministic sequences of kernels, we propose using

$$\mathcal{C}(f_w) = \sum_t \|\delta \mathbf{w}_t\|_2 \|\Phi_t\|_F \quad (22)$$

521 where Φ_t is the tangent feature matrix¹³ at training iteration t , as a heuristic measure of complexity for
 522 neural networks. Following a standard protocol for studying complexity measures, [e.g., 43], Fig. 8
 523 shows its behaviour for MLP on MNIST and VGG19 on CIFAR10 trained with cross entropy loss,
 524 with **(left)** fixed architecture and varying level of corruption in the labels and **(right)** varying hidden
 525 layer size/number of channels up to 4 millions parameters, against other capacity measures proposed
 526 in the recent literature. We observe that it correctly reflects the shape of the generalization gap.

527 5 Related Work

528 **Capacity and Geometry.** In the context of linear models, analysis of the relation between capacity
 529 and feature geometry can be traced back to early work on kernel methods (Schölkopf et al. [53]),
 530 leading to data-dependent error bounds in terms of the eigenvalues of the kernel Gram matrix
 531 (Schölkopf et al. [54]). Recent analysis of the minimum norm interpolators in overparametrized linear
 532 regression emphasized the impact of feature geometry – through the spectrum of the data covariance –
 533 on generalization performance [8, 39].

534 Specifically, these works illustrate the key role of *feature anisotropy*, combined to a high correlation
 535 of the few dominant features with the signal [see also e.g. 13], in the generalization performance
 536 [39]. Both feature anisotropy and high correlation of the dominant features are conditions for a high
 537 alignment between kernel and labels. Our results in this paper emphasizes the key role of the training
 538 dynamics in favouring such conditions.

539 **Generalization Measures.** There has been a large body of work on generalization measures for
 540 neural networks (see Jiang et al. [31] and references therein), some of which theoretically motivated
 541 by norm or margin based bounds (e.g Neyshabur et al. [45], Bartlett et al. [7]). Liang et al. [37]
 542 proposed using the Fisher-Rao norm to measure capacity in a geometrical invariant manner. Our
 543 approach aims at taking into account the geometry along the whole optimization trajectories. Closely
 544 related perspectives in the recent literature are the notion of stiffness [21] and coherent gradients [15],
 545 tied to the structure of tangent kernels for the loss class.

546 **Spectral Bias and Tangent Kernels.** A recent line of work on the so-called *spectral bias* [48, 59],
 547 relying on Fourier analysis, suggested that neural networks prioritize learning the lowest complexity
 548 components of the data during training. In *linearized* regimes where the training dynamics can be
 549 described by a fixed kernel [30, 20, 16], this can be understood in terms of the standard learning bias
 550 along the kernel principal components in linear regression [4, 9, 14]. Several other works [12, 9, 60]
 551 investigated implicit bias of neural networks through a spectral analysis in such regimes. In this paper,
 552 we highlight and discuss *non-linear* effects, in the feature learning regime where the tangent kernel
 553 evolves during training [22, 58].

554 Independent concurrent works highlight alignment phenomena similar to the one we study here
 555 [32, 46]. We offer various complementary empirical insights, and frame the alignment mechanism
 556 from the point of view of implicit regularization.

557 6 Conclusion

558 Through experiments on modern architectures, we highlighted an alignment effect of the tangent
 559 features and their kernel along a small number of task-dependent directions, quantified by centered
 560 kernel alignment. We interpret this phenomenon as combining a *feature selection* mechanism and a
 561 *compression* of the model around the dominant features.

562 We argued that such a dynamical alignment can act as implicit regularization. By extrapolating
 563 Rademacher complexity theory from linear models to learning flows, we introduced a new heuristic
 564 complexity measure for neural networks, and showed that it correlates with the generalization gap
 565 when varying the number of parameters, and when increasing the proportion of corrupted labels.

566 The results of this paper open several avenues for further investigation. The type of complexity
 567 measure we propose suggests a principled way to rescale the geometry in which to perform gradient
 568 descent [56, 44]. Whether a procedure such as SuperNat, which optimizes a preconditioning matrix so

¹³In terms of tangent kernels, $\|\Phi_t\|_F = \sqrt{\text{Tr}\mathbf{K}_t}$ where \mathbf{K}_t is the tangent kernel matrix.

569 as to minimize a generalization bound¹⁴, can produce meaningful practical results for neural networks,
570 remains to be seen.

571 One of the consequences one can expect from alignment effect highlighted here is to encourage
572 learning from a small number of highly predictive features. While this feature selection ability
573 might explain in part the performance of neural networks on a range of supervised tasks, it may also
574 might underpin their notorious sensitivity to spurious correlations [51] and weakness to generalize
575 out-of-distribution [23]. Resolving this tension is a fascinating challenge.

576 A Geometry and Tangent Kernels

577 We describe in more formal detail the notion of geometry we consider in the paper for parametric
578 function classes. Formally, specifying such a geometry relies on a choice a distance measure or metric
579 on the function space, which is then pulled back to parameter space. We will consider general classes
580 of predictors:

$$\mathcal{F} = \{f_{\mathbf{w}} : \mathcal{X} \rightarrow \mathbb{R}^c \mid \mathbf{w} \in \mathcal{W}\}, \quad (23)$$

581 where the parameter space \mathcal{W} is a finite dimensional manifold of dimension P (typically \mathbb{R}^P). For
582 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
583 also be viewed as a scalar function on $\mathcal{X} \times \mathcal{Y}$ where $\mathcal{Y} = \{1 \cdots c\}$ is the set of classes.

584 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
585 distribution. The inclusion $\mathcal{F} \subset L^2(\rho, \mathbb{R}^c)$ equips \mathcal{F} with the L^2 scalar product and corresponding
586 norm:

$$\langle f, g \rangle_{\rho} := \mathbb{E}_{\mathbf{x} \sim \rho} [f(\mathbf{x})^{\top} g(\mathbf{x})], \quad \|f\|_{\rho} := \sqrt{\langle f, f \rangle_{\rho}} \quad (24)$$

587 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} .
588 That is, given $\zeta, \xi \in \mathcal{T}_{\mathbf{w}}\mathcal{W} \cong \mathbb{R}^P$ on the tangent space at \mathbf{w} ,

$$g_{\mathbf{w}}(\zeta, \xi) = \langle \partial_{\zeta} f_{\mathbf{w}}, \partial_{\xi} f_{\mathbf{w}} \rangle_{\rho} \quad (25)$$

589 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
590 of \mathbb{R}^P , the metric is represented by the matrix of gradient second moments:

$$(g_{\mathbf{w}})_{pq} = \mathbb{E}_{\mathbf{x} \sim \rho} \left[\left(\frac{\partial f_{\mathbf{w}}(\mathbf{x})}{\partial w_p} \right)^{\top} \frac{\partial f_{\mathbf{w}}(\mathbf{x})}{\partial w_q} \right] \quad (26)$$

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

$$\|df_{\mathbf{w}}\|_{\rho}^2 = \sum_{p,q=1}^P \left\langle \frac{\partial f_{\mathbf{w}}}{\partial w_p} dw_p, \frac{\partial f_{\mathbf{w}}}{\partial w_q} dw_q \right\rangle_{\rho} = \sum_{p,q=1}^P (g_{\mathbf{w}})_{pq} dw_p dw_q \quad (27)$$

593 This geometry has a dual description in function space in terms of *kernels*. The idea is to view
594 the differential 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
595 (co)tangent space. Thus, in a given basis, the **tangent features** are given by the function derivatives
596 with respect to the parameters

$$\Phi_{w_p}(\mathbf{x})[y] := \frac{\partial f_{\mathbf{w}}(\mathbf{x})[y]}{\partial w_p} \quad (28)$$

597 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 .
598 It defines the so-called **tangent kernel** through the Euclidean dot product in \mathbb{R}^P :

$$k_{\mathbf{w}}(\mathbf{x}, y; \tilde{\mathbf{x}}, y') = \sum_{p=1}^P \frac{\partial f_{\mathbf{w}}(\mathbf{x})[y]}{\partial w_p} \frac{\partial f_{\mathbf{w}}(\tilde{\mathbf{x}})[y']}{\partial w_p} \quad (29)$$

599 Given n input samples $\mathbf{x}_1, \dots, \mathbf{x}_n$, we represent the sample output scores $f_{\mathbf{w}}(\mathbf{x}_i)[y]$ as flattened in
600 a single vector $\mathbf{f}_{\mathbf{w}} \in \mathbb{R}^{nc}$ and the tangent features $\Phi_{w_p}(\mathbf{x}_i)[y]$ as a $nc \times P$ matrix $\Phi_{\mathbf{w}}$. Using this
601 notation, (26) and (29) yield the sample covariance $P \times P$ matrix and kernel (Gram) $nc \times nc$ matrix:

$$\mathbf{G}_{\mathbf{w}} = \Phi_{\mathbf{w}}^{\top} \Phi_{\mathbf{w}}, \quad \mathbf{K}_{\mathbf{w}} = \Phi_{\mathbf{w}} \Phi_{\mathbf{w}}^{\top} \quad (30)$$

¹⁴See the recent work by [57] for further empirical investigations of this problem in the context of linear models.

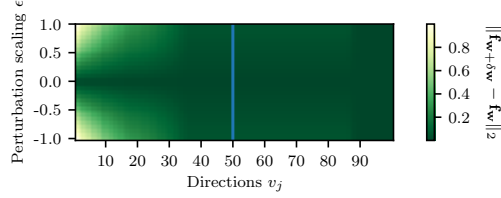


Figure 9: Variations of \mathbf{f}_w (evaluated on a test set) when perturbing the parameters in the directions given by the right singular vectors of the Jacobian (first 50 directions) or in randomly sampled directions (last 50 directions) on a VGG11 network trained for 10 epochs on CIFAR10. We observe that perturbations in most directions have almost no effect, except in those aligned with the top singular vectors.

602 The eigenvalue decompositions of \mathbf{G}_w and \mathbf{K}_w follow from the (SVD) of Φ_w . Assuming $P > nc$,
 603 we can write this SVD by indexing the singular values by a pair $J = (i, y)$ with $i = 1, \dots, n$ and
 604 $y = 1 \dots c$ as $\Phi_w = \sum_{J=1}^{nc} \sqrt{\lambda_J} \mathbf{u}_J \mathbf{v}_J^\top$. Such decompositions summarize the predominant directions
 605 both in parameter and feature space, in the neighborhood of w . Indeed, A small variation δw around
 606 w induces the first order variation $\delta \mathbf{f}_w$ of the function given by

$$\delta \mathbf{f}_w := \Phi_w \delta w = \sum_{J=1}^{nc} \sqrt{\lambda_J} (\mathbf{v}_J^\top \delta w) \mathbf{u}_J \quad (31)$$

607 Fig.9 illustrates this ‘hierarchy’ for a VGG11 network [55] trained for 10 epoches on CIFAR10 [33].
 608 We observe that perturbations in most directions have almost no effect, except in those aligned with
 609 the top singular vectors. This is reflected by a strong anisotropy of tangent kernel spectrum.

610 B Spectral Bias

611 We spell out some more detail for the content of Section 2.

612 B.1 Proof of Lemma 2

613 We consider parameter updates $\delta w_{\text{GD}} := -\eta \nabla_w L$ for gradient descent w.r.t the loss L . Using the
 614 chain rule, we can also write,

$$\delta w_{\text{GD}} = -\eta \Phi_w^\top (\nabla_{\mathbf{f}_w} L) \quad (32)$$

615 **Theorem 5** (Lemma 2 restated). *The gradient descent function updates in first order Taylor*
 616 *approximation, $\delta f_{\text{GD}}(\mathbf{x}) := \langle \delta w_{\text{GD}}, \Phi_w(\mathbf{x}) \rangle$, decompose as,*

$$\delta f_{\text{GD}}(\mathbf{x}) = \sum_{j=1}^n \delta f_j \tilde{u}_{w_j}(\mathbf{x}), \quad \delta f_j = -\eta \lambda_{w_j} (\mathbf{u}_{w_j}^\top \nabla_{\mathbf{f}_w} L) \quad (33)$$

617 *in terms of the kernel principal components \tilde{u}_{w_j} defined by (10).*

618 *Proof.* This follows immediately from (32), the SVD of Φ_w , and the definition (10):

$$\delta f_{\text{GD}}(\mathbf{x}) = -\eta \langle (\nabla_{\mathbf{f}_w} L)^\top \Phi_w, \Phi_w(\mathbf{x}) \rangle = \sum_{j=1}^n \delta f_j \tilde{u}_{w_j}(\mathbf{x}) \quad (34)$$

619 □

620 B.2 The case of linear regression

621 In this case $L = \frac{1}{2} \|\mathbf{f}_w - \mathbf{y}\|^2$ with $f_w = \langle w, \Phi(\mathbf{x}) \rangle$ (setting of Section 4.1), we can make the ‘spectral
 622 bias’ more explicit. A straightforward consequence of (9) is that the linear system governing the
 623 training dynamics in function space decouple in the basis of kernel principal components.

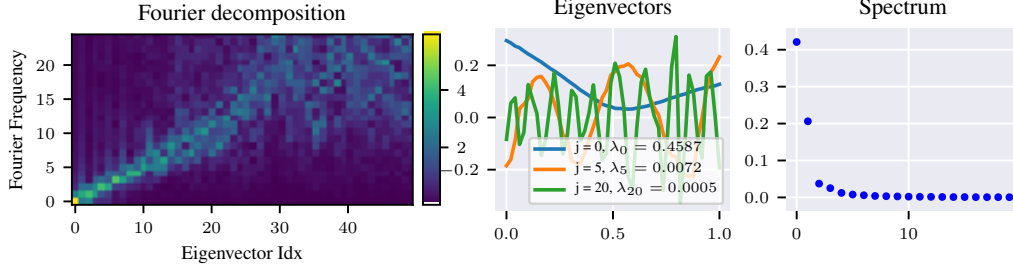


Figure 10: Eigendecomposition of the tangent kernel matrix of a random 6-layer deep 256-unit wide MLP on 1D uniform data (50 equally spaced points in $[0, 1]$). **(left)** Fourier decomposition (y -axis for frequency, colorbar for magnitude) of each eigenvector (x -axis). We observe that eigenvectors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

624 Gradient descent yields the function iterates,

$$f_{\mathbf{w}_t} = f_{\mathbf{w}^*} + (\text{id} - \eta K)^t (f_{\mathbf{w}_0} - f_{\mathbf{w}^*}) \quad (35)$$

625 where id is the identity map and K is the operator acting on functions as $(Kf)(\mathbf{x}) =$
 626 $\sum_{i=1}^n k(\mathbf{x}, \mathbf{x}_i) f(\mathbf{x}_i)$ in terms of the kernel $k(\mathbf{x}, \tilde{\mathbf{x}}) = \langle \Phi(\mathbf{x}), \Phi(\tilde{\mathbf{x}}) \rangle$.

627 *Proof.* The updates (32) induce the functional updates $\delta f_{\text{GD}} = f_{\mathbf{w}_{t+1}} - f_{\mathbf{w}_t}$ given by

$$\delta f_{\text{GD}}(\mathbf{x}) = -\eta \sum_{i=1}^n k(\mathbf{x}, \mathbf{x}_i) (f_{\mathbf{w}_t}(\mathbf{x}_i) - \mathbf{y}_i) \quad (36)$$

628 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. 35 follows by
 629 induction. \square

630 The operator K has eigenvalues $\lambda_1, \dots, \lambda_n$ with eigenfunctions $\tilde{u}_j(\mathbf{x})$ given by (10).

631 *Proof.* We can write $\tilde{u}_j(\mathbf{x}) = \sum_{i=1}^n k(\mathbf{x}, \mathbf{x}_i) u_{ji}$ where $\mathbf{u}_j = [u_{j1} \dots u_{jn}]^\top$ are the eigenvectors of
 632 \mathbf{K} . Observe that $(K\tilde{u}_j)(\mathbf{x}) = \sum_{i=1}^n k(\mathbf{x}, \mathbf{x}_i) (\mathbf{K}\mathbf{u}_j)_i = \sum_{i=1}^n k(\mathbf{x}, \mathbf{x}_i) (\lambda_j u_{ji}) = \lambda_j \tilde{u}_j$. Conversely,
 633 if λ is an eigenvalue of K with eigenfunction \tilde{u} , consider the vector $\mathbf{u} = [\tilde{u}(\mathbf{x}_1) \dots \tilde{u}(\mathbf{x}_n)]^\top$. Since
 634 $\lambda u_i = \tilde{u}(\mathbf{x}_i) = (K\tilde{u})(\mathbf{x}_i) = (\mathbf{K}\mathbf{u})_i$, \mathbf{u} is an eigenvector of \mathbf{K} and λ is one of the λ_j . \square

635 [Spectral Bias for Linear Regression] By initializing $\mathbf{w}_0 = \Phi^\top \alpha_0$ in the span of the features, the
 636 function iterates in Equ.35 uniquely decompose as $f_{\mathbf{w}_t}(\mathbf{x}) = \sum_{j=1}^n f_{jt} \tilde{u}_j(\mathbf{x})$ with

$$f_{jt} - f_j^* = (1 - \eta \lambda_j)^t (f_{j0} - f_j^*) \quad (37)$$

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

638 C Complexity Bounds

639 C.1 Rademacher Complexity

640 Given a family $\mathcal{G} \subset \mathbb{R}^{\mathcal{Z}}$ of real-valued functions on a probability space (\mathcal{Z}, ρ) , the empirical
 641 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 [38]:

$$\widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{G}) = \mathbb{E}_{\sigma \in \{\pm 1\}^n} \left[\sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i g(\mathbf{z}_i) \right], \quad (38)$$

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

644 **C.2 Generalization Bounds**

645 Generalization bounds based on Rademacher complexity are standard [7, 38]. We give here one
646 instance of such a bound, relevant for classification task.

647 **Setup.** We consider a family \mathcal{F} of functions $f_{\mathbf{w}}: \mathcal{X} \rightarrow \mathbb{R}^c$ that output a score or probability $f_{\mathbf{w}}(\mathbf{x})[y]$
648 for each class $y \in \{1 \dots c\}$ (we take $c = 1$ for binary classification). The task is to find a predictor
649 $f_{\mathbf{w}} \in \mathcal{F}$ with small expected classification error, which can be expressed e.g. as

$$L_0(f_{\mathbf{w}}) = \mathbb{P}_{(\mathbf{x}, y) \sim \rho} \{ \mu(f_{\mathbf{w}}(\mathbf{x}), y) < 0 \} \quad (39)$$

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

$$\mu(f(\mathbf{x}), y) = \begin{cases} f(\mathbf{x})y & \text{binary case} \\ f(\mathbf{x})[y] - \max_{y' \neq y} f(\mathbf{x})[y'] & \text{multiclass case} \end{cases} \quad (40)$$

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

$$\ell_{\gamma}(f_{\mathbf{w}}(\mathbf{x}), y) = \phi_{\gamma}(\mu(f_{\mathbf{w}}(\mathbf{x}), y)) \quad (41)$$

652 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
653 $\phi(u) = 1 - u/\gamma$ otherwise. We have the following bound for the expected error (39). With probability
654 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 $f_{\mathbf{w}} \in \mathcal{F}$ [38,
655 Theorems 4.4. and 8.1]:

$$L_0(f_{\mathbf{w}}) \leq \widehat{L}_{\gamma}(f_{\mathbf{w}}) + 2\widehat{\mathcal{R}}_{\mathcal{S}}(\ell_{\gamma}(\mathcal{F}, \cdot)) + 3\sqrt{\frac{\log \frac{2}{\delta}}{2n}} \quad (42)$$

656 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**,

$$\ell_{\gamma}(\mathcal{F}, \cdot) = \{(\mathbf{x}, y) \mapsto \ell_{\gamma}(f_{\mathbf{w}}(\mathbf{x}), y) \mid f_{\mathbf{w}} \in \mathcal{F}\} \quad (43)$$

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

$$\mathcal{R}_{\mathcal{S}}(\ell_{\gamma}(\mathcal{F}, \cdot)) \leq \frac{1}{\gamma} \mathcal{R}_{\mathcal{S}}(\mathcal{F}) \quad (44)$$

658 by Talagrand's contraction lemma [36] (see e.g. Mohri et al. [38, lemma 4.2] for a detailed proof).

659 **C.3 Complexity Bounds: Proofs**

660 We first derive standard bounds for the linear families (70) of scalar functions ($c = 1$):

$$\mathcal{F}_{M_A}^A = \{f_{\mathbf{w}}: \mathbf{x} \mapsto \langle \mathbf{w}, \Phi(\mathbf{x}) \rangle \mid \|\mathbf{w}\|_A \leq M_A\} \quad (45)$$

661 **Theorem 6.** *The empirical Rademacher complexity of $\mathcal{F}_{M_A}^A$ is bounded as,*

$$\widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_{M_A}^A) \leq (M_A/n) \sqrt{\text{Tr} \mathbf{K}_A} \quad (46)$$

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

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

$$\begin{aligned} \sup_{f \in \mathcal{F}_{M_A}^A} \sum_{i=1}^n \sigma_i f(\mathbf{x}_i) &= \sup_{\|\mathbf{w}\|_A \leq M_A} \sum_{i=1}^n \sigma_i \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle \\ &= \sup_{\|A^{\top} \mathbf{w}\|_2 \leq M_A} \sum_{i=1}^n \sigma_i \langle A^{\top} \mathbf{w}, A^{-1} \Phi(\mathbf{x}_i) \rangle \\ &= \sup_{\|\tilde{\mathbf{w}}\|_2 \leq M_A} \langle \tilde{\mathbf{w}}, \sum_{i=1}^n \sigma_i A^{-1} \Phi(\mathbf{x}_i) \rangle \\ &= M_A \left\| \sum_{i=1}^n \sigma_i A^{-1} \Phi(\mathbf{x}_i) \right\|_2 \\ &= M_A \sqrt{\sigma^{\top} \mathbf{K}_A \sigma} \end{aligned} \quad (47)$$

664 From (47) and the definition (38) we obtain:

$$\begin{aligned}\widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_{M_A}^A) &= \frac{M_A}{n} \mathbb{E}_{\sigma} \left[\sqrt{\sigma^{\top} \mathbf{K}_A \sigma} \right] \\ &\leq \frac{M_A}{n} \sqrt{\mathbb{E}_{\sigma} [\sigma^{\top} \mathbf{K}_A \sigma]} \\ &\leq \frac{M_A}{n} \sqrt{\text{Tr} \mathbf{K}_A}\end{aligned}\quad (48)$$

665 where we used Jensen's inequality to pass \mathbb{E}_{σ} under the root, and the properties that $\mathbb{E}[\sigma_i] = 0$ and
666 $\sigma_i^2 = 1$ for all i . \square

667 We now extend the result to the families (16) of learning flows:

$$\mathcal{F}_m^A = \{f_{\mathbf{w}} : \mathbf{x} \mapsto \sum_t \langle \delta \mathbf{w}_t, \Phi(\mathbf{x}) \rangle \mid \|\delta \mathbf{w}_t\|_{A_t} \leq m_t\} \quad (49)$$

668 **Theorem 7** (Theorem 3 restated). *The empirical Rademacher complexity of \mathcal{F}_m^A is bounded as,*

$$\widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_m^A) \leq \sum_t (m_t/n) \sqrt{\text{Tr} \mathbf{K}_{A_t}} \quad (50)$$

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

670 *Proof.* This is simple extension of the previous proof:

$$\begin{aligned}\sup_{f \in \mathcal{F}_m^A} \sum_{i=1}^n \sigma_i f(\mathbf{x}_i) &= \sup_{\|\delta \mathbf{w}_t\|_{A_t} \leq m_t} \sum_{i=1}^n \sigma_i \sum_t \langle \delta \mathbf{w}_t, \Phi(\mathbf{x}_i) \rangle \\ &= \sum_t \sup_{\|\tilde{\delta} \mathbf{w}_t\|_2 \leq m_t} \langle \tilde{\delta} \mathbf{w}_t, \sum_{i=1}^n \sigma_i A_t^{-1} \Phi(\mathbf{x}_i) \rangle \\ &= \sum_t m_t \sqrt{\sigma^{\top} \mathbf{K}_{A_t} \sigma}\end{aligned}\quad (51)$$

671 and we conclude as in (48). \square

672 Finally, we note that the same result can be formulated in terms of an evolving feature map $\Phi_t = A_t^{-1} \Phi$
673 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
674 can also be written as $\delta f_{\mathbf{w}_t}(\mathbf{x}) = \langle \tilde{\delta} \mathbf{w}_t, \Phi_t(\mathbf{x}) \rangle$ where $\tilde{\delta} \mathbf{w}_t = A_t^{\top} \delta \mathbf{w}_t$. The function class (16) can
675 equivalently be written as $\mathcal{F}_m^A = \mathcal{F}_m^{\Phi}$ where Φ denotes a fixed sequence of feature maps, $\Phi = \{\Phi_t\}_t$
676 and

$$\mathcal{F}_m^{\Phi} = \{f_{\mathbf{w}} : \mathbf{x} \mapsto \sum_t \langle \tilde{\delta} \mathbf{w}_t, \Phi_t(\mathbf{x}) \rangle \mid \|\tilde{\delta} \mathbf{w}_t\|_2 \leq m_t\} \quad (52)$$

677 In this formulation, Theorem 3 becomes:

678 **Theorem 3bis.** *The empirical Rademacher complexity of \mathcal{F}_m^{Φ} is bounded as,*

$$\widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_m^{\Phi}) \leq \sum_t (m_t/n) \sqrt{\text{Tr} \mathbf{K}_t} \quad (53)$$

679 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 .

680 C.4 Bounds for Multiclass Classification

681 The generalization bound (42) is based on the **margin loss class** (43). In this section, we show how
682 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 $f_{\mathbf{w}} : \mathcal{X} \rightarrow \mathbb{R}^c$
683 instead. Although the proof is adapted from standard techniques, to our knowledge Lemma C.4 and
684 Theorem 8 below are new results. In what follows, we denote by $\mu_{\mathcal{F}}$ the margin class,

$$\mu_{\mathcal{F}} = \{(\mathbf{x}, y) \mapsto \mu(f_{\mathbf{w}}(\mathbf{x}), y) \mid f_{\mathbf{w}} \in \mathcal{F}\} \quad (54)$$

685 where $\mu(f_{\mathbf{w}}(\mathbf{x}), y)$ is the margin (40). We also define, for each $y \in \{1 \dots c\}$,

$$\mathcal{F}_y = \{\mathbf{x} \mapsto f_{\mathbf{w}}(\mathbf{x})[y] \mid f_{\mathbf{w}} \in \mathcal{F}\}, \quad \mu_{\mathcal{F}, y} = \{\mathbf{x} \mapsto \mu(f_{\mathbf{w}}(\mathbf{x}), y) \mid f_{\mathbf{w}} \in \mathcal{F}\} \quad (55)$$

686 The following inequality holds:

$$\widehat{\mathcal{R}}_{\mathcal{S}}(\ell_{\gamma}(\mathcal{F}, \cdot)) \leq \frac{c}{\gamma} \sum_{y=1}^c \widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_y) \quad (56)$$

687

688 *Proof.* We first follow the first steps of the proof of Mohri et al. [38, Theorem 8.1] to show that

$$\widehat{\mathcal{R}}_{\mathcal{S}}(\ell_{\gamma}(\mathcal{F}, \cdot)) \leq \frac{1}{\gamma} \sum_{y=1}^c \widehat{\mathcal{R}}_{\mathcal{S}}(\mu_{\mathcal{F}, y}) \quad (57)$$

689 We reproduce these steps here for completeness: first, it follows from the $1/\gamma$ -Lipschitzness of the
690 ramp loss ϕ_{γ} in (41) and Talagrand's contraction lemma [38, lemma 4.2] that

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

691 Next, we write

$$\begin{aligned} \widehat{\mathcal{R}}_{\mathcal{S}}(\mu_{\mathcal{F}}) &:= \frac{1}{n} \mathbb{E}_{\sigma} \left[\sup_{f_{\mathbf{w}} \in \mathcal{F}} \sum_{i=1}^n \sigma_i \mu(f_{\mathbf{w}}(\mathbf{x}_i), y_i) \right] \\ &= \frac{1}{n} \mathbb{E}_{\sigma} \left[\sup_{f_{\mathbf{w}} \in \mathcal{F}} \sum_{i=1}^n \sigma_i \sum_{y=1}^c \mu(f_{\mathbf{w}}(\mathbf{x}_i), y) \delta_{y, y_i} \right] \\ &= \frac{1}{n} \sum_{y=1}^c \mathbb{E}_{\sigma} \left[\sup_{f_{\mathbf{w}} \in \mathcal{F}} \sum_{i=1}^n \sigma_i \mu(f_{\mathbf{w}}(\mathbf{x}_i), y) \delta_{y, y_i} \right] \end{aligned} \quad (59)$$

692 where $\delta_{y, y_i} = 1$ if $y = y_i$ and 0 otherwise; the second inequality follows from the sub-additivity of
693 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

$$\begin{aligned} \widehat{\mathcal{R}}_{\mathcal{S}}(\mu_{\mathcal{F}}) &\leq \frac{1}{2n} \sum_{y=1}^c \mathbb{E}_{\sigma} \left[\sup_{f_{\mathbf{w}} \in \mathcal{F}} \sum_{i=1}^n (\epsilon_i \sigma_i) \mu(f_{\mathbf{w}}(\mathbf{x}_i), y) \right] + \frac{1}{2n} \sum_{y=1}^c \mathbb{E}_{\sigma} \left[\sup_{f_{\mathbf{w}} \in \mathcal{F}} \sum_{i=1}^n \sigma_i \mu(f_{\mathbf{w}}(\mathbf{x}_i), y) \right] \\ &= \sum_{y=1}^c \frac{1}{n} \mathbb{E}_{\sigma} \left[\sup_{f_{\mathbf{w}} \in \mathcal{F}} \sum_{i=1}^n \sigma_i \mu(f_{\mathbf{w}}(\mathbf{x}_i), y) \right] \\ &= \sum_{y=1}^c \widehat{\mathcal{R}}_{\mathcal{S}}(\mu_{\mathcal{F}, y}) \end{aligned} \quad (60)$$

694 Together with (58), this leads to (57).

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

$$\begin{aligned} \widehat{\mathcal{R}}_{\mathcal{S}}(\mu_{\mathcal{F}, y}) &= \frac{1}{n} \mathbb{E}_{\sigma} \left[\sup_{f_{\mathbf{w}} \in \mathcal{F}} \sum_{i=1}^n \sigma_i (f_{\mathbf{w}}(\mathbf{x}_i)[y] - \max_{y' \neq y} f_{\mathbf{w}}(\mathbf{x}_i)[y']) \right] \\ &= \widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_y) + \frac{1}{n} \mathbb{E}_{\sigma} \left[\sup_{f_{\mathbf{w}} \in \mathcal{F}} \sum_{i=1}^n (-\sigma_i) \max_{y' \neq y} f_{\mathbf{w}}(\mathbf{x}_i)[y'] \right] \\ &= \widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_y) + \frac{1}{n} \mathbb{E}_{\sigma} \left[\sup_{f_{\mathbf{w}} \in \mathcal{F}} \sum_{i=1}^n \sigma_i \max_{y' \neq y} f_{\mathbf{w}}(\mathbf{x}_i)[y'] \right] \\ &\leq \widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_y) + \widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{G}_y) \end{aligned} \quad (61)$$

696 where $\mathcal{G}_y = \{\max\{f_{y'} : y' \neq y\} \mid f_{y'} \in \mathcal{F}_{y'}\}$. Now Mohri et al. [38, lemma 8.1] show that $\widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{G}_y) \leq$
697 $\sum_{y' \neq y} \widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_{y'})$. This leads to

$$\begin{aligned} \sum_{y=1}^c \widehat{\mathcal{R}}_{\mathcal{S}}(\mu_{\mathcal{F}, y}) &\leq \sum_{y=1}^c \widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_y) + \sum_{y=1}^c \sum_{\substack{y'=1 \\ y' \neq y}}^c \widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_{y'}) \\ &= \sum_{y=1}^c \widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_y) + (c-1) \sum_{y=1}^c \widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_y) \\ &= c \sum_{y=1}^c \widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_y) \end{aligned} \quad (62)$$

698 Substituting in (57) finishes the proof. \square

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

$$\mathcal{F}_{M_A}^A = \{\mathbf{x} \mapsto f_{\mathbf{w}}(\mathbf{x})[y] := \langle \mathbf{w}, \Phi(\mathbf{x})[y] \rangle \mid \|\mathbf{w}\|_A \leq M_A\} \quad (63)$$

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

702 **Theorem 8.** *The emp. Rademacher complexity of the margin loss class $\ell_\gamma(\mathcal{F}_{M_A}^A, \cdot)$ is bounded as,*

$$\widehat{\mathcal{R}}_{\mathcal{S}}(\ell_\gamma(\mathcal{F}_{M_A}^A, \cdot)) \leq (c^{3/2} M_A / \gamma n) \sqrt{\text{Tr} \mathbf{K}_A} \quad (64)$$

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

704 *Proof.* Eq.56, and Theorem 8 applied to each linear family \mathcal{F}_y of (scalar) functions leads to

$$\widehat{\mathcal{R}}_{\mathcal{S}}(\ell_\gamma(\mathcal{F}_{M_A}^A, \cdot)) \leq \frac{c}{\gamma} \sum_{y=1}^c \frac{M_A}{n} \sqrt{\text{Tr} \mathbf{K}_A^{yy}} \quad (65)$$

705 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
706 average $\frac{1}{c} \sum_{y=1}^c$ under the root using Jensen inequality, we conclude:

$$\begin{aligned} \widehat{\mathcal{R}}_{\mathcal{S}}(\ell_\gamma(\mathcal{F}_{M_A}^A, \cdot)) &\leq \frac{c^2 M_A}{\gamma n} \sqrt{\frac{1}{c} \sum_{y=1}^c \text{Tr} \mathbf{K}_A^{yy}} \\ &= \frac{c^{3/2} M_A}{\gamma n} \sqrt{\text{Tr} \mathbf{K}_A} \end{aligned} \quad (66)$$

707

□

708 C.5 Linear models: Which Norm for Measuring Capacity?

709 We consider a family \mathcal{F} of scalar functions $f_{\mathbf{w}}(\mathbf{x}) = \langle \mathbf{w}, \Phi(\mathbf{x}) \rangle$ linearly parametrized by $\mathbf{w} \in \mathbb{R}^P$,
710 where Φ is a fixed mapping of the input space \mathcal{X} into \mathbb{R}^P . Given a training set \mathcal{S} of size n , we denote
711 by $\Phi = [\Phi(\mathbf{x}_1), \dots, \Phi(\mathbf{x}_n)]^\top$ the $n \times P$ feature matrix and by $\mathbf{y} = [y_1 \dots y_n]^\top$ the label vector. We
712 are interested in the ‘overparametrized’ regime: we assume $P \geq n$. We write the SVD of the feature
713 matrix as $\Phi = \sum_{j=1}^n \sqrt{\lambda_j} \mathbf{u}_j \mathbf{v}_j^\top$, where $\lambda_1 \geq \dots \geq \lambda_n$ are ranked in nonincreasing order. We will
714 consider the minimum ℓ^2 norm interpolators [27],

$$\mathbf{w}^* = \Phi^\top \mathbf{K}^{-1} \mathbf{y} = \sum_{j=1}^n \frac{\mathbf{u}_j^\top \mathbf{y}}{\sqrt{\lambda_j}} \mathbf{v}_j \quad (67)$$

715 A standard approach is to measure capacity in terms of the ℓ^2 norm the weight vector. Considering

$$\mathcal{F}_M = \{f_{\mathbf{w}} : x \mapsto \langle \mathbf{w}, \Phi(x) \rangle \mid \|\mathbf{w}\|_2 \leq M\}, \quad (68)$$

716 the Rademacher complexity of \mathcal{F}_M can be bounded as [6, Lemma 22]:

$$\widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_M) \leq (M/n) \|\Phi\|_{\text{F}} \quad (69)$$

717 where $\|\Phi\|_{\text{F}}$ is the Froebenius norm of the feature matrix.¹⁵

718 Is the ℓ^2 norm a good capacity measure, even for algorithms biased towards low ℓ^2 norm solutions?
719 If the distribution of solutions $\mathbf{w}_{\mathcal{S}}^*$, where $\mathcal{S} \sim \rho^n$, is reasonably isotropic, taking the smallest ℓ^2 ball
720 containing them (with high probability) gives an accurate description of the class of trained models.
721 However for very anisotropic distributions, the solutions do not fill any such ball so describing trained
722 models in terms of ℓ^2 balls is wasteful [53]. For the minimum ℓ^2 norm interpolators (67), the solution
723 distribution typically inherits the anisotropy of the features. For example, if $y_i = \bar{y}(\mathbf{x}_i) + \varepsilon_i$ where
724 $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$, the covariance of the solutions with respect to noise is $\text{cov}_{\varepsilon}[\mathbf{w}^*, \mathbf{w}^*] = \sum_j \frac{\sigma^2}{\lambda_j} \mathbf{v}_j \mathbf{v}_j^\top$,
725 which scales as $1/\lambda_j$ along \mathbf{v}_j .

¹⁵Note that $\|\Phi\|_{\text{F}} = \sqrt{\text{Tr} \mathbf{K}}$ where $\mathbf{K} = \Phi \Phi^\top$ is the kernel matrix.

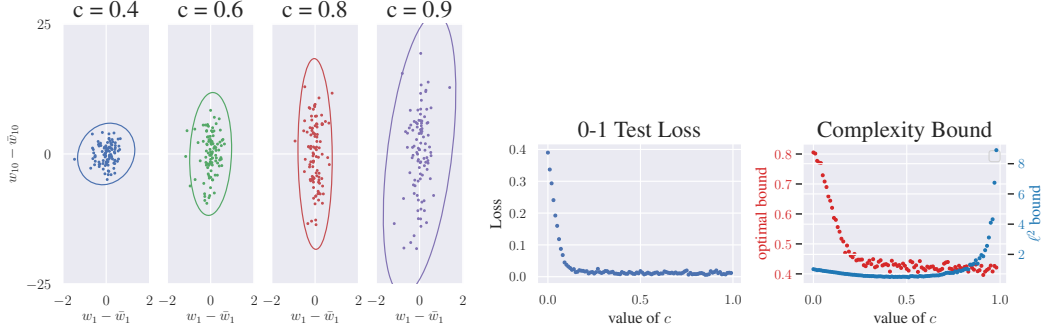


Figure 11: **Left:** 2D projection of the minimum- ℓ^2 -norm interpolators \mathbf{w}_S^* , $\mathcal{S} \sim \rho^n$, for linear models $f_{\mathbf{w}} = \langle \mathbf{w}, \Phi_c \rangle$, as the feature scaling factor varies from 0 (white features) to 1 (original, anisotropic features). For larger c , the solutions scatter in a very anisotropic way. **Right:** Average test classification loss and complexity bounds (69) (blue plot) for the solution vectors \mathbf{w}_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 (71) optimized as in Proposition 9 (red plot) does not suffer from this problem.

726 To visualize this on a simple setting, consider P random Fourier features [49], fit on 1D data \mathbf{x}
727 modelled by N equally spaced points in $[-a, a]$. In this setting, the (true) feature map is represented
728 by a $N \times P$ matrix with SVD $\Phi = \sum_j \sqrt{l_j} \psi_j \varphi_j^\top$. The labels are given by $y(\mathbf{x}) = \text{sign}(\psi_1(\mathbf{x}))$. To
729 highlight the effect of feature anisotropy, we further rescale the singular values as $l_j^c = 1 + c(l_j - 1)$ so
730 as to interpolate between whitened features ($c=0$) and the original ones ($c=1$). We set $P=N=1000$.
731 Fig 11 (left) shows 2D projections in the plane $(\varphi_1, \varphi_{10})$ of (centered) solutions $\mathbf{w}_S^* - \mathbb{E}_S \mathbf{w}_S^*$, for
732 a pool of 100 (sub)samples \mathcal{S} of size $n = 50$, for increasing values of the scaling factor c . As c
733 approaches 1, the solutions begin to scatter in a very anisotropic way in parameter space; as shown in
734 Fig 11 (right), the complexity bound (69) (blue plot) becomes increasingly loose and fails to reflect
735 the shape of the test error.

736 We emphasize that this issue is about the choice of norm and not about complexity-based bounds *per se*.
737 In fact, note that anisotropies can in principle be compensated by a suitable linear reparametrization
738 $\mathbf{w} \mapsto A^\top \mathbf{w}$, $\Phi \mapsto A^{-1} \Phi$. Any such A can be viewed as defining a new norm $\|\mathbf{w}\|_A := \sqrt{\mathbf{w}^\top g_A \mathbf{w}}$
739 induced by the metric $g_A = AA^\top$. The following classes

$$\mathcal{F}_{M_A}^A = \{f_{\mathbf{w}} : \mathbf{x} \mapsto \langle \mathbf{w}, \Phi(\mathbf{x}) \rangle \mid \|\mathbf{w}\|_A \leq M_A\}, \quad (70)$$

740 define a much richer set of complexity classes than (68), represented by ellipsoids of all shapes in
741 parameter space. A direct extension of the standard result (69) yields:

$$\widehat{\mathcal{R}}_{\mathcal{S}}(\mathcal{F}_{M_A}^A) \leq (M_A/n) \|A^{-1} \Phi^\top\|_{\text{F}} \quad (71)$$

742 in terms of the Froebenius norm of the rescaled feature matrix.¹⁶ More meaningful norms than
743 the ℓ^2 norm can be found by optimizing the bound (71) with $M_A = \|\mathbf{w}^*\|_A$, over a given class of
744 reparametrization matrices A . We give an example of this in the following Proposition.

745 **Proposition 9.** Consider the class of reparametrization matrices $A_\nu = \sum_{j=1}^n \sqrt{\nu_j} \mathbf{v}_j \mathbf{v}_j^\top +$
746 $\mathbb{1}_{\text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_n\}^\perp}$, which act as mere rescaling $\lambda_j \rightarrow \lambda_j/\nu_j$ of the singular values of the feature
747 matrix. Any minimizer of (71) for the minimum ℓ^2 -norm interpolator takes the form

$$\nu_j^* = \kappa \frac{\sqrt{\lambda_j}}{|\mathbf{v}_j^\top \mathbf{w}^*|} = \kappa \frac{\lambda_j}{|\mathbf{u}_j^\top \mathbf{y}|} \quad (72)$$

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

749 Note that in the context of Proposition 9, the optimal norm $\|\cdot\|_{A_{\nu^*}}$ depends both on the feature
750 geometry – through the singular values – and on the task – through the labels –. As shown in Fig 1
751 (right, red plot), in the random Fourier feature setting, the corresponding bound has a much nicer
752 behaviour than the standard bound (69) based on the ℓ^2 norm.

¹⁶We also have $\|A^{-1} \Phi^\top\|_{\text{F}} = \sqrt{\text{Tr} \mathbf{K}_A^{-1}}$ where $\mathbf{K}_A = \Phi g_A^{-1} \Phi^\top$ is the rescaled kernel matrix.

753 **D CKA and Spectral Entropy**

754 We make explicit a couple of metrics used in Section 3.

755 **Centered kernel alignment (CKA).** We used CKA [18] to measure the similarity between tangent
 756 features and labels. Given two kernel matrices \mathbf{K} and \mathbf{K}' in $\mathbb{R}^{r \times r}$, it is defined as

$$\text{CKA}(\mathbf{K}, \mathbf{K}') = \frac{\text{Tr}[\mathbf{K}_c \mathbf{K}'_c]}{\|\mathbf{K}_c\|_F \|\mathbf{K}'_c\|_F} \in [0, 1] \quad (73)$$

757 where the c subscript denotes the feature centering operation, i.e. $\mathbf{K}_c = \mathbf{C} \mathbf{K} \mathbf{C}$ where $\mathbf{C} = \mathbf{I}_r - \frac{1}{r} \mathbf{1} \mathbf{1}^T$
 758 is the centering matrix. The normalization by the Froebenius norm makes CKA invariant under
 759 isotropic rescaling.

760 Let $\mathbf{Y} \in \mathbb{R}^{nc}$ be the vector resulting from the concatenation of the one-hot label representations
 761 $\mathbf{Y}_i \in \mathbb{R}^c$ of the n samples. Similarity with the labels is measured through CKA with the rank-one
 762 kernel $\mathbf{K}_Y := \mathbf{Y} \mathbf{Y}^T$,

$$\text{CKA}(\mathbf{K}, \mathbf{K}_Y) = \frac{\mathbf{Y}^T \mathbf{K}_c \mathbf{Y}}{\|\mathbf{K}_c\|_F \|\mathbf{K}_Y\|_F} \quad (74)$$

763 **Effective rank.** We used a notion of effective rank based on **spectral entropy** [50]. Given a kernel
 764 matrix \mathbf{K} with (strictly) positive eigenvalues $\lambda_1, \dots, \lambda_n$, let

$$\mu_j = \lambda_j / \text{Tr} \mathbf{K}, \quad \text{Tr} \mathbf{K} = \sum_{i=1}^n \lambda_i \quad (75)$$

765 be the trace-normalized eigenvalues. The effective rank is defined as [50]:

$$\text{erank} = \exp(H(\mu_1, \dots, \mu_n)) \quad (76)$$

766 where $H(\mu)$ is the Shannon entropy given by

$$H(\mu_1, \dots, \mu_n) = - \sum_{j=1}^n \mu_j \log(\mu_j) \quad (77)$$

767 This effective rank is a real number between 1 and n , upper bounded by $\text{rank}(\mathbf{K})$, which measures
 768 the ‘uniformity’ of the spectrum through the entropy.

769 **E Experiments: Details and Additions**

770 **E.1 Synthetic Experiment of Fig 3**

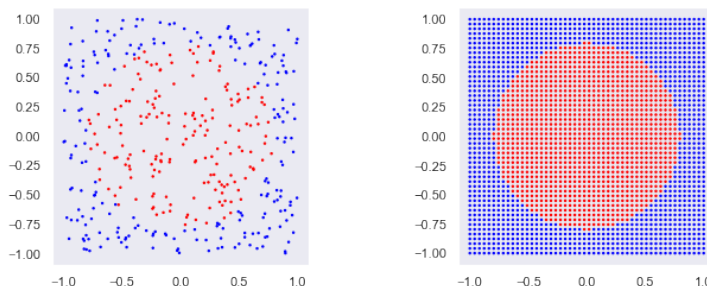


Figure 12: Disk dataset. **Left:** Training set of $n = 500$ points (\mathbf{x}_i, y_i) where $\mathbf{x} \sim \text{Unif}[-1, 1]^2$, $y_i = 1$ if $\|\mathbf{x}_i\|_2 \leq r = \sqrt{2/\pi}$ and -1 otherwise. **Right:** Large test sample (2500 points forming a 50×50 grid) used to evaluate the tangent kernel. In our experiment, we trained a 6-layer deep 256-unit wide MLP by full batch gradient descent of binary cross entropy.

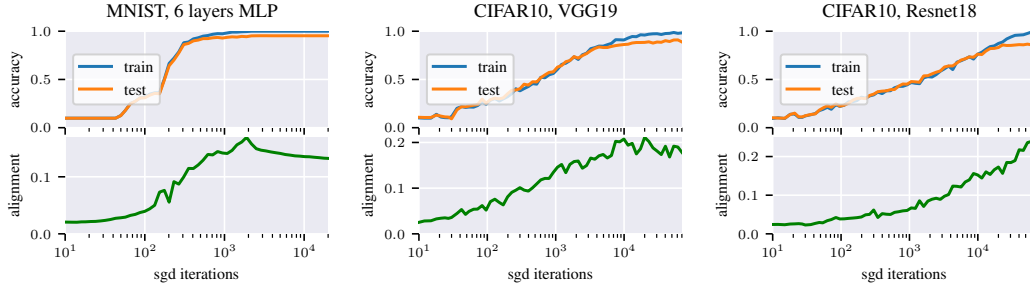


Figure 13: Evolution of the CKA between the tangent kernel and the class label kernel $K_Y = YY^T$ measured on a held-out test set for different architectures: **(left)** 6 layers of 80 hidden units MLP on MNIST **(middle)** VGG19 on CIFAR10 **(right)** Resnet18 on CIFAR10. We observe an increase of the alignment to the target function.

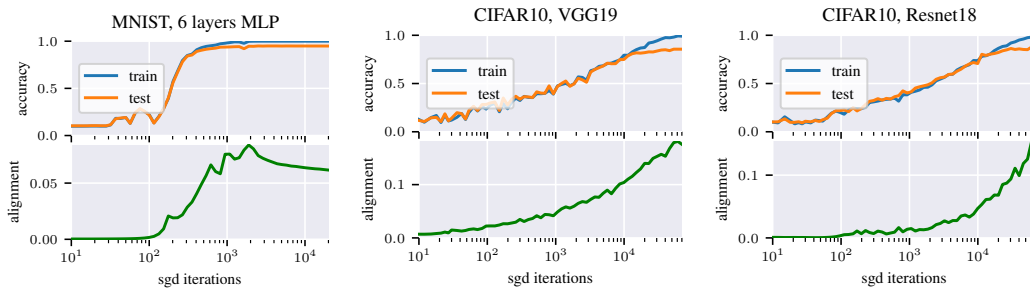


Figure 14: Same as figure 13 but without centering the kernel. Evolution of the uncentered kernel alignment between the tangent kernel and the class label kernel $K_Y = YY^T$ measured on a held-out test set for different architectures: **(left)** 6 layers of 80 hidden units MLP on MNIST **(middle)** VGG19 on CIFAR10 **(right)** Resnet18 on CIFAR10. We observe an increase of the alignment to the target function.

771 **E.2 More alignment plots**

772 **E.3 More plots on spectra**

773 **E.4 Ablation: Effect of depth on alignment**

774 In order to study the influence of the architecture on the alignment effect, we measure the CKA
 775 for different networks and different initialization as we increase the depth. The results in Fig 16
 776 suggest that the alignment effect is magnified as depth increases. We also observe that the ratio of the
 777 maximum alignment between easy and difficult examples is increased with depth, but stays high for a
 778 smaller number of iterations.

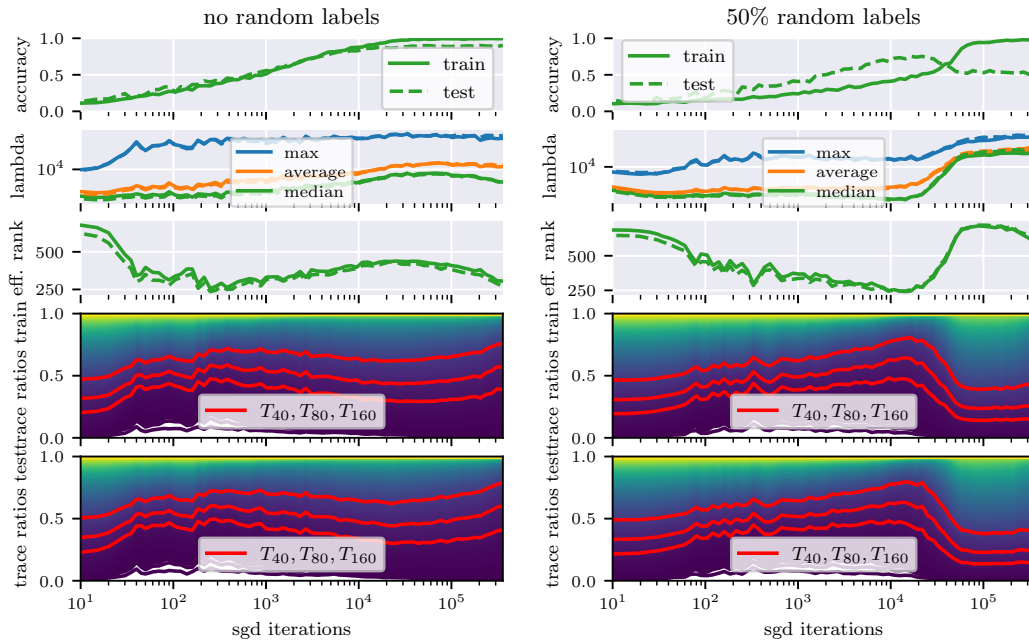


Figure 15: Evolution of tangent kernel spectrum, effective rank and trace ratios of a VGG19 trained by SGD with batch size 100, learning rate 0.003 and momentum 0.9 on dataset **(left)** CIFAR10 and **(right)** CIFAR10 with 50% random labels. We highlight the top 40, 80 and 160 trace ratios in **red**. The small effective rank of the kernel biases the training procedure towards a few top eigenvectors, as can also be observed by remarking that the trace ratio T_{40} account for $\sim 50\%$ of the total trace.

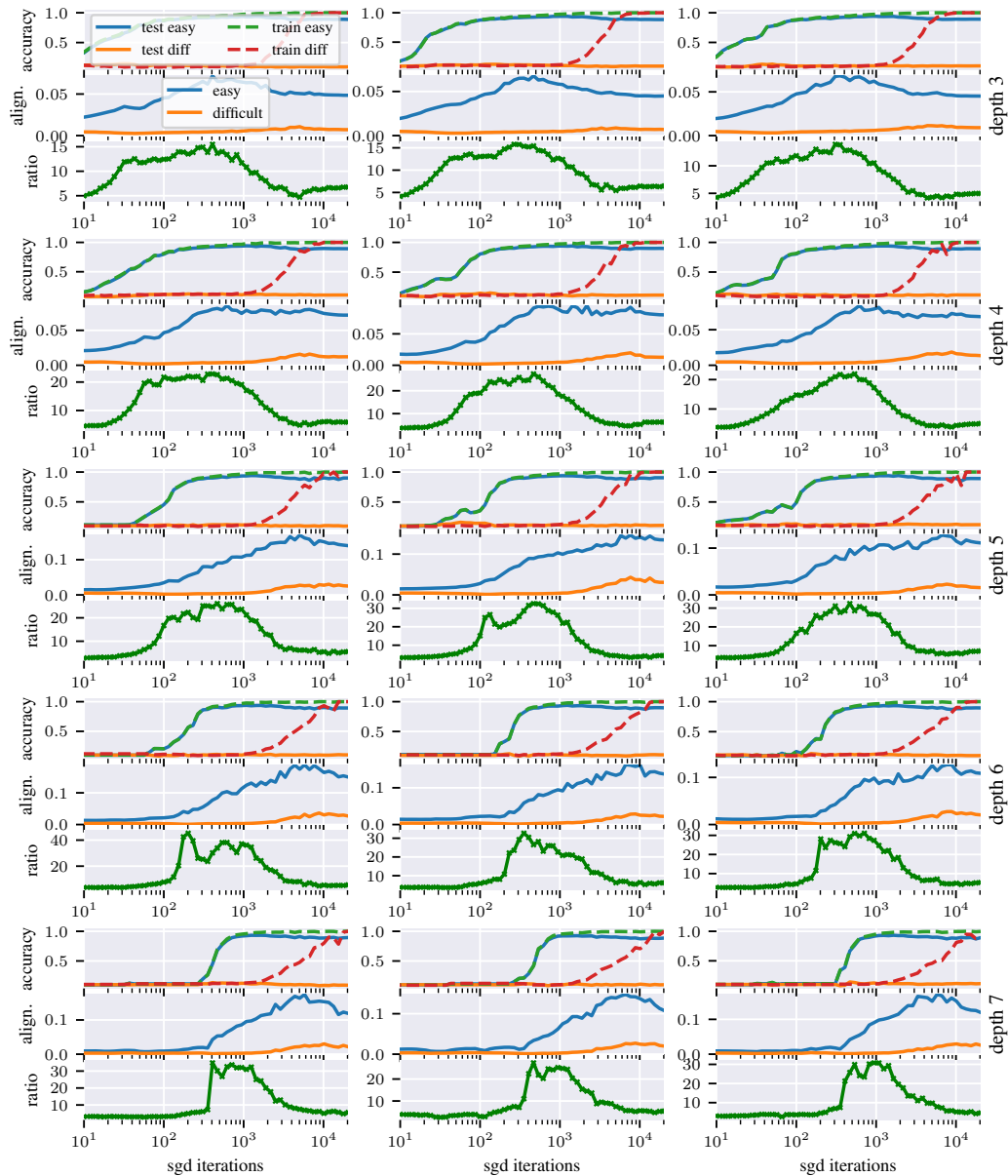


Figure 16: Effect of depth on alignment. 10,000 MNIST examples with 1000 random labels MNIST examples trained with learning rate=0.01, momentum=0.9 and batch size=100 for MLP with hidden layers size 60 and **(in rows)** varying depths **(in columns)** varying random initialization/minibatch sampling. 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.