

Random features and polynomial rules

Fabián Aguirre-López^{1,2,3*}, Silvio Franz¹ and Mauro Pastore^{1,4†}

¹ Université Paris-Saclay, CNRS, LPTMS, 91405 Orsay, France

² Chair of Econophysics and Complex Systems, École polytechnique, 91128 Palaiseau, France

³ LadHyX UMR CNRS 7646, École polytechnique, 91128 Palaiseau, France

⁴ Laboratoire de physique de l'École normale supérieure, CNRS, PSL University, Sorbonne University, Université Paris-Cité, 24 rue Lhomond, 75005 Paris, France

* fabian.aguirre-lopez@polytechnique.edu, † mauro.pastore@phys.ens.fr

Abstract

Random features models play a distinguished role in the theory of deep learning, describing the behavior of neural networks close to their infinite-width limit. In this work, we present a thorough analysis of the generalization performance of random features models for generic supervised learning problems with Gaussian data. Our approach, built with tools from the statistical mechanics of disordered systems, maps the random features model to an equivalent polynomial model, and allows us to plot average generalization curves as functions of the two main control parameters of the problem: the number of random features N and the size P of the training set, both assumed to scale as powers in the input dimension D . Our results extend the case of proportional scaling between N , P and D . They are in accordance with rigorous bounds known for certain particular learning tasks and are in quantitative agreement with numerical experiments performed over many order of magnitudes of N and P . We find good agreement also far from the asymptotic limits where $D \rightarrow \infty$ and at least one between P/D^K , N/D^L remains finite.

Copyright attribution to authors.

This work is a submission to SciPost Physics.

License information to appear upon publication.

Publication information to appear upon publication.

Received Date

Accepted Date

Published Date

1

2 Contents

3	1 Introduction	2
4	1.1 Related works	4
5	2 The model	6
6	3 Generalization error	8
7	4 Kernel learning and polynomial models	9
8	5 Replica calculation	11
9	5.1 Replica symmetric theory	13
10	5.2 Saddle-point equations for quadratic loss	15
11	6 Strongly separated regimes	16

12	7 Effective theory for finite-size random features networks	16
13	8 Conclusions and perspectives	18
14	A Kernel on the Hermite basis	19
15	B Hermite polynomials and Wick products	19
16	C Evaluation of the moments of ν, λ^a	20
17	D Results on random matrix theory	21
18	D.1 Marchenko-Pastur distribution and Stieltjes transformation	21
19	D.2 Spectral density of $C^{\otimes \ell}$	22
20	E Determinant of sum of matrices with orthogonal row spaces	24
21	F Traces over RS matrices	24
22	G Replica-symmetric free energy	25
23	G.1 Measure contribution	25
24	G.2 Pattern contribution	26
25	H Asymptotic limits of the saddle-point equations	26
26	H.1 Case (i)	26
27	H.2 Case (ii)	27
28	H.3 Case (iii)	28
29	I Numerical experiments	28
30	References	29

31
32

33 1 Introduction

34 The connection between deep feed-forward neural networks (DNNs) in the large-width limit
 35 and kernel methods has been well understood in the last years. It has been shown, in a
 36 Bayesian learning perspective, that if the number of units in each hidden layer is taken to
 37 infinity at fixed input dimension and training set size, a DNN becomes a “neural network
 38 Gaussian process” whose kernels can be defined iteratively layer by layer [1–4]. This result
 39 has been recently generalized beyond the infinite-width limit [5–10]. In a dynamical perspec-
 40 tive moreover, it has been shown that wide DNNs trained with gradient-based methods exhibit
 41 a the lazy-training kernel regime [11], evaluated by first order Taylor-expanding the network
 42 with respect to the weights around initialization [12–14].

43 Once a DNN is proven equivalent to a kernel machine, the mechanism by which it realizes
 44 the input-output mapping of the corresponding supervised-learning task is understood: the
 45 input data, which generally speaking are points in \mathbb{R}^D , are mapped with an implicit *feature*
 46 *map* $\psi : \mathbb{R}^D \rightarrow \mathbb{R}^N$ to an N -dimensional space where the classification, or regression, rule is
 47 linear and can be learnt by the read-out layer. The mapping to the feature space is implicit,
 48 in the sense that the learning problem can be solved by a support vector machine (SVM), so
 49 that learning and generalization depend on the features only through the kernel $\tilde{\mathcal{H}}(\mathbf{x}, \mathbf{x}') =$

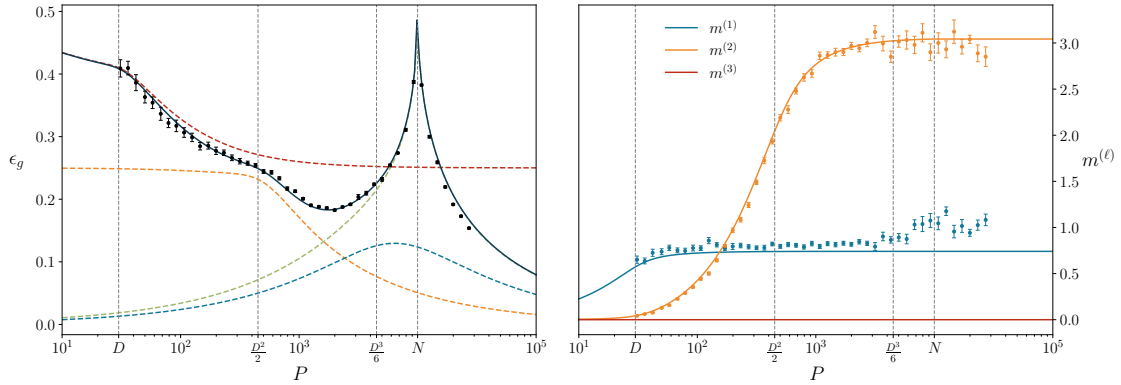


Figure 1: **Left:** generalization error of the RFM on a classification task, as a function of the size of the training set P , for $D = 30$, $N = 10^4$, weights regularization $\zeta = 10^{-8}$, quadratic teacher (balanced: $\tau_1 = \tau_2 = 1/\sqrt{2}$, $\tau_{\ell>2} = 0$) and ELU activation functions; the continuous line is the equivalent polynomial theory devised in Sec. 4, truncated at $L = 3$; dashed lines are the asymptotic theories (see Sec 6 for details) for $N \rightarrow \infty$ and P/D finite (red), $N \rightarrow \infty$ and $P/\binom{D}{2}$ finite (yellow), $N \rightarrow \infty$ and $P/\binom{D}{3}$ finite (blue), $P/\binom{D}{3}$ and N/P finite (green); black points are results from numerical experiments averaged over 50 instances (see Appendix I). The model learns the linear features (first step at $P \sim O(D)$), then learns the quadratic features (second step at $P \sim O(D^2)$), then follows the interpolation peak at $P \sim N$. **Right:** numerical and theoretical teacher-student overlaps – defined in Eq. (35) and (43) – of the linear and quadratic features (the overlap of the cubic features is identically 0 by definition); the parameters of the model are the same as for the left panel.

50 $\sum_{i=1}^N \psi_i(\mathbf{x})\psi_i(\mathbf{x}')/N$ (see, for reference, [15]). Learning curves (generalization error as a
 51 function of the size P of the training set) of kernel machines can be obtained analytically from
 52 a statistical mechanics [16–19] or a mathematical [20–22] perspective. A very interesting
 53 trait of these curves is their staircase shape for $P \sim D^K$: by setting the scaling of the size of the
 54 training set to a certain power K of the input dimension, features of order K can be learnt by
 55 the machine, so that the test error decreases increasing K with subsequent steps.

56 The discovery of the lazy training regime of wide neural networks motivated in the recent
 57 past the study of the *random features model* (RFM) [23,24], a shallow (one-hidden-layer, 1HL)
 58 neural network where the feature map is explicitly parametrized by a fixed random linear
 59 embedding of the input points from \mathbb{R}^D to \mathbb{R}^N , followed by a non-linear activation function. In
 60 this sense, the model mimics the behavior of a neural network in the large-width limit, where
 61 the feature map depends only on initialization and learning is linear.

62 In the present work we study theoretically the generalization performance of the RFM in
 63 the large- D limit, with $P \sim D^K$, $N \sim D^L$. We find, under a quite general teacher/student
 64 setting with a random polynomial teacher and Gaussian i.i.d. input data, that

- 65 • as long as $P \ll N$, the model behaves as an infinite-rank ($N \rightarrow \infty$) kernel machine:
 66 for $P \sim D^K$, features of order K can be learnt, such that the generalization error as a
 67 function of P has a staircase descent (or overfitting peaks if the teacher is less complex)
 68 with steps corresponding to different values of K ;
- 69 • for $P \gg N$ and $N \sim D^L$, the model is equivalent to a degree- L polynomial student: if the
 70 complexity of the teacher is lower than the degree L , the generalization error is equal
 71 to zero, or otherwise, to the minimum error for a degree- L polynomial fitting a more
 72 complex teacher;

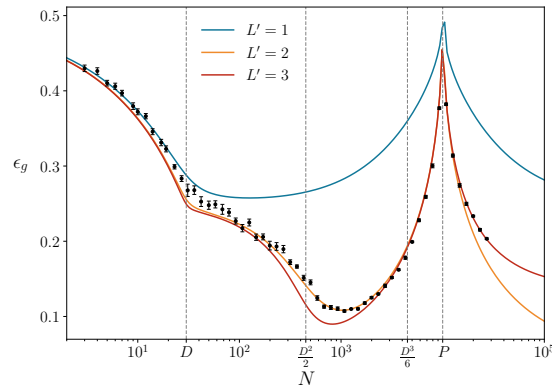


Figure 2: Generalization error of a RFM on a classification task, as a function of the number of hidden units N , for $P = 10^4$ and the rest of the parameters as in Fig. 1; continuous lines are the theories truncated at $L' = 1, 2, 3$ (respectively: blue, yellow, red); numerical points (in black) are nicely interpolating between these curves in the regimes where $N \sim O(D), O(D^2), O(D^3)$, validating Eq. (24), where the truncation L' of the equivalent polynomial theory is fixed at $L \sim \log(N)/\log(D)$.

- for $P \sim N$, an interpolation peak of the generalization error, which depends on the strength of the regularization of the student’s weights, occurs.

This behavior is depicted in Fig. 1. Comparison with numerical experiments shows that our theory, based on the mapping of the RFM to an *equivalent noisy polynomial model*, predicts well the quantitative behavior of the true generalization performance at finite size, over many orders of magnitude.

Our theory, formulated from the point of view of the statistical mechanics of disordered systems, expresses the generalization performance of the RFM in terms of few order parameters with a clear physical interpretation, as overlaps between combinations of the student’s weights and the parameters defining the teacher. In this way, we are offering a complementary take on what is known about RFMs in the computer science community, as we discuss in the following.

1.1 Related works

In this section we give an overview on the previous works that have been of inspiration to our paper, presenting relevant results and differences with our approach.

Random feature models were introduced in [23–26], initially as randomized low-rank approximations of kernels arising in classification or regression problems. Recently, their interest was renewed by the discovery that DNNs behaves as RFMs close to the infinite-width limit, both in a Bayesian learning [1–4] and in a gradient-based learning [11–14] setting. This mapping, which provides one of the few limits where DNNs can be studied with analytical methods, has motivated in the last few years a huge effort to formalize their behavior in terms of expressive power and generalization performance.

In particular, the impressive series of works [14, 27–33] (see [34] for a review) formulates rigorous bounds on the generalization performance of RFMs in different asymptotic regimes. For a non-exhaustive recap of the results (with our notation):

- In [27], the large- D limits where $D^{L+\delta} \leq N \leq D^{L+1-\delta}$ (for small δ) after sending $P \rightarrow \infty$ (underparametrized regime) and $D^{K+\delta} \leq P \leq D^{K+1-\delta}$ after sending $N \rightarrow \infty$ (overparametrized regime) are considered. In the first case the model is found equivalent to

101 degree- L polynomial regression; in the second one, it reduces to (infinite-rank) kernel
 102 regression, which for that number of samples can fit at most a degree- K polynomial in
 103 the inputs, in a way also investigated in literature [16–22].

- 104 • In [29], the limit where both N and P scale linearly with D with their ratio fixed is
 105 considered; the generalization error as a function of the ratio between the number of
 106 hidden units and the size of the training set first decreases for N/P small, then exhibits
 107 a peak at the interpolation threshold $N/P = 1$ and then relaxes again for $N \gg P$ to
 108 the value predicted from the kernel theory with $P \sim D$, coherently with the previous
 109 point. This phenomenology is widely observed in numerical experiments and known in
 110 literature as *double descent* [35] of the generalization error.
- 111 • In [31], the authors push forward the analysis of [27] (that is, P and N scaling poly-
 112 nomially with D) to the regimes where $N \leq P^{1-\delta}$ and $N \geq P^{1+\delta}$. The authors show
 113 indeed that the limiting behavior is given by the smallest of N and P , and they find the
 114 interpolation threshold at $N \sim P$ also in this polynomial scaling.
- 115 • In [33], universality results on training and test error are proven in the $P \sim N$ regime
 116 for a larger class of models, as long as with finite-dimensional outputs, and generic
 117 losses. Indeed, they prove that training and test errors depend on the random features
 118 distribution only through its covariance structure.

119 These papers find bounds to the generalization performance of a RFM with rigorous analytical
 120 methods under quite general assumptions on data distribution and activation functions.

121 A statistical mechanics point of view, complementary to the formal approach discussed
 122 so far, has been formulated in the series of papers [36–42]. Originally aiming at modelling
 123 the role of data structure in machine learning, as in other contemporary approaches [43–50],
 124 the authors obtained in [37] a closed-form expression for the generalization error of RFMs
 125 for regression and classification in the asymptotic regime where $N \sim P \sim D$. Their approach,
 126 based on the replica theory from statistical mechanics [51], can be applied to supervised learn-
 127 ing tasks with generic convex loss functions. Not only their results are supported under mild
 128 hypothesis by analytical proofs [29, 33, 38, 52, 53], but they can predict remarkably well the nu-
 129 merical experiments. Our work extends these results to more general scaling regimes, where
 130 $P \sim D^K$, $N \sim D^L$.

131 One of the main steps in our derivation is the expansion of activation function of the hidden
 132 layer on a polynomial basis, which corresponds to the diagonalization of the kernel (20) on
 133 its eigenbasis (Mercer’s decomposition). This expansion is then truncated to a certain degree
 134 L , corresponding to the integer exponent in the scaling law $N \sim D^L$: similar approximations
 135 appeared recently in [54, 55]. Moreover, while the literature on the double descent behav-
 136 ior of the generalization error is vast and impossible to outline here (see for example [35]),
 137 we mention [56], where the presence of more than one peak in the generalization curve is
 138 remarked: the authors call “linear peak” the one occurring at $P \sim D$ for $N \gg P$, where the
 139 model behaves as a kernel learning the linear features, while for $P \sim N$ there is a “non-linear
 140 peak” due to the non-linearity of the activation function acting as noise and overfitted when
 141 P and N are of the same order; in the present work we show that, as long as $N \gg P$, there is
 142 a peak (or a descent) for each of the regimes $P \sim D^K$.

143 Appeared in parallel with our work, the paper [57] pushes forward the line of research
 144 of [29] from a mathematical perspective, deriving sharp asymptotics for the generalization of
 145 random features ridge regression in the polynomial regime.

Symbol	Definition
D	input space dimension
$N \sim D^L$	feature space dimension
$P \sim D^K$	size of the training set
B	degree of the teacher
n	number of replicas
η_ℓ	$N/\binom{D}{\ell}$
α, β, \dots	indices in input space
i, j, \dots	indices in feature space
μ, ν, \dots	indices spanning the training set
a, b, \dots	indices in replica space
α	multi-index $\{\alpha_1, \dots, \alpha_\ell\}$, $\alpha_1 < \dots < \alpha_\ell$
θ	teacher parameters, $\theta = \{\theta_\alpha^{(\ell)}\}_{\ell=1}^B$
F	$N \times D$ random features matrix
$\mathbf{F}_\alpha, \mathbf{F}_i$	$(F_{i\alpha})_{i=1}^N, (F_{i\alpha})_{\alpha=1}^D$
$\mathbf{F}_\alpha^{\otimes \ell}$	$(F_{i\alpha_1} \dots F_{i\alpha_\ell})_{i=1}^N$
C	FF^\top/D
$C^{\otimes \ell}$	$((C_{ij})_{i,j=1}^N)^\ell \simeq \sum_\alpha \mathbf{F}_\alpha^{\otimes \ell} (\mathbf{F}_\alpha^{\otimes \ell})^\top / \binom{D}{\ell}$
$Q, Q^{(\ell)}, \dots$	$(Q_{ab})_{a,b=1}^n, (Q_{ab}^{(\ell)})_{a,b=1}^n, \dots$

Table 1: Notations used in this paper

146 2 The model

147 We would like to study the generalization performance of the Random Features model in a
 148 teacher/student [58, 59] supervised learning set-up, where the teacher performs an input-
 149 output mapping with various degree of complexity. We summarize in Table 1 the main nota-
 150 tions used in this paper.

151 The input data \mathbf{x} are vectors in \mathbb{R}^D with i.i.d. Gaussian elements, while the labels are
 152 assigned by a polynomial teacher of degree B defined as:

$$y \sim p(y | \nu(\mathbf{x})),$$

$$\nu(\mathbf{x}) = \sum_{\ell=1}^B \frac{\tau_\ell}{\sqrt{\binom{D}{\ell}}} \sum_{\alpha_1 < \dots < \alpha_\ell} \theta_{\alpha_1 \dots \alpha_\ell}^{(\ell)} x_{\alpha_1} \dots x_{\alpha_\ell}, \quad (1)$$

153 where $\theta_\alpha^{(1)}, \theta_{\alpha\beta}^{(2)}, \dots$ are i.i.d. $\mathcal{N}(0, 1)$ parameters collectively denoted as θ , describing the
 154 non-linear decision boundary (diagonal terms, irrelevant for large D , are for simplicity not
 155 included in the sum). Notice that the function $\nu(\mathbf{x})$ coincide with the Hamiltonian of the
 156 “mixed p -spin model” of the statistical physics of the spin-glasses (see, for example, [60]).
 157 The mixture parameters τ_ℓ , weighting the monomials of different degree, are chosen to respect
 158 $\sum_{\ell=1}^B \tau_\ell^2 = 1$. Within this general setting, we will concentrate on the specific simple examples
 159 of a deterministic teacher for binary classification or a noisy teacher for polynomial regression
 160 with variance of the noise Δ , for which Eq. (1) reduces respectively to

$$y \sim \delta[y - \text{sgn } \nu(\mathbf{x})], \quad y \sim \mathcal{N}[y | \nu(\mathbf{x}), \Delta]. \quad (2)$$

161 It has been shown in [16] that a *polynomial* student, defined in the same way as in Eq. (1),
 162 would learn the weights of the teacher in a hierarchical fashion: $O(D^K)$ examples are needed
 163 in order to learn the parameters $\theta^{(\ell)}$ for $\ell \leq K$. However, here the student’s task is to learn

164 the weights of the last layer of a 2-layers NN, $f(\mathbf{x}; \mathbf{w})$, whose first layer realizes a random
 165 embedding of the data in a N -dimensional feature space:

$$f(\mathbf{x}; \mathbf{w}) = \phi[\lambda(\mathbf{x}; \mathbf{w})], \quad (3)$$

$$\lambda(\mathbf{x}; \mathbf{w}) = \frac{1}{\sqrt{N}} \sum_{i=1}^N w_i \sigma\left(\frac{1}{\sqrt{D}} \sum_{\alpha=1}^D F_{i\alpha} x_\alpha\right) \quad (4)$$

166 where F is a $N \times D$ quenched random matrix with i.i.d. standard normal entries, σ is the non-
 167 linear activation function of the hidden layer, $\mathbf{w} \in \mathbb{R}^N$ the student's weight vector and ϕ the
 168 activation function of the last ("readout") layer. It is customary to introduce the pre-activations

$$h_i = \frac{1}{\sqrt{D}} \sum_{\alpha=1}^D F_{i\alpha} x_\alpha, \quad (5)$$

169 which at fixed instance of the random features F , given that we chose x_α i.i.d normal variables,
 170 follow a multivariate Gaussian distribution with covariance

$$C_{ij} = \mathbb{E}_{\mathbf{x}^\mu} [h_i h_j] = \frac{1}{D} \sum_{\alpha=1}^D F_{i\alpha} F_{j\alpha}. \quad (6)$$

171 In our setting with independent random features, C is a Wishart matrix.

172 While our theory is general in the choice of σ that we will suppose it can be expanded in
 173 Hermite polynomials (see Sec. 4). We will test our results for popular choices, such as

$$\sigma(h) = \text{ReLU}(h) = \max(h, 0), \quad (7)$$

$$\sigma(h) = \text{ELU}(h) = \begin{cases} \exp(h) - 1 & \text{if } h < 0, \\ h & \text{if } h \geq 0, \end{cases} \quad (8)$$

174 (respectively, Rectified and Exponential Linear Unit).

175 The training set is made of P input-output pairs, $\mathcal{T} = \{(\mathbf{x}^\mu, y^\mu)\}_{\mu=1}^P$. The student learns by
 176 solving the following optimization problem,

$$\mathbf{w}^* = \underset{\mathbf{w}}{\text{argmin}} \left[\sum_{\mu=1}^P \mathcal{L}[y^\mu, \lambda(\mathbf{x}^\mu; \mathbf{w})] + \frac{\zeta}{2} \|\mathbf{w}\|^2 \right], \quad (9)$$

177 where \mathcal{L} is an opportune convex loss function and ζ controls the regularization of the weights.
 178 The choice of the loss function \mathcal{L} and the readout activation function ϕ in Eq. (3) defines the
 179 specific learning task to perform. To simplify calculations we will mostly look at the cases of
 180 of a pure quadratic loss, reading, both in the case of regression and classification:

$$\mathcal{L}(y, \lambda) = \frac{1}{2} (y - \lambda)^2, \quad (10)$$

181 the use of a regression loss for a classification task dates back to the early days of NNs [59,61].

182 The main aim of this work is the evaluation of the generalization performance of the model,
 183 both for the classification and the regression problems, using a statistical mechanics approach.
 184 From this perspective, the model defines a disordered system with N degrees of freedom \mathbf{w} ,
 185 and quenched disorder given by the realization of the input points \mathbf{x}^μ , the teacher's parameters
 186 θ and the random features F . Our computation will follow the standard path, starting from
 187 the computation partition function at inverse temperature β

$$\mathcal{Z} = \int d\mathbf{w} \exp \left[-\beta \sum_{\mu=1}^P \mathcal{L}[y^\mu, \lambda(\mathbf{x}^\mu; \mathbf{w})] - \frac{\beta\zeta}{2} \|\mathbf{w}\|^2 \right]. \quad (11)$$

188 3 Generalization error

189 In order to quantify how well the student can learn the teacher, we look at the generalization
 190 error, defined as the probability of misclassifying a new sample (in the case of classification)
 191 or as the mean squared error of a new point (in the case of regression). Given a test point
 192 $(\mathbf{x}, y) \sim p_0(\mathbf{x})p(y|\nu(\mathbf{x}))$, both cases can be expressed with the following formula,

$$\epsilon_g(\theta, F, \mathcal{T}) = \int d\mathbf{x} p_0(\mathbf{x}) \int dy p(y|\nu(\mathbf{x})) \frac{1}{4^\kappa} [y - \phi(\lambda(\mathbf{x}; \mathbf{w}_\mathcal{T}^*))]^2, \quad (12)$$

193 where $\kappa = 1$ for binary classification and $\kappa = 0$ for regression. Notice the presence of the
 194 function ϕ in the definition of the generalization error, at variance with the loss function.

195 With (12) we can evaluate the quality of the student NN (3) for a given realization of the
 196 teacher, of the random weights F , and of the dataset \mathcal{T} . In order to get a general view of
 197 the effectiveness of (3), we calculate the average generalization error over all the sources of
 198 randomness. Doing so, we get a function of N , P , and D only,

$$\begin{aligned} \epsilon_g &= \int d\nu d\lambda p(\nu, \lambda) \int dy p(y|\nu) \frac{1}{4^\kappa} [y - \phi(\lambda)]^2 \\ p(\nu, \lambda) &= \mathbb{E} \int d\mathbf{x} p_0(\mathbf{x}) \delta(\nu - \nu(\mathbf{x})) \delta(\lambda - \lambda(\mathbf{x}; \mathbf{w}_\mathcal{T}^*)), \end{aligned} \quad (13)$$

199 where we took $\mathbb{E} = \mathbb{E}_{\theta, F, \mathcal{T}}$.

200 We have written the average generalization error as in Eq. (13) to show that we only need
 201 to know the joint distribution of (ν, λ) to evaluate it. Being \mathbf{x} a test point, and so uncorrelated
 202 from \mathbf{w}^* , we will take the distribution $p(\nu, \lambda)$ as Gaussian: to compute the generalization error
 203 we only need the first and second moments,

$$\begin{aligned} 0 &= \mathbb{E}[\nu], & t^* &= \mathbb{E}[\lambda], \\ \rho &= \mathbb{E}[\nu^2], & m^* &= \mathbb{E}[\nu\lambda], & q^* &= \mathbb{E}[\lambda^2] - t^{*2}. \end{aligned} \quad (14)$$

204 Notice that by definition of the model (*i.e.* the normalization of the mixing parameters τ_ℓ)
 205 ρ is identically equal to 1. In section 5 we will show how to obtain this quantities from a
 206 replica approach. Central limit theorems for sums of non linear functions of Gaussian fields
 207 (the pre-activations (5) at given feature matrix F), of the kind we just used to motivate this
 208 ansatz, have been proven in the past under conditions on the realization of the feature-feature
 209 covariance matrix C and of the vector \mathbf{w}^* [38, 62, 63].

210 For the case of binary classification with $y = \text{sgn}(\nu)$ and $\phi = \text{sgn}$,

$$\begin{aligned} \epsilon_g &= \frac{1}{4} \mathbb{E} [y - \text{sgn}(\lambda)]^2 \\ &= \int_{-\infty}^0 D\nu \left[1 - H\left(\frac{t^* + m^* \nu}{\sqrt{q^* - m^{*2}}}\right) \right] + \int_0^\infty D\nu H\left(\frac{t^* + m^* \nu}{\sqrt{q^* - m^{*2}}}\right), \end{aligned} \quad (15)$$

211 where we use the Gardner's notation [58] $D\nu = \frac{e^{-\nu^2/2}}{\sqrt{2\pi}} d\nu$ and $H(x) = \int_x^\infty Dt$. Notice that
 212 when $t^* = 0$ (that is, when the student is zero-mean) the formula simplifies to

$$\epsilon_g = \frac{1}{\pi} \arccos\left(\frac{m^*}{\sqrt{q^*}}\right). \quad (16)$$

213 For the case of noisy polynomial regression, ($\phi = \text{id}$ and $\Delta = \mathbb{E}[(y - \nu)^2]$) [64, 65],

$$\epsilon_g = \mathbb{E}[y - \lambda]^2 = \rho + \Delta - 2m^* + q^* + t^{*2}. \quad (17)$$

214 These formulas remind the generalization error of a generalized linear model with the same
 215 architecture as the teacher [59]: in that case, $m^*/\sqrt{q^*}$ corresponds to the angle between
 216 the teacher and the student weight vectors. For the RFM, it is not clear *a priori* if we can
 217 interpret $m^*/\sqrt{q^*}$ as a scalar product of the teacher's weight vector and some effective weights
 218 of the student. If this can be done, the RFM could be mapped to an equivalent polynomial
 219 model. In Sec. 4 we will show how to explicitly construct it from \mathbf{w} and F , thus achieving
 220 this mapping. To do so, we need to spend a few words on the connection between RFMs and
 221 kernel machines, in order to explain the truncation of the activation function σ on the basis
 222 of Hermite polynomials, which we will use later on.

223 4 Kernel learning and polynomial models

224 The RFM defined in (3) is a generalized linear model in the learnable parameters \mathbf{w} , so it can
 225 be formulated as a kernel model, as we remind in this section. First of all, because of the
 226 choice of quadratic loss, we can write down the explicit solution to (9),

$$w_i^* = \sum_j \left(\zeta \mathbb{1}_N + \frac{P}{N} \bar{\mathcal{K}} \right)_{ij}^{-1} \frac{1}{\sqrt{N}} \sum_{\mu} y^{\mu} \sigma(h_j^{\mu}), \quad (18)$$

227 where the pre-activations h are given by (5) and the operator

$$\bar{\mathcal{K}}_{ij} = \frac{1}{P} \sum_{\mu} \sigma(h_i^{\mu}) \sigma(h_j^{\mu}) \quad (19)$$

228 defines the kernel in feature space. The properties of the kernel are crucial for the generaliza-
 229 tion performances.

230 While our analysis will be more general, for the purpose of arguing, in this section we
 231 consider the limit $P \rightarrow \infty$. In this case the kernel reduces to

$$\mathcal{K}_{ij} = \mathbb{E}_{\mathbf{x}^{\mu}} [\sigma(h_i^{\mu}) \sigma(h_j^{\mu})]. \quad (20)$$

232 From this formula, it is possible to obtain an explicit formula of the kernel \mathcal{K} as a function of the
 233 covariance matrix of the pre-activations (6). To this aim, as the pre-activations are Gaussian,
 234 it is convenient to expand the activation function on the basis of Hermite polynomials (see
 235 also [27]):

$$\sigma(h_i) = \sum_{\ell=0}^{\infty} \frac{\mu_{\ell}}{\ell!} \text{He}_{\ell}(h_i), \quad (21)$$

236 where He_{ℓ} is the ℓ -th Hermite polynomial and the coefficient μ_{ℓ} are:

$$\mu_{\ell} = \int Dx \text{He}_{\ell}(x) \sigma(x). \quad (22)$$

237 Along these lines, the kernel (20) can be expressed for large D [66, 67] (see App. A for
 238 details) as

$$\mathcal{K}_{ij} = \sum_{\ell=0}^{\infty} \frac{\mu_{\ell}^2}{\ell!} (C_{ij})^{\ell}, \quad (23)$$

239 where C_{ij} , given by (6), is a rank- D Wishart matrix with elements $C_{ii} = 1 + O(D^{-1/2})$ and
 240 $C_{ij} = O(D^{-1/2})$ for $i \neq j$. The matrix with entries $(C_{ij})^{\ell}$, which we denote by $C^{\circ\ell}$, defines an
 241 interesting random matrix ensemble, obtained taking Hadamard (element by element) powers
 242 of the covariance C . A similar ensemble was recently studied in [68].

243 Suppose now the relation between N and D is fixed: $N \sim D^{L+\delta}$ with $0 \leq \delta < 1$. The $N \times N$
 244 matrix $C^{\circ\ell}$ has generically rank equal to $\min\{D^\ell, N\}$ and off-diagonal elements $O(D^{-\ell/2})$. For
 245 $\ell > L$ the matrix is full ranked, the small off-diagonal terms give a vanishing contribution to
 246 eigenvalues and eigenvectors. We can then truncate the expansion substituting $C^{\circ\ell > L}$ by the
 247 identity matrix:

$$\mathcal{K}_{ij} \simeq \sum_{\ell=0}^L \frac{\mu_\ell^2}{\ell!} (C_{ij})^\ell + \mu_{\perp,L}^2 \delta_{ij}, \quad (24)$$

248 where

$$\mu_{\perp,L}^2 = \sum_{\ell=L+1}^{\infty} \frac{\mu_\ell^2}{\ell!}. \quad (25)$$

249 This truncation is proven for $L = 1$ (that is, in the proportional regime $N \sim D$) in [69], and
 250 extended to the case $L > 1$ under generic assumptions on the kernel \mathcal{K} in [31,55]. A convincing
 251 check of this property for moderately large values of N is given by Fig. 2, which shows the
 252 theoretical curves of the generalization error obtained through a truncated effective theory
 253 (that we describe below) at different values of L' , compared with the numerical experiments,
 254 as a function of N ; quantitative agreement is obtained for $L' = L \sim \log N / \log D$, with the
 255 numerical points interpolating nicely the theoretical curves in the various regimes.

256 The analysis above suggests that in the $N \sim D^L$ regime we can represent the RFM as an
 257 effective noisy polynomial student

$$\lambda_{\text{eff}}(\mathbf{x}^\mu; \mathbf{w}) = \mu_0 m^{(0)} + \sum_{\ell=1}^L \frac{\mu_\ell}{\sqrt{D}^\ell} \sum_{\alpha_1, \dots, \alpha_\ell} s_{\alpha_1 \dots \alpha_\ell}^{(\ell)} : x_{\alpha_1}^\mu \dots x_{\alpha_\ell}^\mu : + z^\mu, \quad (26)$$

258 where

- 259 • $m^{(0)} = \sum_i w_i / \sqrt{N}$ is the empirical mean of the vector \mathbf{w} , rescaled by \sqrt{N} ;
- 260 • the student parameters $s_{\alpha_1 \dots \alpha_\ell}^{(\ell)}$ are the scalar product of \mathbf{w} with the “vectors” $\mathbf{F}_{\alpha_1 \dots \alpha_\ell}^{\otimes \ell} / \sqrt{N}$
 261 with components $F_{i\alpha_1} \dots F_{i\alpha_\ell} / \sqrt{N}$ (see Table 1),

$$s_{\alpha_1 \dots \alpha_\ell}^{(\ell)} = \frac{1}{\sqrt{N}} \sum_i w_i F_{i\alpha_1} \dots F_{i\alpha_\ell}. \quad (27)$$

- 262 • we have written the expansion of the Hermite polynomials in terms of the so-called Wick
 263 products of the x 's, routinely used in theoretical physics and defined from the following
 264 generating function (see for example [70]):

$$\begin{aligned} :x_1 \dots x_k: &= \partial_{\lambda_1} \dots \partial_{\lambda_k} G(\boldsymbol{\lambda}; \mathbf{x}) \Big|_{\boldsymbol{\lambda}=0}, \\ G(\boldsymbol{\lambda}; \mathbf{x}) &= \frac{\exp(\boldsymbol{\lambda}^\top \mathbf{x})}{\mathbb{E}[\exp(\boldsymbol{\lambda}^\top \mathbf{x})]} = \exp(\boldsymbol{\lambda}^\top \mathbf{x} - \|\boldsymbol{\lambda}\|^2/2) \end{aligned} \quad (28)$$

265 (see App. B for more details). These quantities have the property $\mathbb{E}[:x_1 \dots x_k:] = 0$.

- 266 • the last term z^μ is a Gaussian noise term with zero mean and variance $\mathbb{E}(z^{\mu 2}(\mathbf{w})) =$
 267 $\mu_{\perp,L}^2 \sum_{i=1}^N w_i^2 / N$ which can be represented as

$$z^\mu = \frac{\mu_{\perp,L}}{\sqrt{N}} \sum_{i=1}^N w_i v_i^\mu, \quad (29)$$

268 in terms of i.i.d. $\mathcal{N}(0, 1)$ variables v_i^μ .

269 Although ultimately the parameters $\mathbf{s}^{(\ell)}$ and \mathbf{z} are functions on the network weights, to
270 enlighten the notation we will not explicitly write the dependence on \mathbf{w} .

271 In (26) we give an effective description of the RFM, mapping it to a polynomial model
272 with correlated weights in presence of a noise term coming from the $\ell > L$ terms in the expansion
273 (21). This is an extension to generic scaling regimes $N \sim D^L$ of the *Gaussian equivalence*
274 *principle* from [38] and related works, to which it reduces when $L = 1$. In the following, we
275 will base our analysis on this representation of λ . This description makes more transparent
276 the meaning of the observables introduced in Sec. 3 and the mechanism by which the RFM
277 learns the teacher's features, as we explain in the following.

278 5 Replica calculation

279 Let us now turn to the analysis of the general case through the replica method. To obtain the
280 generalization error we write the joint probability distribution of ν and λ in Eq. (13) as the
281 zero temperature limit of the equilibrium distribution of a statistical mechanics system, as

$$p(\nu, \lambda) = \lim_{\beta \rightarrow \infty} \mathbb{E} \int \frac{d\mathbf{w}}{\mathcal{Z}} e^{-\beta \sum_{\mu} \mathcal{L}[y^{\mu}, \lambda(\mathbf{x}^{\mu}; \mathbf{w})] - \frac{\beta \zeta}{2} \|\mathbf{w}\|^2} \int d\mathbf{x} p_0(\mathbf{x}) \delta(\nu - \nu(\mathbf{x})) \delta(\lambda - \lambda(\mathbf{x}; \mathbf{w})). \quad (30)$$

282 Through a standard application of the replica trick we rewrite the distribution as

$$p(\nu, \lambda) = \lim_{n \rightarrow 0} \lim_{\beta \rightarrow \infty} \mathbb{E} \int \prod_{a=1}^n d\mathbf{w}^a e^{-\beta \sum_{\mu, a} \mathcal{L}[y^{\mu}, \lambda(\mathbf{x}^{\mu}; \mathbf{w}^a)] - \frac{\beta \zeta}{2} \sum_a \|\mathbf{w}^a\|^2} \times \int d\mathbf{x} p_0(\mathbf{x}) \delta(\nu - \nu(\mathbf{x})) \delta(\lambda - \lambda(\mathbf{x}; \mathbf{w}^1)), \quad (31)$$

283 which can be obtained from the calculation of the n -times replicated partition function

$$Z_n = \mathbb{E}[Z^n] = \int \prod_{a=1}^n d\mathbf{w}^a e^{-\frac{\beta \zeta}{2} \sum_a \|\mathbf{w}^a\|^2} \mathbb{E}_{F, \theta} \left[\mathbb{E}_{\nu, \{\lambda^a\}} \int dy p(y|\nu) e^{-\beta \sum_a \mathcal{L}(y, \lambda^a)} \right]^P. \quad (32)$$

284 In this integral, we treat the distribution of ν and λ^a conditioned by F , θ and \mathbf{w}^a as Gaussian,
285 with moments given by

$$t_a = \mathbb{E}(\lambda_a | F, \theta), \quad M_a = \mathbb{E}(\nu \lambda_a | F, \theta), \quad Q_{ab} = \mathbb{E}(\lambda_a \lambda_b | F, \theta) - t_a t_b. \quad (33)$$

286 from which we can extract the generalization error according to (15), (17). Using the repre-
287 sentation (26) we can decompose these order parameters as (see Appendix C for details)

$$t_a = \mu_0 M_a^{(0)}, \quad M_a = \sum_{\ell=1}^{\min\{L, B\}} \frac{\mu_{\ell} \tau_{\ell}}{\sqrt{\ell!}} M_a^{(\ell)}, \quad Q_{ab} = \mu_{\perp, L}^2 Q_{ab}^{(0)} + \sum_{\ell=1}^L \frac{\mu_{\ell}^2}{\ell!} Q_{ab}^{(\ell)}, \quad (34)$$

288 with the definitions:

$$M_a^{(0)} = \frac{1}{\sqrt{N}} \sum_{i=1}^N w_i^a, \quad M_a^{(\ell)} = \frac{\boldsymbol{\theta}^{(\ell)} \cdot \mathbf{s}_a^{(\ell)}}{\binom{D}{\ell}}, \quad Q_{ab}^{(0)} = \frac{1}{N} \sum_{i=1}^N w_i^a w_i^b, \quad Q_{ab}^{(\ell)} = \frac{1}{N} \sum_{i, j=1}^N w_i^a C_{ij}^{\ell} w_j^b, \quad (35)$$

289 where we are using the notation

$$\boldsymbol{\theta}^{(\ell)} \cdot \mathbf{s}_a^{(\ell)} = \sum_{\alpha} \theta_{\alpha}^{(\ell)} s_{a, \alpha}^{(\ell)} \quad (36)$$

290 (remember that the sum over α is restricted to ordered tuples).

291 Enforcing these definition with delta functions in Fourier representation, and anticipating
 292 saddle point integration for the various M and Q , and their Fourier conjugated parameters
 293 that we denote as \hat{M} and \hat{Q} with the due indexes, we rewrite the partition function as

$$Z_n = e^{PS_P[Q,M]} e^{\frac{N}{2} \sum_{a,b} \hat{Q}_{ab}^{(0)} Q_{ab}^{(0)} + \frac{1}{2} \sum_{\ell,a,b} \binom{D}{\ell} \hat{Q}_{ab}^{(\ell)} Q_{ab}^{(\ell)} + \sum_{\ell,a} \binom{D}{\ell} \hat{M}_a^{(\ell)} M_a^{(\ell)}} \\ \times \mathbb{E}_{F,\theta} \int d\mathbf{w} e^{-\frac{1}{2} \mathbf{w}^\top \left[(\beta \zeta \mathbb{1}_n + Q^{(0)}) \otimes \mathbb{1}_N + \sum_{\ell} \hat{Q}^{(\ell)} \otimes \frac{C^{\otimes \ell}}{\eta_\ell} \right] \mathbf{w} - \sum_{\ell,i,a,\alpha} \hat{M}_a^{(\ell)} w_{i,a}^\alpha F_{i,a}^{\otimes \ell} \theta_a^{(\ell)} / \sqrt{\eta_\ell \binom{D}{\ell}}}, \quad (37)$$

294 where now $\mathbf{w} \in \mathbb{R}^{n \times N}$, the sums over ℓ span $\{1, \dots, L\}$, $\eta_\ell = N / \binom{D}{\ell}$ and

$$S_P[Q, M] = \log \mathbb{E}_{\nu, \{\lambda^a\}} \int dy p(y|\nu) e^{-\beta \sum_a \mathcal{L}(y, \lambda^a)}. \quad (38)$$

295 In writing Eq. (37), we took $\hat{M}_a^{(0)} \rightarrow 0$, as the Fourier conjugate of the mean t_a is suppressed
 296 in the large- N limit [71] (a property that could be checked *a posteriori* from the saddle point
 297 equation for $\hat{M}_a^{(0)}$); moreover, the conventional scalings with N and $\binom{D}{\ell}$ in this equation are
 298 chosen in such a way that the hat variables corresponding to the asymptotic regimes explained
 299 in Sec. 6 have a non-trivial high-dimensional limit.

300 Averaging over θ we obtain:¹

$$Z_n = e^{PS_P[Q,M]} e^{\frac{N}{2} \sum_{a,b} \hat{Q}_{ab}^{(0)} Q_{ab}^{(0)} + \frac{1}{2} \sum_{\ell,a,b} \binom{D}{\ell} \hat{Q}_{ab}^{(\ell)} Q_{ab}^{(\ell)} + \sum_{\ell,a} \binom{D}{\ell} \hat{M}_a^{(\ell)} M_a^{(\ell)}} \\ \times \mathbb{E}_F \int d\mathbf{w} e^{-\frac{1}{2} \mathbf{w}^\top \left[(\beta \zeta \mathbb{1}_n + \hat{Q}^{(0)}) \otimes \mathbb{1}_N + \sum_{\ell} (\hat{Q}^{(\ell)} - \hat{M}^{(\ell)} \hat{M}^{(\ell)\top}) \otimes \frac{C^{\otimes \ell}}{\eta_\ell} \right] \mathbf{w}} \quad (39)$$

301 and integrating over \mathbf{w} ,

$$Z_n = e^{PS_P[Q,M]} e^{\frac{N}{2} \sum_{a,b} \hat{Q}_{ab}^{(0)} Q_{ab}^{(0)} + \frac{1}{2} \sum_{\ell,a,b} \binom{D}{\ell} \hat{Q}_{ab}^{(\ell)} Q_{ab}^{(\ell)} + \sum_{\ell,a} \binom{D}{\ell} \hat{M}_a^{(\ell)} M_a^{(\ell)} - \frac{1}{2} \text{Tr} \log [A^{(0)} \otimes \mathbb{1}_N + \sum_{\ell} B^{(\ell)} \otimes C^{\otimes \ell}]}, \quad (40)$$

302 where traces are taken over replica and feature indices and we introduced for compactness
 303 the $n \times n$ matrices

$$A^{(0)} = \beta \zeta \mathbb{1}_n + \hat{Q}^{(0)}, \quad B^{(\ell)} = (\hat{Q}^{(\ell)} - \hat{M}^{(\ell)} \hat{M}^{(\ell)\top}) / \eta_\ell. \quad (41)$$

304 We notice at this point that, given $N \sim D^{L+\delta}$, for $\ell \leq L$ the matrices $C^{\otimes \ell}$ have rank $r_\ell =$
 305 $O(D^\ell) \ll N$ and have eigenvalues of order $N / \binom{D}{\ell}$. Simple perturbation theory shows that
 306 adding these matrices with coefficients of order 1 only slightly modify the eigenvalues. This is
 307 due to the fact that the row spaces (that is, the complements to their null spaces) corresponding
 308 to the different ℓ are almost orthogonal. In such a situation we approximate the trace-log term
 309 appearing in (40) as

$$\text{Tr} \log \left[A^{(0)} \otimes \mathbb{1}_N + \sum_{\ell=1}^L B^{(\ell)} \otimes C^{\otimes \ell} \right] \simeq N(1-L) \text{Tr} \log(A^{(0)}) + \sum_{\ell=1}^L \text{Tr} \log(A^{(0)} \otimes \mathbb{1}_N + B^{(\ell)} \otimes C^{\otimes \ell}) \quad (42)$$

310 (notice that Tr in $\text{Tr} \log(A^{(0)})$ is over replica indices only). We report a detailed derivation of
 311 Eq. (42) under the hypothesis of orthogonality of the $C^{\otimes \ell}$ row spaces in Appendix E. Notice
 312 that we could have gotten to the same result decomposing the vectors \mathbf{w} on the row spaces
 313 of the $C^{\otimes \ell}$ supposed orthogonal. This decomposition clearly shows the hierarchical nature of
 314 learning.

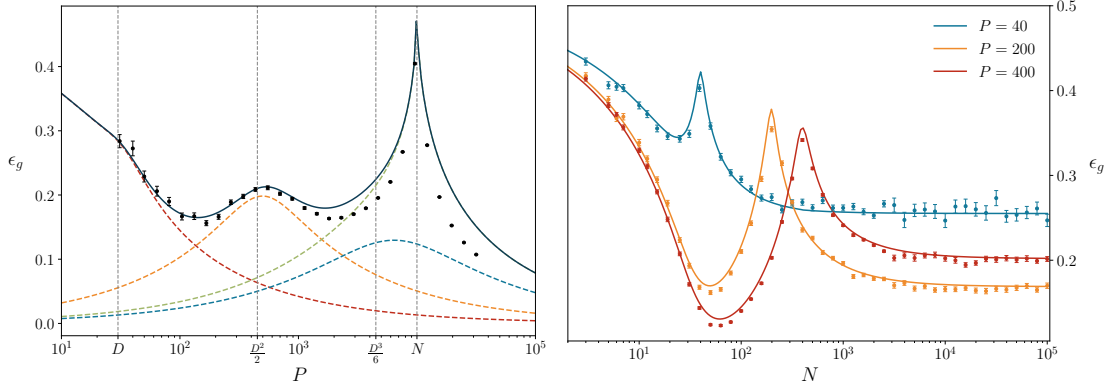


Figure 3: **Left:** generalization error of the RFM on a classification task, as a function of the size of the training set P , for $D = 30$, $N = 10^4$, weights regularization $\zeta = 10^{-8}$, linear teacher ($\tau_1 = 1$, $\tau_{\ell>1} = 0$) and ELU activation functions; the continuous line is the mean-field theory truncated at $L = 3$; dashed lines are the asymptotic theories for P/D finite and $L > 1$ (red), $P/\binom{D}{2}$ finite and $L > 2$ (yellow), $P/\binom{D}{3}$ finite and $L > 3$ (blue), $P/\binom{D}{3}$ finite and $L = 3$ (green); black points are results from numerical experiments averaged over 50 instances (see Appendix I). The model learns the linear features (first step at $P \sim O(D)$), then overfits the quadratic features before learning they are zero (peak at $P \sim O(D^2)$), then follows the interpolation peak $P \sim N$. Notice how the accordance between the mean-field theory and the experiment is only qualitative around the last peak. **Right:** Generalization error on classification for a linear teacher, as a function of the number of random features N , for different amounts of data P ($D = 30$, $\zeta = 10^{-4}$, see Appendix I). The optimal amount of hidden units, for which ϵ_g is minimal, shifts from overparametrization to underparametrization, as it is visible in the curves for $P = 40$ and $P = 200, 400$. At fixed value of N , not always more data means better generalization: after the interpolation peak, the order between the red ($P = 400$) and yellow ($P = 200$) curves is reversed (point of view complementary to the plot in the left panel, where, at fixed N , the error can increase with P). The curves as functions of N are obtained by gluing together the theories truncated at the corresponding L .

315 5.1 Replica symmetric theory

316 In order to complete the evaluation of the partition function, we need to specify the form of
 317 the replica parameters. In this paper we use the replica symmetry (RS) ansatz

$$Q_{ab}^{(\ell)} = \frac{\chi^{(\ell)}}{\beta} \delta_{ab} + q^{(\ell)}, \quad M_a^{(\ell)} = m^{(\ell)}, \quad t_a = t. \quad (43)$$

318 Notice that the diagonal elements of the matrix $Q^{(\ell)}$ are $Q_{aa}^{(\ell)} = \frac{\chi^{(\ell)}}{\beta} + q^{(\ell)}$. We anticipate the
 319 scaling with β of the variables χ : the quantities $Q_{aa}^{(\ell)}$ measures the variance of the variables λ ,
 320 tending to zero for $\beta \rightarrow \infty$. This implies the following form for the conjugate order parameters
 321 in the RS:

$$\hat{Q}_{ab}^{(\ell)} = \beta \hat{\chi}^{(\ell)} \delta_{ab} - \beta^2 \hat{q}^{(\ell)}, \quad \hat{M}_a^{(\ell)} = -\beta \hat{m}^{(\ell)}. \quad (44)$$

¹For the sake of simplicity, to write Eq. (39) we collected a common $C^{\otimes \ell}$ between the terms $\hat{Q}^{(\ell)}$ and $\hat{M}^{(\ell)} \hat{M}^{(\ell)\top}$, even though the average over the teacher gives instead a term $\sum_a \mathbf{F}_a^{\otimes \ell} (\mathbf{F}_a^{\otimes \ell})^\top / \binom{D}{\ell}$, with ordered indices a 's, in front of $\hat{M}^{(\ell)} \hat{M}^{(\ell)\top}$. See discussion around Eq. (47).

322 Exploiting the explicit parametrization of the RS matrices, we can perform the traces over
323 replica indices in Eq. (42), to get (see Appendix F)

$$\begin{aligned} \text{Tr} \log(A \otimes \mathbb{1}_N + B \otimes C^{\otimes \ell}) &= nN \log(\beta \hat{\chi}^{(\ell)}) + n \text{Tr} \log(\gamma_\ell \mathbb{1} + C^{\otimes \ell}) \\ &\quad - n\beta \eta_\ell \frac{\hat{q}^{(0)}}{\hat{\chi}^{(\ell)}} \text{Tr}(\gamma_\ell \mathbb{1} + C^{\otimes \ell})^{-1} - n\beta \frac{\hat{q}^{(\ell)} + (\hat{m}^{(\ell)})^2}{\hat{\chi}^{(\ell)}} \text{Tr}[C^{\otimes \ell}(\gamma_\ell \mathbb{1} + C^{\otimes \ell})^{-1}], \end{aligned} \quad (45)$$

324 where we introduced the parameter

$$\gamma_\ell = \eta_\ell \frac{(\zeta + \hat{\chi}^{(0)})}{\hat{\chi}^{(\ell)}} \quad (46)$$

325 and the remaining traces are over feature indices only.

326 We need now to evaluate the traces in feature indexes. In order to proceed, we make
327 make at this point a crucial approximation, and treat $C^{\otimes \ell}$ as a Wishart matrix with parameter
328 $\eta_\ell = N/\binom{D}{\ell}$. This amounts essentially in approximating $C^{\otimes \ell}$, by

$$C_{ij}^{\otimes \ell} = \frac{\ell!}{D^\ell} \sum_{\alpha_1 < \dots < \alpha_\ell} F_{\alpha_1}^i F_{\alpha_1}^j \dots F_{\alpha_\ell}^i F_{\alpha_\ell}^j \quad (47)$$

329 i.e. in neglecting the terms with equal indexes α in the sum that defines $C^{\otimes \ell}$. While this
330 approximation can be fully justified in the regimes where $N, D \rightarrow \infty$ with N/D^L finite, as we
331 will see, it turns out to be an excellent approximation even for moderately large values of the
332 parameters (see Sec. 6 and Appendix D).

333 Using the properties of Wishart matrices (see Appendix D), we can write that, for large N ,

$$\frac{1}{N} \text{Tr}(\gamma_\ell \mathbb{1} + C^{\otimes \ell})^{-1} \approx g_\ell(-\gamma_\ell), \quad (48)$$

334 where g_ℓ is the Stieltjes transformation of the Marchenko-Pastur distribution with ratio $\eta_\ell =$
335 $N/\binom{D}{\ell}$:

$$g_\ell(z) = \frac{1 - z - \eta_\ell - \sqrt{(1 - z - \eta_\ell)^2 - 4z\eta_\ell}}{2z\eta_\ell}. \quad (49)$$

336 Re-arranging terms we get, for large β ,

$$Z_n \sim e^{PS_p + NS_M}, \quad (50)$$

337 where

$$\begin{aligned} \frac{1}{\beta n} S_M &= - \sum_{\ell=1}^{\min\{L, B\}} \frac{m^{(\ell)} \hat{m}^{(\ell)}}{\eta_\ell} + \frac{1}{2} \sum_{\ell=0}^L \frac{q^{(\ell)} \hat{\chi}^{(\ell)} - \chi^{(\ell)} \hat{q}^{(\ell)}}{\eta_\ell} + \frac{(1-L)}{2} \frac{\hat{q}^{(0)}}{\zeta + \hat{\chi}^{(0)}} \\ &\quad + \frac{1}{2} \sum_{\ell=1}^L \eta_\ell \frac{\hat{q}^{(0)}}{\hat{\chi}^{(\ell)}} g_\ell(-\gamma_\ell) + \frac{1}{2} \sum_{\ell=1}^L \frac{\hat{q}^{(\ell)} + (\hat{m}^{(\ell)})^2}{\hat{\chi}^{(\ell)}} [1 - \gamma_\ell g(-\gamma_\ell)] \end{aligned} \quad (51)$$

338 and, for the quadratic loss 10,

$$\frac{1}{\beta n} S_P = \frac{2m^* \langle y \nu \rangle - q^* - \langle (t^* - y)^2 \rangle}{2(1 + \chi^*)}, \quad (52)$$

339 where $\langle \cdot \rangle = \int dy D\nu p(y|\nu)(\cdot)$ is the average over the teacher distribution (1) and

$$\begin{aligned} m^* &= \sum_{\ell=1}^{\min\{L, B\}} \frac{\tau_\ell \mu_\ell}{\sqrt{\ell!}} m^{(\ell)}, \quad t^* = \mu_0 m^{(0)}, \\ \chi^* &= \mu_\perp^2 \chi^{(0)} + \sum_{\ell=1}^L \frac{\mu_\ell^2}{\ell!} \chi^{(\ell)}, \quad q^* = \mu_\perp^2 q^{(0)} + \sum_{\ell=1}^L \frac{\mu_\ell^2}{\ell!} q^{(\ell)}. \end{aligned} \quad (53)$$

340 A detailed derivation of the terms S_M and S_P , with the form of S_P valid for generic loss func-
 341 tions, is reported in Appendix G.

342 Eq. (53) gives the RS version of Eq. (34): these quantities are precisely the ones appearing
 343 in Eq. (14), giving the low-order statistics of the distribution used to evaluate the generaliza-
 344 tion error. Once their value is known from the saddle point equations implicit in the derivation
 345 of the partition function, they can be used to obtain the generalization curves reported in this
 346 paper.

347 5.2 Saddle-point equations for quadratic loss

348 The free energy in Eq. (50) has to be evaluated at the saddle point with respect to all the
 349 RS order parameters and their Fourier conjugates. The resulting equations, which we report
 350 here for the case of quadratic loss function. Remark however that only the equations where
 351 P appears explicitly depend on the form of the loss. The equations can be solved in steps.
 352 First, a set of $2L + 2$ nonlinear equations is used to determine the variables $\chi^{(0)}, \dots, \chi^{(L)}$ and
 353 $\hat{\chi}^{(0)}, \dots, \hat{\chi}^{(L)}$:

$$\begin{aligned} \hat{\chi}^{(0)} &= \frac{P}{N} \frac{\mu_{\perp}^2}{1 + \chi^*}, & \chi^{(0)} &= \frac{1 - \sum_{\ell=1}^L [1 - \gamma_{\ell} g_{\ell}(-\gamma_{\ell})]}{\hat{\chi}^{(0)} + \zeta}, \\ \hat{\chi}^{(\ell)} &= \frac{P}{\binom{D}{\ell}} \frac{\mu_{\ell}^2}{\ell!} \frac{1}{1 + \chi^*}, & \chi^{(\ell)} &= \frac{N}{\binom{D}{\ell}} \frac{1 - \gamma_{\ell} g_{\ell}(-\gamma_{\ell})}{\hat{\chi}^{(\ell)}}. \end{aligned} \quad (54)$$

354 From the solution of Eq. (54), we can fully determine $m^{(\ell)}, \hat{m}^{(\ell)}$ according to

$$m^{(0)} = \frac{\langle y \rangle}{\mu_0}, \quad m^{(\ell)} = \chi^{(\ell)} \hat{m}^{(\ell)}, \quad \hat{m}^{(\ell)} = \frac{P}{\binom{D}{\ell}} \frac{\mu_{\ell} \tau_{\ell}}{\sqrt{\ell!}} \frac{\langle y \nu \rangle}{1 + \chi^*}. \quad (55)$$

355 With all the previous values we can determine the rest of the variables through the following
 356 set of linear equations:

$$\begin{aligned} q^{(0)} &= \frac{\hat{q}^{(0)}}{(\zeta + \hat{\chi}^{(0)})^2} \left(1 - \sum_{\ell=1}^L [1 - \gamma_{\ell}^2 g'_{\ell}(-\gamma_{\ell})] \right) + \sum_{\ell=1}^L \frac{\hat{m}^{(\ell)2} + \hat{q}^{(\ell)}}{(\zeta + \hat{\chi}^{(0)}) \hat{\chi}^{(\ell)}} [\gamma_{\ell} g_{\ell}(-\gamma_{\ell}) - \gamma_{\ell}^2 g'_{\ell}(-\gamma_{\ell})], \\ q^{(\ell)} &= \frac{N}{\binom{D}{\ell}} \frac{\hat{q}^{(0)}}{(\zeta + \hat{\chi}^{(0)}) \hat{\chi}^{(\ell)}} [\gamma_{\ell} g_{\ell}(-\gamma_{\ell}) - \gamma_{\ell}^2 g'_{\ell}(-\gamma_{\ell})] \\ &\quad + \frac{N}{\binom{D}{\ell}} \frac{\hat{m}^{(\ell)2} + \hat{q}^{(\ell)}}{\hat{\chi}^{(\ell)2}} [1 + \gamma_{\ell}^2 g'_{\ell}(-\gamma_{\ell}) - 2\gamma_{\ell} g_{\ell}(-\gamma_{\ell})], \\ \hat{q}^{(0)} &= \frac{P}{N} \mu_{\perp}^2 \frac{\langle (\mu_0 m^{(0)} - y)^2 \rangle - 2\langle y \nu \rangle m^* + q^*}{(1 + \chi^*)^2}, \\ \hat{q}^{(\ell)} &= \frac{P}{\binom{D}{\ell}} \frac{\mu_{\ell}^2}{\ell!} \frac{\langle (\mu_0 m^{(0)} - y)^2 \rangle - 2\langle y \nu \rangle m^* + q^*}{(1 + \chi^*)^2}. \end{aligned} \quad (56)$$

357 Notice that, because of the conventional scalings we chose for the hat variables starting from
 358 Eq. (37) and for the definition of γ_{ℓ} , these equations give $O(1)$ results for the order parameters
 359 m, χ, q .

360 By numerically integrating Eq. (54), (55), (56), we obtain the theoretical curves for the
 361 generalization error in Eq. (16) and for the order parameters we report in this paper. We com-
 362 pare the result with numerical simulations: despite its asymptotic nature and the hypothesis
 363 of row space orthogonality, our theory works reasonably well even if D is not large. The results
 364 are shown in Fig. 1, 2 ($D = 30$ in this case), where the generalization error is quantitatively
 365 predicted by the theory both when varying P and N .

366 6 Strongly separated regimes

367 Our analysis relies on two cornerstones: (1) The possibility of taking the saddle point on the
 368 replica parameters (2) Treating the traces of $\log(\gamma_\ell + C^{\circ\ell})$ as if $C^{\circ\ell}$ were Wishart matrices.
 369 These assumptions can be justified if $P, N, D \rightarrow \infty$. Depending on the relation between the
 370 three parameters one is led to consider the following different asymptotic regimes:

371 (i) $N, P, D \rightarrow \infty, P/N \rightarrow 0, P/D^K$ finite; (this includes the case $N \sim D^L$ with $L > K$).

372 (ii) $N, P, D \rightarrow \infty, N/D^L$ finite, P/N finite;

373 (iii) $N, P, D \rightarrow \infty, P/N \rightarrow \infty, N/D^L$ finite; (this includes the case $P \sim D^K$ with $K > L$).

374 In all these cases we need to specify the relation of K and L with B , the maximum degree of
 375 the teacher polynomial. This will be done in the following of this section.

376 In order to understand these regimes, we need to evaluate terms of the kind

$$k_\ell = \text{Tr} \log(a\mathbb{1} + bC^{\circ\ell}), \quad C_{ij}^{\circ\ell} = \left(\frac{1}{D} \sum_{\alpha} F_{i\alpha} F_{j\alpha} \right)^\ell \quad (57)$$

377 in three situations (a) $D^\ell \gg N$; (b) $D^\ell \ll N$; (c) $D^\ell \sim N$. Notice that in all cases, while
 378 the diagonal elements are $C_{ii}^{\circ\ell} = 1 + O(\sqrt{1/D^\ell})$, the off-diagonal elements $C_{i \neq j}^{\circ\ell}$ are of the
 379 order $D^{-\ell/2}$. In case (a), $D^\ell \gg N$, apart for a negligible number of possible eigenvalue of
 380 order $N/D^{\ell/2}$, all the other eigenvalues are $\lambda = 1 + O(\sqrt{N/D^\ell})$, and to the leading order we
 381 simply have $k_\ell = N \log(a + b)$. If we are in the opposite situation, (b), $D^\ell \ll N$, we have only
 382 $O(D^\ell)$ non-zero eigenvalues, roughly equal to $\ell! N/D^\ell + O(\sqrt{N/D^\ell})$, and to the leading order
 383 $k_\ell = N \log(a)$. The interesting case is (c) $N = O(D^\ell)$: we have here D^ℓ eigenvalues of order 1
 384 that contribute to k_ℓ . The leading contribution can be understood writing

$$C_{ij}^{\circ\ell} = \frac{\ell!}{D^\ell} \sum_{\alpha} F_{i,\alpha}^{\otimes \ell} F_{j,\alpha}^{\otimes \ell} + \text{terms with less } \alpha\text{'s} \quad (58)$$

385 where the sum includes the terms where the α 's in the multi-index α are ordered, coherently
 386 with our definition in Table 1. This leading term is a matrix of rank $\min\{N, \binom{D}{\ell}\}$. Other terms
 387 with smaller number of indexes in the sum lead to matrices of lower rank r (with $r/N \rightarrow$
 388 0). Moreover, due to the randomness of the F , the row spaces of these term are effectively
 389 orthogonal to the leading one. We conclude that we can compute k_ℓ as if $C^{\circ\ell}$ were a Wishart
 390 matrix of parameter $\eta_\ell = N/\binom{D}{\ell}$. The explicit formula is given in eq. (D.12), and both limits
 391 $\eta_\ell \rightarrow 0$ and $\eta_\ell \rightarrow \infty$ agree with the previous analysis of cases (a) and (b) respectively. We
 392 show in appendix D that approximating $C^{\circ\ell}$ as a Wishart matrix gives good results also for
 393 moderately large values of N and D .

394 In all our three cases, most of the order parameters go to trivial limits, while only the ones
 395 corresponding to the selected scaling regime converge to non-trivial values. We report the
 396 corresponding equations in Appendix H. In this way, we are able to plot the dashed lines in
 397 Fig. 1 and 3.

398 7 Effective theory for finite-size random features networks

399 In the last sections we devised a theory able to capture the relevant phenomenology of general-
 400 ization in RFMs at finite values of input dimension, hidden layer width and size of the training
 401 set. Indeed, even though the asymptotic approximation leading to the system of saddle-point

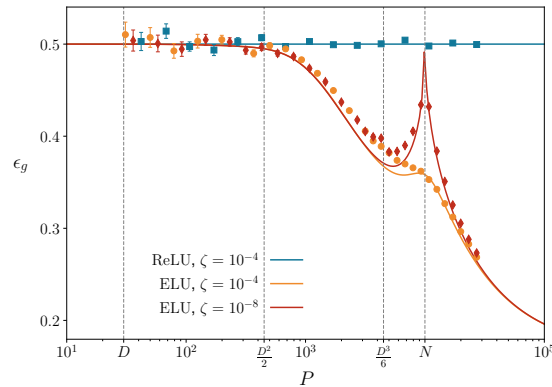


Figure 4: Generalization error vs P ($D = 30$, $N = 10^4$) on classification for a *purely cubic teacher* ($\tau_3 = 1$); in blue, polynomial theory and numerical experiments for ReLU activation function (7): in this case, $\mu_3 = 0$ and the model cannot learn the cubic features, so the error remains $1/2$; in yellow and red (respectively, for $\zeta = 10^{-4}$, 10^{-8}), the case of ELU (8), for which $\mu_3 \neq 0$ and the model can learn the cubic features.

402 equations (54), (55), (56) is justified only for N large and N/D^L finite, the curves obtained by
 403 fixing the values of N , P and D at finite values are in accordance with numerical simulations
 404 over several orders of magnitudes of the control parameters. This occurs thanks to the fact
 405 that we kept into account quantities that scale differently with D , as $N/\binom{D}{\ell}$ or $P/\binom{D}{\ell}$, that are
 406 formally zero or infinity in the asymptotic regimes presented in Sec. 6.

407 By developing a theory from Eq. (26), we show that the RFM is in essence equivalent to a
 408 polynomial model: the student tries to tune its weights through the combinations $\mathbf{s}^{(\ell)}$ defined
 409 in (27) to fit the corresponding coefficients $\theta^{(\ell)}$ of the teacher. This interpretation is also
 410 confirmed in the numerical experiments: see Fig. 1 (right) for the behavior of the teacher-
 411 student overlaps $m^{(\ell)}$ in the case of a quadratic teacher.

412 However, a crucial difference from a purely polynomial setting arises: the degree of the
 413 equivalent polynomial model is controlled by the scaling L of the random features, and higher
 414 order terms in the expansion of the kernel \mathcal{K} on the Hermite basis act as noise, given by
 415 Eq. (29). This eventually produces the interpolation peak in the generalization error at $N \sim P$,
 416 which would not be present for a vanilla polynomial student (see Fig. 1 and 3): in this regime,
 417 the model is overfitting the effective noise. In terms of the order parameters, overlaps of
 418 different orders are coupled by an additional set of parameters $\chi^{(0)}$, $q^{(0)}$, related to the noise
 419 term in the equivalent polynomial model.

420 In summary, the learning of features of a certain order is possible as long as the number
 421 of parameters N is enough: the scaling $L \sim \log N / \log D$ controls the learning process through
 422 the truncation of the kernel (24). At the same time, P also plays an important role: if $K \sim$
 423 $\log P / \log D$ is smaller than L , the model only learns as a K -degree polynomial; on the other
 424 hand, if $K > L$, the model learns as a L -degree polynomial.

425 By choosing a polynomial teacher of arbitrary degree B , we are able to explore to some
 426 extent the interplay between the complexity of the data and the one of the neural network.
 427 In the case where the teacher is less complex than the network, we can see that overfitting
 428 can occur and that overparametrization is not always optimal. This can be seen in Fig. 3.
 429 In the case of a linear teacher, if the amount of data P is $O(D)$, an overparametrized network
 430 generalizes better. However, as soon as P hits the quadratic regime, but is still far from enabling
 431 the network to realize that there is no quadratic feature, then overparametrization leads to
 432 overfitting and therefore the optimal N is less than P .

433 Interestingly, in order for the model to learn features of order ℓ , the activation function σ
 434 must have a non-zero Hermite coefficient μ_ℓ in Eq. (21). This can be seen from our theory
 435 by the fact that in the total teacher-student overlap m^* in Eq. (53) the single entry $m^{(\ell)}$ is
 436 weighted by the corresponding coefficient. This theoretical prediction was tested by using a
 437 cubic teacher and two different students, one with ReLU activation function and the other one
 438 with ELU: the ReLU one, which has no third order term in the Hermite basis ($\mu_3 = 0$) could
 439 not learn the teacher, while the ELU one, that does have a nonzero component ($\mu_3 \neq 0$), was
 440 able to (see Fig. 4).

441 8 Conclusions and perspectives

442 The approach we have explored so far provides a way to analytically evaluate the general-
 443 ization performance of a RFM in the limit of large input dimension D , in the scaling regimes
 444 $N \sim D^L, P \sim D^K$.

445 We considered a teacher-student setting, where a shallow random features student is re-
 446 quired to fit a polynomial teacher. The student network learns as an equivalent polynomial
 447 model with effective noise. We showed this property by expanding the kernel in feature space
 448 on a convenient basis (21).

449 The resulting theory is effective, in the sense that it is formulated in terms of a few collective
 450 order parameters (the teacher-student overlaps $m^{(\ell)}$, the student-student overlaps $q^{(\ell)}, \chi^{(\ell)}$)
 451 with a clear physical interpretation and whose values are fixed via a variational principle,
 452 as explained in Sec. 5. To perform the calculation we neglect the correlations between the
 453 student's coefficients, assuming orthogonality between the row spaces of the components $C^{\otimes \ell}$
 454 of the kernel.

455 We find quantitative agreement with numerical simulations, except close to the interpo-
 456 lation peak at $N \sim P$ in some cases (see Fig. 3, left, where this effect is more apparent).
 457 Nevertheless, even then the effective theory gives a good qualitative picture, predicting the
 458 location and the shape of the peak. See also Fig. 1, right, depicting how the teacher-student
 459 overlaps of already learned features become noisy in the interpolation regime. A precise finite-
 460 size analysis of this effect, to address the gap between theory and numerics in this regime, is
 461 left for future work.

462 One possible direction to continue this work is to consider how close is the learning of a
 463 fully-trained network to this model. The role of the variables $\mathbf{s}^{(\ell)}$ could play a similar role even
 464 if the values for $F_{i\alpha}$ are also learned, at least close to the lazy regime. However, what is the
 465 fate of row space orthogonality of the kernel components, which is ultimately responsible for
 466 the staircase behavior of the generalization error, for networks that are trained end-to-end in
 467 a feature learning regime?

468 Moreover, it would be interesting to extend our analysis to deeper models [10, 72] in dif-
 469 ferent scaling regimes of the dimensions. Even if the RFM, whatever the activation function
 470 of the last layer, is essentially bounded by a polynomial model, the precise shape of the kernel
 471 in cases where a deeper architecture is involved could help understanding to some extent the
 472 feature learning regimes of realistic models, in view of the discussion above.

473 Acknowledgements

474 The authors would like to thank Pietro Rotondo, Rosalba Pacelli, Bruno Loureiro, Valentina
 475 Ros, the QBio group at ENS for discussions and suggestions. MP and FAL are grateful to the
 476 organizers and speakers of the Statistical Physics of Deep Learning summer school held in June

477 2022 in Como, where the idea was in part conceived.

478 **Funding information** The authors have been supported by a grant from the Simons Founda-
479 tion (grant No. 454941, S. Franz), thanks to which most of this work was performed at
480 LPTMS (CNRS, Université Paris-Saclay).

481 FAL conducted part of this research within the Econophysics & Complex Systems Research
482 Chair, under the aegis of the Fondation du Risque, the Fondation de l'École polytechnique, the
483 École polytechnique and Capital Fund Management.

484 A Kernel on the Hermite basis

485 In this section we report the steps needed to obtain the expression of the feature-feature kernel
486 in Sec. 4. The kernel to evaluate is defined as

$$\begin{aligned} \mathcal{K}_{ii} &= \mathbb{E}_{h_i}[\sigma(h_i)^2] = \int \frac{du}{\sqrt{2\pi C_{ii}}} e^{-\frac{u^2}{2C_{ii}}} \sigma(u)^2 \\ \mathcal{K}_{ij} &= \mathbb{E}_{h_i, h_j}[\sigma(h_i)\sigma(h_j)] = \int \frac{du dv}{2\pi\sqrt{\det \bar{C}}} e^{-\frac{1}{2}(u,v)\bar{C}^{-1}(u,v)^\top} \sigma(u)\sigma(v) \quad i \neq j \end{aligned} \quad (\text{A.1})$$

487 where

$$\bar{C} = \begin{pmatrix} C_{ii} & C_{ij} \\ C_{ij} & C_{jj} \end{pmatrix}. \quad (\text{A.2})$$

488 Using the fact that $C_{ii} \simeq C_{jj} \simeq 1$, this kernel can be written as a series of separable kernels
489 exploiting Mehler's formula [66, 67], that we report here for convenience:

$$\frac{1}{2\pi\sqrt{1-c^2}} e^{-\frac{1}{2}(u,v)\begin{pmatrix} 1 & c \\ c & 1 \end{pmatrix}^{-1}(u,v)^\top} = \frac{e^{-\frac{u^2}{2}}}{\sqrt{2\pi}} \frac{e^{-\frac{v^2}{2}}}{\sqrt{2\pi}} \sum_{\ell=0}^{\infty} \frac{c^\ell}{\ell!} \text{He}_\ell(u)\text{He}_\ell(v), \quad (\text{A.3})$$

490 from which we find Eq. (23) using the fact that, by orthogonality of the Hermite polynomials,

$$\mathcal{K}_{ii} = \sum_{\ell=0}^{\infty} \frac{\mu_\ell^2}{\ell!}. \quad (\text{A.4})$$

491 Mehler's formula, which dates back to 1866, can be viewed as an example of Mercer's decom-
492 position [15].

493 B Hermite polynomials and Wick products

494 For completeness, we show in this section that, asymptotically for D large,

$$\text{He}_\ell(h_i) \simeq \sum_{\alpha_1, \dots, \alpha_\ell} \frac{F_{i\alpha_1} \cdots F_{i\alpha_\ell}}{\sqrt{D}^\ell} :x_{\alpha_1} \cdots x_{\alpha_\ell}:, \quad (\text{B.1})$$

495 for $\ell \geq 1$. The equivalence follows from the generating function of the Hermite polynomials,

$$\text{He}_\ell(h_i) = \frac{d^\ell}{dt^\ell} \exp(th_i - t^2/2) \Big|_{t=0}, \quad (\text{B.2})$$

496 with $h_i = \sum_\alpha F_{i\alpha} x_\alpha / \sqrt{D}$. Defining

$$\lambda_\alpha = t \frac{F_{i\alpha}}{\sqrt{D}}, \quad (\text{B.3})$$

497 we have, for D large,

$$\sum_{\alpha} \lambda_{\alpha}^2 \approx t^2, \quad \sum_{\alpha} \frac{F_{i\alpha} \lambda_{\alpha}}{\sqrt{D}} \approx t, \quad \sum_{\alpha} \frac{F_{i\alpha}}{\sqrt{D}} \frac{\partial}{\partial \lambda_{\alpha}} \approx \frac{d}{dt}, \quad (\text{B.4})$$

498 where we used repeatedly $\sum_{\alpha} (F_{i\alpha})^2 / D \simeq 1$. The thesis follows from comparison with Eq. (28).
 499 Notice that, in the simpler case of a single standard Gaussian variable x , the identity $\text{He}_{\ell}(x) =$
 500 $:x^{\ell}:$ is exact and trivially follows from the definition of the Wick power.

501 C Evaluation of the moments of ν, λ^a

502 We assume that the variables $(\nu, \{\lambda^a\})$ are normally distributed with mean and covariance

$$\mathbb{E}_{\mathbf{x}}[(\nu, \{\lambda^a\})] = (0, \{t^a\}), \quad \text{cov}_{\mathbf{x}}[(\nu, \{\lambda^a\})] = \begin{pmatrix} \rho & M^{\top} \\ M & Q \end{pmatrix}, \quad (\text{C.1})$$

503 where

$$\begin{aligned} t_a &= \mathbb{E}_{\mathbf{x}}[\lambda^a] = \sum_{i=1}^N \frac{w_i^a}{\sqrt{N}} \mathbb{E}_{h_i}[\sigma(h_i)], \\ \rho &= \mathbb{E}_{\mathbf{x}}[\nu^2] - \mathbb{E}_{\mathbf{x}}[\nu]^2 = \sum_{\ell=1}^B \tau_{\ell}^2 \frac{\|\theta^{(\ell)}\|^2}{\binom{D}{\ell}}, \\ M_a &= \mathbb{E}_{\mathbf{x}}[\nu \lambda^a] = \sum_{i,\ell} \frac{w_i^a \tau_{\ell}}{\sqrt{N \binom{D}{\ell}}} \sum_{\alpha_1 < \dots < \alpha_{\ell}} \theta_{\alpha_1 \dots \alpha_{\ell}}^{(\ell)} \mathbb{E}_{\mathbf{x}}[x_{\alpha_1} \dots x_{\alpha_{\ell}} \sigma(h_i)], \\ Q_{ab} &= \mathbb{E}_{\mathbf{x}}[\lambda^a \lambda^b] - t^a t^b = \sum_{i,j=1}^N \frac{w_i^a w_j^b}{N} \mathbb{E}_{h_i, h_j}[\sigma(h_i) \sigma(h_j)] - t^a t^b, \end{aligned} \quad (\text{C.2})$$

504 To proceed, we make the following steps, starting from the expansion of the activation
 505 function on the Hermite basis, Eq. (21). For t_a we simply observe that $\mathbb{E}_{h_i}[\sigma(h_i)] = \mu_0$. For
 506 ρ we use the fact that \mathbf{x} is distributed as a standard normal random vector. To deal with Q_{ab}
 507 we introduce the truncation of (24). Finally, for M_a we write explicitly

$$\sum_{\alpha_1 < \dots < \alpha_k} \theta_{\alpha_1 \dots \alpha_k}^{(k)} \mathbb{E}_{\mathbf{x}}[x_{\alpha_1} \dots x_{\alpha_k} \sigma(h_i)] = \sum_{\alpha_1 < \dots < \alpha_k} \theta_{\alpha_1 \dots \alpha_k}^{(k)} \sum_{\ell=0}^{\infty} \frac{\mu_{\ell}}{\ell!} \mathbb{E}_{\mathbf{x}}[x_{\alpha_1} \dots x_{\alpha_k} \text{He}_{\ell}(h_i)] \quad (\text{C.3})$$

508 and we perform Wick's contractions in order to evaluate the expected value, exploiting the
 509 mapping to Wick's product explained in Appednix B. As the indices α of the teacher are strictly
 510 ordered, they must be paired only with the ones in the Wick product, leaving only the term
 511 $\ell = k$ in the sum over ℓ . The number of possible contractions is $k!$, so the result is

$$\begin{aligned} t_a &= \frac{\mu_0}{\sqrt{N}} \sum_{i=1}^N w_i^a, \\ M_a &= \sum_i \frac{w_i^a}{\sqrt{N}} \sum_{\ell=1}^B \frac{\tau_{\ell}}{\binom{D}{\ell} \sqrt{\ell!}} \sum_{\alpha} \theta_{\alpha}^{(\ell)} F_{i,\alpha}^{\otimes \ell}, \\ Q_{ab} &= \frac{1}{N} \sum_{i,j=1}^N w_i^a w_j^b \left(\delta_{ij} \mu_{\perp,L}^2 + \sum_{\ell=1}^L \frac{\mu_{\ell}^2}{\ell!} (C_{ij})^{\ell} \right), \end{aligned} \quad (\text{C.4})$$

512 from which Eq. (34) follows.

513 D Results on random matrix theory

514 D.1 Marchenko-Pastur distribution and Stieltjes transformation

515 In this section, we remind some textbook results in Random Matrix Theory we used in the
516 main text, for the reader's convenience. First of all, random matrices of the form

$$C = FF^\top/D, \quad (\text{D.1})$$

517 where F is a $N \times D$ random matrix with i.i.d. entries $F_{i\alpha}$ such that $\mathbb{E}[F_{i\alpha}] = 0$, $\mathbb{E}[(F_{i\alpha})^2] = \sigma^2$,
518 define the Wishart (or Wishart-Laguerre) ensemble. For large N and D , parameter $\eta \equiv N/D$
519 finite, their spectral density follows the Marchenko-Pastur (MP) distribution,

$$\rho_{\text{MP}}(\lambda) = \begin{cases} (1 - 1/\eta) \delta(\lambda) + \rho_{\text{bulk}}(\lambda/\sigma^2)/\sigma^2 & \text{if } \eta > 1, \\ \rho_{\text{bulk}}(\lambda/\sigma^2)/\sigma^2 & \text{if } \eta \leq 1, \end{cases} \quad (\text{D.2})$$

520 with

$$\rho_{\text{bulk}}(\lambda) = \frac{\sqrt{(\lambda_+ - \lambda)(\lambda - \lambda_-)}}{2\pi\eta\lambda}, \quad \lambda_{\pm} = (1 \pm \sqrt{\eta})^2 \quad (\text{D.3})$$

521 with support in $\lambda_- \leq \lambda \leq \lambda_+$.

522 The MP distribution can be obtained with standard methods [73, 74]. The determinant of
523 the resolvent can be evaluated as follows:

$$\mathbb{E} \left[\det \left(\gamma \mathbb{1}_N + \frac{FF^\top}{D} \right) \right]^{-\frac{1}{2}} = \mathbb{E} \int \frac{d\mathbf{x}}{(2\pi)^{\frac{N}{2}}} e^{-\frac{1}{2} \mathbf{x}^\top (\gamma \mathbb{1}_N + \frac{FF^\top}{D}) \mathbf{x}}. \quad (\text{D.4})$$

524 By Gaussian linearization,

$$\mathbb{E} \int \frac{d\mathbf{y}}{(2\pi)^{\frac{D}{2}}} \frac{d\mathbf{x}}{(2\pi)^{\frac{N}{2}}} e^{-\frac{\|\mathbf{y}\|^2}{2} - \frac{\gamma}{2} \|\mathbf{x}\|^2 + i\mathbf{x}^\top \frac{F}{\sqrt{D}} \mathbf{y}} \quad (\text{D.5})$$

525 The average over F gives

$$\int \frac{d\mathbf{y}}{(2\pi)^{\frac{D}{2}}} \frac{d\mathbf{x}}{(2\pi)^{\frac{N}{2}}} e^{-\frac{\|\mathbf{y}\|^2}{2} - \frac{\gamma}{2} \|\mathbf{x}\|^2 - \frac{1}{2D} \|\mathbf{x}\|^2 \|\mathbf{y}\|^2}. \quad (\text{D.6})$$

526 Integrating over \mathbf{y} ,

$$\int \frac{d\mathbf{x}}{(2\pi)^{\frac{N}{2}}} e^{-\frac{\gamma}{2} \|\mathbf{x}\|^2 - \frac{D}{2} \log(1 + \|\mathbf{x}\|^2/D)}. \quad (\text{D.7})$$

527 Inserting $r = \|\mathbf{x}\|^2/N$ with a Dirac delta, we can integrate over \mathbf{x} :

$$\int \frac{dr d\hat{r}}{4\pi} e^{\frac{iN\hat{r}r}{2} - \frac{N}{2} \log(i\hat{r}) - \frac{N}{2} \gamma r - \frac{N}{2\eta} \log(1 + \eta r)}. \quad (\text{D.8})$$

528 The integral over the Fourier variable \hat{r} can be solved via asymptotic integration, the saddle-
529 point being in $\hat{r} = -ir^{-1}$:

$$\int dr e^{\frac{N}{2} [1 + \log(r) - \gamma r - \frac{1}{\eta} \log(1 + \eta r)]} \quad (\text{D.9})$$

530 The saddle point equation in r gives

$$\frac{1}{r} - \gamma - \frac{1}{1 + \eta r} = 0 \quad (\text{D.10})$$

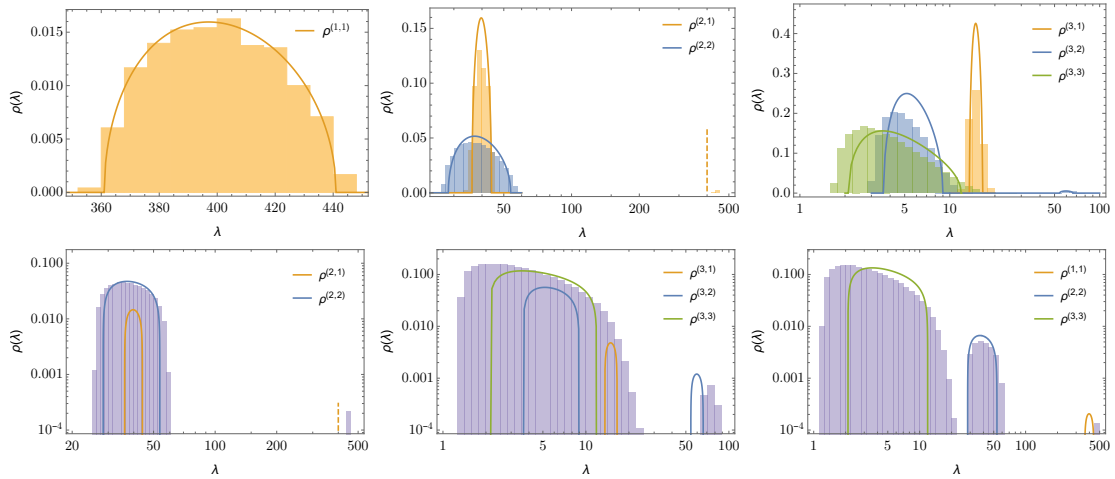


Figure 5: **Top row** – empirical (30 instances, $D = 20$, $N = D^3$) vs. analytical (MP) distributions of the non-zero eigenvalues of the matrices defined in Sec. D.2: $C^{(1,1)}$ (left), $C^{(2,1)}/D$, $C^{(2,2)}$ (center), $C^{(3,1)}/D^2$, $3C^{(3,2)}/D$, $C^{(3,3)}$ (right). **Bottom row** – comparison of the analytical curves with the empirical distribution (notice the log scale on the axes) of $C^{\circ 2}$ (left), $C^{\circ 3}$ (center) and $C^{\circ 1} + C^{\circ 2} + C^{\circ 3}$ (right); analytical curves in the bottom row are rescaled in such a way that the sum of the densities in each panel is normalized.

531 with solutions

$$r_{\pm} = \frac{\eta - \gamma - 1 \pm \sqrt{(\eta - \gamma - 1)^2 + 4\eta\gamma}}{2\eta\gamma}. \quad (\text{D.11})$$

532 The correct branch can be proven to be $r = r_+$. From this analysis, the relation

$$\frac{1}{N} \mathbb{E} \text{Tr} \log(\gamma \mathbb{1} + C) = -(1 - \gamma r) - \log(r) + \frac{1}{\eta} \log(1 + \eta r) \quad (\text{D.12})$$

533 follows. Deriving with respect to γ ,

$$\frac{1}{N} \mathbb{E} \text{Tr}(\gamma \mathbb{1} + C)^{-1} = r(\gamma). \quad (\text{D.13})$$

534 By definition of Stieltjes transform, $r(\gamma) = g(-\gamma)$, which gives Eq. (49).

535 D.2 Spectral density of $C^{\circ \ell}$

536 In this Appendix we discuss the spectral density of the matrices $C^{\circ \ell}$, to clarify the kind of
537 approximation we used in the main text. We are interested to the large N computation of the
538 following traces:

$$a_{\ell} = \frac{1}{N} \text{Tr}(\gamma_{\ell} \mathbb{1} + C^{\circ \ell})^{-1}, \quad b_{\ell} = \frac{1}{N} \text{Tr} C^{\circ \ell} (\gamma_{\ell} \mathbb{1} + C^{\circ \ell})^{-1} \quad (\text{D.14})$$

539 under the hypothesis that $\eta_L = N/\binom{D}{L}$ remain finite. We anticipate that γ_{ℓ} given by (46) either
540 remain finite (if P/N remains finite) or tends to infinity (if $P/N \rightarrow \infty$) in that limit. As we
541 have already discussed, for $\ell > L$, the matrix $C^{\circ \ell}$ is fully ranked, with diagonal elements close
542 to one and off-diagonal elements of order $D^{-\ell/2}$: all eigenvalues will be equal to one up to
543 a negligible correction. For that reason we could neglect off-diagonal terms for $\ell > L$ and
544 $a_{\ell} \approx b_{\ell} \approx (1 + \gamma_{\ell})^{-1}$. For $\ell < L$ conversely, the matrix has rank D^{ℓ} at most, and it is easy to see
545 that its max eigenvalue cannot be larger than $N \max_i \left(\frac{1}{D} \sum_{\alpha} F_{i,\alpha}^2 \right)^{\ell} = N(1 + O(\sqrt{\log(N)/D}))$.²

² $\lambda_{\max} = \max_{\|v\|_2=1} \frac{1}{D^{\ell}} \sum_{\alpha_1, \dots, \alpha_{\ell}} \left(\sum_i v_i F_{i,\alpha_1} \dots F_{i,\alpha_{\ell}} \right)^2 \leq N \sum_i v_i^2 \left(\frac{1}{D} \sum_{\alpha} F_{i,\alpha}^2 \right)^{\ell}$

546 We get therefore

$$\frac{1}{N} \left((N - D^\ell) / \gamma_\ell + D^\ell / (\gamma_\ell + N) \right) \leq a_\ell \leq \frac{1}{\gamma_\ell}, \quad 0 \leq b_\ell \leq \frac{D^\ell}{N} \frac{N}{\gamma_\ell + N}. \quad (\text{D.15})$$

547 It remains to be discussed the only non trivial case: $\ell = L$. In that case, we can decompose the
548 matrix $C^{\odot L}$ as a Wishart matrix with rank $\min\{N, \binom{D}{L}\}$ and parameter η_L , plus a contribution
549 with rank at most D^{L-1} which for reasoning similar to the previous case, do not contribute to
550 a_L and b_L in the thermodynamic limit.

551 We would like now to show, that even for moderate values of N and D , the neglect of all
552 non-Wishart contribution provides an excellent approximation to the spectrum. To fix ideas,
553 let us consider $L = 3$ ($N \sim D^3$), so that we consider the matrices

$$C^{\odot 1} = C^{(1,1)}, \quad C^{\odot 2} = \frac{1}{D} C^{(2,1)} + C^{(2,2)}, \quad C^{\odot 3} = \frac{1}{D^2} C^{(3,1)} + \frac{3}{D} C^{(3,2)} + C^{(3,3)}, \quad (\text{D.16})$$

554 where (we use the label (ℓ, k) , where ℓ is the corresponding exponent in $C^{\odot \ell}$, and k the number
555 of different summation indices)

$$\begin{aligned} C_{ij}^{(1,1)} &= \frac{1}{D} \sum_{\alpha} F_{i\alpha} F_{j\alpha} = C_{ij}, \\ C_{ij}^{(2,1)} &= \frac{1}{D} \sum_{\alpha} F_{i\alpha}^2 F_{j\alpha}^2, \\ C_{ij}^{(2,2)} &= \frac{2}{D^2} \sum_{\alpha < \beta} F_{i\alpha} F_{i\beta} F_{j\alpha} F_{j\beta}, \\ C_{ij}^{(3,1)} &= \frac{1}{D} \sum_{\alpha} F_{i\alpha}^3 F_{j\alpha}^3, \\ C_{ij}^{(3,2)} &= \frac{1}{D^2} \sum_{\alpha \neq \beta} F_{i\alpha}^2 F_{i\beta} F_{j\alpha}^2 F_{j\beta}, \\ C_{ij}^{(3,3)} &= \frac{6}{D^3} \sum_{\alpha < \beta < \gamma} F_{i\alpha} F_{i\beta} F_{i\gamma} F_{j\alpha} F_{j\beta} F_{j\gamma}. \end{aligned} \quad (\text{D.17})$$

556 We can say the following on the matrices $C^{(\ell,k)}$ when N, D are both (generically) large:

557 • $C^{(1,1)} = C$ has a Marchenko-Pastur (MP) spectrum with parameter $\eta_1 = N/D$ and $\sigma^2 =$
558 1 , with D bulk eigenvalues $\lambda = N/D + O(\sqrt{N/D})$ (and $N - D$ zero eigenvalues).

559 • $C^{(2,1)}$ can be written as

$$C_{ij}^{(2,1)} \simeq 1 + \frac{1}{D} \sum_{\alpha} (\Delta_{i\alpha} \Delta_{j\alpha}), \quad (\text{D.18})$$

560 where $\Delta_{i\alpha} = F_{i\alpha}^2 - \mathbb{E}[F_{i\alpha}^2] = F_{i\alpha}^2 - 1$. Notice that $\mathbb{E}[\Delta_{i\alpha}^2] = 2$. From this, it follows
561 that $C^{(2,1)}$ has an MP spectrum with parameter η_1 and $\sigma^2 = 2$, with D bulk eigenvalues
562 $O(\sigma^2 \eta_1)$, plus an additional outlier eigenvalue of order N (due to the finite mean);
563 however, in $C^{\odot 2}$ this matrix is scaled by an additional factor of $1/D$, so it contributes to
564 the sum with D eigenvalues $O(2N/D^2)$ and an outlier $O(N/D)$.

565 • $C^{(2,2)}$ has an MP spectrum with parameter $\eta_2 = 2N/D^2$ and $\sigma^2 = 1$, with $D^2/2$ bulk
566 eigenvalues $O(\eta_2)$.

567 • $C^{(3,1)}$ has an MP spectrum with parameter η_1 and $\sigma^2 = 15$, with D bulk eigenvalues
568 $O(\eta_1)$; however, in $C^{\odot 3}$ this matrix is scaled by an additional factor of $1/D^2$, so it con-
569 tributes to the sum with D eigenvalues $O(N/D^3)$.

570 • $C^{(3,2)}$ can be written as

$$C^{(3,2)} \simeq \frac{1}{D^2} \sum_{\alpha \neq \beta} \Delta_{i\alpha} F_{i\beta} \Delta_{j\alpha} F_{j\beta} + \frac{1}{D} \sum_{\alpha} F_{i\alpha} F_{j\alpha}. \quad (\text{D.19})$$

571 The first addendum (notice that the double sum is not symmetric) has an MP spectrum
572 with parameter N/D^2 and $\sigma^2 = 2$, with D^2 eigenvalues $O(2N/D^2)$, while the second
573 addendum is C ; however, in $C^{\odot 3}$ they are both scaled by a factor $3/D$, so they contribute
574 to the sum with D^2 eigenvalues $O(6N/D^3)$ and with D eigenvalues $O(3N/D^2)$.

575 • $C^{(3,3)}$ has an MP spectrum with parameter $\eta_3 = 6N/D^3$ and $\sigma^2 = 1$, with $D^3/6$ bulk
576 eigenvalues $O(\alpha_3)$.

577 This heuristics is compared with numerical results in Fig. 5, which shows a remarkable ac-
578 cordance. In the main text, we took the approximation $C^{\odot \ell} \simeq C^{(\ell, \ell)}$, and considered the row
579 spaces of $C^{\odot \ell}$ for different ℓ as orthogonal: in Fig. 5, bottom right, we show how the spectrum
580 of a sum of the full matrices $C^{\odot \ell}$ is reasonably approximated by the sum of the (analytical)
581 spectra of the corresponding $C^{(\ell, \ell)}$ matrices, validating our approach.

582 E Determinant of sum of matrices with orthogonal row spaces

583 In this section we derive Eq. (42). Let us take the $N \times N$ matrix given by

$$K = a\mathbb{1} + \sum_{\ell=1}^L b_{\ell} C_{\ell}, \quad (\text{E.1})$$

584 where the matrices C_{ℓ} are such that $\text{rank}(C_{\ell}) = r_{\ell}$, $\sum_{\ell} r_{\ell} \leq N$ and their row spaces \mathcal{R}_{ℓ} (that
585 is, the orthogonal complements to their null spaces) are mutually orthogonal ($\mathcal{R}_{\ell} \perp \mathcal{R}_k$ for
586 $k \neq \ell$). Then,

$$\det K = a^{N - \sum_{\ell} r_{\ell}} \prod_{\ell} \det_{\parallel}^{(\ell)}(a\mathbb{1} + b_{\ell} C_{\ell}), \quad (\text{E.2})$$

587 where $\det_{\parallel}^{(\ell)}(\cdot)$ is the determinant restricted to the row space of C_{ℓ} :

$$\det_{\parallel}^{(\ell)}(a\mathbb{1} + b_{\ell} C_{\ell}) = \prod_{\alpha=1}^{r_{\ell}} (a + b_{\ell} \lambda_{\alpha}), \quad (\text{E.3})$$

588 with λ_{α} the non-zero eigenvalues of C_{ℓ} . Eq. (E.2) can be proven by noticing that, if $\{\mathbf{e}_{\ell}^{\alpha}\}_{\alpha=1}^{r_{\ell}}$
589 is a basis of \mathcal{R}_{ℓ} and $\{\mathbf{e}_{\perp}^{\alpha}\}_{\alpha=1}^{N - \sum_{\ell} r_{\ell}}$ a basis of $(\bigcup_{\ell} \mathcal{R}_{\ell})^{\perp}$, the set $(\bigcup_{\ell} \{\mathbf{e}_{\ell}^{\alpha}\}) \cup \{\mathbf{e}_{\perp}^{\alpha}\}$ is a basis of \mathbb{R}^N
590 in which the matrix K is in block-diagonal form. Moreover, from Eq. (E.3)

$$\det_{\parallel}^{(\ell)}(a\mathbb{1} + b_{\ell} C_{\ell}) = \det(a\mathbb{1} + b_{\ell} C_{\ell}) a^{-(N - r_{\ell})}, \quad (\text{E.4})$$

591 so we can conclude that

$$\det K = a^{N(1-L)} \prod_{\ell} \det(a\mathbb{1} + b_{\ell} C_{\ell}). \quad (\text{E.5})$$

592 F Traces over RS matrices

593 In this section we derive Eq. (45). We need to evaluate

$$\text{Tr} \log (A \otimes \mathbb{1}_N + B \otimes C^{\odot \ell}), \quad (\text{F.1})$$

594 where A, B are RS $n \times n$ matrices. We can write

$$A \otimes \mathbb{1}_N + B \otimes C^{\otimes \ell} = (B \otimes \mathbb{1}_N)(B^{-1}A \oplus C^{\otimes \ell}), \quad (\text{F.2})$$

595 where the Kronecker sum is defined as

$$B^{-1}A \oplus C^{\otimes \ell} = B^{-1}A \otimes \mathbb{1}_N + \mathbb{1}_n \otimes C^{\otimes \ell}. \quad (\text{F.3})$$

596 The eigenvalues of a Kronecker sum are the sums of the eigenvalues of the addenda. Calling
597 σ_a the eigenvalues of $B^{-1}A$ and λ_i the eigenvalues of $C^{\otimes \ell}$, this means that

$$\log \det(B^{-1}A \oplus C^{\otimes \ell}) = \sum_{a,i} \log(\sigma_a + \lambda_i). \quad (\text{F.4})$$

598 Given that $B^{-1}A$ is RS, it has 2 different eigenvalues, σ with multiplicity $n-1$ and $\sigma + n\tilde{\sigma}$ with
599 multiplicity 1, so that for small n

$$\log \det(B^{-1}A \oplus C^{\otimes \ell}) = n \sum_i \log(\sigma + \lambda_i) + n \sum_i \frac{\tilde{\sigma}}{\sigma + \lambda_i}. \quad (\text{F.5})$$

600 In total we get

$$\text{Tr} \log(A \otimes \mathbb{1}_N + B \otimes C^{\otimes \ell}) = nN \log b + nN \frac{\tilde{b}}{b} + n \sum_i \log(\sigma + \lambda_i) + n \sum_i \frac{\tilde{\sigma}}{\sigma + \lambda_i}. \quad (\text{F.6})$$

601 Using the RS algebra, we know that $\sigma = a/b$, $\tilde{\sigma} = (b\tilde{a} - a\tilde{b})/b^2$, so that

$$\begin{aligned} \text{Tr} \log(A \otimes \mathbb{1}_N + B \otimes C^{\otimes \ell}) &= n \text{Tr} \log(a\mathbb{1} + bC^{\otimes \ell}) + n\tilde{a} \text{Tr}(a\mathbb{1} + bC^{\otimes \ell})^{-1} \\ &\quad + n\tilde{b} \text{Tr}[C^{\otimes \ell}(a\mathbb{1} + bC^{\otimes \ell})^{-1}]. \end{aligned} \quad (\text{F.7})$$

602 It only remains to find $a, \tilde{a}, b, \tilde{b}$:

$$a = \beta(\zeta + \hat{\chi}^{(0)}), \quad \tilde{a} = -\beta^2 \hat{q}^{(0)}, \quad b = \beta \hat{\chi}^{(\ell)}/\eta_\ell, \quad \tilde{b} = -\beta^2[\hat{q}^{(\ell)} + (\hat{m}^{(\ell)})^2]/\eta_\ell. \quad (\text{F.8})$$

603 We define $\gamma_\ell = a/b = \eta_\ell(\zeta + \hat{\chi}^{(0)})/\hat{\chi}^{(\ell)}$ to get Eq. (45).

604 G Replica-symmetric free energy

605 In this section we report the main steps to obtain the terms S_M and S_P in Eq. (51) and (52),
606 that is the measure and pattern contributions to the free energy.

607 G.1 Measure contribution

608 By plugging the RS ansatz (43), (44) and Eq. (45) in Eq. (40), we readily obtain

$$\begin{aligned} S_M &= -n\beta \sum_{\ell=1}^L \frac{m^{(\ell)} \hat{m}^{(\ell)}}{\eta_\ell} + \frac{n}{2} \sum_{\ell=0}^L \frac{1}{\eta_\ell} [\chi^{(\ell)} \hat{\chi}^{(\ell)} + \beta(q^{(\ell)} \hat{\chi}^{(\ell)} - \chi^{(\ell)} \hat{q}^{(\ell)})] \\ &\quad - \frac{n}{2} \log(\beta(\zeta + \hat{\chi}^{(0)})) + \frac{\beta n(1-L)}{2} \frac{\hat{q}^{(0)}}{\zeta + \hat{\chi}^{(0)}} - \frac{n}{2N} \sum_{\ell=1}^L \text{Tr} \log(\mathbb{1} + C^{\otimes \ell}/\gamma_\ell) \\ &\quad + \frac{\beta n}{2N} \sum_{\ell=1}^L \eta_\ell \frac{\hat{q}^{(0)}}{\hat{\chi}^{(\ell)}} \text{Tr}(\gamma_\ell \mathbb{1} + C^{\otimes \ell})^{-1} + \frac{\beta n}{2N} \sum_{\ell=1}^L \frac{\hat{q}^{(\ell)} + (\hat{m}^{(\ell)})^2}{\hat{\chi}^{(\ell)}} \text{Tr}[C^{\otimes \ell}(\gamma_\ell \mathbb{1} + C^{\otimes \ell})^{-1}]. \end{aligned} \quad (\text{G.1})$$

609 We obtain Eq. (51) by keeping the leading order terms for β large and using Eq. (48).

610 G.2 Pattern contribution

611 S_P is a function only of the order parameters:

$$S_P = \log \left[\int d\nu \prod_{a=1}^n d\lambda^a p(\nu, \{\lambda^a\}) \int dy p(y|\nu) e^{-\beta \sum_a \mathcal{L}(y, \lambda^a)} \right], \quad (\text{G.2})$$

$$p(\nu, \{\lambda^a\}) = \mathcal{N} \left((\nu, \{\lambda^a\}) \mid (0, \{t^a\}), \begin{pmatrix} 1 & M^\top \\ M & Q \end{pmatrix} \right).$$

612 With the RS ansatz and for small n ,

$$S_P = \log \left[\int dy d\nu \prod_{a=1}^n d\lambda^a p(y|\nu) e^{-\frac{\nu^2}{2} + \beta \frac{m^* \nu}{\chi^*} \sum_a \lambda^a - \frac{\beta}{2\chi^*} \sum_a \lambda_a^2 - \beta \sum_a \mathcal{L}(y, \lambda^a + t^*) - \beta^2 \frac{m^{*2} - q^*}{2\chi^{*2}} \sum_{a,b} \lambda^a \lambda^b} \right]$$

$$- \frac{n}{2} \log(2\pi) - \frac{1}{2} \log \det \begin{pmatrix} 1 & M^\top \\ M & Q \end{pmatrix}. \quad (\text{G.3})$$

613 To factorize the integral over replicas we use the Hubbard-Stratonovich transformation

$$e^{-\beta^2 \frac{m^{*2} - q^*}{2\chi^{*2}} \sum_{a,b} \lambda^a \lambda^b} = \mathbb{E}_\xi e^{\beta \frac{\sqrt{q^* - m^{*2}}}{\chi^*} \sum_a \lambda^a \xi}, \quad (\text{G.4})$$

614 obtaining, to leading order in n ,

$$S_P = -\frac{n}{2} \log \frac{\chi^*}{\beta} - \frac{n\beta}{2} \frac{q^*}{\chi^*} + n \mathbb{E}_\xi \int dy D\nu p(y|\nu)$$

$$\times \log \int d\lambda e^{\beta \left(\sqrt{q^* - m^{*2}} \xi + m^* \nu \right) \frac{\lambda}{\chi^*} - \frac{\beta \lambda^2}{2\chi^*} - \beta \mathcal{L}(y, \lambda + t^*)}. \quad (\text{G.5})$$

615 For our choice of loss (10) and for β large, we obtain Eq. (52).

616 H Asymptotic limits of the saddle-point equations

617 The system of saddle-point equations can be studied in different asymptotic limits, as we anticipated in Sec. 6:

- 619 (i) $N, P, D \rightarrow \infty, P/N \rightarrow 0, P/D^K$ finite;
- 620 (ii) $N, P, D \rightarrow \infty, N/D^L$ finite, P/N finite;
- 621 (iii) $N, P, D \rightarrow \infty, P/N \rightarrow \infty, N/D^L$ finite.

622 H.1 Case (i)

623 In the limit where N scales faster to infinity than P , Eq. (54) reduces to

$$\hat{\chi}^{(0)} \rightarrow 0, \quad \chi^{(0)} \rightarrow \frac{1}{\zeta},$$

$$\hat{\chi}^{(\ell)} \rightarrow \begin{cases} \infty & \text{for } \ell < K, \\ \frac{P}{\binom{D}{K}} \frac{\mu_K^2}{K!(1+\chi^*)} & \text{for } \ell = K, \\ 0 & \text{for } \ell > K, \end{cases} \quad \chi^{(\ell)} \rightarrow \begin{cases} 0 & \text{for } \ell < K, \\ \frac{1}{\hat{\chi}^{(K)} + \zeta} & \text{for } \ell = K, \\ \frac{1}{\zeta} & \text{for } \ell > K, \end{cases} \quad (\text{H.1})$$

624 where we used the asymptotic results for the Stieltjes transformation of the Marchenko-Pastur
625 distribution,

$$1 - \gamma_\ell g(-\gamma_\ell; \eta_\ell) \sim \begin{cases} \frac{1}{\eta_\ell} & \text{for } \ell < K, \\ \frac{1}{\eta_K + \gamma_K} & \text{for } \ell = K, \\ \frac{1}{\gamma_\ell} & \text{for } \ell > K. \end{cases} \quad (\text{H.2})$$

626 Notice that now, consistently,

$$\chi^* = \frac{\mu_{\perp,K}^2}{\zeta} + \frac{\mu_K^2}{K!} \chi^{(K)}, \quad (\text{H.3})$$

627 because $\mu_{\perp,L}^2$ recombines with the terms coming from $K < \ell \leq L$ to give $\mu_{\perp,K}^2$. Eq. (55) reduces
628 to

$$m^{(0)} = \frac{\langle y \rangle}{\mu_0}$$

$$\hat{m}^{(\ell)} \rightarrow \begin{cases} \infty & \text{for } \ell < K, \\ \frac{P}{\binom{D}{K}} \frac{\mu_K \tau_K}{\sqrt{K!}} \frac{\langle y \nu \rangle}{1 + \chi^*}, & \text{for } \ell = K, \\ 0 & \text{for } \ell > K, \end{cases} \quad m^{(\ell)} \rightarrow \begin{cases} \sqrt{\ell!} \frac{\tau_\ell}{\mu_\ell} \langle y \nu \rangle & \text{for } \ell < K, \\ \sqrt{K!} \frac{\tau_K}{\mu_K} \langle y \nu \rangle (1 - \zeta \chi^{(K)}) & \text{for } \ell = K, \\ 0 & \text{for } \ell > K, \end{cases} \quad (\text{H.4})$$

629 while Eq. (56) becomes

$$\hat{q}^{(0)} \rightarrow 0, \quad q^{(0)} \rightarrow 0$$

$$\hat{q}^{(\ell)} \rightarrow \begin{cases} \infty & \text{for } \ell < K, \\ \frac{P}{\binom{D}{K}} \frac{\mu_K^2}{K!} \frac{(\mu_0 m^{(0)} - y)^2 - 2 \langle y \nu \rangle m^* + q^*}{(1 + \chi^*)^2} & \text{for } \ell = K, \\ 0 & \text{for } \ell > K, \end{cases} \quad q^{(\ell)} \rightarrow \begin{cases} \ell! \frac{\tau_\ell^2}{\mu_\ell^2} \langle y \nu \rangle^2 & \text{for } \ell < K, \\ \frac{(\hat{m}^{(K)})^2 + \hat{q}^{(K)}}{(\hat{\chi}^{(K)} + \zeta)^2} & \text{for } \ell = K, \\ 0 & \text{for } \ell > K, \end{cases} \quad (\text{H.5})$$

630 where now

$$q^* = \langle y \nu \rangle^2 \sum_{\ell=1}^{K-1} \tau_\ell^2 + \frac{\mu_K^2}{K!} q^{(K)}, \quad m^* = \langle y \nu \rangle \sum_{\ell=1}^{K-1} \tau_\ell^2 + \frac{\mu_K \tau_K}{\sqrt{K!}} m^{(K)}. \quad (\text{H.6})$$

631 H.2 Case (ii)

632 In the limit where both P and N scale in the the same way, $N \sim P \sim O(D^L)$, we have, for

633 $0 < \ell < L$,

$$\begin{aligned} \hat{\chi}^{(\ell)} &\rightarrow \infty, & \hat{m}^{(\ell)} &\rightarrow \infty, & \hat{q}^{(\ell)} &\rightarrow \infty, \\ \chi^{(\ell)} &\rightarrow 0, & m^{(\ell)} &\rightarrow \sqrt{\ell!} \frac{\tau_\ell}{\mu_\ell} \langle y \nu \rangle, & q^{(\ell)} &\rightarrow \ell! \frac{\tau_\ell^2}{\mu_\ell^2} \langle y \nu \rangle^2. \end{aligned} \quad (\text{H.7})$$

634 For the other parameters we need to solve the equations for χ

$$\begin{aligned} \hat{\chi}^{(0)} &= \frac{P}{N} \frac{\mu_{\perp,L}^2}{1 + \chi^*}, & \chi^{(0)} &= \frac{\gamma_L g_L(-\gamma_L)}{\hat{\chi}^{(0)} + \zeta}, \\ \hat{\chi}^{(L)} &= \frac{P}{\binom{D}{L} L!} \frac{\mu_L^2}{1 + \chi^*}, & \chi^{(L)} &= \frac{N}{\binom{D}{L}} \frac{1 - \gamma_L g_L(-\gamma_L)}{\hat{\chi}^{(L)}}, \end{aligned} \quad (\text{H.8})$$

635 for m ,

$$m^{(0)} = \langle y \rangle / \mu_0, \quad m^{(L)} = \chi^{(L)} \hat{m}^{(L)}, \quad \hat{m}^{(L)} = \frac{P}{\binom{D}{L}} \frac{\mu_L \tau_L}{\sqrt{L!}} \frac{\langle y \nu \rangle}{1 + \chi^*}, \quad (\text{H.9})$$

636 and for q

$$\begin{aligned}
\hat{q}^{(0)} &= \frac{P}{N} \mu_{\perp,L}^2 \frac{\langle (\mu_0 m^{(0)} - y)^2 \rangle - 2\langle y \nu \rangle m^* + q^*}{(1 + \chi^*)^2}, \\
\hat{q}^{(L)} &= \frac{P}{\binom{D}{L}} \frac{\mu_L^2 \langle (\mu_0 m^{(0)} - y)^2 \rangle - 2\langle y \nu \rangle m^* + q^*}{L! (1 + \chi^*)^2}, \\
q^{(0)} &= \frac{\hat{q}^{(0)}}{(\zeta + \hat{\chi}^{(0)})^2} \gamma_L^2 g_L'(-\gamma_L) + \frac{\hat{m}^{(L)2} + \hat{q}^{(L)}}{(\zeta + \hat{\chi}^{(0)}) \hat{\chi}^{(L)}} [\gamma_L g_L(-\gamma_L) - \gamma_L^2 g_L'(-\gamma_L)], \\
q^{(L)} &= \frac{N}{\binom{D}{L}} \frac{\hat{q}^{(0)}}{(\zeta + \hat{\chi}^{(0)}) \hat{\chi}^{(L)}} [\gamma_L g_L(-\gamma_L) - \gamma_L^2 g_L'(-\gamma_L)] \\
&\quad + \frac{N}{\binom{D}{L}} \frac{\hat{m}^{(L)2} + \hat{q}^{(L)}}{\hat{\chi}^{(L)2}} [1 + \gamma_L^2 g_L'(-\gamma_L) - 2\gamma_L g_L(-\gamma_L)].
\end{aligned} \tag{H.10}$$

637 The values χ^* , m^* and q^* are consistent with their definition. At variance with case (i), $\chi^{(0)}$ and
638 $q^{(0)}$ have non-trivial values, responsible for the interpolation peak appearing in this regime.
639 Notice that their value is controlled explicitly by the regularizer ζ : the lower it is, the sharper
640 is the peak. Moreover, the spectral function relative to the active component, g_L , also gives a
641 non-trivial contribution.

642 H.3 Case (iii)

643 In the limit where P is scaling faster than N to infinity, we have that for all $0 < \ell < L$ the
644 order parameters behave as in Eq. (H.7), meaning that the degree- L student learns perfectly
645 all the terms of the teacher of degree less than L , as the amount of training data P is effectively
646 infinite. In this case

$$\gamma_L = \frac{L! \mu_{\perp,L}^2}{\mu_L^2} \tag{H.11}$$

647 and we have $\chi^{(L)}, \hat{\chi}^{(L)} \rightarrow 0$; $\hat{q}^{(0)}, \hat{q}^{(L)} \rightarrow \infty$ and

$$\begin{aligned}
m^{(L)} &= \eta_L \langle y \nu \rangle \sqrt{L!} \frac{\tau_L}{\mu_L} (1 - \gamma_L g_L(-\gamma_L)), \\
q^{(0)} &= \eta_L \langle y \nu \rangle^2 \frac{\tau_L^2}{\mu_{\perp,L}^2} [\gamma_L g_L(-\gamma_L) - \gamma_L^2 g_L'(-\gamma_L)], \\
q^{(L)} &= \eta_L \langle y \nu \rangle^2 L! \frac{\tau_L^2}{\mu_L^2} [1 + \gamma_L^2 g_L'(-\gamma_L) - 2\gamma_L g_L(-\gamma_L)].
\end{aligned} \tag{H.12}$$

648 I Numerical experiments

649 All numerical experiments were done in Python using JAX, [75], to generate the synthetic ran-
650 dom data, and scikit, [76], to optimize the parameters. The optimizer has a simple analytic
651 form given by (18). Nevertheless, it is potentially inefficient to implement the formula naively,
652 as it would require the inversion of a very large matrix. Since we used very large values of N
653 and P , we performed the ridge regression with the function `sklearn.linear_model.Ridge`.
654 In this way we could explore regimes of N, P up to order D^3 .

655 Almost all numerical experiments were performed with $D = 30$. In most of the simulations
656 we sampled 50 times for each combination of N, P, D . For the right panel of Figure 3 we used a
657 larger number of samples since in that case both $D = 30$ and $P = 40 \sim 400$ were small, hence

658 the generalization error had higher variability. For $N < 3000$ we used 500, 200, 300 samples
659 respectively for $P = 40, 200, 400$. For $N > 3000$ we used 100, 100, 50 samples respectively for
660 $P = 40, 200, 400$.

661 A GitHub repository collecting the code needed to reproduce the figures of this paper (both
662 numerical experiments and theoretical curves from the integration of the saddle-point equa-
663 tions) can be found at [77].

664 References

- 665 [1] R. M. Neal, *Priors for Infinite Networks*, pp. 29–53, Springer New York, New York, NY,
666 ISBN 978-1-4612-0745-0, doi:10.1007/978-1-4612-0745-0_2 (1996).
- 667 [2] C. Williams, *Computing with infinite networks*, In M. Mozer, M. Jordan
668 and T. Petsche, eds., *Advances in Neural Information Processing Systems*,
669 vol. 9. MIT Press (1996), [https://proceedings.neurips.cc/paper/1996/file/
670 ae5e3ce40e0404a45ecacaaf05e5f735-Paper.pdf](https://proceedings.neurips.cc/paper/1996/file/ae5e3ce40e0404a45ecacaaf05e5f735-Paper.pdf).
- 671 [3] J. Lee, J. Sohl-dickstein, J. Pennington, R. Novak, S. Schoenholz and Y. Bahri, *Deep neural*
672 *networks as Gaussian processes*, In *International Conference on Learning Representations*
673 (2018), <https://openreview.net/forum?id=B1EA-M-OZ>.
- 674 [4] A. G. de G. Matthews, J. Hron, M. Rowland, R. E. Turner and Z. Ghahramani, *Gaussian*
675 *process behaviour in wide deep neural networks*, In *International Conference on Learning*
676 *Representations* (2018), <https://openreview.net/forum?id=H1-nGgWC->.
- 677 [5] G. Naveh and Z. Ringel, *A self consistent theory of Gaussian processes captures feature*
678 *learning effects in finite CNNs*, In M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang and
679 J. W. Vaughan, eds., *Advances in Neural Information Processing Systems*, vol. 34, pp.
680 21352–21364. Curran Associates, Inc. (2021), [https://proceedings.neurips.cc/paper_
681 files/paper/2021/file/b24d21019de5e59da180f1661904f49a-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/b24d21019de5e59da180f1661904f49a-Paper.pdf).
- 682 [6] S. Ariosto, R. Pacelli, F. Ginelli, M. Gherardi and P. Rotondo, *Universal mean-field up-*
683 *per bound for the generalization gap of deep neural networks*, *Phys. Rev. E* **105**, 064309
684 (2022), doi:10.1103/PhysRevE.105.064309, <https://arxiv.org/abs/2201.11022>.
- 685 [7] R. Pacelli, S. Ariosto, M. Pastore, F. Ginelli, M. Gherardi and P. Rotondo, *A statistical*
686 *mechanics framework for Bayesian deep neural networks beyond the infinite-width limit*,
687 *Nature Machine Intelligence* **5**(12), 1497 (2023), doi:10.1038/s42256-023-00767-6.
- 688 [8] A. Atanasov, B. Bordelon, S. Sainathan and C. Pehlevan, *The onset of variance-limited*
689 *behavior for networks in the lazy and rich regimes*, In *The Eleventh International Conference*
690 *on Learning Representations* (2023), <https://openreview.net/forum?id=JLinxPOVTh7>.
- 691 [9] I. Seroussi, G. Naveh and Z. Ringel, *Separation of scales and a thermodynamic descrip-*
692 *tion of feature learning in some CNNs*, *Nature Communications* **14**(1), 908 (2023),
693 doi:10.1038/s41467-023-36361-y, <https://arxiv.org/abs/2112.15383>.
- 694 [10] H. Cui, F. Krzakala and L. Zdeborova, *Bayes-optimal learning of deep random networks*
695 *of extensive-width*, In A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato and
696 J. Scarlett, eds., *Proceedings of the 40th International Conference on Machine Learning*,
697 vol. 202 of *Proceedings of Machine Learning Research*, pp. 6468–6521. PMLR (2023),
698 <https://proceedings.mlr.press/v202/cui23b.html>.

- 699 [11] L. Chizat, E. Oyallon and F. Bach, *On lazy training in differentiable pro-*
700 *gramming*, In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc,
701 E. Fox and R. Garnett, eds., *Advances in Neural Information Processing Systems*,
702 vol. 32. Curran Associates, Inc. (2019), [https://proceedings.neurips.cc/paper/2019/](https://proceedings.neurips.cc/paper/2019/file/ae614c557843b1df326cb29c57225459-Paper.pdf)
703 [file/ae614c557843b1df326cb29c57225459-Paper.pdf](https://proceedings.neurips.cc/paper/2019/file/ae614c557843b1df326cb29c57225459-Paper.pdf).
- 704 [12] A. Jacot, F. Gabriel and C. Hongler, *Neural tangent kernel: Convergence and gener-*
705 *alization in neural networks*, In S. Bengio, H. Wallach, H. Larochelle, K. Grauman,
706 N. Cesa-Bianchi and R. Garnett, eds., *Advances in Neural Information Processing Sys-*
707 *tems*, vol. 31. Curran Associates, Inc. (2018), [https://proceedings.neurips.cc/paper/](https://proceedings.neurips.cc/paper/2018/file/5a4be1fa34e62bb8a6ec6b91d2462f5a-Paper.pdf)
708 [2018/file/5a4be1fa34e62bb8a6ec6b91d2462f5a-Paper.pdf](https://proceedings.neurips.cc/paper/2018/file/5a4be1fa34e62bb8a6ec6b91d2462f5a-Paper.pdf).
- 709 [13] A. Bietti and J. Mairal, *On the inductive bias of neural tangent kernels*, In
710 H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox and R. Gar-
711 nett, eds., *Advances in Neural Information Processing Systems*, vol. 32. Cur-
712 ran Associates, Inc. (2019), [https://proceedings.neurips.cc/paper/2019/file/](https://proceedings.neurips.cc/paper/2019/file/c4ef9c39b300931b69a36fb3dbb8d60e-Paper.pdf)
713 [c4ef9c39b300931b69a36fb3dbb8d60e-Paper.pdf](https://proceedings.neurips.cc/paper/2019/file/c4ef9c39b300931b69a36fb3dbb8d60e-Paper.pdf).
- 714 [14] A. Montanari and Y. Zhong, *The interpolation phase transition in neural networks: Mem-*
715 *orization and generalization under lazy training*, *The Annals of Statistics* **50**(5), 2816
716 (2022), doi:[10.1214/22-AOS2211](https://doi.org/10.1214/22-AOS2211), <https://arxiv.org/abs/2007.12826>.
- 717 [15] N. Cristianini and J. Shawe-Taylor, *An Introduction to Support Vector Ma-*
718 *chines and Other Kernel-based Learning Methods*, Cambridge University Press,
719 doi:[10.1017/CBO9780511801389](https://doi.org/10.1017/CBO9780511801389) (2000).
- 720 [16] H. Yoon and J.-H. Oh, *Learning of higher-order perceptrons with tunable complexities*, *Jour-*
721 *nal of Physics A: Mathematical and General* **31**(38), 7771 (1998), doi:[10.1088/0305-](https://doi.org/10.1088/0305-4470/31/38/012)
722 [4470/31/38/012](https://doi.org/10.1088/0305-4470/31/38/012).
- 723 [17] R. Dietrich, M. Opper and H. Sompolinsky, *Statistical mechanics of support vector net-*
724 *works*, *Phys. Rev. Lett.* **82**, 2975 (1999), doi:[10.1103/PhysRevLett.82.2975](https://doi.org/10.1103/PhysRevLett.82.2975).
- 725 [18] B. Bordelon, A. Canatar and C. Pehlevan, *Spectrum dependent learning curves in kernel*
726 *regression and wide neural networks*, In H. D. III and A. Singh, eds., *Proceedings of the*
727 *37th International Conference on Machine Learning*, vol. 119 of *Proceedings of Machine*
728 *Learning Research*, pp. 1024–1034. PMLR (2020), [http://proceedings.mlr.press/v119/](http://proceedings.mlr.press/v119/bordelon20a/bordelon20a.pdf)
729 [bordelon20a/bordelon20a.pdf](http://proceedings.mlr.press/v119/bordelon20a/bordelon20a.pdf).
- 730 [19] A. Canatar, B. Bordelon and C. Pehlevan, *Spectral bias and task-model alignment explain*
731 *generalization in kernel regression and infinitely wide neural networks*, *Nature Communi-*
732 *cations* **12**(1), 2914 (2021), doi:[10.1038/s41467-021-23103-1](https://doi.org/10.1038/s41467-021-23103-1), [https://arxiv.org/abs/](https://arxiv.org/abs/2006.13198)
733 [2006.13198](https://arxiv.org/abs/2006.13198).
- 734 [20] T. Misiakiewicz, *Spectrum of inner-product kernel matrices in the polynomial regime and*
735 *multiple descent phenomenon in kernel ridge regression*, doi:[10.48550/ARXIV.2204.10425](https://doi.org/10.48550/ARXIV.2204.10425)
736 (2022).
- 737 [21] H. Hu and Y. M. Lu, *Sharp asymptotics of kernel ridge regression beyond the linear regime*,
738 doi:[10.48550/ARXIV.2205.06798](https://doi.org/10.48550/ARXIV.2205.06798) (2022).
- 739 [22] L. Xiao, H. Hu, T. Misiakiewicz, Y. M. Lu and J. Pennington, *Precise learning curves*
740 *and higher-order scaling limits for dot-product kernel regression*, *Journal of Statistical*
741 *Mechanics: Theory and Experiment* **2023**(11), 114005 (2023), doi:[10.1088/1742-](https://doi.org/10.1088/1742-5468/ad01b7)
742 [5468/ad01b7](https://doi.org/10.1088/1742-5468/ad01b7).

- 743 [23] A. Rahimi and B. Recht, *Random features for large-scale kernel machines*, In J. Platt,
744 D. Koller, Y. Singer and S. Roweis, eds., *Advances in Neural Information Processing Sys-*
745 *tems*, vol. 20. Curran Associates, Inc. (2007), [https://proceedings.neurips.cc/paper/](https://proceedings.neurips.cc/paper/2007/file/013a006f03dbc5392effeb8f18fda755-Paper.pdf)
746 [2007/file/013a006f03dbc5392effeb8f18fda755-Paper.pdf](https://proceedings.neurips.cc/paper/2007/file/013a006f03dbc5392effeb8f18fda755-Paper.pdf).
- 747 [24] M.-F. Balcan, A. Blum and S. Vempala, *Kernels as features: On kernels, margins, and*
748 *low-dimensional mappings*, *Machine Learning* **65**(1), 79 (2006), doi:[10.1007/s10994-](https://doi.org/10.1007/s10994-006-7550-1)
749 [006-7550-1](https://doi.org/10.1007/s10994-006-7550-1).
- 750 [25] A. Rahimi and B. Recht, *Weighted sums of random kitchen sinks: Replac-*
751 *ing minimization with randomization in learning*, In D. Koller, D. Schuurmans,
752 Y. Bengio and L. Bottou, eds., *Advances in Neural Information Processing Systems*,
753 vol. 21. Curran Associates, Inc. (2008), [https://proceedings.neurips.cc/paper/2008/](https://proceedings.neurips.cc/paper/2008/file/0efe32849d230d7f53049ddc4a4b0c60-Paper.pdf)
754 [file/0efe32849d230d7f53049ddc4a4b0c60-Paper.pdf](https://proceedings.neurips.cc/paper/2008/file/0efe32849d230d7f53049ddc4a4b0c60-Paper.pdf).
- 755 [26] A. Rahimi and B. Recht, *Uniform approximation of functions with random bases*, In *2008*
756 *46th Annual Allerton Conference on Communication, Control, and Computing*, pp. 555–
757 561, doi:[10.1109/ALLERTON.2008.4797607](https://doi.org/10.1109/ALLERTON.2008.4797607) (2008).
- 758 [27] B. Ghorbani, S. Mei, T. Misiakiewicz and A. Montanari, *Linearized two-layers neural net-*
759 *works in high dimension*, *The Annals of Statistics* **49**(2), 1029 (2021), doi:[10.1214/20-](https://doi.org/10.1214/20-AOS1990)
760 [AOS1990](https://doi.org/10.1214/20-AOS1990), <https://arxiv.org/abs/1904.12191>.
- 761 [28] B. Ghorbani, S. Mei, T. Misiakiewicz and A. Montanari, *Limitations of lazy training of*
762 *two-layers neural network*, In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-
763 Buc, E. Fox and R. Garnett, eds., *Advances in Neural Information Processing Systems*,
764 vol. 32. Curran Associates, Inc. (2019), [https://proceedings.neurips.cc/paper/2019/](https://proceedings.neurips.cc/paper/2019/file/c133fb1bb634af68c5088f3438848bfd-Paper.pdf)
765 [file/c133fb1bb634af68c5088f3438848bfd-Paper.pdf](https://proceedings.neurips.cc/paper/2019/file/c133fb1bb634af68c5088f3438848bfd-Paper.pdf).
- 766 [29] S. Mei and A. Montanari, *The generalization error of random features regression: Precise*
767 *asymptotics and the double descent curve*, *Communications on Pure and Applied Mathe-*
768 *matics* **75**(4), 667 (2022), doi:[10.1002/cpa.22008](https://doi.org/10.1002/cpa.22008), <https://arxiv.org/abs/1908.05355>.
- 769 [30] B. Ghorbani, S. Mei, T. Misiakiewicz and A. Montanari, *When do neural networks out-*
770 *perform kernel methods?*, In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan and
771 H. Lin, eds., *Advances in Neural Information Processing Systems*, vol. 33, pp. 14820–
772 14830. Curran Associates, Inc. (2020), [https://proceedings.neurips.cc/paper/2020/](https://proceedings.neurips.cc/paper/2020/file/a9df2255ad642b923d95503b9a7958d8-Paper.pdf)
773 [file/a9df2255ad642b923d95503b9a7958d8-Paper.pdf](https://proceedings.neurips.cc/paper/2020/file/a9df2255ad642b923d95503b9a7958d8-Paper.pdf).
- 774 [31] S. Mei, T. Misiakiewicz and A. Montanari, *Generalization error of random feature and*
775 *kernel methods: Hypercontractivity and kernel matrix concentration*, *Applied and Com-*
776 *putational Harmonic Analysis* **59**, 3 (2022), doi:[10.1016/j.acha.2021.12.003](https://doi.org/10.1016/j.acha.2021.12.003), Special
777 Issue on Harmonic Analysis and Machine Learning, <https://arxiv.org/abs/2101.10588>.
- 778 [32] S. Mei, T. Misiakiewicz and A. Montanari, *Learning with invariances in random features*
779 *and kernel models*, In M. Belkin and S. Kpotufe, eds., *Proceedings of Thirty Fourth Con-*
780 *ference on Learning Theory*, vol. 134 of *Proceedings of Machine Learning Research*, pp.
781 3351–3418. PMLR (2021), <http://proceedings.mlr.press/v134/mei21a/mei21a.pdf>.
- 782 [33] A. Montanari and B. N. Saeed, *Universality of empirical risk minimization*, In P.-L. Loh
783 and M. Raginsky, eds., *Proceedings of Thirty Fifth Conference on Learning Theory*, vol.
784 178 of *Proceedings of Machine Learning Research*, pp. 4310–4312. PMLR (2022), [https:](https://proceedings.mlr.press/v178/montanari22a.html)
785 [//proceedings.mlr.press/v178/montanari22a.html](https://proceedings.mlr.press/v178/montanari22a.html).

- 786 [34] P. L. Bartlett, A. Montanari and A. Rakhlin, *Deep learning: a statistical viewpoint*, Acta
787 Numerica **30**, 87–201 (2021), doi:[10.1017/S0962492921000027](https://doi.org/10.1017/S0962492921000027).
- 788 [35] M. Belkin, D. Hsu, S. Ma and S. Mandal, *Reconciling modern machine-learning practice
789 and the classical bias–variance trade-off*, Proceedings of the National Academy of Sciences
790 **116**(32), 15849 (2019), doi:[10.1073/pnas.1903070116](https://doi.org/10.1073/pnas.1903070116).
- 791 [36] S. Goldt, M. Mézard, F. Krzakala and L. Zdeborová, *Modeling the influence of data structure
792 on learning in neural networks: The hidden manifold model*, Phys. Rev. X **10**, 041044
793 (2020), doi:[10.1103/PhysRevX.10.041044](https://doi.org/10.1103/PhysRevX.10.041044), <https://arxiv.org/abs/1909.11500>.
- 794 [37] F. Gerace, B. Loureiro, F. Krzakala, M. Mézard and L. Zdeborová, *Generalisation error
795 in learning with random features and the hidden manifold model*, Journal of Statisti-
796 cal Mechanics: Theory and Experiment **2021**(12), 124013 (2021), doi:[10.1088/1742-
797 5468/ac3ae6](https://doi.org/10.1088/1742-5468/ac3ae6), <https://arxiv.org/abs/2002.09339>.
- 798 [38] S. Goldt, B. Loureiro, G. Reeves, F. Krzakala, M. Mezard and L. Zdeborová, *The Gaussian
799 equivalence of generative models for learning with shallow neural networks*, In J. Bruna,
800 J. Hesthaven and L. Zdeborová, eds., *Proceedings of the 2nd Mathematical and Scientific
801 Machine Learning Conference*, vol. 145 of *Proceedings of Machine Learning Research*, pp.
802 426–471. PMLR (2022), <https://proceedings.mlr.press/v145/goldt22a/goldt22a.pdf>.
- 803 [39] B. Loureiro, C. Gerbelot, H. Cui, S. Goldt, F. Krzakala, M. Mezard and L. Zde-
804 borová, *Learning curves of generic features maps for realistic datasets with a teacher-
805 student model*, In M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang and J. W.
806 Vaughan, eds., *Advances in Neural Information Processing Systems*, vol. 34, pp. 18137–
807 18151. Curran Associates, Inc. (2021), [https://proceedings.neurips.cc/paper/2021/
808 file/9704a4fc48ae88598dcbdcdf57f3fdef-Paper.pdf](https://proceedings.neurips.cc/paper/2021/file/9704a4fc48ae88598dcbdcdf57f3fdef-Paper.pdf).
- 809 [40] M. Refinetti, S. Goldt, F. Krzakala and L. Zdeborová, *Classifying high-dimensional Gaus-
810 sian mixtures: Where kernel methods fail and neural networks succeed*, In M. Meila and
811 T. Zhang, eds., *Proceedings of the 38th International Conference on Machine Learning*,
812 vol. 139 of *Proceedings of Machine Learning Research*, pp. 8936–8947. PMLR (2021),
813 <http://proceedings.mlr.press/v139/refinetti21b/refinetti21b.pdf>.
- 814 [41] H. Cui, B. Loureiro, F. Krzakala and L. Zdeborová, *Generalization error rates in kernel re-
815 gression: The crossover from the noiseless to noisy regime*, In A. Beygelzimer, Y. Dauphin,
816 P. Liang and J. W. Vaughan, eds., *Advances in Neural Information Processing Systems*
817 (2021), https://openreview.net/forum?id=Da_EHrAcfwd.
- 818 [42] D. Schröder, H. Cui, D. Dmitriev and B. Loureiro, *Deterministic equivalent and error
819 universality of deep random features learning*, In A. Krause, E. Brunskill, K. Cho, B. En-
820 gelhardt, S. Sabato and J. Scarlett, eds., *Proceedings of the 40th International Conference
821 on Machine Learning*, vol. 202 of *Proceedings of Machine Learning Research*, pp. 30285–
822 30320. PMLR (2023), [https://proceedings.mlr.press/v202/schroder23a/schroder23a.
823 pdf](https://proceedings.mlr.press/v202/schroder23a/schroder23a.pdf).
- 824 [43] S. Chung, D. D. Lee and H. Sompolinsky, *Linear readout of object manifolds*, Phys. Rev.
825 E **93**, 060301 (2016), doi:[10.1103/PhysRevE.93.060301](https://doi.org/10.1103/PhysRevE.93.060301), [https://arxiv.org/abs/1512.
826 01834](https://arxiv.org/abs/1512.01834).
- 827 [44] S. Chung, D. D. Lee and H. Sompolinsky, *Classification and geometry of general perceptual
828 manifolds*, Phys. Rev. X **8**, 031003 (2018), doi:[10.1103/PhysRevX.8.031003](https://doi.org/10.1103/PhysRevX.8.031003), [https://
829 arxiv.org/abs/1710.06487](https://arxiv.org/abs/1710.06487).

- 830 [45] F. Borra, M. C. Lagomarsino, P. Rotondo and M. Gherardi, *Generalization from correlated*
831 *sets of patterns in the perceptron*, Journal of Physics A: Mathematical and Theoretical
832 **52**(38), 384004 (2019), doi:[10.1088/1751-8121/ab3709](https://doi.org/10.1088/1751-8121/ab3709), [https://arxiv.org/abs/1903.](https://arxiv.org/abs/1903.06818)
833 [06818](https://arxiv.org/abs/1903.06818).
- 834 [46] P. Rotondo, M. C. Lagomarsino and M. Gherardi, *Counting the learnable func-*
835 *tions of geometrically structured data*, Phys. Rev. Res. **2**, 023169 (2020),
836 doi:[10.1103/PhysRevResearch.2.023169](https://doi.org/10.1103/PhysRevResearch.2.023169), <https://arxiv.org/abs/1903.12021>.
- 837 [47] P. Rotondo, M. Pastore and M. Gherardi, *Beyond the storage capacity:*
838 *Data-driven satisfiability transition*, Phys. Rev. Lett. **125**, 120601 (2020),
839 doi:[10.1103/PhysRevLett.125.120601](https://doi.org/10.1103/PhysRevLett.125.120601), <https://arxiv.org/abs/2005.09992>.
- 840 [48] M. Pastore, P. Rotondo, V. Erba and M. Gherardi, *Statistical learning theory of structured*
841 *data*, Phys. Rev. E **102**, 032119 (2020), doi:[10.1103/PhysRevE.102.032119](https://doi.org/10.1103/PhysRevE.102.032119), [https:](https://arxiv.org/abs/2005.10002)
842 [//arxiv.org/abs/2005.10002](https://arxiv.org/abs/2005.10002).
- 843 [49] M. Pastore, *Critical properties of the SAT/UNSAT transitions in the classification problem*
844 *of structured data*, Journal of Statistical Mechanics: Theory and Experiment **2021**(11),
845 113301 (2021), doi:[10.1088/1742-5468/ac312b](https://doi.org/10.1088/1742-5468/ac312b), <https://arxiv.org/abs/2109.08502>.
- 846 [50] M. Gherardi, *Solvable model for the linear separability of structured data*, Entropy **23**(3)
847 (2021), doi:[10.3390/e23030305](https://doi.org/10.3390/e23030305).
- 848 [51] M. Mezard, G. Parisi and M. Virasoro, *Spin Glass Theory and Beyond*, World Scientific,
849 doi:[10.1142/0271](https://doi.org/10.1142/0271) (1986).
- 850 [52] O. Dhifallah and Y. M. Lu, *A precise performance analysis of learning with random features*
851 (2020), <https://arxiv.org/abs/2008.11904>.
- 852 [53] H. Hu and Y. M. Lu, *Universality laws for high-dimensional learning with ran-*
853 *dom features*, IEEE Transactions on Information Theory **69**(3), 1932 (2023),
854 doi:[10.1109/TIT.2022.3217698](https://doi.org/10.1109/TIT.2022.3217698), <https://arxiv.org/abs/2009.07669>.
- 855 [54] Z. Wang and Y. Zhu, *Overparameterized random feature regression with nearly orthogonal*
856 *data*, In F. Ruiz, J. Dy and J.-W. van de Meent, eds., *Proceedings of The 26th Interna-*
857 *tional Conference on Artificial Intelligence and Statistics*, vol. 206 of *Proceedings of Machine*
858 *Learning Research*, pp. 8463–8493. PMLR (2023), [https://proceedings.mlr.press/v206/](https://proceedings.mlr.press/v206/wang23m/wang23m.pdf)
859 [wang23m/wang23m.pdf](https://proceedings.mlr.press/v206/wang23m/wang23m.pdf).
- 860 [55] Y. M. Lu and H.-T. Yau, *An equivalence principle for the spectrum of random inner-product*
861 *kernel matrices with polynomial scalings* (2023), <https://arxiv.org/abs/2205.06308>.
- 862 [56] S. d'Ascoli, L. Sagun and G. Biroli, *Triple descent and the two kinds of overfitting: where*
863 *& why do they appear?*, In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan
864 and H. Lin, eds., *Advances in Neural Information Processing Systems*, vol. 33, pp. 3058–
865 3069. Curran Associates, Inc. (2020), [https://proceedings.neurips.cc/paper/2020/file/](https://proceedings.neurips.cc/paper/2020/file/1fd09c5f59a8ff35d499c0ee25a1d47e-Paper.pdf)
866 [1fd09c5f59a8ff35d499c0ee25a1d47e-Paper.pdf](https://proceedings.neurips.cc/paper/2020/file/1fd09c5f59a8ff35d499c0ee25a1d47e-Paper.pdf).
- 867 [57] H. Hu, Y. M. Lu and T. Misiakiewicz, *Asymptotics of random feature regression beyond the*
868 *linear scaling regime* (2024), <https://arxiv.org/abs/2403.08160>.
- 869 [58] E. Gardner and B. Derrida, *Three unfinished works on the optimal storage capacity*
870 *of networks*, Journal of Physics A: Mathematical and General **22**(12), 1983 (1989),
871 doi:[10.1088/0305-4470/22/12/004](https://doi.org/10.1088/0305-4470/22/12/004).

- 872 [59] A. Engel and C. Van den Broeck, *Statistical Mechanics of Learning*, Cambridge University
873 Press, doi:[10.1017/CBO9781139164542](https://doi.org/10.1017/CBO9781139164542) (2001).
- 874 [60] G. Folena, S. Franz and F. Ricci-Tersenghi, *Rethinking mean-field glassy dynamics and its
875 relation with the energy landscape: The surprising case of the spherical mixed p-spin model*,
876 Phys. Rev. X **10**, 031045 (2020), doi:[10.1103/PhysRevX.10.031045](https://doi.org/10.1103/PhysRevX.10.031045).
- 877 [61] B. Widrow and M. E. Hoff, *Adaptive switching circuits*, In *1960 IRE WESCON Convention
878 Record, Part 4*, pp. 96–104 (1960).
- 879 [62] P. Breuer and P. Major, *Central limit theorems for non-linear functionals of Gaussian fields*,
880 Journal of Multivariate Analysis **13**(3), 425 (1983), doi:[10.1016/0047-259X\(83\)90019-
881 2](https://doi.org/10.1016/0047-259X(83)90019-2).
- 882 [63] J.-M. Bardet and D. Surgailis, *Moment bounds and central limit theorems for
883 Gaussian subordinated arrays*, Journal of Multivariate Analysis **114**, 457 (2013),
884 doi:<https://doi.org/10.1016/j.jmva.2012.08.002>.
- 885 [64] A. Mozeika, M. Sheikh, F. Aguirre-Lopez, F. Antenucci and A. C. C. Coolen, *Exact results
886 on high-dimensional linear regression via statistical physics*, Phys. Rev. E **103**, 042142
887 (2021), doi:[10.1103/PhysRevE.103.042142](https://doi.org/10.1103/PhysRevE.103.042142).
- 888 [65] A. C. C. Coolen, M. Sheikh, A. Mozeika, F. Aguirre-Lopez and F. Antenucci, *Replica analysis
889 of overfitting in generalized linear regression models*, Journal of Physics A: Mathematical
890 and Theoretical **53**(36), 365001 (2020), doi:[10.1088/1751-8121/aba028](https://doi.org/10.1088/1751-8121/aba028).
- 891 [66] W. F. Kibble, *An extension of a theorem of Mehler's on Hermite polynomials*, Math-
892 ematical Proceedings of the Cambridge Philosophical Society **41**(1), 12–15 (1945),
893 doi:[10.1017/S0305004100022313](https://doi.org/10.1017/S0305004100022313).
- 894 [67] T. Liang and H. Tran-Bach, *Mehler's formula, branching process, and compositional kernels
895 of deep neural networks*, Journal of the American Statistical Association **117**(539), 1324
896 (2022), doi:[10.1080/01621459.2020.1853547](https://doi.org/10.1080/01621459.2020.1853547), <https://arxiv.org/abs/2004.04767>.
- 897 [68] J. Bryson, R. Vershynin and H. Zhao, *Marchenko–Pastur law with relaxed indepen-
898 dence conditions*, Random Matrices: Theory and Applications **10**(04), 2150040 (2021),
899 doi:[10.1142/S2010326321500404](https://doi.org/10.1142/S2010326321500404), <https://arxiv.org/abs/1912.12724>.
- 900 [69] N. E. Karoui, *The spectrum of kernel random matrices*, The Annals of Statistics **38**(1), 1
901 (2010), doi:[10.1214/08-AOS648](https://doi.org/10.1214/08-AOS648).
- 902 [70] J. Glimm and A. Jaffe, *Quantum Physics: A Functional Integral Point of View*, Springer-
903 Verlag, doi:[10.1007/978-1-4612-5158-3](https://doi.org/10.1007/978-1-4612-5158-3) (1987).
- 904 [71] E. Gardner, *The space of interactions in neural network models*, Journal of Physics A:
905 Mathematical and General **21**(1), 257 (1988), doi:[10.1088/0305-4470/21/1/030](https://doi.org/10.1088/0305-4470/21/1/030).
- 906 [72] D. Bosch, A. Panahi and B. Hassibi, *Precise asymptotic analysis of deep random feature
907 models*, In G. Neu and L. Rosasco, eds., *Proceedings of Thirty Sixth Conference on Learn-
908 ing Theory*, vol. 195 of *Proceedings of Machine Learning Research*, pp. 4132–4179. PMLR
909 (2023), <https://proceedings.mlr.press/v195/bosch23a/bosch23a.pdf>.
- 910 [73] M. Potters and J.-P. Bouchaud, *A First Course in Random Matrix Theory: for Physicists,
911 Engineers and Data Scientists*, Cambridge University Press, doi:[10.1017/9781108768900
912 \(2020\)](https://doi.org/10.1017/9781108768900).

- 913 [74] G. Livan, M. Novaes and P. Vivo, *Introduction to random matrices theory and practice*,
914 Monograph Award **63**, 54 (2018).
- 915 [75] J. Bradbury, R. Frostig, P. Hawkins, M. J. Johnson, C. Leary, D. Maclaurin, G. Necula,
916 A. Paszke, J. VanderPlas, S. Wanderman-Milne and Q. Zhang, *JAX: composable transfor-*
917 *mations of Python+NumPy programs* (2018).
- 918 [76] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel,
919 P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos *et al.*, *Scikit-learn: Machine*
920 *learning in Python*, Journal of Machine Learning Research **12**(85), 2825 (2011), [http:](http://jmlr.org/papers/v12/pedregosa11a.html)
921 [://jmlr.org/papers/v12/pedregosa11a.html](http://jmlr.org/papers/v12/pedregosa11a.html).
- 922 [77] *GitHub repository to reproduce the figures in this paper*, [https://github.com/](https://github.com/MauroPastore/RandomFeatures/)
923 [MauroPastore/RandomFeatures/](https://github.com/MauroPastore/RandomFeatures/).