

Random features and polynomial rules

Fabián Aguirre-López^{1,2,3*}, Silvio Franz¹ and Mauro Pastore^{1,4,5†}

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

⁵ The Abdus Salam International Centre for Theoretical Physics, Strada Costiera 11, 34151 Trieste, Italy

* fabian.aguirre-lopez@polytechnique.edu, † mpastore@ictp.it

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	12
9	5.1 Replica symmetric theory	14

10	5.2 Saddle-point equations for quadratic loss	16
11	6 Strongly separated regimes	17
12	7 Effective theory for finite-size random features networks	18
13	8 Conclusions and perspectives	20
14	A Kernel on the Hermite basis	21
15	B Hermite polynomials and Wick products	22
16	C Evaluation of the moments of γ, λ^a	22
17	D Results on random matrix theory	23
18	D.1 Marchenko-Pastur distribution and Stieltjes transformation	23
19	D.2 Spectral density of $C^{\otimes \ell}$	24
20	E Determinant of sum of matrices with orthogonal row spaces	26
21	F Traces over RS matrices	27
22	G Replica-symmetric free energy	28
23	G.1 Measure contribution	28
24	G.2 Pattern contribution	28
25	H Asymptotic limits of the saddle-point equations	29
26	H.1 Case (i)	29
27	H.2 Case (ii)	30
28	H.3 Case (iii)	31
29	I Numerical experiments	31
30	References	31

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 has
 39 been recently generalized beyond the infinite-width limit [5–10]. In a dynamical perspective
 40 moreover, it has been shown that wide DNNs trained with gradient-based methods exhibit the
 41 lazy-training kernel regime [11], evaluated by first order Taylor-expanding the network with
 42 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*

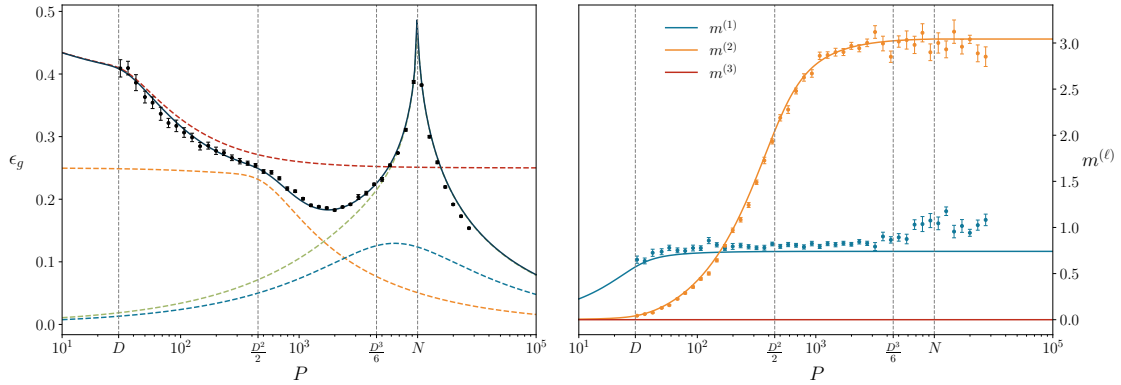


Figure 1: **Left:** generalization error of the RFM on a classification task, as a function of the size of the training set P , for $D = 30$, $N = 10^4$, weights regularization $\zeta = 10^{-8}$, *quadratic teacher* (balanced: $\tau_1 = \tau_2 = 1/\sqrt{2}$, $\tau_{\ell>2} = 0$) and ELU activation functions (defined in Eq. (8) below); 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. (37) and (45) – 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.

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 $\tilde{\mathcal{H}}(\mathbf{x}, \mathbf{x}') =$
 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 for empirical risk minimization, with $P \sim D^K$, $N \sim D^L$. We find, under a quite
 64 general teacher/student setting with a random polynomial teacher and Gaussian i.i.d. input
 65 data, that

- 66 • as long as $P \ll N$, the model behaves as an infinite-rank ($N \rightarrow \infty$) kernel machine:
 67 for $P \sim D^K$, features of order K can be learnt, such that the generalization error as a
 68 function of P has a staircase descent (or overfitting peaks if the teacher is less complex)

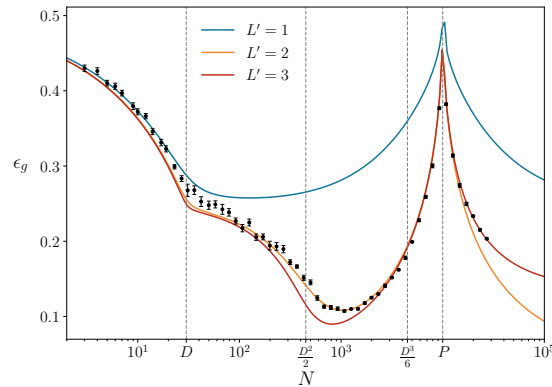


Figure 2: Generalization error of a RFM on a classification task, as a function of the number of hidden units N , for $P = 10^4$ and the rest of the parameters as in Fig. 1; continuous lines are the theories truncated at $L' = 1, 2, 3$ (respectively: blue, yellow, red); numerical points (in black) are nicely interpolating between these curves in the regimes where $N \sim O(D), O(D^2), O(D^3)$, validating Eq. (25), where the truncation L' of the equivalent polynomial theory is fixed at $L \sim \log(N)/\log(D)$.

69 with steps corresponding to different values of K ;

- 70 • for $P \gg N$ and $N \sim D^L$, the model is equivalent to a degree- L polynomial student: if the
- 71 complexity of the teacher is lower than the degree L , the generalization error is equal
- 72 to zero, or otherwise, to the minimum error for a degree- L polynomial fitting a more
- 73 complex teacher;
- 74 • for $P \sim N$, an interpolation peak of the generalization error, which depends on the
- 75 strength of the regularization of the student’s weights, occurs.

76 This behavior is depicted in Fig. 1. Comparison with numerical experiments shows that our

77 theory, based on the mapping of the RFM to an *equivalent noisy polynomial model*, predicts

78 well the quantitative behavior of the true generalization performance at finite size, over many

79 orders of magnitude.

80 Our theory, formulated from the point of view of the statistical mechanics of disordered

81 systems, expresses the generalization performance of the RFM in terms of few order param-

82 eters with a clear physical interpretation, as overlaps between combinations of the student’s

83 weights and the parameters defining the teacher. In this way, we are offering a complementary

84 take on what is known about RFMs in the computer science community, as we discuss in the

85 following.

86 1.1 Related works

87 In this section we give an overview on the previous works that have been of inspiration to our

88 paper, presenting relevant results and differences with our approach.

89 Random feature models were introduced in [23–26], initially as randomized low-rank ap-

90 proximations of kernels arising in classification or regression problems. Recently, their interest

91 was renewed by the discovery that DNNs behaves as RFMs close to the infinite-width limit, both

92 in a Bayesian learning [1–4] and in a gradient-based learning [11–14] setting. This mapping,

93 which provides one of the few limits where DNNs can be studied with analytical methods, has

94 motivated in the last few years a huge effort to formalize their behavior in terms of expressive

95 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 degree- L polynomial regression; in the second one, it reduces to (infinite-rank) kernel regression, which for that number of samples can fit at most a degree- K polynomial in the inputs, in a way also investigated in literature [16–22].
- In [29], the limit where both N and P scale linearly with D with their ratio fixed is considered; the generalization error as a function of the ratio between the number of hidden units and the size of the training set first decreases for N/P small, then exhibits a peak at the interpolation threshold $N/P = 1$ and then relaxes again for $N \gg P$ to the value predicted from the kernel theory with $P \sim D$, coherently with the previous point. This phenomenology is widely observed in numerical experiments and known in literature as *double descent* [35] of the generalization error.
- In [31], the authors push forward the analysis of [27] (that is, P and N scaling polynomially with D) to the regimes where $N \leq P^{1-\delta}$ and $N \geq P^{1+\delta}$. The authors show indeed that the limiting behavior is given by the smallest of N and P , and they find the interpolation threshold at $N \sim P$ also in this polynomial scaling.
- In [33], universality results on training and test error are proven in the $P \sim N$ regime for a larger class of models, as long as with finite dimensional outputs, and generic losses. Indeed, they prove that training and test errors depend on the random features distribution only through its covariance structure.

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

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

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

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
$C^{\otimes \ell}$	$((C_{ij})_{i,j=1}^\ell)^N \simeq \sum_\alpha \mathbf{F}_\alpha^{\otimes \ell} (\mathbf{F}_\alpha^{\otimes \ell})^\top / \binom{D}{\ell}$
$Q, Q^{(\ell)}, \dots$	$(Q_{ab})_{a,b=1}^n, (Q_{ab}^{(\ell)})_{a,b=1}^n, \dots$

Table 1: Notations used in this paper

144 Appeared in parallel with our work, the paper [57] pushes forward the line of research
 145 of [29] from a mathematical perspective, deriving sharp asymptotics for the generalization
 146 of random features ridge regression in the polynomial regime. The even more recent [58]
 147 bounds the test error of random features ridge regression with a dimension-free (that is, for
 148 arbitrary input dimension D) non-asymptotic (depending explicitly on N and P , converging to
 149 the test error when at least one of them is large) deterministic equivalent, depending only on
 150 the feature map eigenvalues through a set of self-consistent equations. The mapping of our
 151 approach to [57, 58] is left for future work.

152 2 The model

153 We would like to study the generalization performance of the Random Features model in a
 154 teacher/student [59, 60] supervised learning set-up, where the teacher performs an input-
 155 output mapping with various degree of complexity. We summarize in Table 1 the main nota-
 156 tions used in this paper.

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

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

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

159 where $\theta_\alpha^{(1)}, \theta_{\alpha\beta}^{(2)}, \dots$ are i.i.d. $\mathcal{N}(0, 1)$ parameters collectively denoted as θ , describing the
 160 non-linear decision boundary (diagonal terms, irrelevant for large D , are for simplicity not
 161 included in the sum). Notice that the function $\nu(\mathbf{x})$ coincide with the Hamiltonian of the
 162 “mixed p -spin model” of the statistical physics of the spin-glasses (see, for example, [61]).

163 The mixture parameters τ_ℓ , weighting the monomials of different degree, are chosen to respect
 164 $\sum_{\ell=1}^B \tau_\ell^2 = 1$. Within this general setting, we will concentrate on the specific simple examples
 165 of a deterministic teacher for binary classification or a noisy teacher for polynomial regression
 166 with variance of the noise Δ , for which Eq. (1) reduces respectively to

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

167 It has been shown in [16] that a *polynomial* student, defined in the same way as in Eq. (1),
 168 would learn the weights of the teacher in a hierarchical fashion: $O(D^K)$ examples are needed
 169 in order to learn the parameters $\theta^{(\ell)}$ for $\ell \leq K$. However, here the student's task is to learn
 170 the weights of the last layer of a 2-layers NN, $f(\mathbf{x}; \mathbf{w})$, whose first layer realizes a random
 171 embedding of the data in a N -dimensional feature space:

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

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

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

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

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

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

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

178 While our theory is general in the choice of σ (as long as it can be expanded on the basis
 179 of Hermite polynomials – see Sec. 4), we will test our results for popular choices, such as

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

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

180 (respectively, Rectified and Exponential Linear Unit).

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

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

183 where \mathcal{L} is an opportune convex loss function and ζ controls the regularization of the weights.
 184 Notice how the solution of this optimization problem is an implicit function of the training set
 185 \mathcal{T} , the parameters of the teacher θ and the random features F , that is $\mathbf{w}^* = \mathbf{w}^*(\mathcal{T}, \theta, F)$; we
 186 will omit this dependence to lighten notations.

187 The choice of the loss function \mathcal{L} and the readout activation function ϕ in Eq. (3) defines
 188 the specific learning task to perform. The approach we present in the following can be followed

189 for any choice of \mathcal{L} , as long as the optimization problem (9) is convex (to justify the Replica
190 Symmetric ansatz, see below); this is true in particular if $\mathcal{L}(y, \lambda)$ is convex as a function of
191 λ , as the student's weights \mathbf{w} enter linearly in the definition of $\lambda(\mathbf{x}; \mathbf{w})$. However, to simplify
192 formulas, we will report in the main text only the case of a pure quadratic loss, reading, both
193 in the case of regression and classification:

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

194 The use of a regression loss for a classification task (λ instead of $\phi(\lambda)$ even when $\phi = \text{sgn}$)
195 is not unusual in practical cases (e.g. the `linear_model.RidgeClassifier` class in the
196 Scikit-learn library for Python [62]) and dates back to the early days of NNs [60, 63].

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

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

203 3 Generalization error

204 In order to quantify how well the student can learn the teacher, we look at the generalization
205 error, defined as the probability of misclassifying a new sample (in the case of classification)
206 or as the mean squared error of a new point (in the case of regression). Given a test point
207 $(\mathbf{x}, y) \sim p_0(\mathbf{x})p(y|\nu(\mathbf{x}))$, both cases can be expressed with the following formula,

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

208 where $\kappa = 1$ for binary classification and $\kappa = 0$ for regression. Notice the presence of the
209 function ϕ in the definition of the generalization error, at variance with the loss function (10).

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

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

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

215 We have written the average generalization error as in Eq. (13) to show that we only
216 need to know the joint distribution of (ν, λ) to evaluate it. Since \mathbf{x} is a test point, and is
217 thus uncorrelated with \mathbf{w}^* , we will take the distribution $p(\nu, \lambda)$ as Gaussian: to compute the
218 generalization error we only need the first and second moments,

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

219 Notice that by definition of the model (*i.e.* the normalization of the mixing parameters τ_ℓ)
 220 ρ is identically equal to 1. In section 5 we will show how to obtain these quantities from a
 221 replica approach. Stating formally hypotheses on \mathbf{w}^* , F and the functions $\nu(\mathbf{x})$, $\lambda(\mathbf{x}; \mathbf{w})$ in
 222 order to justify this ansatz is beyond the scope of this paper: we will check *a posteriori* its
 223 validity with numerical experiments. Central limit theorems for sums of non linear functions
 224 of Gaussian fields (the pre-activations (5) at given feature matrix F), of the kind we just used
 225 to motivate this ansatz, have been proven in the past under rather technical conditions on the
 226 realization of the feature-feature covariance matrix C and of the vector \mathbf{w}^* [33, 38, 53, 64, 65].
 227 The interested reader can find a sketch of proof in [36], Appendix A.2, where the moments of
 228 the variables λ are evaluated and the leading order diagrams identified as the Gaussian ones.

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

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

230 where we use the Gardner notation [66] $D\nu = \frac{e^{-\nu^2/2}}{\sqrt{2\pi}} d\nu$ and $H(x) = \int_x^{\infty} Dt$. Notice that when
 231 $t^* = 0$ (that is, when the student is zero-mean) the formula simplifies to

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

232 For the case of noisy polynomial regression, ($\phi = \text{id}$ and $\Delta = \mathbb{E}[(y - \nu)^2]$) [67, 68],

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

233 These formulas remind the generalization error of a generalized linear model with the same
 234 architecture as the teacher [60]: in that case, $m^*/\sqrt{q^*}$ corresponds to the angle between
 235 the teacher and the student weight vectors. For the RFM, it is not clear *a priori* if we can
 236 interpret $m^*/\sqrt{q^*}$ as a scalar product of the teacher's weight vector and some effective weights
 237 of the student. If this can be done, the RFM could be mapped to an equivalent polynomial
 238 model. In Sec. 4 we will show how to explicitly construct it from \mathbf{w} and F , thus achieving
 239 this mapping. To do so, we need to spend a few words on the connection between RFMs and
 240 kernel machines, in order to explain the truncation of the activation function σ on the basis
 241 of Hermite polynomials, which we will use later on.

242 4 Kernel learning and polynomial models

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

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

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

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

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

249 While our analysis will be more general, in this section we consider the limit $P \rightarrow \infty$, for
250 the purpose of arguing.¹ In this case the empirical kernel reduces to

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

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

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

255 where He_ℓ is the ℓ -th Hermite polynomial and the coefficient μ_ℓ are:

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

256 Along these lines, the kernel (20) can be expressed for large D [69, 70] (see App. A for
257 details) as

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

258 where C_{ij} , given by (6), is a rank- D Wishart matrix with elements $C_{ii} = 1 + O(D^{-1/2})$ and
259 $C_{ij} = O(D^{-1/2})$ for $i \neq j$. The matrix with entries $(C_{ij})^\ell$, which we denote by $C^{\circ\ell}$, defines an
260 interesting random matrix ensemble, obtained taking Hadamard (element by element) powers
261 of the covariance C . A similar ensemble was recently studied in [71].

262 Suppose now the relation between N and D is fixed: $N \sim D^{L+\delta}$ with $0 \leq \delta < 1$. The
263 $N \times N$ matrix $C^{\circ\ell}$ has generically rank equal to $\min\{D^\ell/\ell!, N\}$ (neglecting possible smaller
264 contributions to the rank coming from outliers, see Sec. 6 where we discuss more in detail
265 the properties of these matrices) and off-diagonal elements $O(D^{-\ell/2})$. For $\ell > L$ the matrix
266 is full ranked, the small off-diagonal terms give a vanishing contribution to eigenvalues and
267 eigenvectors. In other words, when D^ℓ is scaling faster than N to infinity, we can take the large
268 D limit *before* the large N one in the combination

$$(C_{ij})^\ell = \delta_{ij}[1 + \ell O(D^{-1/2})] + (1 - \delta_{ij})O(D^{-\ell/2}) \underset{D \text{ large, } D^\ell \gg N}{\simeq} \delta_{ij}, \quad (24)$$

269 in the same way as the Wishart matrix $C_{ij} = \delta_{ij}[1 + O(D^{-1/2})] + (1 - \delta_{ij})O(D^{-1/2})$ concentrates
270 around δ_{ij} for $D \gg N$ (the Marchenko-Pastur distribution, providing the asymptotic distribu-
271 tion of the spectrum of C , concentrates around 1 for $N/D \rightarrow 0$). We can thus truncate the
272 expansion substituting $C^{\circ\ell > L}$ by the identity matrix:

$$\mathcal{K}_{ij} \simeq \sum_{\ell=0}^L \frac{\mu_\ell^2}{\ell!} (C_{ij})^\ell + \mu_{\perp,L}^2 \delta_{ij}, \quad (25)$$

273 where

$$\mu_{\perp,L}^2 = \sum_{\ell=L+1}^{\infty} \frac{\mu_\ell^2}{\ell!} = \mathbb{E}_x[\sigma(x)^2] - \sum_{\ell=0}^L \frac{\mu_\ell^2}{\ell!} \quad (26)$$

¹We do so in this section to introduce the kernel \mathcal{K} as a limit of the empirical kernel $\bar{\mathcal{K}}$; in the replica approach in Sec. 5 the expectation over the data will be taken explicitly and the kernel will appear naturally without taking the $P \rightarrow \infty$ limit from the start.

274 for $x \sim \mathcal{N}(0, 1)$.

275 This truncation is proven for $L = 1$ (that is, in the proportional regime $N \sim D$) in [72], and
 276 extended to the case $L > 1$ under generic assumptions on the kernel \mathcal{K} in [31,55]. A convincing
 277 check of this property for moderately large values of N is given by Fig. 2, which shows the
 278 theoretical curves of the generalization error obtained through a truncated effective theory
 279 (that we describe below) at different values of L' , compared with the numerical experiments,
 280 as a function of N ; quantitative agreement is obtained for $L' = L \sim \log N / \log D$, with the
 281 numerical points interpolating nicely the theoretical curves in the various regimes.

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

$$\lambda_{\text{eff}}(\mathbf{x}^\mu; \mathbf{w}) = \mu_0 m^{(0)} + \sum_{\ell=1}^L \frac{\mu_\ell}{\sqrt{D}^\ell} \sum_{\alpha_1, \dots, \alpha_\ell} s_{\alpha_1 \dots \alpha_\ell}^{(\ell)} : x_{\alpha_1}^\mu \cdots x_{\alpha_\ell}^\mu : + z^\mu, \quad (27)$$

284 where

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

$$s_{\alpha_1 \dots \alpha_\ell}^{(\ell)} = \frac{1}{\sqrt{N}} \sum_i w_i F_{i\alpha_1} \cdots F_{i\alpha_\ell}. \quad (28)$$

- 288 • we have written the expansion of the Hermite polynomials in terms of the so-called Wick
 289 products of the x 's, routinely used in theoretical physics and defined from the following
 290 generating function (see for example [73]):

$$\begin{aligned} :x_1 \cdots x_k: &= \partial_{\lambda_1} \cdots \partial_{\lambda_k} G(\boldsymbol{\lambda}; \mathbf{x}) \Big|_{\boldsymbol{\lambda}=0}, \\ G(\boldsymbol{\lambda}; \mathbf{x}) &= \frac{\exp(\boldsymbol{\lambda}^\top \mathbf{x})}{\mathbb{E}[\exp(\boldsymbol{\lambda}^\top \mathbf{x})]} = \exp(\boldsymbol{\lambda}^\top \mathbf{x} - \|\boldsymbol{\lambda}\|^2/2) \end{aligned} \quad (29)$$

291 These quantities have the property $\mathbb{E}[:x_1 \cdots x_k:] = 0$. The mapping

$$\text{He}_\ell(h_i) \simeq \sum_{\alpha_1, \dots, \alpha_\ell} \frac{F_{i\alpha_1} \cdots F_{i\alpha_\ell}}{\sqrt{D}^\ell} :x_{\alpha_1} \cdots x_{\alpha_\ell}:, \quad (30)$$

292 which is true for D large and which we used to write λ_{eff} in terms of Wick products
 293 starting from (21), is proven in App. B.

- 294 • the last term z^μ is a Gaussian noise term with zero mean and variance $\mathbb{E}(z^{\mu 2}(\mathbf{w})) =$
 295 $\mu_{\perp, L}^2 \sum_{i=1}^N w_i^2 / N$ which can be represented as

$$z^\mu = \frac{\mu_{\perp, L}}{\sqrt{N}} \sum_{i=1}^N w_i v_i^\mu, \quad (31)$$

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

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

299 In (27) we give an effective description of the RFM, mapping it to a polynomial model
 300 with correlated weights in presence of a noise term coming from the $\ell > L$ terms in the expan-
 301 sion (21). The mapping is motivated by the fact that $\lambda_{\text{eff}}(\mathbf{x}^\mu; \mathbf{w})$ defined in this way, admits

302 as second moment $\mathbf{w}^\top \mathcal{K} \mathbf{w} / N$ at given F and \mathbf{w} , with the kernel truncated according to (25);
 303 we show this explicitly for the replicated version of λ in Appendix C, together with the covari-
 304 ance structure with the polynomial $\nu(\mathbf{x})$ defining the teacher. This is an extension to generic
 305 scaling regimes $N \sim D^L$ of the *Gaussian equivalence principle* from [38] and related works, to
 306 which it reduces when $L = 1$. In the following, we will base our analysis on this representation
 307 of λ . This description makes more transparent the meaning of the observables introduced in
 308 Sec. 3 and the mechanism by which the RFM learns the teacher's features, as we explain in
 309 the following.

310 5 Replica calculation

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

$$p(\nu, \lambda) = \lim_{\beta \rightarrow \infty} \mathbb{E} \int d\mathbf{w} \frac{1}{\mathcal{Z}} e^{-\beta \sum_{\mu} \mathcal{L}[y^{\mu}, \lambda(\mathbf{x}^{\mu}; \mathbf{w})] - \frac{\beta \zeta}{2} \|\mathbf{w}\|^2} \int d\mathbf{x} p_0(\mathbf{x}) \delta(\nu - \nu(\mathbf{x})) \delta(\lambda - \lambda(\mathbf{x}; \mathbf{w})). \quad (32)$$

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

$$p(\nu, \lambda) = \lim_{n \rightarrow 0} \lim_{\beta \rightarrow \infty} \mathbb{E} \int \prod_{a=1}^n d\mathbf{w}^a e^{-\beta \sum_{\mu, a} \mathcal{L}[y^{\mu}, \lambda(\mathbf{x}^{\mu}; \mathbf{w}^a)] - \frac{\beta \zeta}{2} \sum_a \|\mathbf{w}^a\|^2} \times \int d\mathbf{x} p_0(\mathbf{x}) \delta(\nu - \nu(\mathbf{x})) \delta(\lambda - \lambda(\mathbf{x}; \mathbf{w}^1)), \quad (33)$$

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

$$Z_n = \mathbb{E}[\mathcal{Z}^n] = \int \prod_{a=1}^n d\mathbf{w}^a e^{-\frac{\beta \zeta}{2} \sum_a \|\mathbf{w}^a\|^2} \mathbb{E}_{F, \theta} \left[\mathbb{E}_{\nu, \{\lambda^a\}} \int dy p(y|\nu) e^{-\beta \sum_a \mathcal{L}(y, \lambda^a)} \right]^P. \quad (34)$$

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

$$t_a = \mathbb{E}(\lambda_a | F, \theta), \quad M_a = \mathbb{E}(\nu \lambda_a | F, \theta), \quad Q_{ab} = \mathbb{E}(\lambda_a \lambda_b | F, \theta) - t_a t_b. \quad (35)$$

318 from which we can extract the generalization error according to (15), (17). Using the repre-
 319 sentation (27) we can decompose these order parameters as (see Appendix C for details)

$$t_a = \mu_0 M_a^{(0)}, \quad M_a = \sum_{\ell=1}^{\min\{L, B\}} \frac{\mu_{\ell} \tau_{\ell}}{\sqrt{\ell!}} M_a^{(\ell)}, \quad Q_{ab} = \mu_{\perp, L}^2 Q_{ab}^{(0)} + \sum_{\ell=1}^L \frac{\mu_{\ell}^2}{\ell!} Q_{ab}^{(\ell)}, \quad (36)$$

320 with the definitions:

$$M_a^{(0)} = \frac{1}{\sqrt{N}} \sum_{i=1}^N w_i^a, \quad M_a^{(\ell)} = \frac{\boldsymbol{\theta}^{(\ell)} \cdot \mathbf{s}_a^{(\ell)}}{\binom{D}{\ell}}, \quad Q_{ab}^{(0)} = \frac{1}{N} \sum_{i=1}^N w_i^a w_i^b, \quad Q_{ab}^{(\ell)} = \frac{1}{N} \sum_{i, j=1}^N w_i^a C_{ij}^{\ell} w_j^b, \quad (37)$$

321 where we are using the notation

$$\boldsymbol{\theta}^{(\ell)} \cdot \mathbf{s}_a^{(\ell)} = \sum_{\alpha} \theta_{\alpha}^{(\ell)} s_{a, \alpha}^{(\ell)} \quad (38)$$

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

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

$$Z_n = e^{PS_p[Q,M]} e^{\frac{N}{2} \sum_{a,b} \hat{Q}_{ab}^{(0)} Q_{ab}^{(0)} + \frac{1}{2} \sum_{\ell,a,b} \binom{D}{\ell} \hat{Q}_{ab}^{(\ell)} Q_{ab}^{(\ell)} + \sum_{\ell,a} \binom{D}{\ell} \hat{M}_a^{(\ell)} M_a^{(\ell)}} \\ \times \mathbb{E}_{F,\theta} \int d\mathbf{w} e^{-\frac{1}{2} \mathbf{w}^\top [(\beta \zeta \mathbb{1}_n + \hat{Q}^{(0)}) \otimes \mathbb{1}_N + \sum_{\ell} \hat{Q}^{(\ell)} \otimes \frac{C^{\otimes \ell}}{\eta_\ell}] \mathbf{w} - \sum_{\ell,i,a,\alpha} \hat{M}_a^{(\ell)} w_i^\alpha F_{i,\alpha}^{\otimes \ell} \theta_a^{(\ell)} / \sqrt{\eta_\ell \binom{D}{\ell}}}, \quad (39)$$

326 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

$$S_p[Q, M] = \log \mathbb{E}_{\nu, \{\lambda^a\}} \int dy p(y|\nu) e^{-\beta \sum_a \mathcal{L}(y, \lambda^a)}. \quad (40)$$

327 In writing Eq. (39), we took $\hat{M}_a^{(0)} \rightarrow 0$, as the Fourier conjugate of the mean t_a is suppressed
328 in the large- N limit [66] (a property that could be checked *a posteriori* from the saddle point
329 equation for $\hat{M}_a^{(0)}$);² moreover, the conventional scalings with N and $\binom{D}{\ell}$ in this equation are
330 chosen in such a way that the hat variables corresponding to the asymptotic regimes explained
331 in Sec. 6 have a non-trivial high-dimensional limit.

332 Averaging over θ we obtain:³

$$Z_n = e^{PS_p[Q,M]} e^{\frac{N}{2} \sum_{a,b} \hat{Q}_{ab}^{(0)} Q_{ab}^{(0)} + \frac{1}{2} \sum_{\ell,a,b} \binom{D}{\ell} \hat{Q}_{ab}^{(\ell)} Q_{ab}^{(\ell)} + \sum_{\ell,a} \binom{D}{\ell} \hat{M}_a^{(\ell)} M_a^{(\ell)}} \\ \times \mathbb{E}_F \int d\mathbf{w} e^{-\frac{1}{2} \mathbf{w}^\top [(\beta \zeta \mathbb{1}_n + \hat{Q}^{(0)}) \otimes \mathbb{1}_N + \sum_{\ell} (\hat{Q}^{(\ell)} - \hat{M}^{(\ell)} \hat{M}^{(\ell)\top}) \otimes \frac{C^{\otimes \ell}}{\eta_\ell}] \mathbf{w}} \quad (41)$$

333 and integrating over \mathbf{w} ,

$$Z_n = e^{PS_p[Q,M]} e^{\frac{N}{2} \sum_{\ell,a,b} \hat{Q}_{ab}^{(\ell)} Q_{ab}^{(\ell)} + \frac{1}{2} \sum_{\ell,a,b} \binom{D}{\ell} \hat{Q}_{ab}^{(\ell)} Q_{ab}^{(\ell)} + \sum_{\ell,a} \binom{D}{\ell} \hat{M}_a^{(\ell)} M_a^{(\ell)} - \frac{1}{2} \text{Tr} \log [A^{(0)} \otimes \mathbb{1}_N + \sum_{\ell} B^{(\ell)} \otimes C^{\otimes \ell}]}, \quad (42)$$

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

$$A^{(0)} = \beta \zeta \mathbb{1}_n + \hat{Q}^{(0)}, \quad B^{(\ell)} = (\hat{Q}^{(\ell)} - \hat{M}^{(\ell)} \hat{M}^{(\ell)\top}) / \eta_\ell. \quad (43)$$

336 We notice at this point that, given $N \sim D^{L+\delta}$, for $\ell \leq L$ the matrices $C^{\otimes \ell}$ have rank $r_\ell =$
337 $O(D^\ell) \ll N$ and have eigenvalues of order $N / \binom{D}{\ell}$. Simple perturbation theory shows that
338 adding these matrices with coefficients of order 1 only slightly modify the eigenvalues. This is
339 due to the fact that the row spaces (that is, the complements to their null spaces) corresponding
340 to the different ℓ are almost orthogonal (we postpone a throughout discussion on this point
341 to Sec. 6, where we collect and motivate the assumptions we are using on the matrices $C^{\otimes \ell}$).
342 In such a situation we approximate the trace-log term appearing in (42) as

$$\text{Tr} \log \left[A^{(0)} \otimes \mathbb{1}_N + \sum_{\ell=1}^L B^{(\ell)} \otimes C^{\otimes \ell} \right] \simeq N(1-L) \text{Tr} \log(A^{(0)}) + \sum_{\ell=1}^L \text{Tr} \log(A^{(0)} \otimes \mathbb{1}_N + B^{(\ell)} \otimes C^{\otimes \ell}) \quad (44)$$

²The terms depending on $\hat{M}^{(0)}$ are given by

$$S_{\hat{M}^{(0)}} = \frac{M^{(0)\top} \hat{M}^{(0)}}{\sqrt{N}} + \frac{1}{2} \hat{M}^{(0)\top} \frac{1}{N} \sum_{ij} \left[(\beta \zeta \mathbb{1}_n + \hat{Q}^{(0)}) \otimes \mathbb{1}_N + \sum_{\ell} (\hat{Q}^{(\ell)} - \hat{M}^{(\ell)} \hat{M}^{(\ell)\top}) \otimes \frac{C^{\otimes \ell}}{\eta_\ell} \right]_{ij}^{-1} \hat{M}^{(0)},$$

so that the saddle point equation for $\hat{M}^{(0)}$ gives $\hat{M}^{(0)} = O(1/\sqrt{N})$.

³For the sake of simplicity, to write Eq. (41) 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. (49).

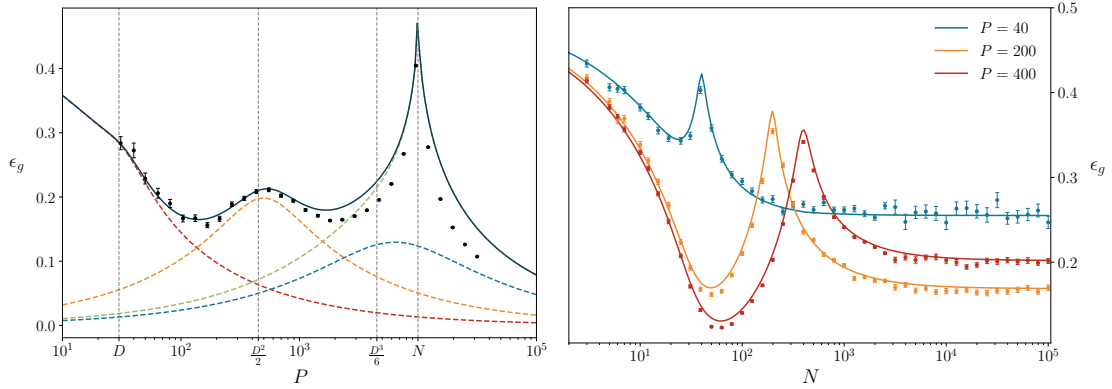


Figure 3: **Left:** generalization error of the RFM on a classification task, as a function of the size of the training set P , for $D = 30$, $N = 10^4$, weights regularization $\zeta = 10^{-8}$, linear teacher ($\tau_1 = 1$, $\tau_{\ell>1} = 0$) and ELU activation functions; the continuous line is the 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 .

343 (notice that Tr in $\text{Tr} \log(A^{(0)})$ is over replica indices only). We report a detailed derivation of
 344 Eq. (44) under the hypothesis of orthogonality of the $C^{\otimes \ell}$ row spaces in Appendix E. Notice
 345 that we could have gotten to the same result decomposing the vectors \mathbf{w} on the row spaces
 346 of the $C^{\otimes \ell}$ supposed orthogonal. This decomposition clearly shows the hierarchical nature of
 347 learning.

348 5.1 Replica symmetric theory

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

$$Q_{ab}^{(\ell)} = \frac{\chi^{(\ell)}}{\beta} \delta_{ab} + q^{(\ell)}, \quad M_a^{(\ell)} = m^{(\ell)}, \quad t_a = t. \quad (45)$$

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

354 order parameters in the RS:

$$\hat{Q}_{ab}^{(\ell)} = \beta \hat{\chi}^{(\ell)} \delta_{ab} - \beta^2 \hat{q}^{(\ell)}, \quad \hat{M}_a^{(\ell)} = -\beta \hat{m}^{(\ell)}. \quad (46)$$

355 Exploiting the explicit parametrization of the RS matrices, we can perform the traces over
356 replica indices in Eq. (44), to get (see Appendix F)

$$\begin{aligned} \text{Tr} \log(A \otimes \mathbb{1}_N + B \otimes C^{\otimes \ell}) &= nN \log(\beta \hat{\chi}^{(\ell)}) + n \text{Tr} \log(\gamma_\ell \mathbb{1} + C^{\otimes \ell}) \\ &- n\beta \eta_\ell \frac{\hat{q}^{(0)}}{\hat{\chi}^{(\ell)}} \text{Tr}(\gamma_\ell \mathbb{1} + C^{\otimes \ell})^{-1} - n\beta \frac{\hat{q}^{(\ell)} + (\hat{m}^{(\ell)})^2}{\hat{\chi}^{(\ell)}} \text{Tr}[C^{\otimes \ell}(\gamma_\ell \mathbb{1} + C^{\otimes \ell})^{-1}], \end{aligned} \quad (47)$$

357 where we introduced the parameter

$$\gamma_\ell = \eta_\ell \frac{(\zeta + \hat{\chi}^{(0)})}{\hat{\chi}^{(\ell)}} \quad (48)$$

358 and the remaining traces are over feature indices only.

359 We need now to evaluate the traces in feature indices. In order to proceed, we make at this
360 point a crucial approximation, and treat $C^{\otimes \ell}$ as a matrix with a Merchenko-Pastur spectrum
361 with parameter $\eta_\ell = N/\binom{D}{\ell}$. This amounts essentially in approximating $C^{\otimes \ell}$, by

$$C_{ij}^{\otimes \ell} = \frac{\ell!}{D^\ell} \sum_{\alpha_1 < \dots < \alpha_\ell} F_{\alpha_1}^i F_{\alpha_1}^j \dots F_{\alpha_\ell}^i F_{\alpha_\ell}^j \quad (49)$$

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

366 Using the properties of the resolvent of large random matrices (see Appendix D), we can
367 write that, for large N ,

$$\frac{1}{N} \text{Tr}(\gamma_\ell \mathbb{1} + C^{\otimes \ell})^{-1} \approx g_\ell(-\gamma_\ell), \quad (50)$$

368 where g_ℓ is the Stieltjes transformation of the Marchenko-Pastur distribution with ratio $\eta_\ell =$
369 $N/\binom{D}{\ell}$:

$$g_\ell(z) = \frac{1 - z - \eta_\ell - \sqrt{(1 - z - \eta_\ell)^2 - 4z\eta_\ell}}{2z\eta_\ell}. \quad (51)$$

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

$$Z_n \sim e^{PS_p + NS_M}, \quad (52)$$

371 where

$$\begin{aligned} \frac{1}{\beta n} S_M &= - \sum_{\ell=1}^{\min\{L, B\}} \frac{m^{(\ell)} \hat{m}^{(\ell)}}{\eta_\ell} + \frac{1}{2} \sum_{\ell=0}^L \frac{q^{(\ell)} \hat{\chi}^{(\ell)} - \chi^{(\ell)} \hat{q}^{(\ell)}}{\eta_\ell} + \frac{(1-L)}{2} \frac{\hat{q}^{(0)}}{\zeta + \hat{\chi}^{(0)}} \\ &+ \frac{1}{2} \sum_{\ell=1}^L \eta_\ell \frac{\hat{q}^{(0)}}{\hat{\chi}^{(\ell)}} g_\ell(-\gamma_\ell) + \frac{1}{2} \sum_{\ell=1}^L \frac{\hat{q}^{(\ell)} + (\hat{m}^{(\ell)})^2}{\hat{\chi}^{(\ell)}} [1 - \gamma_\ell g(-\gamma_\ell)] \end{aligned} \quad (53)$$

372 and, for the quadratic loss (10),

$$\frac{1}{\beta n} S_p = \frac{2m^* \langle y v \rangle - q^* - \langle (t^* - y)^2 \rangle}{2(1 + \chi^*)}, \quad (54)$$

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

$$\begin{aligned}
 m^* &= \sum_{\ell=1}^{\min\{L,B\}} \frac{\tau_\ell \mu_\ell}{\sqrt{\ell!}} m^{(\ell)}, & t^* &= \mu_0 m^{(0)}, \\
 \chi^* &= \mu_\perp^2 \chi^{(0)} + \sum_{\ell=1}^L \frac{\mu_\ell^2}{\ell!} \chi^{(\ell)}, & q^* &= \mu_\perp^2 q^{(0)} + \sum_{\ell=1}^L \frac{\mu_\ell^2}{\ell!} q^{(\ell)}.
 \end{aligned}
 \tag{55}$$

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

376 Eq. (55) gives the RS version of Eq. (36): these quantities are precisely the ones appearing
 377 in Eq. (14), giving the low-order statistics of the distribution used to evaluate the generaliza-
 378 tion error. Once their value is known from the saddle point equations implicit in the derivation
 379 of the partition function, they can be used to obtain the generalization curves reported in this
 380 paper.

381 5.2 Saddle-point equations for quadratic loss

382 The free energy in Eq. (52) has to be evaluated at the saddle point with respect to all the RS
 383 order parameters and their Fourier conjugates. We report here the resulting equations, in the
 384 special case of quadratic loss function (10). Remark however that only the equations where
 385 P appears explicitly depend on the form of the loss, and have to be modified for other choices
 386 (see Appendix G.2). The equations can be solved in steps. First, a set of $2L + 2$ nonlinear
 387 equations is used to determine the variables $\chi^{(0)}, \dots, \chi^{(L)}$ and $\hat{\chi}^{(0)}, \dots, \hat{\chi}^{(L)}$:

$$\begin{aligned}
 \hat{\chi}^{(0)} &= \frac{P}{N} \frac{\mu_\perp^2}{1 + \chi^*}, & \chi^{(0)} &= \frac{1 - \sum_{\ell=1}^L [1 - \gamma_\ell g_\ell(-\gamma_\ell)]}{\hat{\chi}^{(0)} + \zeta}, \\
 \hat{\chi}^{(\ell)} &= \frac{P}{\binom{D}{\ell}} \frac{\mu_\ell^2}{\ell!} \frac{1}{1 + \chi^*}, & \chi^{(\ell)} &= \frac{N}{\binom{D}{\ell}} \frac{1 - \gamma_\ell g_\ell(-\gamma_\ell)}{\hat{\chi}^{(\ell)}}.
 \end{aligned}
 \tag{56}$$

388 From the solution of Eq. (56), we can fully determine $m^{(\ell)}, \hat{m}^{(\ell)}$ according to

$$m^{(0)} = \frac{\langle y \rangle}{\mu_0}, \quad m^{(\ell)} = \chi^{(\ell)} \hat{m}^{(\ell)}, \quad \hat{m}^{(\ell)} = \frac{P}{\binom{D}{\ell}} \frac{\mu_\ell \tau_\ell}{\sqrt{\ell!}} \frac{\langle y \nu \rangle}{1 + \chi^*}.
 \tag{57}$$

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

$$\begin{aligned}
 q^{(0)} &= \frac{\hat{q}^{(0)}}{(\zeta + \hat{\chi}^{(0)})^2} \left(1 - \sum_{\ell=1}^L [1 - \gamma_\ell^2 g_\ell'(-\gamma_\ell)] \right) + \sum_{\ell=1}^L \frac{\hat{m}^{(\ell)2} + \hat{q}^{(\ell)}}{(\zeta + \hat{\chi}^{(0)}) \hat{\chi}^{(\ell)}} [\gamma_\ell g_\ell(-\gamma_\ell) - \gamma_\ell^2 g_\ell'(-\gamma_\ell)], \\
 q^{(\ell)} &= \frac{N}{\binom{D}{\ell}} \frac{\hat{q}^{(0)}}{(\zeta + \hat{\chi}^{(0)}) \hat{\chi}^{(\ell)}} [\gamma_\ell g_\ell(-\gamma_\ell) - \gamma_\ell^2 g_\ell'(-\gamma_\ell)] \\
 &\quad + \frac{N}{\binom{D}{\ell}} \frac{\hat{m}^{(\ell)2} + \hat{q}^{(\ell)}}{\hat{\chi}^{(\ell)2}} [1 + \gamma_\ell^2 g_\ell'(-\gamma_\ell) - 2\gamma_\ell g_\ell(-\gamma_\ell)], \\
 \hat{q}^{(0)} &= \frac{P}{N} \mu_\perp^2 \frac{\langle (\mu_0 m^{(0)} - y)^2 \rangle - 2\langle y \nu \rangle m^* + q^*}{(1 + \chi^*)^2}, \\
 \hat{q}^{(\ell)} &= \frac{P}{\binom{D}{\ell}} \frac{\mu_\ell^2}{\ell!} \frac{\langle (\mu_0 m^{(0)} - y)^2 \rangle - 2\langle y \nu \rangle m^* + q^*}{(1 + \chi^*)^2}.
 \end{aligned}
 \tag{58}$$

391 Notice that, because of the conventional scalings we chose for the hat variables starting from
 392 Eq. (39) and for the definition of γ_ℓ , these equations give $O(1)$ results for the order parameters
 393 m, χ, q .

394 By numerically integrating Eq. (56), (57), (58), we obtain the theoretical curves for the
 395 generalization error in Eq. (16) and for the order parameters we report in this paper. We com-
 396 pare the result with numerical simulations: despite its asymptotic nature and the hypothesis
 397 of row space orthogonality, our theory works reasonably well even if D is not large. The results
 398 are shown in Fig. 1, 2 ($D = 30$ in this case), where the generalization error is quantitatively
 399 predicted by the theory both when varying P and N .

400 6 Strongly separated regimes

401 Our analysis relies on a number of assumptions:

- 402 1. the Gaussian ansatz on the distribution of $(\nu, \{\lambda^a\}_a)$ at given F, θ and \mathbf{w}^a in the repli-
 403 cated partition function (39);
- 404 2. the truncation of the kernel \mathcal{K} at order L , based on a concentration property of the
 405 matrices $C^{\circ\ell}$;
- 406 3. the fact that the row spaces of the matrices $C^{\circ\ell}$ and $C^{\circ k}$ are orthogonal for $\ell \neq k$, in
 407 order to factorize their contribution to the partition function;
- 408 4. the possibility of taking $C^{\circ\ell}$ as matrices with a spectrum asymptotically described by the
 409 Marchenko-Pastur distribution with aspect ratio $N/(D^\ell/\ell!)$;
- 410 5. the Replica Symmetric ansatz for the overlap matrices describing the teacher-student
 411 distribution;
- 412 6. the possibility of taking the saddle point on the replica parameters for large N , consid-
 413 ering only the leading order in N, P before fixing their relative scaling with D .

414 Some of these assumptions have been already discussed in the previous sections. In the fol-
 415 lowing, we revise and motivate the assumptions on the matrices $C^{\circ\ell}$, namely 2-4, that can
 416 be justified if $P, N, D \rightarrow \infty$ (see Appendix D.2 for more details). Depending on the rela-
 417 tion between the three parameters one is led to consider the following different asymptotic
 418 regimes:

- 419 (i) $N, P, D \rightarrow \infty, P/N \rightarrow 0, P/D^K$ finite; (this includes the case $N \sim D^L$ with $L > K$).
- 420 (ii) $N, P, D \rightarrow \infty, N/D^L$ finite, P/N finite;
- 421 (iii) $N, P, D \rightarrow \infty, P/N \rightarrow \infty, N/D^L$ finite; (this includes the case $P \sim D^K$ with $K > L$).

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

$$k_\ell = \text{Tr} \log(a\mathbb{1} + bC^{\circ\ell}), \quad C_{ij}^{\circ\ell} = \left(\frac{1}{D} \sum_\alpha F_{i\alpha} F_{j\alpha} \right)^\ell \quad (59)$$

423 in three situations (a) $D^\ell \gg N$; (b) $D^\ell \ll N$; (c) $D^\ell \sim N$. Notice that in all cases, while
 424 the diagonal elements are $C_{ii}^{\circ\ell} = 1 + \ell O(\sqrt{1/D})$, the off-diagonal elements $C_{i \neq j}^{\circ\ell}$ are of the
 425 order $D^{-\ell/2}$. In case (a), $D^\ell \gg N$, apart for a negligible number of possible eigenvalue of
 426 order $N/D^{\ell/2}$, all the other eigenvalues are $\lambda = 1 + O(\sqrt{N/D^\ell})$, and to the leading order we

427 simply have $k_\ell = N \log(a + b)$. If we are in the opposite situation, (b), $D^\ell \ll N$, we have only
 428 $O(D^\ell)$ non-zero eigenvalues, roughly equal to $\ell! N/D^\ell + O(\sqrt{N/D^\ell})$, and to the leading order
 429 $k_\ell = N \log(a)$. The interesting case is (c) $N = O(D^\ell)$: we have here D^ℓ eigenvalues of order 1
 430 that contribute to k_ℓ . The leading contribution can be understood writing

$$C_{ij}^{\otimes \ell} = \frac{\ell!}{D^\ell} \sum_{\alpha} F_{i,\alpha}^{\otimes \ell} F_{j,\alpha}^{\otimes \ell} + \text{terms with less different } \alpha\text{'s} \quad (60)$$

431 where the sum includes the terms where the α 's in the multi-index α are ordered, coherently
 432 with our definition in Table 1. This leading term is a matrix of rank $\min\{N, D^\ell/\ell!\}$: the $D^\ell/\ell!$
 433 vectors $\mathbf{F}_\alpha^{\otimes \ell}$ are approximately orthogonal in \mathbb{R}^N , as

$$\mathbf{F}_\alpha^{\otimes \ell} \cdot \mathbf{F}_\beta^{\otimes \ell} = \sum_{i=1}^N F_{i\alpha_1} F_{i\beta_1} \cdots F_{i\alpha_\ell} F_{i\beta_\ell} = N \delta_{\alpha_1\beta_1} \cdots \delta_{\alpha_\ell\beta_\ell} + O(N^{1/2}) \quad (61)$$

434 when α and β are ordered, by law of large numbers, so that the sum of outer products
 435 $\sum_{\alpha} \mathbf{F}_\alpha^{\otimes \ell} (\mathbf{F}_\alpha^{\otimes \ell})^\top$ has rank $D^\ell/\ell!$ as long as $N > D^\ell/\ell!$; if $N < D^\ell/\ell!$, this $N \times N$ matrix is full
 436 rank. Other terms with smaller number of indices in the sum lead to matrices of lower rank r
 437 (with $r/N \rightarrow 0$). Moreover, due to the randomness of the F , the row spaces of these term are
 438 effectively orthogonal to the leading one. To understand this, take for example the case $\ell = 2$
 439 and $N = O(D^2)$: the leading order term of the matrix $C^{\otimes 2}$ has eigenvectors approximately
 440 equal to the vectors $(F_{i\beta_1} F_{i\beta_2})_i$ for $\beta_1 < \beta_2$, as

$$\sum_j \left(\frac{2}{D^2} \sum_{\alpha_1 < \alpha_2} F_{i\alpha_1} F_{i\alpha_2} F_{j\alpha_1} F_{j\alpha_2} \right) F_{j\beta_1} F_{j\beta_2} = \frac{2N}{D^2} F_{i\beta_1} F_{i\beta_2} + O(N^{1/2}/D^2). \quad (62)$$

441 When we apply to this vector the next-to-leading order term of the matrix $C^{\otimes 2}$ we find

$$\sum_j \left(\frac{1}{D^2} \sum_{\alpha} F_{i\alpha}^2 F_{j\alpha}^2 \right) F_{j\beta_1} F_{j\beta_2} = O(N^{1/2}/D^2), \quad (63)$$

442 because the indices β_1 and β_2 are different and one among them remains unpaired. In this
 443 way we can say that the vectors $(F_{i\beta_1} F_{i\beta_2})_i$ are in the null space of the terms we are discarding
 444 in (60). With similar arguments, one can show that the leading terms of $C^{\otimes \ell}$ and $C^{\otimes k}$ have
 445 approximately orthogonal row spaces when $k \neq \ell$ and the scaling of N with D is fixed. We
 446 conclude that we can compute k_ℓ as if $C^{\otimes \ell}$ were a Wishart matrix with aspect ratio $\eta_\ell = N/\binom{D}{\ell}$.
 447 The explicit formula is given in eq. (D.12), and both limits $\eta_\ell \rightarrow 0$ and $\eta_\ell \rightarrow \infty$ agree
 448 with the previous analysis of cases (a) and (b) respectively. We show in Appendix D.2 that
 449 approximating $C^{\otimes \ell}$ as a Wishart matrix gives good results also for moderately large values of
 450 N and D .

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

455 7 Effective theory for finite-size random features networks

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

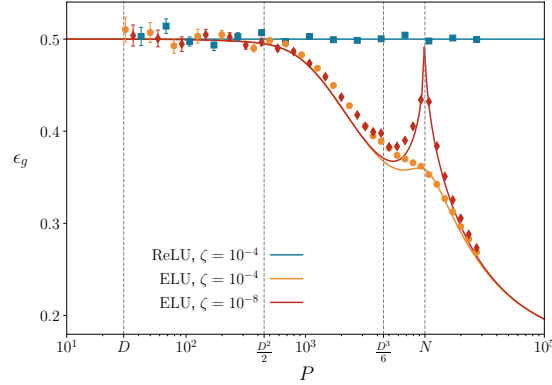


Figure 4: Generalization error vs P ($D = 30$, $N = 10^4$) on classification for a *purely cubic teacher* ($\tau_3 = 1$); in blue, polynomial theory and numerical experiments for ReLU activation function (7): in this case, $\mu_3 = 0$ and the model cannot learn the cubic features, so the error remains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

459 equations (56), (57), (58) is justified only for N large and N/D^L finite, the curves obtained by
 460 fixing the values of N , P and D at finite values are in accordance with numerical simulations
 461 over several orders of magnitudes of the control parameters. This occurs thanks to the fact
 462 that we kept into account quantities that scale differently with D , as $N/\binom{D}{\ell}$ or $P/\binom{D}{\ell}$, that are
 463 formally zero or infinity in the asymptotic regimes presented in Sec. 6.

464 By developing a theory from Eq. (27), we show that the RFM is in essence equivalent to a
 465 polynomial model: the student tries to tune its weights through the combinations $\mathbf{s}^{(\ell)}$ defined
 466 in (28) to fit the corresponding coefficients $\theta^{(\ell)}$ of the teacher. This interpretation is also
 467 confirmed in the numerical experiments: see Fig. 1 (right) for the behavior of the teacher-
 468 student overlaps $m^{(\ell)}$ in the case of a quadratic teacher.

469 However, a crucial difference from a purely polynomial setting arises: the degree of the
 470 equivalent polynomial model is controlled by the scaling L of the random features, and higher
 471 order terms in the expansion of the kernel \mathcal{K} on the Hermite basis act as noise, given by
 472 Eq. (31). This eventually produces the interpolation peak in the generalization error at $N \sim P$,
 473 which would not be present for a vanilla polynomial student (see Fig. 1 and 3): in this regime,
 474 the model is using the effective noise to overfit the teacher. In terms of the order parameters,
 475 overlaps of different orders are coupled by an additional set of parameters $\chi^{(0)}$, $q^{(0)}$, related
 476 to the noise term in the equivalent polynomial model.

477 In summary, the learning of features of a certain order is possible as long as the number
 478 of parameters N is enough: the scaling $L \sim \log N / \log D$ controls the learning process through
 479 the truncation of the kernel (25). At the same time, P also plays an important role: if $K \sim$
 480 $\log P / \log D$ is smaller than L , the model only learns as a K -degree polynomial; on the other
 481 hand, if $K > L$, the model learns as a L -degree polynomial.

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

490 Interestingly, in order for the model to learn features of order ℓ , the activation function σ
 491 must have a non-zero Hermite coefficient μ_ℓ in Eq. (21). This can be seen from our theory
 492 by the fact that in the total teacher-student overlap m^* in Eq. (55) the single entry $m^{(\ell)}$ is
 493 weighted by the corresponding coefficient. This theoretical prediction was tested by using a
 494 cubic teacher and two different students, one with ReLU activation function and the other one
 495 with ELU: the ReLU one, which has no third order term in the Hermite basis ($\mu_3 = 0$) could
 496 not learn the teacher, while the ELU one, that does have a nonzero component ($\mu_3 \neq 0$), was
 497 able to (see Fig. 4).

498 8 Conclusions and perspectives

499 The approach we have explored so far provides a way to analytically evaluate the general-
 500 ization performance of a RFM in the limit of large input dimension D , in the scaling regimes
 501 $N \sim D^L, P \sim D^K$.

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

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

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

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

525 Moreover, it would be interesting to extend our analysis to deeper models [10, 74] in differ-
 526 ent scaling regimes of the dimensions. Even if the RFM, whatever the activation function of the
 527 last layer, is essentially bounded by a polynomial model, the precise shape of the kernel in cases
 528 where a deeper architecture is involved could help understanding to some extent the feature
 529 learning regimes of realistic models, in view of the discussion above. Our approach can also
 530 be extended beyond the case of unstructured input data, following for example [36, 43–50, 75]
 531 and, in particular, [76–78]: in those cases, we expect the intrinsic dimension of the data to
 532 play a role similar to the parameter D used here, possibly determining the order of features
 533 that a RFM can learn at given N and P .

534 Finally, we mention how the replica approach we adopted here can be applied to non-
 535 convex optimization problems, at the cost of choosing a more complicated ansatz for the
 536 overlap matrices, accounting for replica symmetry breaking. Even in those cases, the replica

537 symmetric treatment we provided can be applied as a qualitative approximation, often quan-
 538 titatively correct in the teacher-student setting (that is, whenever a low-energy configuration
 539 of \mathbf{w} is planted by a teacher in the energy landscape defined by the loss (9), effectively con-
 540 vexifying even an *a priori* non-convex problem, *i.e.* setting the problem in a replica symmetric
 541 region of a generically non-convex phase diagram – see, for example, [79, 80], studying the
 542 perceptron with hinge loss in the random labels vs. teacher-student settings).

543 Acknowledgements

544 The authors would like to thank Pietro Rotondo, Rosalba Pacelli, Bruno Loureiro, Valentina
 545 Ros, the QBio group at ENS for discussions and suggestions. MP and FAL are grateful to the
 546 organizers and speakers of the Statistical Physics of Deep Learning summer school held in June
 547 2022 in Como, where the idea was in part conceived.

548 **Funding information** The authors have been supported by a grant from the Simons Foun-
 549 dation (grant No. 454941, S. Franz), thanks to which most of this work was performed at
 550 LPTMS (CNRS, Université Paris-Saclay).

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

554 A Kernel on the Hermite basis

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

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

557 where

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

558 Using the fact that $C_{ii} \simeq C_{jj} \simeq 1$, this kernel can be written as a series of separable kernels
 559 exploiting Mehler's formula [69, 70], that we report here for convenience:

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

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

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

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

563 B Hermite polynomials and Wick products

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

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

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

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

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

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

567 we have, for D large,

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

568 where we used repeatedly $\sum_\alpha (F_{i\alpha})^2 / D \simeq 1$. The thesis follows from comparison with Eq. (29).
 569 Notice that, in the simpler case of a single standard Gaussian variable x , the identity $\text{He}_\ell(x) =$
 570 $:x^\ell :$ is exact and trivially follows from the definition of the Wick power.

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

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

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

573 where

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

574 To proceed, we make the following steps, starting from the expansion of the activation
 575 function on the Hermite basis, Eq. (21). For t_a we simply observe that $\mathbb{E}_{h_i}[\sigma(h_i)] = \mu_0$. For
 576 ρ we use the fact that \mathbf{x} is distributed as a standard normal random vector. To deal with Q_{ab}
 577 we introduce the truncation of (25). Finally, for M_a we write explicitly

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

578 and we perform Wick's contractions in order to evaluate the expected value, exploiting the
 579 mapping to Wick's product explained in Appednix B. As the indices α of the teacher are strictly
 580 ordered, they must be paired only with the ones in the Wick product, leaving only the term
 581 $\ell = k$ in the sum over ℓ . The number of possible contractions is $k!$, so the result is

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

582 from which Eq. (36) follows.

583 D Results on random matrix theory

584 D.1 Marchenko-Pastur distribution and Stieltjes transformation

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

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

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

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

590 with

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

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

592 The MP distribution can be obtained with standard methods [81, 82]. The determinant of
 593 the resolvent can be evaluated as follows:

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

594 By Gaussian linearization,

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

595 The average over F gives

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

596 Integrating over \mathbf{y} ,

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

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

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

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

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

600 The saddle point equation in r gives

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

601 with solutions

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

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

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

603 follows. Deriving with respect to γ ,

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

604 By definition of Stieltjes transformaiton, $r(\gamma) = g(-\gamma)$, which gives Eq. (51).

605 D.2 Spectral density of $C^{\otimes \ell}$

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

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

609 under the hypothesis that $\eta_L = N/\binom{D}{L}$ remain finite. We anticipate that γ_{ℓ} given by (48) either
610 remain finite (if P/N remains finite) or tends to infinity (if $P/N \rightarrow \infty$) in that limit. As we
611 have already discussed, for $\ell > L$, the matrix $C^{\otimes \ell}$ is fully ranked, with diagonal elements close
612 to one and off-diagonal elements of order $D^{-\ell/2}$: all eigenvalues will be equal to one up to
613 a negligible correction. For that reason we could neglect off-diagonal terms for $\ell > L$ and
614 $a_{\ell} \approx b_{\ell} \approx (1 + \gamma_{\ell})^{-1}$. For $\ell < L$ conversely, the matrix has rank D^{ℓ} at most, and it is easy to see
615 that its max eigenvalue cannot be larger than $N \max_i \left(\frac{1}{D} \sum_{\alpha} F_{i,\alpha}^2 \right)^{\ell} = N(1 + O(\sqrt{\log(N)/D}))$.⁴
616 We get therefore

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

⁴ $\lambda_{\max} = \max_{\nu_i \nu_j^2=1} \frac{1}{D^{\ell}} \sum_{\alpha_1, \dots, \alpha_{\ell}} \left(\sum_i \nu_i F_{i,\alpha_1} \dots F_{i,\alpha_{\ell}} \right)^2 \leq N \sum_i \nu_i^2 \left(\frac{1}{D} \sum_{\alpha} F_{i,\alpha}^2 \right)^{\ell}$

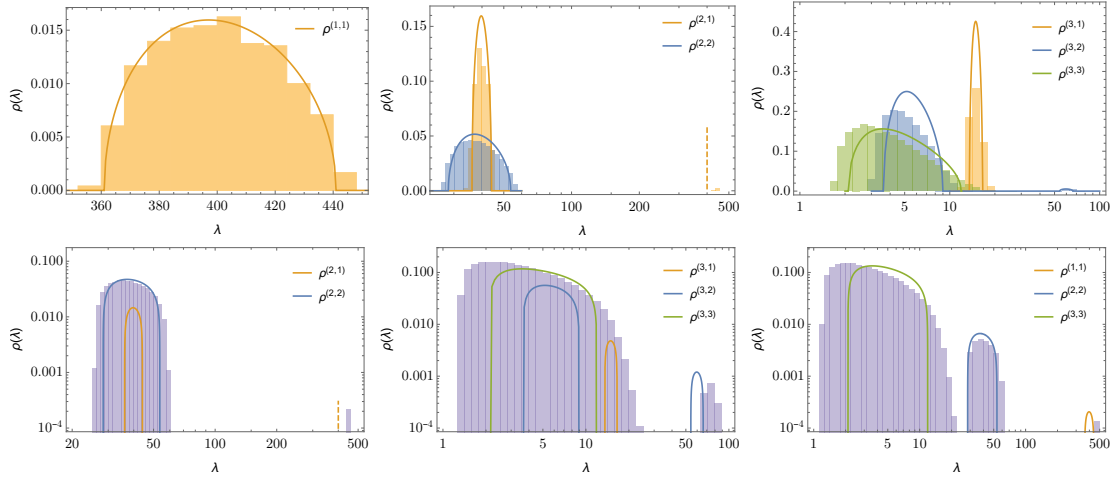


Figure 5: **Top row** – empirical (30 instances, $D = 20$, $N = D^3$) vs. analytical (MP) distributions of the non-zero eigenvalues of the matrices defined in Sec. D.2: $C^{(1,1)}$ (left), $C^{(2,1)}/D$, $C^{(2,2)}$ (center), $C^{(3,1)}/D^2$, $3C^{(3,2)}/D$, $C^{(3,3)}$ (right). **Bottom row** – comparison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

617 It remains to be discussed the only non trivial case: $\ell = L$ In that case, we can decompose
 618 the matrix $C^{\odot L}$ as a matrix with rank $\min\{N, \binom{D}{L}\}$ and spectrum asymptotically distributed
 619 according to the Marchenko-Pastur law with parameter η_L , plus a contribution with rank at
 620 most D^{L-1} which for reasoning similar to the previous case, do not contribute to a_L and b_L in
 621 the thermodynamic limit.

622 We would like now to show, that even for moderate values of N and D , neglecting all the
 623 subleading contributions provides an excellent approximation to the spectrum. To fix ideas,
 624 let us consider $L = 3$ ($N \sim D^3$), so that we consider the matrices

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

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

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

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

628 • $C^{(1,1)} = C$ has a Marchenko-Pastur (MP) spectrum with parameter $\eta_1 = N/D$ and $\sigma^2 =$
629 1 , with D bulk eigenvalues $\lambda = N/D + O(\sqrt{N/D})$ (and $N - D$ zero eigenvalues).

630 • $C^{(2,1)}$ can be written as

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

631 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
632 that $C^{(2,1)}$ has an MP spectrum with parameter η_1 and $\sigma^2 = 2$, with D bulk eigenvalues
633 $O(\sigma^2 \eta_1)$, plus an additional outlier eigenvalue of order N (due to the finite mean);
634 however, in $C^{\odot 2}$ this matrix is scaled by an additional factor of $1/D$, so it contributes to
635 the sum with D eigenvalues $O(2N/D^2)$ and an outlier $O(N/D)$.

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

638 • $C^{(3,1)}$ has an MP spectrum with parameter η_1 and $\sigma^2 = 15$, with D bulk eigenvalues
639 $O(\eta_1)$; however, in $C^{\odot 3}$ this matrix is scaled by an additional factor of $1/D^2$, so it con-
640 tributes to the sum with D eigenvalues $O(N/D^3)$.

641 • $C^{(3,2)}$ can be written as

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

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

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

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

653 E Determinant of sum of matrices with orthogonal row spaces

654 In this section we derive Eq. (44). Let us take the $N \times N$ matrix given by

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

655 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
656 is, the orthogonal complements to their null spaces) are mutually orthogonal ($\mathcal{R}_{\ell} \perp \mathcal{R}_k$ for
657 $k \neq \ell$). Then,

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

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

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

659 with λ_α the non-zero eigenvalues of C_ℓ . Eq. (E.2) can be proven by noticing that, if $\{\mathbf{e}_\ell^\alpha\}_{\alpha=1}^{r_\ell}$
 660 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
 661 in which the matrix K is in block-diagonal form. Moreover, from Eq. (E.3)

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

662 so we can conclude that

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

663 F Traces over RS matrices

664 In this section we derive Eq. (47). We need to evaluate

$$\text{Tr} \log(A \otimes \mathbb{1}_N + B \otimes C^{\otimes \ell}), \quad (\text{F1})$$

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

$$A \otimes \mathbb{1}_N + B \otimes C^{\otimes \ell} = (B \otimes \mathbb{1}_N) (B^{-1}A \oplus C^{\otimes \ell}), \quad (\text{F2})$$

666 where the Kronecker sum is defined as

$$B^{-1}A \oplus C^{\otimes \ell} = B^{-1}A \otimes \mathbb{1}_N + \mathbb{1}_n \otimes C^{\otimes \ell}. \quad (\text{F3})$$

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

$$\log \det(B^{-1}A \oplus C^{\otimes \ell}) = \sum_{a,i} \log(\sigma_a + \lambda_i). \quad (\text{F4})$$

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

$$\log \det(B^{-1}A \oplus C^{\otimes \ell}) = n \sum_i \log(\sigma + \lambda_i) + n \sum_i \frac{\tilde{\sigma}}{\sigma + \lambda_i}. \quad (\text{F5})$$

671 In total we get

$$\text{Tr} \log(A \otimes \mathbb{1}_N + B \otimes C^{\otimes \ell}) = nN \log b + nN \frac{\tilde{b}}{b} + n \sum_i \log(\sigma + \lambda_i) + n \sum_i \frac{\tilde{\sigma}}{\sigma + \lambda_i}. \quad (\text{F6})$$

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

$$\begin{aligned} \text{Tr} \log(A \otimes \mathbb{1}_N + B \otimes C^{\otimes \ell}) &= n \text{Tr} \log(a\mathbb{1} + bC^{\otimes \ell}) + n\tilde{a} \text{Tr}(a\mathbb{1} + bC^{\otimes \ell})^{-1} \\ &\quad + n\tilde{b} \text{Tr}[C^{\otimes \ell}(a\mathbb{1} + bC^{\otimes \ell})^{-1}]. \end{aligned} \quad (\text{F7})$$

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

$$a = \beta(\zeta + \hat{\chi}^{(0)}), \quad \tilde{a} = -\beta^2 \hat{q}^{(0)}, \quad b = \beta \hat{\chi}^{(\ell)}/\eta_\ell, \quad \tilde{b} = -\beta^2 [\hat{q}^{(\ell)} + (\hat{m}^{(\ell)})^2]/\eta_\ell. \quad (\text{F8})$$

674 We define $\gamma_\ell = a/b = \eta_\ell(\zeta + \hat{\chi}^{(0)})/\hat{\chi}^{(\ell)}$ to get Eq. (47).

675 G Replica-symmetric free energy

676 In this section we report the main steps to obtain the terms S_M and S_P in Eq. (53) and (54),
677 that is the measure and pattern contributions to the free energy.

678 G.1 Measure contribution

679 By plugging the RS ansatz (45), (46) and Eq. (47) in Eq. (42), we readily obtain

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

680 We obtain Eq. (53) by keeping the leading order terms for β large and using Eq. (50).

681 G.2 Pattern contribution

682 S_P is a function only of the order parameters:

$$\begin{aligned}
S_P = & \log \left[\int d\nu \prod_{a=1}^n d\lambda^a p(\nu, \{\lambda^a\}) \int dy p(y|\nu) e^{-\beta \sum_a \mathcal{L}(y, \lambda^a)} \right], \\
p(\nu, \{\lambda^a\}) = & \mathcal{N} \left((\nu, \{\lambda^a\}) \mid (0, \{t^a\}), \begin{pmatrix} 1 & M^\top \\ M & Q \end{pmatrix} \right).
\end{aligned} \tag{G.2}$$

683 With the RS ansatz and for small n ,

$$\begin{aligned}
S_P = & \log \left[\int dy d\nu \prod_{a=1}^n d\lambda^a p(y|\nu) e^{-\frac{\nu^2}{2} + \beta \frac{m^* \nu}{\chi^*} \sum_a \lambda^a - \frac{\beta}{2\chi^*} \sum_a \lambda_a^2 - \beta \sum_a \mathcal{L}(y, \lambda^a + t^*) - \beta^2 \frac{m^{*2} - q^*}{2\chi^{*2}} \sum_{a,b} \lambda^a \lambda^b} \right] \\
& - \frac{n}{2} \log(2\pi) - \frac{1}{2} \log \det \begin{pmatrix} 1 & M^\top \\ M & Q \end{pmatrix}.
\end{aligned} \tag{G.3}$$

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

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

685 obtaining, to leading order in n ,

$$\begin{aligned}
S_P = & -\frac{n}{2} \log \frac{\chi^*}{\beta} - \frac{n\beta}{2} \frac{q^*}{\chi^*} + n \mathbb{E}_\xi \int dy D\nu p(y|\nu) \\
& \times \log \int d\lambda e^{\beta \left(\sqrt{q^* - m^{*2}} \xi + m^* \nu \right) \frac{\lambda}{\chi^*} - \frac{\beta \lambda^2}{2\chi^*} - \beta \mathcal{L}(y, \lambda + t^*)}.
\end{aligned} \tag{G.5}$$

686 For our choice of loss (10) and for β large, we obtain Eq. (54).

687 For a generic choice of loss \mathcal{L} , the integral in λ in Eq. (G.5) can still be evaluated asymptotically for large β . The saddle point in λ is given by

$$\lambda^* = \underset{\lambda}{\text{argmin}} \left[\frac{\lambda^2}{2\chi^*} + \mathcal{L}(y, \lambda + t^*) - \frac{\sqrt{q^* - m^{*2}} \xi + m^* \nu}{\chi^*} \lambda \right], \tag{G.6}$$

689 that is by the solution of the stationary equation

$$\lambda + \chi^* \frac{\partial}{\partial \lambda} \mathcal{L}(y, \lambda + t^*) = \sqrt{q^* - m^{*2}} \xi + m^* \nu. \quad (\text{G.7})$$

690 For any choice of \mathcal{L} , this equation gives the value of λ^* as a function of y , ν , ξ and the order
691 parameters. By substituting this value in (G.5) we obtain a generalized form of S_p valid for any
692 loss. By differentiating with respect to the order parameters, we obtain saddle point equations
693 valid for any loss, generalizing the ones for the hat variables reported in Sec. 5.2.

694 H Asymptotic limits of the saddle-point equations

695 The system of saddle-point equations can be studied in different asymptotic limits, as we an-
696 ticipated in Sec. 6:

- 697 (i) $N, P, D \rightarrow \infty, P/N \rightarrow 0, P/D^K$ finite;
698 (ii) $N, P, D \rightarrow \infty, N/D^L$ finite, P/N finite;
699 (iii) $N, P, D \rightarrow \infty, P/N \rightarrow \infty, N/D^L$ finite.

700 H.1 Case (i)

701 In the limit where N scales faster to infinity than P , Eq. (56) reduces to

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

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

702 where we used the asymptotic results for the Stieltjes transformation of the Marchenko-Pastur
703 distribution,

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

704 Notice that now, consistently,

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

705 because $\mu_{\perp, L}^2$ recombines with the terms coming from $K < \ell \leq L$ to give $\mu_{\perp, K}^2$. Eq. (57) reduces
706 to

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

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

707 while Eq. (58) becomes

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

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

708 where now

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

709 H.2 Case (ii)

710 In the limit where both P and N scale in the the same way, $N \sim P \sim O(D^L)$, we have, for
711 $0 < \ell < L$,

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

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

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

713 for m ,

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

714 and for q

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

715 The values χ^* , m^* and q^* are consistent with their definition. At variance with case (i), $\chi^{(0)}$ and
716 $q^{(0)}$ have non-trivial values, responsible for the interpolation peak appearing in this regime.
717 Notice that their value is controlled explicitly by the regularizer ζ : the lower it is, the sharper
718 is the peak. Moreover, the spectral function relative to the active component, g_L , also gives a
719 non-trivial contribution.

720 H.3 Case (iii)

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

$$\gamma_L = \frac{L! \mu_{\perp, L}^2}{\mu_L^2} \quad (\text{H.11})$$

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

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

726 I Numerical experiments

727 All numerical experiments were done in Python using JAX, [83], to generate the synthetic ran-
 728 dom data, and scikit, [62], to optimize the parameters. The optimizer has a simple analytic
 729 form given by (18). Nevertheless, it is potentially inefficient to implement the formula naively,
 730 as it would require the inversion of a very large matrix. Since we used very large values of N
 731 and P , we performed the ridge regression with the function `sklearn.linear_model.Ridge`.
 732 In this way we could explore regimes of N, P up to order D^3 .

733 Almost all numerical experiments were performed with $D = 30$. In most of the simulations
 734 we sampled 50 times for each combination of N, P, D . For the right panel of Figure 3 we used a
 735 larger number of samples since in that case both $D = 30$ and $P = 40 \sim 400$ were small, hence
 736 the generalization error had higher variability. For $N < 3000$ we used 500, 200, 300 samples
 737 respectively for $P = 40, 200, 400$. For $N > 3000$ we used 100, 100, 50 samples respectively for
 738 $P = 40, 200, 400$.

739 A GitHub repository collecting the code needed to reproduce the figures of this paper (both
 740 numerical experiments and theoretical curves from the integration of the saddle-point equa-
 741 tions) can be found at [84].

742 References

- 743 [1] R. M. Neal, *Priors for Infinite Networks*, pp. 29–53, Springer New York, New York, NY,
 744 ISBN 978-1-4612-0745-0, doi:10.1007/978-1-4612-0745-0_2 (1996).
- 745 [2] C. Williams, *Computing with infinite networks*, In M. Mozer, M. Jor-
 746 dan and T. Petsche, eds., *Advances in Neural Information Processing Systems*,
 747 vol. 9. MIT Press (1996), <https://proceedings.neurips.cc/paper/1996/file/ae5e3ce40e0404a45ecacaaf05e5f735-Paper.pdf>.
 748
- 749 [3] J. Lee, J. Sohl-dickstein, J. Pennington, R. Novak, S. Schoenholz and Y. Bahri, *Deep neural*
 750 *networks as Gaussian processes*, In *International Conference on Learning Representations*
 751 (2018), <https://openreview.net/forum?id=B1EA-M-OZ>.

- 752 [4] A. G. de G. Matthews, J. Hron, M. Rowland, R. E. Turner and Z. Ghahramani, *Gaussian*
753 *process behaviour in wide deep neural networks*, In *International Conference on Learning*
754 *Representations* (2018), <https://openreview.net/forum?id=H1-nGgWC->.
- 755 [5] G. Naveh and Z. Ringel, *A self consistent theory of Gaussian processes captures feature*
756 *learning effects in finite CNNs*, In M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang and
757 J. W. Vaughan, eds., *Advances in Neural Information Processing Systems*, vol. 34, pp.
758 21352–21364. Curran Associates, Inc. (2021), [https://proceedings.neurips.cc/paper_](https://proceedings.neurips.cc/paper_files/paper/2021/file/b24d21019de5e59da180f1661904f49a-Paper.pdf)
759 [files/paper/2021/file/b24d21019de5e59da180f1661904f49a-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/b24d21019de5e59da180f1661904f49a-Paper.pdf).
- 760 [6] S. Ariosto, R. Pacelli, F. Ginelli, M. Gherardi and P. Rotondo, *Universal mean-field up-*
761 *per bound for the generalization gap of deep neural networks*, *Phys. Rev. E* **105**, 064309
762 (2022), doi:[10.1103/PhysRevE.105.064309](https://doi.org/10.1103/PhysRevE.105.064309), <https://arxiv.org/abs/2201.11022>.
- 763 [7] R. Pacelli, S. Ariosto, M. Pastore, F. Ginelli, M. Gherardi and P. Rotondo, *A statistical*
764 *mechanics framework for Bayesian deep neural networks beyond the infinite-width limit*,
765 *Nature Machine Intelligence* **5**(12), 1497 (2023), doi:[10.1038/s42256-023-00767-6](https://doi.org/10.1038/s42256-023-00767-6).
- 766 [8] A. Atanasov, B. Bordelon, S. Sainathan and C. Pehlevan, *The onset of variance-limited*
767 *behavior for networks in the lazy and rich regimes*, In *The Eleventh International Conference*
768 *on Learning Representations* (2023), <https://openreview.net/forum?id=JLINxPOVTh7>.
- 769 [9] I. Seroussi, G. Naveh and Z. Ringel, *Separation of scales and a thermodynamic descrip-*
770 *tion of feature learning in some CNNs*, *Nature Communications* **14**(1), 908 (2023),
771 doi:[10.1038/s41467-023-36361-y](https://doi.org/10.1038/s41467-023-36361-y), <https://arxiv.org/abs/2112.15383>.
- 772 [10] H. Cui, F. Krzakala and L. Zdeborova, *Bayes-optimal learning of deep random networks*
773 *of extensive-width*, In A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato and
774 J. Scarlett, eds., *Proceedings of the 40th International Conference on Machine Learning*,
775 vol. 202 of *Proceedings of Machine Learning Research*, pp. 6468–6521. PMLR (2023),
776 <https://proceedings.mlr.press/v202/cui23b.html>.
- 777 [11] L. Chizat, E. Oyallon and F. Bach, *On lazy training in differentiable pro-*
778 *gramming*, In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc,
779 E. Fox and R. Garnett, eds., *Advances in Neural Information Processing Systems*,
780 vol. 32. Curran Associates, Inc. (2019), [https://proceedings.neurips.cc/paper/2019/](https://proceedings.neurips.cc/paper/2019/file/ae614c557843b1df326cb29c57225459-Paper.pdf)
781 [file/ae614c557843b1df326cb29c57225459-Paper.pdf](https://proceedings.neurips.cc/paper/2019/file/ae614c557843b1df326cb29c57225459-Paper.pdf).
- 782 [12] A. Jacot, F. Gabriel and C. Hongler, *Neural tangent kernel: Convergence and gener-*
783 *alization in neural networks*, In S. Bengio, H. Wallach, H. Larochelle, K. Grauman,
784 N. Cesa-Bianchi and R. Garnett, eds., *Advances in Neural Information Processing Sys-*
785 *tems*, vol. 31. Curran Associates, Inc. (2018), [https://proceedings.neurips.cc/paper/](https://proceedings.neurips.cc/paper/2018/file/5a4be1fa34e62bb8a6ec6b91d2462f5a-Paper.pdf)
786 [2018/file/5a4be1fa34e62bb8a6ec6b91d2462f5a-Paper.pdf](https://proceedings.neurips.cc/paper/2018/file/5a4be1fa34e62bb8a6ec6b91d2462f5a-Paper.pdf).
- 787 [13] A. Bietti and J. Mairal, *On the inductive bias of neural tangent kernels*, In
788 H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox and R. Gar-
789 nett, eds., *Advances in Neural Information Processing Systems*, vol. 32. Cur-
790 ran Associates, Inc. (2019), [https://proceedings.neurips.cc/paper/2019/file/](https://proceedings.neurips.cc/paper/2019/file/c4ef9c39b300931b69a36fb3dbb8d60e-Paper.pdf)
791 [c4ef9c39b300931b69a36fb3dbb8d60e-Paper.pdf](https://proceedings.neurips.cc/paper/2019/file/c4ef9c39b300931b69a36fb3dbb8d60e-Paper.pdf).
- 792 [14] A. Montanari and Y. Zhong, *The interpolation phase transition in neural networks: Mem-*
793 *orization and generalization under lazy training*, *The Annals of Statistics* **50**(5), 2816
794 (2022), doi:[10.1214/22-AOS2211](https://doi.org/10.1214/22-AOS2211), <https://arxiv.org/abs/2007.12826>.

- 795 [15] N. Cristianini and J. Shawe-Taylor, *An Introduction to Support Vector Ma-*
796 *chines and Other Kernel-based Learning Methods*, Cambridge University Press,
797 doi:[10.1017/CBO9780511801389](https://doi.org/10.1017/CBO9780511801389) (2000).
- 798 [16] H. Yoon and J.-H. Oh, *Learning of higher-order perceptrons with tunable complexities*, Jour-
799 *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)
800 [4470/31/38/012](https://doi.org/10.1088/0305-4470/31/38/012).
- 801 [17] R. Dietrich, M. Opper and H. Sompolinsky, *Statistical mechanics of support vector net-*
802 *works*, Phys. Rev. Lett. **82**, 2975 (1999), doi:[10.1103/PhysRevLett.82.2975](https://doi.org/10.1103/PhysRevLett.82.2975).
- 803 [18] B. Bordelon, A. Canatar and C. Pehlevan, *Spectrum dependent learning curves in kernel*
804 *regression and wide neural networks*, In H. D. III and A. Singh, eds., *Proceedings of the*
805 *37th International Conference on Machine Learning*, vol. 119 of *Proceedings of Machine*
806 *Learning Research*, pp. 1024–1034. PMLR (2020), [http://proceedings.mlr.press/v119/](http://proceedings.mlr.press/v119/bordelon20a/bordelon20a.pdf)
807 [bordelon20a/bordelon20a.pdf](http://proceedings.mlr.press/v119/bordelon20a/bordelon20a.pdf).
- 808 [19] A. Canatar, B. Bordelon and C. Pehlevan, *Spectral bias and task-model alignment explain*
809 *generalization in kernel regression and infinitely wide neural networks*, Nature Communi-
810 *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)
811 [2006.13198](https://arxiv.org/abs/2006.13198).
- 812 [20] T. Misiakiewicz, *Spectrum of inner-product kernel matrices in the polynomial regime and*
813 *multiple descent phenomenon in kernel ridge regression*, doi:[10.48550/ARXIV.2204.10425](https://doi.org/10.48550/ARXIV.2204.10425)
814 (2022).
- 815 [21] H. Hu and Y. M. Lu, *Sharp asymptotics of kernel ridge regression beyond the linear regime*,
816 doi:[10.48550/ARXIV.2205.06798](https://doi.org/10.48550/ARXIV.2205.06798) (2022).
- 817 [22] L. Xiao, H. Hu, T. Misiakiewicz, Y. M. Lu and J. Pennington, *Precise learning curves*
818 *and higher-order scaling limits for dot-product kernel regression*, Journal of Statistical
819 *Mechanics: Theory and Experiment* **2023**(11), 114005 (2023), doi:[10.1088/1742-](https://doi.org/10.1088/1742-5468/ad01b7)
820 [5468/ad01b7](https://doi.org/10.1088/1742-5468/ad01b7).
- 821 [23] A. Rahimi and B. Recht, *Random features for large-scale kernel machines*, In J. Platt,
822 D. Koller, Y. Singer and S. Roweis, eds., *Advances in Neural Information Processing Sys-*
823 *tems*, vol. 20. Curran Associates, Inc. (2007), [https://proceedings.neurips.cc/paper/](https://proceedings.neurips.cc/paper/2007/file/013a006f03dbc5392effeb8f18fda755-Paper.pdf)
824 [2007/file/013a006f03dbc5392effeb8f18fda755-Paper.pdf](https://proceedings.neurips.cc/paper/2007/file/013a006f03dbc5392effeb8f18fda755-Paper.pdf).
- 825 [24] M.-F. Balcan, A. Blum and S. Vempala, *Kernels as features: On kernels, margins, and*
826 *low-dimensional mappings*, Machine Learning **65**(1), 79 (2006), doi:[10.1007/s10994-](https://doi.org/10.1007/s10994-006-7550-1)
827 [006-7550-1](https://doi.org/10.1007/s10994-006-7550-1).
- 828 [25] A. Rahimi and B. Recht, *Weighted sums of random kitchen sinks: Replac-*
829 *ing minimization with randomization in learning*, In D. Koller, D. Schuurmans,
830 Y. Bengio and L. Bottou, eds., *Advances in Neural Information Processing Systems*,
831 vol. 21. Curran Associates, Inc. (2008), [https://proceedings.neurips.cc/paper/2008/](https://proceedings.neurips.cc/paper/2008/file/0efe32849d230d7f53049ddc4a4b0c60-Paper.pdf)
832 [file/0efe32849d230d7f53049ddc4a4b0c60-Paper.pdf](https://proceedings.neurips.cc/paper/2008/file/0efe32849d230d7f53049ddc4a4b0c60-Paper.pdf).
- 833 [26] A. Rahimi and B. Recht, *Uniform approximation of functions with random bases*, In *2008*
834 *46th Annual Allerton Conference on Communication, Control, and Computing*, pp. 555–
835 561, doi:[10.1109/ALLERTON.2008.4797607](https://doi.org/10.1109/ALLERTON.2008.4797607) (2008).
- 836 [27] B. Ghorbani, S. Mei, T. Misiakiewicz and A. Montanari, *Linearized two-layers neural net-*
837 *works in high dimension*, The Annals of Statistics **49**(2), 1029 (2021), doi:[10.1214/20-](https://doi.org/10.1214/20-AOS1990)
838 [AOS1990](https://doi.org/10.1214/20-AOS1990), <https://arxiv.org/abs/1904.12191>.

- 839 [28] B. Ghorbani, S. Mei, T. Misiakiewicz and A. Montanari, *Limitations of lazy training of*
840 *two-layers neural network*, In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-
841 Buc, E. Fox and R. Garnett, eds., *Advances in Neural Information Processing Systems*,
842 vol. 32. Curran Associates, Inc. (2019), [https://proceedings.neurips.cc/paper/2019/](https://proceedings.neurips.cc/paper/2019/file/c133fb1bb634af68c5088f3438848bfd-Paper.pdf)
843 [file/c133fb1bb634af68c5088f3438848bfd-Paper.pdf](https://proceedings.neurips.cc/paper/2019/file/c133fb1bb634af68c5088f3438848bfd-Paper.pdf).
- 844 [29] S. Mei and A. Montanari, *The generalization error of random features regression: Precise*
845 *asymptotics and the double descent curve*, *Communications on Pure and Applied Mathe-*
846 *matics* **75**(4), 667 (2022), doi:[10.1002/cpa.22008](https://doi.org/10.1002/cpa.22008), <https://arxiv.org/abs/1908.05355>.
- 847 [30] B. Ghorbani, S. Mei, T. Misiakiewicz and A. Montanari, *When do neural networks out-*
848 *perform kernel methods?*, In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan and
849 H. Lin, eds., *Advances in Neural Information Processing Systems*, vol. 33, pp. 14820–
850 14830. Curran Associates, Inc. (2020), [https://proceedings.neurips.cc/paper/2020/](https://proceedings.neurips.cc/paper/2020/file/a9df2255ad642b923d95503b9a7958d8-Paper.pdf)
851 [file/a9df2255ad642b923d95503b9a7958d8-Paper.pdf](https://proceedings.neurips.cc/paper/2020/file/a9df2255ad642b923d95503b9a7958d8-Paper.pdf).
- 852 [31] S. Mei, T. Misiakiewicz and A. Montanari, *Generalization error of random feature and*
853 *kernel methods: Hypercontractivity and kernel matrix concentration*, *Applied and Com-*
854 *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
855 Issue on Harmonic Analysis and Machine Learning, <https://arxiv.org/abs/2101.10588>.
- 856 [32] S. Mei, T. Misiakiewicz and A. Montanari, *Learning with invariances in random features*
857 *and kernel models*, In M. Belkin and S. Kpotufe, eds., *Proceedings of Thirty Fourth Con-*
858 *ference on Learning Theory*, vol. 134 of *Proceedings of Machine Learning Research*, pp.
859 3351–3418. PMLR (2021), <http://proceedings.mlr.press/v134/mei21a/mei21a.pdf>.
- 860 [33] A. Montanari and B. N. Saeed, *Universality of empirical risk minimization*, In P.-L. Loh
861 and M. Raginsky, eds., *Proceedings of Thirty Fifth Conference on Learning Theory*, vol.
862 178 of *Proceedings of Machine Learning Research*, pp. 4310–4312. PMLR (2022), [https:](https://arxiv.org/abs/2202.08832)
863 [//arxiv.org/abs/2202.08832](https://arxiv.org/abs/2202.08832).
- 864 [34] P. L. Bartlett, A. Montanari and A. Rakhlin, *Deep learning: a statistical viewpoint*, *Acta*
865 *Numerica* **30**, 87–201 (2021), doi:[10.1017/S0962492921000027](https://doi.org/10.1017/S0962492921000027).
- 866 [35] M. Belkin, D. Hsu, S. Ma and S. Mandal, *Reconciling modern machine-learning practice*
867 *and the classical bias–variance trade-off*, *Proceedings of the National Academy of Sciences*
868 **116**(32), 15849 (2019), doi:[10.1073/pnas.1903070116](https://doi.org/10.1073/pnas.1903070116).
- 869 [36] S. Goldt, M. Mézard, F. Krzakala and L. Zdeborová, *Modeling the influence of data structure*
870 *on learning in neural networks: The hidden manifold model*, *Phys. Rev. X* **10**, 041044
871 (2020), doi:[10.1103/PhysRevX.10.041044](https://doi.org/10.1103/PhysRevX.10.041044), <https://arxiv.org/abs/1909.11500>.
- 872 [37] F. Gerace, B. Loureiro, F. Krzakala, M. Mézard and L. Zdeborová, *Generalisation error*
873 *in learning with random features and the hidden manifold model*, *Journal of Statisti-*
874 *cal Mechanics: Theory and Experiment* **2021**(12), 124013 (2021), doi:[10.1088/1742-](https://doi.org/10.1088/1742-5468/ac3ae6)
875 [5468/ac3ae6](https://doi.org/10.1088/1742-5468/ac3ae6), <https://arxiv.org/abs/2002.09339>.
- 876 [38] S. Goldt, B. Loureiro, G. Reeves, F. Krzakala, M. Mezard and L. Zdeborová, *The Gaussian*
877 *equivalence of generative models for learning with shallow neural networks*, In J. Bruna,
878 J. Hesthaven and L. Zdeborová, eds., *Proceedings of the 2nd Mathematical and Scientific*
879 *Machine Learning Conference*, vol. 145 of *Proceedings of Machine Learning Research*, pp.
880 426–471. PMLR (2022), <https://proceedings.mlr.press/v145/goldt22a/goldt22a.pdf>.

- 881 [39] B. Loureiro, C. Gerbelot, H. Cui, S. Goldt, F. Krzakala, M. Mezard and L. Zde-
882 borová, *Learning curves of generic features maps for realistic datasets with a teacher-*
883 *student model*, In M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang and J. W.
884 Vaughan, eds., *Advances in Neural Information Processing Systems*, vol. 34, pp. 18137–
885 18151. Curran Associates, Inc. (2021), [https://proceedings.neurips.cc/paper/2021/](https://proceedings.neurips.cc/paper/2021/file/9704a4fc48ae88598dcbdcdf57f3fdef-Paper.pdf)
886 [file/9704a4fc48ae88598dcbdcdf57f3fdef-Paper.pdf](https://proceedings.neurips.cc/paper/2021/file/9704a4fc48ae88598dcbdcdf57f3fdef-Paper.pdf).
- 887 [40] M. Refinetti, S. Goldt, F. Krzakala and L. Zdeborová, *Classifying high-dimensional Gaus-*
888 *sian mixtures: Where kernel methods fail and neural networks succeed*, In M. Meila and
889 T. Zhang, eds., *Proceedings of the 38th International Conference on Machine Learning*,
890 vol. 139 of *Proceedings of Machine Learning Research*, pp. 8936–8947. PMLR (2021),
891 <http://proceedings.mlr.press/v139/refinetti21b/refinetti21b.pdf>.
- 892 [41] H. Cui, B. Loureiro, F. Krzakala and L. Zdeborová, *Generalization error rates in kernel re-*
893 *gression: The crossover from the noiseless to noisy regime*, In A. Beygelzimer, Y. Dauphin,
894 P. Liang and J. W. Vaughan, eds., *Advances in Neural Information Processing Systems*
895 (2021), https://openreview.net/forum?id=Da_EHrAcfwd.
- 896 [42] D. Schröder, H. Cui, D. Dmitriev and B. Loureiro, *Deterministic equivalent and error*
897 *universality of deep random features learning*, In A. Krause, E. Brunskill, K. Cho, B. En-
898 gelhardt, S. Sabato and J. Scarlett, eds., *Proceedings of the 40th International Conference*
899 *on Machine Learning*, vol. 202 of *Proceedings of Machine Learning Research*, pp. 30285–
900 30320. PMLR (2023), [https://proceedings.mlr.press/v202/schroder23a/schroder23a.](https://proceedings.mlr.press/v202/schroder23a/schroder23a.pdf)
901 [pdf](https://proceedings.mlr.press/v202/schroder23a/schroder23a.pdf).
- 902 [43] S. Chung, D. D. Lee and H. Sompolinsky, *Linear readout of object manifolds*, *Phys. Rev.*
903 *E* **93**, 060301 (2016), doi:[10.1103/PhysRevE.93.060301](https://doi.org/10.1103/PhysRevE.93.060301), [https://arxiv.org/abs/1512.](https://arxiv.org/abs/1512.01834)
904 [01834](https://arxiv.org/abs/1512.01834).
- 905 [44] S. Chung, D. D. Lee and H. Sompolinsky, *Classification and geometry of general perceptual*
906 *manifolds*, *Phys. Rev. X* **8**, 031003 (2018), doi:[10.1103/PhysRevX.8.031003](https://doi.org/10.1103/PhysRevX.8.031003), [https:](https://arxiv.org/abs/1710.06487)
907 [//arxiv.org/abs/1710.06487](https://arxiv.org/abs/1710.06487).
- 908 [45] F. Borra, M. C. Lagomarsino, P. Rotondo and M. Gherardi, *Generalization from correlated*
909 *sets of patterns in the perceptron*, *Journal of Physics A: Mathematical and Theoretical*
910 **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)
911 [06818](https://arxiv.org/abs/1903.06818).
- 912 [46] P. Rotondo, M. C. Lagomarsino and M. Gherardi, *Counting the learnable func-*
913 *tions of geometrically structured data*, *Phys. Rev. Res.* **2**, 023169 (2020),
914 doi:[10.1103/PhysRevResearch.2.023169](https://doi.org/10.1103/PhysRevResearch.2.023169), <https://arxiv.org/abs/1903.12021>.
- 915 [47] P. Rotondo, M. Pastore and M. Gherardi, *Beyond the storage capacity:*
916 *Data-driven satisfiability transition*, *Phys. Rev. Lett.* **125**, 120601 (2020),
917 doi:[10.1103/PhysRevLett.125.120601](https://doi.org/10.1103/PhysRevLett.125.120601), <https://arxiv.org/abs/2005.09992>.
- 918 [48] M. Pastore, P. Rotondo, V. Erba and M. Gherardi, *Statistical learning theory of structured*
919 *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)
920 [//arxiv.org/abs/2005.10002](https://arxiv.org/abs/2005.10002).
- 921 [49] M. Pastore, *Critical properties of the SAT/UNSAT transitions in the classification problem*
922 *of structured data*, *Journal of Statistical Mechanics: Theory and Experiment* **2021**(11),
923 113301 (2021), doi:[10.1088/1742-5468/ac312b](https://doi.org/10.1088/1742-5468/ac312b), <https://arxiv.org/abs/2109.08502>.

- 924 [50] M. Gherardi, *Solvable model for the linear separability of structured data*, Entropy **23**(3)
925 (2021), doi:[10.3390/e23030305](https://doi.org/10.3390/e23030305).
- 926 [51] M. Mezard, G. Parisi and M. Virasoro, *Spin Glass Theory and Beyond*, World Scientific,
927 doi:[10.1142/0271](https://doi.org/10.1142/0271) (1986).
- 928 [52] O. Dhifallah and Y. M. Lu, *A precise performance analysis of learning with random features*
929 (2020), <https://arxiv.org/abs/2008.11904>.
- 930 [53] H. Hu and Y. M. Lu, *Universality laws for high-dimensional learning with ran-*
931 *dom features*, IEEE Transactions on Information Theory **69**(3), 1932 (2023),
932 doi:[10.1109/TIT.2022.3217698](https://doi.org/10.1109/TIT.2022.3217698), <https://arxiv.org/abs/2009.07669>.
- 933 [54] Z. Wang and Y. Zhu, *Overparameterized random feature regression with nearly orthogonal*
934 *data*, In F. Ruiz, J. Dy and J.-W. van de Meent, eds., *Proceedings of The 26th Interna-*
935 *tional Conference on Artificial Intelligence and Statistics*, vol. 206 of *Proceedings of Machine*
936 *Learning Research*, pp. 8463–8493. PMLR (2023), [https://proceedings.mlr.press/v206/](https://proceedings.mlr.press/v206/wang23m/wang23m.pdf)
937 [wang23m/wang23m.pdf](https://proceedings.mlr.press/v206/wang23m/wang23m.pdf).
- 938 [55] Y. M. Lu and H.-T. Yau, *An equivalence principle for the spectrum of random inner-product*
939 *kernel matrices with polynomial scalings* (2023), <https://arxiv.org/abs/2205.06308>.
- 940 [56] S. d'Ascoli, L. Sagun and G. Biroli, *Triple descent and the two kinds of overfitting: where*
941 *& why do they appear?*, In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan
942 and H. Lin, eds., *Advances in Neural Information Processing Systems*, vol. 33, pp. 3058–
943 3069. Curran Associates, Inc. (2020), [https://proceedings.neurips.cc/paper/2020/file/](https://proceedings.neurips.cc/paper/2020/file/1fd09c5f59a8ff35d499c0ee25a1d47e-Paper.pdf)
944 [1fd09c5f59a8ff35d499c0ee25a1d47e-Paper.pdf](https://proceedings.neurips.cc/paper/2020/file/1fd09c5f59a8ff35d499c0ee25a1d47e-Paper.pdf).
- 945 [57] H. Hu, Y. M. Lu and T. Misiakiewicz, *Asymptotics of random feature regression beyond the*
946 *linear scaling regime* (2024), <https://arxiv.org/abs/2403.08160>.
- 947 [58] L. Defilippis, B. Loureiro and T. Misiakiewicz, *Dimension-free deterministic equivalents for*
948 *random feature regression* (2024), <https://arxiv.org/abs/2405.15699>.
- 949 [59] E. Gardner and B. Derrida, *Three unfinished works on the optimal storage capacity*
950 *of networks*, Journal of Physics A: Mathematical and General **22**(12), 1983 (1989),
951 doi:[10.1088/0305-4470/22/12/004](https://doi.org/10.1088/0305-4470/22/12/004).
- 952 [60] A. Engel and C. Van den Broeck, *Statistical Mechanics of Learning*, Cambridge University
953 Press, doi:[10.1017/CBO9781139164542](https://doi.org/10.1017/CBO9781139164542) (2001).
- 954 [61] G. Folena, S. Franz and F. Ricci-Tersenghi, *Rethinking mean-field glassy dynamics and its*
955 *relation with the energy landscape: The surprising case of the spherical mixed p-spin model*,
956 Phys. Rev. X **10**, 031045 (2020), doi:[10.1103/PhysRevX.10.031045](https://doi.org/10.1103/PhysRevX.10.031045).
- 957 [62] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel,
958 P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos *et al.*, *Scikit-learn: Machine*
959 *learning in Python*, Journal of Machine Learning Research **12**(85), 2825 (2011), [http:](http://jmlr.org/papers/v12/pedregosa11a.html)
960 [//jmlr.org/papers/v12/pedregosa11a.html](http://jmlr.org/papers/v12/pedregosa11a.html).
- 961 [63] B. Widrow and M. E. Hoff, *Adaptive switching circuits*, In *1960 IRE WESCON Convention*
962 *Record, Part 4*, pp. 96–104 (1960).
- 963 [64] P. Breuer and P. Major, *Central limit theorems for non-linear functionals of Gaussian fields*,
964 Journal of Multivariate Analysis **13**(3), 425 (1983), doi:[10.1016/0047-259X\(83\)90019-](https://doi.org/10.1016/0047-259X(83)90019-2)
965 [2](https://doi.org/10.1016/0047-259X(83)90019-2).

- 966 [65] J.-M. Bardet and D. Surgailis, *Moment bounds and central limit theorems for*
967 *Gaussian subordinated arrays*, *Journal of Multivariate Analysis* **114**, 457 (2013),
968 doi:<https://doi.org/10.1016/j.jmva.2012.08.002>.
- 969 [66] E. Gardner, *The space of interactions in neural network models*, *Journal of Physics A:*
970 *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).
- 971 [67] A. Mozeika, M. Sheikh, F. Aguirre-Lopez, F. Antenucci and A. C. C. Coolen, *Exact results*
972 *on high-dimensional linear regression via statistical physics*, *Phys. Rev. E* **103**, 042142
973 (2021), doi:[10.1103/PhysRevE.103.042142](https://doi.org/10.1103/PhysRevE.103.042142).
- 974 [68] A. C. C. Coolen, M. Sheikh, A. Mozeika, F. Aguirre-Lopez and F. Antenucci, *Replica analysis*
975 *of overfitting in generalized linear regression models*, *Journal of Physics A: Mathematical*
976 *and Theoretical* **53**(36), 365001 (2020), doi:[10.1088/1751-8121/aba028](https://doi.org/10.1088/1751-8121/aba028).
- 977 [69] W. F. Kibble, *An extension of a theorem of Mehler's on Hermite polynomials*, *Math-*
978 *ematical Proceedings of the Cambridge Philosophical Society* **41**(1), 12–15 (1945),
979 doi:[10.1017/S0305004100022313](https://doi.org/10.1017/S0305004100022313).
- 980 [70] T. Liang and H. Tran-Bach, *Mehler's formula, branching process, and compositional kernels*
981 *of deep neural networks*, *Journal of the American Statistical Association* **117**(539), 1324
982 (2022), doi:[10.1080/01621459.2020.1853547](https://doi.org/10.1080/01621459.2020.1853547), <https://arxiv.org/abs/2004.04767>.
- 983 [71] J. Bryson, R. Vershynin and H. Zhao, *Marchenko–Pastur law with relaxed indepen-*
984 *dence conditions*, *Random Matrices: Theory and Applications* **10**(04), 2150040 (2021),
985 doi:[10.1142/S2010326321500404](https://doi.org/10.1142/S2010326321500404), <https://arxiv.org/abs/1912.12724>.
- 986 [72] N. E. Karoui, *The spectrum of kernel random matrices*, *The Annals of Statistics* **38**(1), 1
987 (2010), doi:[10.1214/08-AOS648](https://doi.org/10.1214/08-AOS648).
- 988 [73] J. Glimm and A. Jaffe, *Quantum Physics: A Functional Integral Point of View*, Springer-
989 Verlag, doi:[10.1007/978-1-4612-5158-3](https://doi.org/10.1007/978-1-4612-5158-3) (1987).
- 990 [74] D. Bosch, A. Panahi and B. Hassibi, *Precise asymptotic analysis of deep random feature*
991 *models*, In G. Neu and L. Rosasco, eds., *Proceedings of Thirty Sixth Conference on Learn-*
992 *ing Theory*, vol. 195 of *Proceedings of Machine Learning Research*, pp. 4132–4179. PMLR
993 (2023), <https://proceedings.mlr.press/v195/bosch23a/bosch23a.pdf>.
- 994 [75] F. Cagnetta, L. Petrini, U. M. Tomasini, A. Favero and M. Wyart, *How deep neural networks*
995 *learn compositional data: The random hierarchy model*, *Phys. Rev. X* **14**, 031001 (2024),
996 doi:[10.1103/PhysRevX.14.031001](https://doi.org/10.1103/PhysRevX.14.031001).
- 997 [76] J. A. Zavatone-Veth and C. Pehlevan, *Learning curves for deep structured gaussian feature*
998 *models* (2023), <https://arxiv.org/abs/2303.00564>.
- 999 [77] D. Schröder, D. Dmitriev, H. Cui and B. Loureiro, *Asymptotics of learning with deep struc-*
1000 *tured (random) features* (2024), <https://arxiv.org/abs/2402.13999>.
- 1001 [78] A. Atanasov, J. A. Zavatone-Veth and C. Pehlevan, *Risk and cross validation in ridge*
1002 *regression with correlated samples* (2024), <https://arxiv.org/abs/2408.04607>.
- 1003 [79] S. Franz, A. Sclocchi and P. Urbani, *Critical jammed phase of the linear perceptron*, *Phys.*
1004 *Rev. Lett.* **123**, 115702 (2019), doi:[10.1103/PhysRevLett.123.115702](https://doi.org/10.1103/PhysRevLett.123.115702).
- 1005 [80] E. Loffredo, M. Pastore, S. Cocco and R. Monasson, *Restoring balance: principled un-*
1006 *der/oversampling of data for optimal classification*, In *Forty-first International Conference*
1007 *on Machine Learning* (2024).

- 1008 [81] M. Potters and J.-P. Bouchaud, *A First Course in Random Matrix Theory: for Physicists,*
1009 *Engineers and Data Scientists*, Cambridge University Press, doi:[10.1017/9781108768900](https://doi.org/10.1017/9781108768900)
1010 (2020).
- 1011 [82] G. Livan, M. Novaes and P. Vivo, *Introduction to random matrices theory and practice*,
1012 *Monograph Award* **63**, 54 (2018).
- 1013 [83] J. Bradbury, R. Frostig, P. Hawkins, M. J. Johnson, C. Leary, D. Maclaurin, G. Necula,
1014 A. Paszke, J. VanderPlas, S. Wanderman-Milne and Q. Zhang, *JAX: composable transfor-*
1015 *mations of Python+NumPy programs* (2018).
- 1016 [84] *GitHub repository to reproduce the figures in this paper*, [https://github.com/](https://github.com/MauroPastore/RandomFeatures/)
1017 [MauroPastore/RandomFeatures/](https://github.com/MauroPastore/RandomFeatures/).