

A statistical mechanics framework for Bayesian deep neural networks beyond the infinite-width limit

Original

A statistical mechanics framework for Bayesian deep neural networks beyond the infinite-width limit / Pacelli, R.; Ariosto, S.; Pastore, M.; Ginelli, F.; Gherardi, M.; Rotondo, P.. - In: NATURE MACHINE INTELLIGENCE. - ISSN 2522-5839. - (In corso di stampa).

Availability:

This version is available at: 11583/2983565 since: 2023-11-03T08:00:34Z

Publisher:

Springer Nature

Published

DOI:

Terms of use:

This article is made available under terms and conditions as specified in the corresponding bibliographic description in the repository

Publisher copyright

(Article begins on next page)

A statistical mechanics framework for Bayesian deep neural networks beyond the infinite-width limit

R. Pacelli,^{1,2} S. Ariosto,^{3,4} M. Pastore,^{5,6} F. Ginelli,^{3,4} M. Gherardi,^{7,4} and P. Rotondo^{8,4}

¹*Dipartimento di Scienza Applicata e Tecnologia, Politecnico di Torino, 10129 Torino, Italy*

²*Artificial Intelligence Lab, Bocconi University, 20136 Milano, Italy*

³*Dipartimento di Scienza e Alta Tecnologia and Center for Nonlinear and Complex Systems, Università degli Studi dell'Insubria, Via Valleggio 11, 22100 Como, Italy*

⁴*I.N.F.N. Sezione di Milano, Via Celoria 16, 20133 Milano, Italy*

⁵*Université Paris-Saclay, CNRS, LPTMS, 91405 Orsay, France*

⁶*Laboratoire de physique de l'École normale supérieure, CNRS, PSL University, Sorbonne University, Université Paris-Cité, 24 rue Lhomond, 75005 Paris, France*

⁷*Università degli Studi di Milano, Via Celoria 16, 20133 Milano, Italy*

⁸*Dipartimento di Scienze Matematiche, Fisiche e Informatiche, Università degli Studi di Parma, Parco Area delle Scienze, 7/A 43124 Parma, Italy*

Decades-long literature testifies to the success of statistical mechanics at clarifying fundamental aspects of deep learning. Yet the ultimate goal remains elusive: we lack a complete theoretical framework to predict practically relevant scores, such as the train and test accuracy, from knowledge of the training data. Huge simplifications arise in the infinite-width limit, where the number of units N_ℓ in each hidden layer ($\ell = 1, \dots, L$, being L the finite depth of the network) far exceeds the number P of training examples. This idealisation, however, blatantly departs from the reality of deep learning practice, where training sets are larger than the widths of the networks. Here, we show one way to overcome these limitations. The partition function for fully-connected architectures, which encodes information about the trained models, can be evaluated analytically with the toolset of statistical mechanics. The computation holds in the “thermodynamic limit” where both N_ℓ and P are large and their ratio $\alpha_\ell = P/N_\ell$, which vanishes in the infinite-width limit, is now finite and generic. This advance allows us to obtain (i) a closed formula for the generalisation error associated to a regression task in a one-hidden layer network with finite α_1 ; (ii) an approximate expression of the partition function for deep architectures (technically, via an “effective action” that depends on a finite number of “order parameters”); (iii) a link between deep neural networks in the proportional asymptotic limit and Student’s t processes; (iv) a simple criterion to predict whether finite-width networks (with ReLU activation) achieve better test accuracy than infinite-width ones. As exemplified by these results, our theory provides a starting point to tackle the problem of generalisation in realistic regimes of deep learning.

I. INTRODUCTION

The rise of deep learning, driven by advances in computing technology and foreshadowed by decades of research, has outpaced our ability to develop a solid theoretical foundation [1–3].

Filling the gaps in our understanding of deep learning on a fundamental level is a long-time collective effort involving several communities. Statistical physics achieved far-reaching results in this regard, and remains a well-spring of fresh perspectives and breakthroughs [4–12]. One notable recent advance was obtained by considering the infinite-width limit, where the number of training data P is fixed and the size of the hidden layers is taken to infinity. The observation that such deep models are equivalent to Gaussian processes (GPs) [13–21] established a connection between deep learning and kernel methods [22], and provided a statistical physics description of this regime [8, 23, 24].

However, there is agreement that a more complete theory should address deep learning beyond the infinite-width limit [25–30]: in fact, realistic neural networks operate in a qualitatively different regime, where the number of training examples exceeds the width of the largest

layer. Modeling the finite-width regime in the thermodynamic limit, where the number of degrees of freedom diverges and the tools of statistical mechanics are most effective, amounts to taking the asymptotic limit where both the size of the training set P and the number of units in each hidden layer N_ℓ are taken to infinity with their ratios fixed, as we consider in the present work:

$$P, N_\ell \rightarrow \infty, \quad \alpha_\ell = \frac{P}{N_\ell} \text{ finite} \quad \forall \ell = 1, \dots, L \quad (1)$$

with L being the (finite) depth of the network (the scaling of P with the input size N_0 , which deserves special care, is discussed in section III F of the Methods). This choice guarantees that such networks work in the over-parametrised regime.

Another fruitful line of research, in the direction of overcoming the limitations of the infinite-width limit, sacrifices the non-linear nature of the network by considering a deep *linear* input-output mapping: even if the resulting architecture lacks the expressive power [31–38] of the same model with non-linearities, the multi-layer structure maintains the learning problem non-convex, while amenable to analytical investigation [39, 40]. Very recently, Li and Sompolinsky [6] proposed a method to

analytically evaluate properties of deep finite-width linear networks (e.g., their generalisation error) trained on a generic fixed training set. However, the more relevant case of generic non-linear DNNs remains an open problem, despite some recent notable attempts to address it [4, 27–30, 41–43].

In statistical mechanics, the partition function is the central object encoding the properties of the system in the thermodynamic limit. In this work we address the analytical computation of the partition function of a fully-connected, multi-layer, non-linear neural network, as a function of the training set in the asymptotic limit defined in (1). Technically, the computation amounts to integrating out an extensive number of degrees of freedom (the weights of the network), thus landing on an expression that involves only a finite number (proportional to the depth L) of integrals, to be evaluated by the saddle-point method. In the one-hidden-layer (1HL) case, the only key approximation is justified by a generalised central limit theorem (the Breuer-Major (BM) theorem [44–46]).

In the general case of an architecture with L hidden fully-connected layers, we show that the distribution of the pre-activations at each layer ℓ is a mixture of Gaussians that depends on ℓ parameters. Notably, the back-propagating integration performed in [6] is not a viable option as soon as non-linearities are added to the model. We introduce a forward-propagating method to carry out nested integrations starting from the input layer. This result depends on an assumption that is similar, at least in spirit, to the Gaussian equivalence principle employed for random and generic feature models [47–53].

From these developments, we are able to obtain quantitative predictions for the generalisation error of the network below the interpolation threshold. Moreover, our results have an intriguing interpretation from the point of view of stochastic processes: we show that in the case of finite α_ℓ the GP arising in the infinite-width limit of Bayesian neural networks [16] should be generalised to a Student’s t stochastic process [54].

As a first application of the theory, we establish a simple criterion (equivalent to the one found in the linear case [6, 55] and in finite- P perturbation theory [56, 57]) to predict whether it is convenient, in terms of generalisation performance, to employ a finite-width deep neural network over its infinite-width version.

Problem setting - We consider a supervised learning problem with training set $\mathcal{T}_P = \{\mathbf{x}^\mu, y^\mu\}_{\mu=1}^P$, where each $\mathbf{x}^\mu \in \mathbb{R}^{N_0}$ and the corresponding labels $y^\mu \in \mathbb{R}$. The architecture is a deep neural network $f_{\text{DNN}}(\mathbf{x})$ with $(L-1)$ fully-connected hidden (FC) layers and a final linear readout layer as defined in (23). We analyse regression

problems with a quadratic loss function:

$$\mathcal{L} = \frac{1}{2} \sum_{\mu=1}^P [y^\mu - f_{\text{DNN}}(\mathbf{x}^\mu)]^2 + \mathcal{L}_{\text{reg}}, \quad (2)$$

$$\mathcal{L}_{\text{reg}} = \frac{\lambda_L}{2\beta} \sum_{i_L=1}^{N_L} v_{i_L}^2 + \frac{1}{2\beta} \sum_{\ell=0}^{L-1} \lambda^{(\ell)} \|W^{(\ell)}\|^2, \quad (3)$$

where L^2 regularisations have been added for each layer to the loss function, $\|\cdot\|$ is the standard Frobenius norm defined for the weights matrices $W^{(\ell)}$, and $\beta = 1/T$ is the inverse temperature parameter.

As a standard practice in statistical mechanics of deep learning, we define the partition function of the problem as:

$$Z = \int \mathcal{D}\theta e^{-\beta \mathcal{L}(\theta)}. \quad (4)$$

where the symbol $\int \mathcal{D}\theta$ indicates the collective integration over the weights of the network, $\theta = \{W^{(\ell)}, v\}$. This choice enforces minimization of the training error for $\beta \rightarrow \infty$. We notice that scaling \mathcal{L}_{reg} by $1/\beta$ has a natural Bayesian learning interpretation: the Gibbs probability $P_\beta(\theta) = Z^{-1} e^{-\beta \mathcal{L}(\theta)}$ associated with the partition function in equation (4) is the posterior distribution of the weights after training, whereas the Gaussian regularization is a prior equivalent to assuming that weights at initialization have been drawn from a Gaussian distribution

In this framework, the average test error over a new (unseen) example (\mathbf{x}^0, y^0) is given by:

$$\langle \epsilon_g(\mathbf{x}^0, y^0) \rangle = \int \mathcal{D}\theta [y^0 - f_{\text{DNN}}(\mathbf{x}^0)]^2 \frac{e^{-\beta \mathcal{L}(\theta)}}{Z}. \quad (5)$$

II. RESULTS

A. Asymptotic effective action for one-hidden-layer neural networks in the Bayesian setting

In the case of 1HL architectures, we are able to reduce the partition function (4) to the following two-variables integral in the thermodynamic limit described in (1):

$$Z = \int dQ \int d\bar{Q} \exp \left[-\frac{N_1}{2} S(Q, \bar{Q}) \right] \quad (6)$$

where we have defined an effective action S given by:

$$S = -Q\bar{Q} + \log(1+Q) + \frac{\alpha_1}{P} \text{Tr} \log \beta \left[\frac{1}{\beta} + \frac{\bar{Q}K}{\lambda_1} \right] + \frac{\alpha_1}{P} y^\top \left[\frac{1}{\beta} + \frac{\bar{Q}K}{\lambda_1} \right]^{-1} y \quad (7)$$

and we have introduced a vectorial notation for the output $y^\top = (y^1, y^2, \dots, y^P)$. The $P \times P$, input-dependent

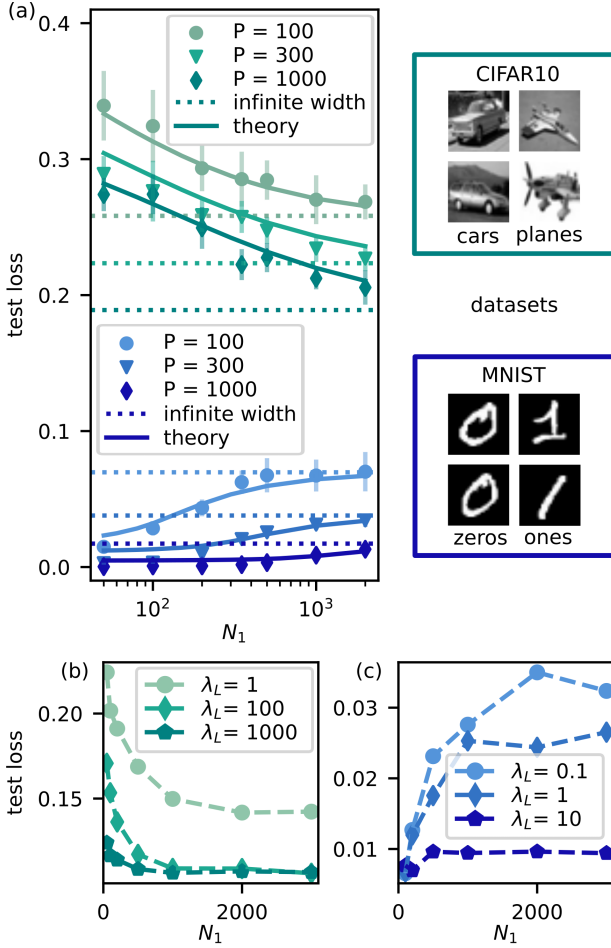


FIG. 1. (a) Learning curves of 1HL architectures with Erf activation (trained with a discretised Langevin dynamics, see also Methods) as a function of the hidden layer size N_1 for two regression tasks on the CIFAR10 (above) and MNIST (below) datasets. Zero/one labels have been chosen in both cases and the images of the CIFAR10 dataset have been gray-scaled and down-scaled to $N_0 = 28 \times 28$. The experimental test loss at different values of the trainset size P (points with error bars indicating one standard deviation) are compared with the theory computed from equation (8) (solid lines). (b,c) Experimental learning curves as a function of N_1 for increasing values of the Gaussian prior of the last layer λ_1 . Dashed lines connecting the points are shown to guide the eye. The nets are trained on $P = 3000$ examples from the CIFAR10 dataset in (b) and $P = 500$ examples from MNIST in (c). Two qualitative predictions of the theory at zero temperature are checked: (i) the generalisation loss should decrease for any N_1 when λ_1 grows; (ii) the dependence of the learning curves on N_1 disappears in the large- λ_1 limit, since the bias is constant (error bars are within points).

kernel K/λ_1 is the neural network Gaussian process (NNGP) kernel [16] arising in the infinite-width limit and its precise definition in terms of the input covariance matrix (rescaled by the Gaussian prior of the first layer λ_0) $C_{\mu\nu} = \mathbf{x}^\mu \cdot \mathbf{x}^\nu / (\lambda_0 N_0)$ is given in the Methods, equation (40). Note also that equation (7) holds for zero-

mean activation functions, that is functions whose average over a centered Gaussian is zero (see equation (39); an effective action for the generic finite-mean case is reported in the supplemental material [58], Sec. IV and that for many reasonable non-linearities and input data distributions the derivation goes through at least in the regime $P = O(N_0)$ (we discuss this key technical point in the Methods). This is the first main result of our work and we conjecture it is exact since the only key Gaussian approximation that we perform is justified by the Breuer-Major theorem [44], as shown in the Methods.

In the supplemental material [58], we obtain a number of additional results that did not enter here for space limitations: (i) a re-derivation of the effective action in the case of linear activation function, valid at fixed P, N_1, N_0 , together with a comparison with the results given in [6, 59]; (ii) a specific derivation of the effective action for quadratic activation function, which makes no use of the Breuer-Major theorem; (iii) the generalisation of the effective action in equation (7) to the case of multiple (but finite) outputs.

We can now solve equation (7) using the saddle-point method, since $N_1 \rightarrow \infty$, which amounts to finding the solutions Q^*, \bar{Q}^* of the system of equations $\partial_Q S = 0$, $\partial_{\bar{Q}} S = 0$ (the infinite-width limit is re-obtained for $\alpha_1 \rightarrow 0$ and corresponds to the particular solution $Q^* = 0$, $\bar{Q}^* = 1$). In the zero-temperature limit, we can find the analytical solution of the saddle-point equations (see Methods). A straightforward computation shows that the generalisation error is given in terms of the usual bias-variance decomposition:

$$\begin{aligned} \langle \epsilon_g(\mathbf{x}^0, y^0) \rangle &= (y^0 - \Gamma_1)^2 + \sigma_1^2, \\ \Gamma_1 &= \sum_{\mu, \nu} \kappa_\mu(\mathbf{x}^0) K_{\mu\nu}^{-1} y_\nu, \\ \sigma_1^2 &= \frac{\bar{Q}^*}{\lambda_1} \left[\kappa_0(\mathbf{x}^0) - \sum_{\mu, \nu} \kappa_\mu(\mathbf{x}^0) K_{\mu\nu}^{-1} \kappa_\nu(\mathbf{x}^0) \right], \end{aligned} \quad (8)$$

where $\kappa_\mu(\mathbf{x}^0)$, $\kappa_0(\mathbf{x}^0)$ can be computed from the functional definition of the NNGP kernel using the new unseen input \mathbf{x}^0 , as shown in the Methods.

We can directly employ equation (8) to obtain testable predictions for the generalisation error of finite-width 1HL architectures trained in the Bayesian learning setting, as we do in panel (a) of Fig. 1 for two specific regression tasks defined on the CIFAR10 and MNIST datasets (details on the numerical experiments are provided in the Methods section III H). It turns out that the generalisation curves for the two regression tasks are monotonically increasing (decreasing) as a function of N_1 depending on the fact that the observable $y^\top (K/\lambda_1)^{-1} y/P$ is smaller (larger) than one. The importance of this quantity in controlling the generalisation performance has been already noted in linear networks [6, 55] as well as in direct perturbation theory at finite P for non-linear networks [56, 57].

We also point out two semi-quantitative predictions

for the general behavior of the generalisation error, just by looking at the dependence of equation (8) on the size of the hidden layer N_1 and on the Gaussian prior of the last layer λ_1 . At $T = 0$, the bias is constant as a function of N_1 (as explicitly observed also in the linear case in Ref. [6]) and of λ_1 . On the contrary, the variance depends on N_1 and decreases as $1/\sqrt{\lambda_1}$ in the large- λ_1 limit. These observations lead to the following two testable predictions: (i) increasing the magnitude of the Gaussian prior λ_1 should systematically improve the generalisation performance at any N_1 ; (ii) for large λ_1 the dependence on N_1 of the generalisation error should disappear (see also the numerical experiments performed in panel (b) in Fig. 1).

B. Link between Student's t -processes and shallow neural networks in the proportional limit

In obtaining the results reported in Sec. II A, our theory can be formulated as a statement on the probability distribution of the output variables

$$s^\mu \equiv \frac{1}{\sqrt{N_1}} \sum_{i_1=1}^{N_1} v_{i_1} \sigma(h_{i_1}^\mu), \quad (9)$$

where $h \sim \mathcal{N}(0, C \otimes \mathbb{1}_{N_1})$, $v \sim \mathcal{N}(0, \lambda_1^{-1} \mathbb{1}_{N_1})$. Proceeding as in the derivation of the partition function presented in Methods, the p.d.f. of these variables can be written as a re-weighted Fourier transform,

$$P(s|\mathcal{T}_P) = \frac{e^{-\frac{\beta}{2} \sum_\mu (y^\mu - s^\mu)^2}}{Z} \int \prod_\mu \frac{d\bar{s}^\mu}{2\pi} e^{i\bar{s}^\top s} \Xi(\bar{s}), \quad (10)$$

of the function

$$\Xi(\bar{s}) = \left(1 + \frac{1}{\lambda_1 N_1} \sum_{\mu, \nu} \bar{s}^\mu K_{\mu\nu}(C) \bar{s}^\nu \right)^{-\frac{N_1}{2}}. \quad (11)$$

It is straightforward to notice that as long as $N_1 \rightarrow \infty$ and $N_1 \gg P$, the dependence on N_1 disappears and we get:

$$\Xi(\bar{s}) \rightarrow e^{-\frac{1}{2\lambda_1} \sum_{\mu, \nu} \bar{s}^\mu K_{\mu\nu}(C) \bar{s}^\nu}. \quad (12)$$

This quantity has a very natural interpretation in view of the NNGP literature. Indeed, for N_1 large and P finite, the variables (9) are jointly multivariate Gaussian distributed according to the central limit theorem, as noted for example in [16]: this limit corresponds indeed to the RHS of our equation (12) and is the cornerstone of the mapping of an infinite-width Bayesian neural network to a GP. This is however no more the case when P is comparable to N_1 : equation (11), derived exploiting the Gaussian equivalence based on the BM theorem in the

proportional asymptotic limit $P/N_1 \sim O(1)$, is suggesting that the variables \bar{s}^μ are distributed according to a multivariate Student's t -distribution [54, 60–62].

The need of considering Student's t -processes as a generalisation of NNGPs has been noted already in the case of different priors on the distribution of the last layer's weights [63]. Non-Gaussianity of the posterior in a form similar to that of Eq. (11) has appeared also in [64–66]. The reason why this kind of process arises in the case we are considering here can be understood with an heuristic argument: when N_1 and P are of the same order, we cannot take the limit $N_1 \rightarrow \infty$ before $P \rightarrow \infty$, and so we need to use the empirical covariance of the output variables s^μ instead of their true one in estimating their probability distribution. A more precise characterization of these neural network Student's t -processes (NNTPs) and the regime where they arise represent interesting topics for future work.

C. Asymptotic effective action for deep neural networks in the Bayesian setting

In the generic case of a deep fully-connected architecture with a finite number of layers L and zero-mean activation function, we express the partition function in terms of a $2L$ -dimensional integral (see Methods):

$$Z_{\text{DNN}} = \int \prod_{\ell=1}^L dQ_\ell d\bar{Q}_\ell e^{-\frac{N_L}{2} S_{\text{DNN}}(\{Q_\ell, \bar{Q}_\ell\})}, \quad (13)$$

where the effective action is given by:

$$\begin{aligned} S_{\text{DNN}} = & \sum_{\ell=1}^L \frac{\alpha_L}{\alpha_\ell} [-Q_\ell \bar{Q}_\ell + \log(1 + Q_\ell)] \\ & + \frac{\alpha_L}{P} \text{Tr} \log \beta \left(\frac{\mathbb{1}}{\beta} + K_L^{(R)}(\{\bar{Q}_\ell\}) \right) \\ & + \frac{\alpha_L}{P} y^T \left(\frac{\mathbb{1}}{\beta} + K_L^{(R)}(\{\bar{Q}_\ell\}) \right)^{-1} y \end{aligned} \quad (14)$$

and we have introduced a renormalised kernel $K^{(R)}$ that generalises the recurrence relation for the L -layer NNGP kernel as:

$$K_\ell^{(R)}(\{\bar{Q}_\ell\}) = \bar{Q}_\ell / \lambda_\ell K \circ [K_{\ell-1}^{(R)}(\{\bar{Q}_\ell\})], \quad K_0^{(R)} = C, \quad (15)$$

where C is the covariance matrix of the inputs defined above and we stress that each $K_\ell^{(R)}$ depends on the variables $\bar{Q}_1, \dots, \bar{Q}_{\ell-1}$ only. For completeness, we notice that the recurrence relation for the infinite-width kernel K_L is given by equation (15) with $\bar{Q}_\ell = 1 \forall \ell = 1, \dots, L$.

This action shares the same structure as the one found in section II A for the special case of 1HL, with the difference that for L hidden layers, the recursive nature of the derivation introduces additional order parameters that are nested in the definition of the kernel K_L . Furthermore, since our derivation applies to layers of arbitrary

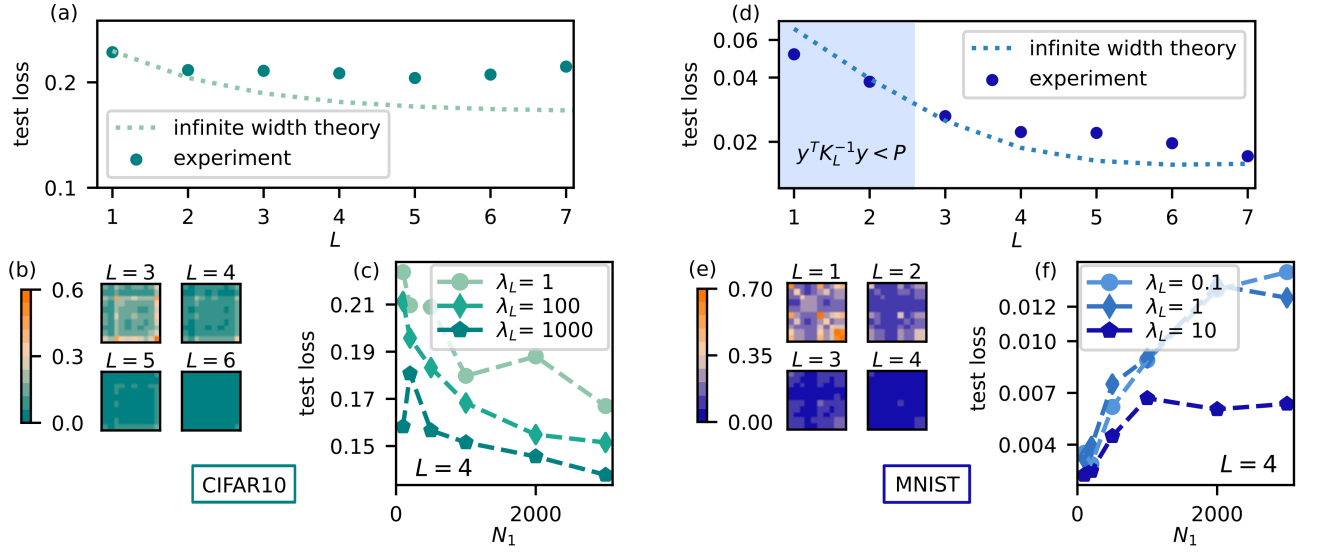


FIG. 2. (a,d) Test loss of a L -HL neural network with ReLU activation, as a function of the depth L , for $P = 100$. The net is trained on a regression task in the small α regime ($\alpha = 0.1$), close to the infinite-width limit. The finite-width network can outperform the infinite-width prediction only when $s_L < 1$ (shaded area), i.e. only for the MNIST task and for depth $L < 3$. (b,e) Visualisation of the entries of the infinite-width NNGP kernel at different layers of the network. The ReLU NNGP kernel converges to zero after repeated iterations. This generates almost vanishing eigenvalues that makes s_L eventually always larger than one. (c,f) Test loss of a 4-HL network trained on $P = 1000$ examples with different regularisation strengths (with $N_\ell = N = 1000$). While increasing the magnitude of the Gaussian prior of the last layer still improves generalisation for all N , it is not clear anymore (as it was for 1HL networks) that the curve at large λ_L is a constant as a function of N . The dashed line is shown to guide the eye (errors bars are within points).

size N_ℓ , the action also depends on the aspect ratios α_ℓ . In the supplemental material [58], we derive a series of additional results: (i) we generalise this effective action for finite-mean activation functions; (ii) we show how to recover the linear case in the isotropic limit $\alpha_\ell = \alpha \forall \ell = 1, \dots, L$; (iii) using (i) we show how to correct the heuristic theory for ReLU activation presented in Ref. [6].

The computation of the generalisation error over a new example (\mathbf{x}^0, y^0) gives:

$$\langle \epsilon_g(\mathbf{x}^0, y^0) \rangle = (y^0 - \Gamma_L)^2 + \sigma_L^2 \quad (16)$$

where

$$\Gamma_L = \sum_{\mu\nu} \kappa_{L\mu}^{(R)} \left(\frac{1}{\beta} + K_L^{(R)}(\{\bar{Q}_\ell\}) \right)_{\mu\nu}^{-1} y_\nu, \quad (17)$$

$$\sigma_L^2 = \kappa_{L0}^{(R)} - \sum_{\mu\nu} \kappa_{L\mu}^{(R)} \left(\frac{1}{\beta} + K_L^{(R)}(\{\bar{Q}_\ell\}) \right)_{\mu\nu}^{-1} \kappa_{L\nu}^{(R)} \quad (18)$$

and $\kappa_{L\mu}^{(R)}$, $\kappa_{L0}^{(R)}$ are recursive kernels computed from the recurrence given in equation (15) using the input \mathbf{x}^0 in the initial conditions.

Note that also in this case we can perform the same scaling analysis of the dependence of the generalisation error on the Gaussian prior in the last layer λ_L (in the zero temperature limit). It turns out that the bias does not depend on it, whereas the variance σ_L^2 approaches

zero as $1/\sqrt{\lambda_L}$ as λ_L is taken to infinity. This means that also in the case of finite depth $L > 1$, training at large values of the Gaussian prior of the last layer should improve generalisation at any aspect ratio of the network. We confirm this general observation with numerical experiments in panels (c) and (f) of Fig. 2. However, differently from the 1HL case, we observe that the bias does depend on the aspect ratio even in the zero-temperature limit and we cannot expect anymore that the dependence on the aspect ratios of the networks α_ℓ disappears in the $\lambda_L \rightarrow \infty$ limit.

We can obtain another prediction of the theory at L layers (that again confirms previous results on linear networks and perturbative calculations for non-linear networks [6, 55–57]) by considering the effective action for ReLU activation. A straightforward Taylor expansion around the infinite-width limit $\alpha_\ell = \alpha = 0 \forall \ell = 1, \dots, L$ shows that the first correction to the test loss $\Delta\epsilon_g$ is proportional to:

$$\Delta\epsilon_g \propto \alpha \left(\frac{1}{P} y^T K_L^{-1} y - 1 \right). \quad (19)$$

where K_L is the solution of recurrence in equation (15) for $\bar{Q}_\ell = 1 \forall \ell = 1, \dots, L$ and ReLU activation. This means that there exists a simple scalar observable that determines whether the finite-width deep neural network will outperform its infinite-width counterpart that gen-

eralises the one found at 1HL:

$$s_L = \frac{1}{P} y^T K_L^{-1} y. \quad (20)$$

In particular, we expect the finite-width network to outperform its infinite-width counterpart whenever $s_L < 1$. In panel (a) and (c) of Fig. 2 we check this prediction for deep architectures with ReLU activation on the same regression tasks employed in the 1HL case. Notice that s_L quickly diverges to infinity as the number of hidden layers L grows. The reason for this is simply that the ReLU NNGP kernel K_L develops at least one zero eigenvalue as $L \rightarrow \infty$. This ultimately occurs because each element of the matrix K_L converges to zero as L grows (see panel (b) and (c) of Fig. 2), as one can easily check by looking at the explicit recurrence relation for the NNGP ReLU kernel. [16, 67]. We note that this singularity can be equivalently thought as the fixed point of the discrete dynamical map defined by the recurrence relation for the NNGP kernel and therefore it might be worth investigating the relation between the generalisation performance in our asymptotic limit and the line of work on the edge of chaos in random neural networks [68, 69].

Equation (19) provides an additional link with Student's t inference. In fact, the same criterion has been found by Tracey and Wolpert [70] in the study of Bayesian optimization with Student's t -processes. Here the authors show that the value of s_L determines whether the Student's t -process they consider has a larger/smaller variance than the corresponding GP with the same kernel.

III. DISCUSSION

In our work we have described a strategy to investigate the statistical mechanics of deep neural networks beyond the infinite-width limit, that is in the finite asymptotic regime $P, N_\ell \rightarrow \infty$ at $\alpha_\ell = P/N_\ell > 0$ as opposed to the infinite-width $\alpha_\ell = 0$. In the 1HL case, we conjecture that our evaluation is exact in the above thermodynamic limit. As such, we do not expect any additional corrections to the result, at least in the asymptotic regime. In particular, we have found a closed expression for the generalisation error that in principle provides a Bayesian estimator of the generalisation capabilities of fully-connected architectures for any given empirical dataset, provided that the chosen architecture is capable of perfectly fitting the trainset.

For the case of finite depth $L > 1$ networks, it should be possible, at least in principle, to take systematically into account non-Gaussian corrections to the saddle-point action to check whether these are relevant or not for the theory at finite width, since the assumptions we made in deriving the results are clear [71] (see also Methods).

From the mathematical perspective, we find the link with Student's t -processes very promising. The precise

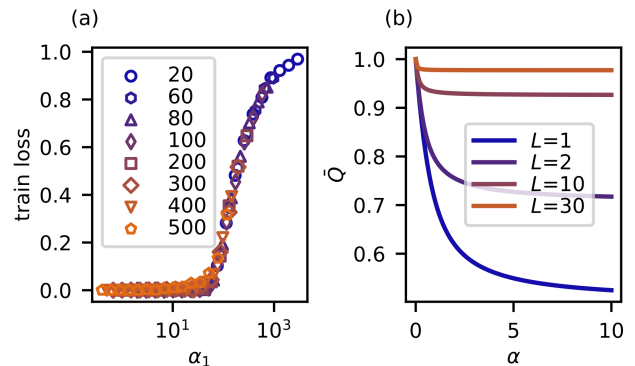


FIG. 3. (Left panel) Training loss of different one-hidden layer architectures trained on a completely random task (i.e. both the inputs $\mathbf{x} \in \mathbb{R}^{N_0}$ with $N_0 = 50$ and the scalar outputs y are i.i.d. random variables sampled from a normal distribution with zero mean and unit variance) as a function of α_1 . At the moment we can not capture this universal phenomenon with our theory, which only describes the overparametrised limit where the training error is exactly zero. (Right panel) The numerical evaluation of the solution \bar{Q} is shown in the case of ReLU activation function and isotropic network $\alpha_\ell = \alpha \forall \ell$, for different depths L . The parameter \bar{Q} quickly approaches -1 for all α , suggesting that also DNNs in the asymptotic regime converge to a kernel limit in the sequential limit where the depth L is taken to infinity after P, N .

characterization of this mapping and its limits of validity represent a research line for future investigation.

Notably, our theory predicts that a kernel limit should also appear in the asymptotic regime as the depth L approaches infinity. This could be checked, for instance, considering the isotropic limit $\alpha_\ell = \alpha \forall \ell$ and ReLU activation. Here one can numerically solve the saddle-point equation for \bar{Q} at large L and verify that $\bar{Q} \rightarrow 1$ for all α , as shown in panel (b) Fig. 3. As such, also in this limit we expect an equivalence with a kernel theory with kernel given by $K_\infty(C)$. Note that from our framework it is clear that we are taking the depth L infinity only after P, N . As such we are not making claims about the challenging simultaneous limit $L, N \rightarrow \infty$ at fixed L/N , as done for instance in Refs. [28, 29, 59, 71].

It is fair to stress that our theory only describes the equilibrium regime of zero train loss, so that our analysis should not apply in the regime $P/N_1 \gg 1$. Interestingly, numerical simulations performed with 1HL architectures of varying width and random training labels show that the train loss follows a universal behavior [55, 72] w.r.t. α_1 (see Fig. 3, panel (a)) also in the regime where the DNN is not capable to perfectly fit the data. It would be desirable to develop a theory that also captures this phase.

Another interesting aspect to understand is the degree to which this mean-field static analysis can be extended beyond equilibrium in order to assess the full training dynamics; such a theory would indeed make it possible to investigate the performance of the many (often heuristics) learning algorithms employed to train deep neural

networks.

We conclude by pointing out that it would be interesting to compare our theory at fixed data with the data-averaged cases studied in [55, 73] and to extend our results to convolutional layers, as done in the infinite-width case in Ref. [18].

METHODS

A. Setting of the learning problem and notation

We consider deep neural networks $f_{\text{DNN}}(\mathbf{x})$ with L fully-connected hidden layers, where the pre-activations of each layer $h_{i_\ell}^{(\ell)}$ ($i_\ell = 1, \dots, N_\ell$; $\ell = 1, \dots, L$) are given recursively as a non-linear function of the pre-activations at the previous layer $h_{i_{\ell-1}}^{(\ell-1)}$ ($i_{\ell-1} = 1, \dots, N_{\ell-1}$):

$$h_{i_\ell}^{(\ell)} = \frac{1}{\sqrt{N_{\ell-1}}} \sum_{i_{\ell-1}=1}^{N_{\ell-1}} W_{i_\ell i_{\ell-1}}^{(\ell)} \sigma(h_{i_{\ell-1}}^{(\ell-1)}) + b_{i_\ell}^{(\ell)}, \quad (21)$$

$$h_{i_1}^{(1)} = \frac{1}{\sqrt{N_0}} \sum_{i_0=1}^{N_0} W_{i_1 i_0}^{(1)} x_{i_0} + b_{i_1}^{(1)} \quad (22)$$

where $W^{(\ell)}$ and $b^{(\ell)}$ are respectively the weights and the biases of the ℓ -th layer, whereas the input layer has dimension N_0 (the input data dimension). σ is a non-linear activation function and it is common to each layer. We add one last readout layer and we define the function implemented by the deep neural network as:

$$f_{\text{DNN}}(\mathbf{x}) = \frac{1}{\sqrt{N_L}} \sum_{i_L=1}^{N_L} v_{i_L} \sigma[h_{i_L}^{(L)}(\mathbf{x})], \quad (23)$$

where \mathbf{v} is the vector of weights of the last layer.

The average training error at a given inverse temperature β is given by:

$$\langle \epsilon_t \rangle = \frac{1}{P} \int \mathcal{D}\theta [\mathcal{L}(\theta) - \mathcal{L}_{\text{reg}}(\theta)] \frac{e^{-\beta \mathcal{L}(\theta)}}{Z}, \quad (24)$$

Training and test errors (as defined in equation (5)) represent two special observables, but more in general, for an arbitrary observable O we have:

$$\langle O \rangle = \int \mathcal{D}\theta O(\theta) \frac{e^{-\beta \mathcal{L}(\theta)}}{Z}. \quad (25)$$

B. The Breuer-Major theorem as a justification for the Gaussian equivalence in shallow networks

The Breuer-Major theorem and its extensions deal with the following sequence of random variables:

$$S_N = \frac{1}{\sqrt{N}} \sum_{i=1}^N c_i F(x_i) \quad N \geq 1. \quad (26)$$

Clearly, if the distribution of the vector $\mathbf{x} = (x_1, \dots, x_N)$ is factorized over its coordinates, i.e. $p(\mathbf{x}) = \prod_i p(x_i)$ and $F(x) = x$, the random variable $S = \lim_{N \rightarrow \infty} S_N$ is normal distributed as long as the mean $\mathbb{E}(x_i) = 0$, the variance $\mathbb{E}(x_i^2)$ is finite and the c_i 's satisfy the so-called Lindeberg's condition. This is also true whenever F is a well-behaved non-linearity.

The Breuer-Major theorem essentially extends this result to generic GPs, providing sufficient conditions on the covariance matrix of the GP and on the non-linearity F that guarantee convergence of S_N to the normal distribution. We report here the modern statement of the theorem given in Ref. [45].

We first consider a stationary (unidimensional) GP $x = (x_k)_{k \in \mathbb{Z}}$. Stationarity—which is not essential and will be replaced by a weaker condition in the following—amounts to require that the covariance of the process $C_{ij} = \mathbb{E}(x_i x_j)$ is a function of the difference $i - j$, i.e. $C_{ij} = C(i - j)$. The only technical condition to be imposed on the non-linear function F is to have well-defined Hermite rank R . The Hermite rank is the smallest positive integer that appears in the decomposition of F over the Hermite polynomials:

$$F(x) = \sum_{k=R}^{\infty} f_k \text{He}_k(x), \quad (27)$$

where $\text{He}_k(x)$ is the k -th Hermite polynomial and f_k the coefficient of the expansion. For many reasonable activation functions F , $R = 1$.

Theorem 1 (Breuer and Major, 1983) Let

$x = (x_k)_{k \in \mathbb{Z}}$ be a stationary unidimensional GP with covariance $C(i - j)$. Let $\mathbb{E}[F(x_1)] = 0$ and $\mathbb{E}[F^2(x_1)] < \infty$ and assume that the function F has Hermite rank $R \geq 1$. Suppose that:

$$\sum_{j \in \mathbb{Z}} |C_{1j}|^R < \infty. \quad (28)$$

Then $\sigma^2 := \mathbb{E}[F(x_1)^2] + 2 \sum_{j=1}^{\infty} \mathbb{E}[F(x_1)F(x_j)]$ is finite. Moreover, one has that the sequence of random variables

$$S_N = \frac{1}{\sqrt{N}} \sum_{i=1}^N F(x_i) \quad N \geq 1 \quad (29)$$

converges in distribution to $\mathcal{N}(0, \sigma^2)$, i.e. to a Gaussian distribution with zero mean and variance σ^2 .

For our scopes we will need a slightly stronger statement than the one just mentioned: (i) in our calculation the covariance will not be stationary and (ii) we will need to consider a more general sequence of nonlinear functions $c_i F(x_i)$, such that each term of the sum (29) is weighted by a factor $c_i \neq 1$.

It has been shown, already in the original reference [44], that the hypothesis of stationarity can be weakened

and replaced with a requirement of uniform convergence of the elements of the covariance, namely:

$$\sum_{j \in \mathbb{Z}} |C_{ij}|^R < B_0 \quad \forall i \in \mathbb{Z}, \quad (30)$$

where B_0 is a positive finite constant. Extension (ii) has been addressed more recently in [46] under mild technical assumptions.

C. Sketch of the calculation of the effective action in the Bayesian setup for one-hidden layer fully-connected neural networks

We now discuss the salient aspects of the calculation. The starting point is the following partition function:

$$Z = \int \prod_{i_1}^{N_1} dv_{i_1} \prod_{i_1, i_0}^{N_1, N_0} dw_{i_1 i_0} \exp \left\{ -\frac{\lambda_1}{2} \sum_{i_1} v_{i_1}^2 - \frac{\lambda_0}{2} \|w\|^2 - \frac{\beta}{2} \sum_{\mu} \left[y^{\mu} - \frac{1}{\sqrt{N_1}} \sum_{i_1} v_{i_1} \sigma \left(\sum_{i_0} \frac{w_{i_1 i_0} x_{i_0}^{\mu}}{\sqrt{N_0}} \right) \right]^2 \right\}. \quad (31)$$

where $w = W^{(1)}$ and we took $b^{(1)} = 0$ without loss generality¹. The first step is to decouple the weights of the different layers in the loss function. This can be done including standard identities built over two families of Dirac deltas, one for the pre-activations of the hidden layer and one for the output of the network:

$$1 = \int \prod_{\mu}^P ds^{\mu} \delta \left[s^{\mu} - \frac{1}{\sqrt{N_1}} \sum_{i_1} v_{i_1} \sigma(h_{i_1}^{\mu}) \right], \quad (32)$$

$$1 = \int \prod_{\mu}^P \prod_{i_1}^{N_1} dh_{i_1}^{\mu} \delta \left(h_{i_1}^{\mu} - \frac{1}{\sqrt{N_0}} \sum_{i_0} w_{i_1 i_0} x_{i_0}^{\mu} \right). \quad (33)$$

By using a standard Fourier representation of these deltas, which introduces the conjugate variables $\bar{h}_{i_1}^{\mu}$ and \bar{s}^{μ} , we can perform the gaussian integrals on the internal and external weights:

$$Z = \int \prod_{\mu}^P \frac{ds^{\mu} d\bar{s}^{\mu}}{2\pi} e^{-\frac{\beta}{2} \sum_{\mu} (y^{\mu} - s^{\mu})^2 + i \sum_{\mu}^P s^{\mu} \bar{s}^{\mu}} \times \left\{ \int \prod_{\mu}^P \frac{dh^{\mu} d\bar{h}^{\mu}}{2\pi} e^{i \sum_{\mu}^P h^{\mu} \bar{h}^{\mu} - \frac{1}{2\lambda_1 N_1} [\sum_{\mu}^P \bar{s}^{\mu} \sigma(h^{\mu})]^2} \times e^{-\frac{1}{2\lambda_0 N_0} \sum_{i_0}^{N_0} (\sum_{\mu}^P \bar{h}^{\mu} x_{i_0}^{\mu})^2} \right\}^{N_1}, \quad (34)$$

where we used the fact that the integrals on $h_{i_1}^{\mu}$ and $\bar{h}_{i_1}^{\mu}$ can be factorized on the index i_1 . The integral over the \bar{h}^{μ} is Gaussian and can be solved:

$$\int \prod_{\mu}^P \frac{d\bar{h}^{\mu}}{2\pi} e^{i \sum_{\mu}^P h^{\mu} \bar{h}^{\mu} - \frac{1}{2\lambda_0 N_0} \sum_{i_0}^{N_0} (\sum_{\mu}^P \bar{h}^{\mu} x_{i_0}^{\mu})^2} = P_1(\{h^{\mu}\}), \quad (35)$$

where

$$P_1(\{h^{\mu}\}) = \frac{e^{-\frac{1}{2} \sum_{\mu, \nu}^P h^{\mu} C_{\mu, \nu}^{-1} h^{\nu}}}{\sqrt{(2\pi)^P \det C}}, \quad C_{\mu \nu} = \frac{1}{\lambda_0 N_0} \sum_{i_0}^{N_0} x_{i_0}^{\mu} x_{i_0}^{\nu}. \quad (36)$$

This last step requires the covariance matrix C to be invertible. Note that this is false as soon as $P > N_0$, but adding a small diagonal term to C solves the issue. One can explicitly check that the final result does not depend on this extra regularization.

To deal with the integral over h^{μ} we can include a further Dirac delta identity for the random variable $q = 1/\sqrt{\lambda_1 N_1} \sum_{\mu}^P \bar{s}^{\mu} \sigma(h^{\mu})$. This leaves us with the problem of finding the probability density $P(q)$. In the limit defined in (1), this is exactly the same setting of the Breuer-Major theorems [44–46]. As such, it is sufficient that both the (regularized) covariance C and the activation function σ satisfy the hypotheses of the theorem to guarantee that the probability distribution $P(q)$ converges in distribution to a Gaussian:

$$P(q) = \int d^P h P_1(\{h^{\mu}\}) \delta \left[q - \frac{1}{\sqrt{\lambda_1 N_1}} \sum_{\mu}^P \bar{s}^{\mu} \sigma(h^{\mu}) \right], \quad (37)$$

$$P(q) \rightarrow \mathcal{N}_q(0, Q).$$

with variance

$$Q(\bar{s}, C) = \frac{1}{\lambda_1 N_1} \sum_{\mu, \nu}^P \bar{s}^{\mu} \left[\int d^P h P_1(\{h^{\rho}\}) \sigma(h^{\mu}) \sigma(h^{\nu}) \right] \bar{s}^{\nu} = \frac{1}{\lambda_1 N_1} \sum_{\mu, \nu}^P \bar{s}^{\mu} K_{\mu \nu}(C) \bar{s}^{\nu}. \quad (38)$$

One can show that there exist special configurations \bar{s} in the domain of integration for which we are not allowed to invoke a Gaussian equivalence (see for instance our discussion at the end of Sec. II in the supplemental material [58]). In our derivation, we are assuming that the contribution of these special configurations to the effective action is negligible in the thermodynamic limit. Here we have also assumed that the variable q has zero mean, a condition verified as long as

$$\int d^P h P_1(\{h^{\mu}\}) \sigma(h^{\nu}) = 0, \quad (39)$$

that is whenever σ is zero-mean; for a more general derivation, relevant for finite-mean activation functions such as ReLU, see the supplemental material [58], Sec. IV.

¹ One can map a system with non-zero biases in a zero-bias one increasing by one the dimensions of the input and of the activations at each layer. The original biases are then trivially mapped in the extra weights of the augmented system.

Each element of the kernel matrix $K_{\mu\nu}(C)$ can be easily reduced from a P -dimensional integral to a simpler two-dimensional one:

$$K_{\mu\nu}(C) = \int \frac{dt_1 dt_2}{\sqrt{(2\pi)^2 \det \tilde{C}}} e^{-\frac{1}{2} \mathbf{t}^T \tilde{C}^{-1} \mathbf{t}} \sigma(t_1) \sigma(t_2), \quad (40)$$

where $\mathbf{t} = (t_1, t_2)^T$ and

$$\tilde{C} = \begin{pmatrix} C_{\mu\mu} & C_{\mu\nu} \\ C_{\mu\nu} & C_{\nu\nu} \end{pmatrix}. \quad (41)$$

is the reduced 2×2 input covariance matrix. It is worth pointing out that the kernel we find here is the so-called *neural network Gaussian process* (NNGP) kernel. It differs from the neural tangent kernel (NTK) that is found in the infinite-width limit of networks trained under gradient descent [74]. The fact that the infinite-width limit of a Bayesian neural network differs from the one obtained from gradient descent is indeed known and discussed in literature [21].

Now we can integrate over the variable q and obtain:

$$\left[\int \frac{dq e^{-\frac{q^2}{2} - \frac{q^2}{2Q(\bar{s}, C)}}}{\sqrt{2\pi Q(\bar{s}, C)}} \right]^{\frac{N_1}{2}} = [Q(\bar{s}, C) + 1]^{-\frac{N_1}{2}}. \quad (42)$$

In the general case of finite $\alpha_1 = P/N_1$, we are only left with the integrals in s^μ and \bar{s}^μ . To solve them it is convenient to introduce one final Dirac delta identity:

$$1 = \int dQ \delta \left[Q - \frac{1}{\lambda_1 N_1} \sum_{\mu, \nu} \bar{s}^\mu K(C)_{\mu\nu} \bar{s}^\nu \right], \quad (43)$$

where $Q \geq -1$ is now an integration variable and not a function of \bar{s} , so that we have removed the explicit dependence on $\sqrt{Q(\bar{s}, C) + 1}$ in the partition function. Finally, the integrals in s^μ and \bar{s}^μ are Gaussian once another integral representation of the delta via a conjugate variable \bar{Q} is inserted. This allows us to get the final effective action obtained in equation (7).

D. Exact solution of the saddle-point equations in the zero temperature limit

The saddle-point equations obtained from (7) considerably simplify in the zero temperature limit ($\beta \rightarrow \infty$). In particular, using the fact that the kernel K has only positive eigenvalues (in the asymptotic regime α_1, α_0 finite), we get:

$$\bar{Q} = \frac{1}{1 + Q}, \quad Q = +\frac{\alpha_1}{\bar{Q}} - \frac{\alpha_1}{\bar{Q}^2} \frac{1}{P} y^T \left(\frac{K}{\lambda_1} \right)^{-1} y. \quad (44)$$

Given the condition $Q \geq -1$, the unique exact solution for \bar{Q} is positive and reads:

$$\bar{Q}^* = \frac{\sqrt{(\alpha_1 - 1)^2 + 4\alpha_1 \frac{1}{P} y^T \left(\frac{K}{\lambda_1} \right)^{-1} y} - (\alpha_1 - 1)}{2}. \quad (45)$$

E. Predictors statistics

The main observable we are interested in is the generalisation error (5). We can proceed along the same lines of the calculation performed in Sec. III C introducing, other than the variables s^μ, h_i^μ defined by (32), (33), additional variables s^0, h_i^0 that describe output and pre-activations of the new test example. We thus get:

$$\begin{aligned} \langle \epsilon_g(\mathbf{x}^0, y^0) \rangle &= \frac{1}{Z} \int \frac{ds^0 d\bar{s}^0}{2\pi} \int \prod_{\mu=1}^P \frac{ds^\mu d\bar{s}^\mu}{2\pi} (y^0 - s^0)^2 \\ &\times e^{-\frac{\beta}{2} \sum_{\mu=1}^P (y^\mu - s^\mu)^2 + i \sum_{\mu=1}^P s^\mu \bar{s}^\mu + i s^0 \bar{s}^0} \\ &\times \left[1 + \frac{1}{\lambda_1 N_1} \left(\sum_{\mu, \nu=1}^P \bar{s}^\mu K_{\mu\nu} \bar{s}^\nu \right. \right. \\ &\quad \left. \left. + 2\bar{s}^0 \sum_{\mu=1}^P \bar{s}^\mu \kappa_\mu(\mathbf{x}^0) + (\bar{s}^0)^2 \kappa_0(\mathbf{x}^0) \right) \right]^{-\frac{N_1}{2}}, \end{aligned} \quad (46)$$

where κ_μ and κ_0 are respectively the train-test and the test-test kernel integrals defined as in (38) when the covariance matrix involves the test input, namely:

$$\kappa_\mu = \int \frac{dt_1 dt_2}{\sqrt{(2\pi)^2 \det \tilde{C}_\mu}} e^{-\frac{1}{2} \mathbf{t}^T \tilde{C}_\mu^{-1} \mathbf{t}} \sigma(t_1) \sigma(t_2), \quad (47)$$

$$\kappa_0 = \int \frac{dt}{\sqrt{2\pi C_{00}}} e^{-\frac{t^2}{2C_{00}}} \sigma(t)^2, \quad (48)$$

where

$$\begin{aligned} \tilde{C}_\mu &= \begin{pmatrix} C_{\mu\mu} & C_{\mu 0} \\ C_{\mu 0} & C_{00} \end{pmatrix}, \quad C_{\mu 0} = \frac{1}{\lambda_0 N_0} \sum_{i_0}^{N_0} x_{i_0}^\mu x_{i_0}^0, \\ C_{00} &= \frac{1}{\lambda_0 N_0} \sum_{i_0}^{N_0} (x_{i_0}^0)^2. \end{aligned} \quad (49)$$

Now we can introduce the order parameters Q and \bar{Q} via equation (43) and their Fourier representation and perform the integration over all the s^μ, \bar{s}^μ and over the s^0 . Doing so yields a single integral in s^0 and integrals on Q and \bar{Q} .

$$\begin{aligned} \langle \epsilon_g(\mathbf{x}^0, y^0) \rangle &= \frac{1}{Z} \int \frac{dQ d\bar{Q}}{2\pi} e^{-\frac{N_1}{2} S(Q, \bar{Q})} \\ &\times \int \frac{ds^0 (y^0 - s^0)^2}{\sqrt{2\pi \sigma^2}} e^{-\frac{(s^0 + \Gamma_1)^2}{2\sigma_1^2}}, \end{aligned} \quad (50)$$

with

$$\begin{aligned}\Gamma_1 &= \frac{\bar{Q}}{\lambda_1} \sum_{\mu\nu} \kappa_\mu(\mathbf{x}^0) \left(\frac{1}{\beta} + \frac{\bar{Q}}{\lambda_1} K \right)_{\mu\nu}^{-1} y_\nu, \\ \sigma_1^2 &= \frac{\bar{Q}}{\lambda_1} \left[\kappa_0(\mathbf{x}^0) - \frac{\bar{Q}}{\lambda_1} \sum_{\mu\nu} \kappa_\mu(\mathbf{x}^0) \left(\frac{1}{\beta} + \frac{\bar{Q}}{\lambda_1} K \right)_{\mu\nu}^{-1} \kappa_\nu(\mathbf{x}^0) \right]\end{aligned}\quad (51)$$

We can then unfold the easy integrals in s^0 and evaluate the result on the saddle point solution. The generalisation error is expressed in terms of Γ_1 and σ_1^2 as in equation (8). Taking the $\beta \rightarrow \infty$ limit in equations (51) yields the expressions in (8).

F. Constraints on the scaling of the size of the dataset P with the input dimension N_0

In this section we address the additional constraints to the thermodynamic scaling ($P, N_1 \rightarrow \infty$ with $\alpha_1 = P/N_1$ finite) that may come from the hypotheses of the Breuer-Major on the covariance matrix C . The only stringent condition to verify is equation (30), that is

$$\sum_{\mu=1}^P |C_{\mu\nu}|^R < B_0 \quad \forall \nu = 1, \dots, P, \quad (52)$$

where B_0 is a given finite constant and R the Hermite rank of the activation function σ . In the case of inputs \mathbf{x} with i.i.d. standard Gaussian coordinates, $C_{\mu\nu}$ is a Wishart random matrix with off-diagonal entries of order $1/\sqrt{N_0}$ and random signs: after taking the absolute value, the sum in Eq (52) is of order $P(N_0)^{-R/2}$. Note that this provides an infinite class of activation functions (those with Hermite rank $R \geq 2$) where we can safely work at least at finite $\alpha_0 = P/N_0$. For activation functions with Hermite rank $R = 1$ (such as Erf or ReLU) we cannot provide such a guarantee by only looking at the hypothesis of the Breuer-Major theorem. It is also worth stressing that, given any odd (non-odd) activation function $\sigma(x)$ with Hermite rank $R = 1$, it is easy to engineer a new reasonable activation function with Hermite rank $R = 3$ ($R = 2$), just by replacing the old activation function with a new one $\sigma_1(x) = \sigma(x) - g_1 x$, where the coefficient $g_1 = \langle \sigma(x) \text{He}_1(x) \rangle$ and the average is over a normal distribution of zero mean and unit variance.

We observe that there is at least one case of activation function with $R = 1$ where the derivation goes through at finite α_0 , i.e. the linear function $\sigma(x) = x$ (in this case we can obtain the result at finite P, N_1, N_0 , as done also in Ref. [59]). In the supplemental material [58], Sec. II, we examine the specific case of quadratic activation $\sigma(x) = x + x^2$ (that has $R = 1$), deriving the final effective action without employing the BM theorem. As in the linear case, this derivation goes through at finite α_0 . We are

thus led to think that the scaling $P = O(\sqrt{N_0})$ suggested for $R = 1$ is overly-pessimistic.

G. Generalisation to deep neural networks with a finite number of hidden layers $L > 1$ and zero-mean activation

In the same spirit of the 1HL calculation, we introduce L sets of auxiliary variables $h_{i_\ell}^\mu$ (where $i_\ell = 1, \dots, N_\ell$) that are equal to the pre-activations at each layer. The strategy to perform the calculation is to show that the probability distribution of the preactivations at each layer $P_\ell(\{h_{i_\ell}^\mu\})$ can be computed recursively, starting from the input layer. We notice that this is conceptually different from the backpropagating kernel renormalisation group introduced in Ref. [6]. It is still a kernel renormalisation group, but forward-propagating, and represents a generalisation to NNTPs of the kernel recurrence arising in NNGPs [16]. In practice, our approach amounts to a systematic, layer-by-layer description of the pre-activation statistics by the Student's t distribution that we have shown to appear in the 1HL case. This can be seen as a quantitative correction to the standard Gaussian statistics that is recovered in the infinite width limit. At the moment we are not able to re-derive the same result using the backpropagating method introduced in [6].

Let us start by integrating the weights of the first layer. This defines a probability distribution over the pre-activations of the first layer via:

$$\begin{aligned}P_1(\{h_{i_1}^\mu\}) &= \int \mathcal{D}W^{(1)} \prod_{i_1, \mu} \delta \left(h_{i_1}^\mu - \frac{1}{\sqrt{N_0}} \sum_{i_0=1}^{N_0} W_{i_1 i_0}^{(1)} x_{i_0}^\mu \right) \\ &= \prod_{i_1=1}^{N_1} \frac{e^{-\frac{1}{2} \sum_{\mu\nu} h_{i_1}^\mu C_{\mu\nu}^{-1} h_{i_1}^\nu}}{\sqrt{(2\pi)^P \det C}}.\end{aligned}\quad (53)$$

where C is defined in (36). This result is straightforward and it is valid for any N_0, P and N_1 , since the prior for the weights is gaussian. At the second layer we have:

$$\begin{aligned}P_2(\{h_{i_2}^\mu\}) &= \int \mathcal{D}W^{(2)} \mathcal{D}h_1 P_1(\{h_{i_1}^\mu\}) \\ &\times \prod_{i_2, \mu} \delta \left(h_{i_2}^\mu - \frac{1}{\sqrt{N_1}} \sum_{i_1=1}^{N_1} W_{i_2 i_1}^{(2)} \sigma(h_{i_1}^\mu) \right)\end{aligned}\quad (54)$$

We now introduce conjugate variables $\bar{h}_{i_2}^\mu$ to the activation of the second layer and the calculation proceeds as in the case of 1HL architectures. To make analytical progress we need to make two fundamental approximations: (i) assuming that the set of random variables $q_{i_2} = 1/(\sqrt{N_1} \lambda_1) \sum_{\mu} \bar{h}_{i_2}^\mu \sigma(h_{i_1}^\mu)$, where $\mathbf{h} \sim \mathcal{N}(0, C)$, is Gaussian-distributed; (ii) neglecting correlations between different pre-activations of the second hidden layer.

In conclusion we get:

$$P_2(\{h_{i_2}^\mu\}) = \int dQ_1 d\bar{Q}_1 e^{-\frac{N_1}{2}(-Q_1\bar{Q}_1 + \log(1+Q_1))} \times \prod_{i_2=1}^{N_2} \frac{e^{-\frac{1}{2}\sum_{\mu\nu} h_{i_2}^\mu (\bar{Q}_1 K(C)/\lambda_1)^{-1}_{\mu\nu} h_{i_2}^\nu}}{\sqrt{(2\pi)^P \det(\bar{Q}_1 K(C)/\lambda_1)}}, \quad (55)$$

where $K(C)$ is defined by equation (38). Notice that except for the integration over the two variables Q_1 and \bar{Q}_1 , this is the same as the probability distribution of the 1HL system (36) if we replace C with $\bar{Q}_1 K(C)/\lambda_1$. This reasoning can be easily iterated across layers and gives:

$$P_L(\{h_{i_L}^\mu\}) = \int \prod_{\ell=1}^{L-1} dQ_\ell d\bar{Q}_\ell e^{-\sum_{\ell=1}^{L-1} \frac{N_\ell}{2}[-Q_\ell\bar{Q}_\ell + \log(1+Q_\ell)]} \times \prod_{i_L=1}^{N_L} \frac{e^{-\frac{1}{2}\sum_{\mu\nu} h_{i_L}^\mu (K_{L-1}^{(R)}(\{\bar{Q}_\ell\}))^{-1}_{\mu\nu} h_{i_L}^\nu}}{\sqrt{(2\pi)^P \det(K_{L-1}^{(R)}(\{\bar{Q}_\ell\}))}}, \quad (56)$$

where $K_\ell^{(R)}(\{\bar{Q}_\ell\})$ is a renormalised kernel that satisfies the recurrence relation in equation (15).

The computation of the generalisation error over a new example (\mathbf{x}^0, y^0) gives:

$$\langle \epsilon_g(\mathbf{x}^0, y^0) \rangle = (y^0 - \Gamma_L)^2 + \sigma_L^2 \quad (57)$$

where Γ_L and σ_L^2 are defined respectively in Eqs. (17) and (18). Note that $\kappa_{L\mu}^{(R)}, \kappa_{L0}^{(R)}$ are recursive kernels that generalise the train-test and test-test kernels (47)-(48). They are defined starting from equation (15) where the kernel K is now evaluated with the covariance matrix C involving train-test or test-test points. Note that L -hidden layers generalisation error is found replacing the 1HL kernel with its recursive generalisation (15).

H. Numerical experiments

1. Network architectures

We perform numerical experiments with deep fully-connected architectures trained on two regression tasks in computer vision. In particular we use the 0 and 1 classes of the MNIST and CIFAR10 datasets, which for the latter correspond to the labels “cars” and “planes”. Examples from CIFAR10 are coarse grained to $N_0 = 28 \times 28$ pixels and converted to grayscale.

To test our theory in the zero-mean activation function case, we used the Erf function, for which the NNGP kernel can be computed analytically [75]:

$$K_{\mu\nu}^{\text{Erf}}(C) = \frac{2}{\pi} \arcsin \left(\frac{2C_{\mu\nu}}{\sqrt{(1+2C_{\mu\mu})(1+2C_{\nu\nu})}} \right). \quad (58)$$

In Fig. 2 we train networks with $\sigma = \text{ReLU}$. The kernel can be computed analytically also in this case [67] and reads:

$$K_{\mu\nu}^{\text{ReLU}}(C) = \sqrt{C_{\mu\mu}C_{\nu\nu}} \kappa \left(\frac{C_{\mu\nu}}{\sqrt{C_{\mu\mu}C_{\nu\nu}}} \right), \quad (59)$$

$$\kappa(x) = \frac{1}{2\pi} \left[x(\pi - \arccos(x)) + \sqrt{1-x^2} \right].$$

2. Sampling from the Bayesian posterior

To ensure convergence of the posterior weights distribution to the Gibbs ensemble, we train our networks using a discretised Langevin dynamics, similarly to what is done in [4, 6]. At each training step t the parameters $\theta = \{W^\ell, v\}$ are updated according to:

$$\theta(t+1) = \theta(t) - \eta \nabla_\theta \mathcal{L}(\theta(t)) + \sqrt{2T\eta} \epsilon(t) \quad (60)$$

where $T = 1/\beta$ is the temperature, η is the learning rate, $\epsilon(t)$ is a white Gaussian noise vector with entries drawn from a standard normal distribution, and the loss is the one defined in equation (3). We employ $T = \eta = 10^{-3}$ throughout all the experiments. This is sufficient to approximate the $T = 0$ dynamics in the regime we are considering. This dynamics requires $10^5/10^6$ steps to reach thermalisation, depending on the sizes of the dataset and network. We extract the generalisation loss within a single run: after the train error has reached its minimum and the test loss is thermalised, we average test loss values every $10^3/10^4$ epochs (depending again on the magnitude of P, N_ℓ). For the sake of completeness, we report the best test accuracy achieved on both datasets by 1HL architectures: 0.86 on CIFAR10 with $P = 3000$ and $\lambda_1 = 1000$, 0.999 on MNIST with the same Gaussian prior and $P = 1000$. The train accuracy is always 1. Additional comments on the technical issues encountered in simulating the Bayesian dynamics are discussed in the Supplemental material [58] in Sec. V.

DATA AVAILABILITY

The CIFAR10 and MNIST datasets that we used for all our experiments are publicly available online, respectively at <https://www.cs.toronto.edu/~kriz/cifar.html> and <http://yann.lecun.com/exdb/mnist/>.

CODE AVAILABILITY

The code used to perform experiments, compute theory predictions and analyze data is available at: https://github.com/rpacelli/FC_deep_bayesian_networks [76].

- [1] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning* (MIT Press, 2016).
- [2] A. Engel and C. Van den Broeck, *Statistical Mechanics of Learning* (Cambridge University Press, 2001).
- [3] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals, Understanding deep learning (still) requires rethinking generalization, *Commun. ACM* **64**, 107–115 (2021).
- [4] I. Seroussi, G. Naveh, and Z. Ringel, Separation of scales and a thermodynamic description of feature learning in some cnns, *Nature Communications* **14**, 908 (2023).
- [5] A. J. Wakhloo, T. J. Sussman, and S. Chung, Linear classification of neural manifolds with correlated variability, *Phys. Rev. Lett.* **131**, 027301 (2023).
- [6] Q. Li and H. Sompolinsky, Statistical mechanics of deep linear neural networks: The backpropagating kernel renormalization, *Phys. Rev. X* **11**, 031059 (2021).
- [7] C. Baldassi, C. Lauditi, E. M. Malatesta, R. Pacelli, G. Perugini, and R. Zecchina, Learning through atypical phase transitions in overparameterized neural networks, *Phys. Rev. E* **106**, 014116 (2022).
- [8] A. Canatar, B. Bordonon, and C. Pehlevan, Spectral bias and task-model alignment explain generalization in kernel regression and infinitely wide neural networks, *Nature communications* **12**, 1 (2021).
- [9] A. Mozeika, B. Li, and D. Saad, Space of functions computed by deep-layered machines, *Phys. Rev. Lett.* **125**, 168301 (2020).
- [10] S. Goldt, M. Mézard, F. Krzakala, and L. Zdeborová, Modeling the influence of data structure on learning in neural networks: The hidden manifold model, *Phys. Rev. X* **10**, 041044 (2020).
- [11] Y. Bahri, J. Kadmon, J. Pennington, S. S. Schoenholz, J. Sohl-Dickstein, and S. Ganguli, Statistical mechanics of deep learning, *Annual Review of Condensed Matter Physics* **11**, 501 (2020).
- [12] B. Li and D. Saad, Exploring the function space of deep-learning machines, *Phys. Rev. Lett.* **120**, 248301 (2018).
- [13] R. M. Neal, Priors for infinite networks, in *Bayesian Learning for Neural Networks* (Springer New York, New York, NY, 1996) pp. 29–53.
- [14] C. Williams, Computing with infinite networks, in *Advances in Neural Information Processing Systems*, Vol. 9, edited by M. Mozer, M. Jordan, and T. Petsche (MIT Press, 1996).
- [15] A. G. de G. Matthews, J. Hron, M. Rowland, R. E. Turner, and Z. Ghahramani, Gaussian process behaviour in wide deep neural networks, in *International Conference on Learning Representations* (2018).
- [16] J. Lee, J. Sohl-dickstein, J. Pennington, R. Novak, S. Schoenholz, and Y. Bahri, Deep neural networks as gaussian processes, in *International Conference on Learning Representations* (2018).
- [17] A. Garriga-Alonso, C. E. Rasmussen, and L. Aitchison, Deep convolutional networks as shallow gaussian processes, in *International Conference on Learning Representations* (2019).
- [18] R. Novak, L. Xiao, Y. Bahri, J. Lee, G. Yang, D. A. Abolafia, J. Pennington, and J. Sohl-dickstein, Bayesian deep convolutional networks with many channels are gaussian processes, in *International Conference on Learning Representations* (2019).
- [19] A. Jacot, F. Gabriel, and C. Hongler, Neural tangent kernel: Convergence and generalization in neural networks, in *Advances in Neural Information Processing Systems*, Vol. 31, edited by S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (Curran Associates, Inc., 2018).
- [20] L. Chizat, E. Oyallon, and F. Bach, On lazy training in differentiable programming, in *Advances in Neural Information Processing Systems*, Vol. 32, edited by H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Curran Associates, Inc., 2019).
- [21] J. Lee, L. Xiao, S. Schoenholz, Y. Bahri, R. Novak, J. Sohl-Dickstein, and J. Pennington, Wide neural networks of any depth evolve as linear models under gradient descent, in *Advances in Neural Information Processing Systems*, Vol. 32, edited by H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Curran Associates, Inc., 2019).
- [22] C. Cortes and V. Vapnik, Support-vector networks, *Machine Learning* **20**, 273 (1995).
- [23] B. Bordonon, A. Canatar, and C. Pehlevan, Spectrum dependent learning curves in kernel regression and wide neural networks, in *Proceedings of the 37th International Conference on Machine Learning*, Proceedings of Machine Learning Research, Vol. 119, edited by H. D. III and A. Singh (PMLR, 2020) pp. 1024–1034.
- [24] R. Dietrich, M. Oppen, and H. Sompolinsky, Statistical mechanics of support vector networks, *Phys. Rev. Lett.* **82**, 2975 (1999).
- [25] M. Seleznova and G. Kutyniok, Neural tangent kernel beyond the infinite-width limit: Effects of depth and initialization, arXiv preprint arXiv:2206.00553 (2022).
- [26] N. Vyas, Y. Bansal, and N. Preetum, Limitations of the ntk for understanding generalization in deep learning, arXiv preprint arXiv:2206.10012 (2022).
- [27] J. M. Antognini, Finite size corrections for neural network gaussian processes (2019), arXiv:1908.10030 [cs.LG].
- [28] S. Yaida, Non-Gaussian processes and neural networks at finite widths, in *Proceedings of The First Mathematical and Scientific Machine Learning Conference*, Proceedings of Machine Learning Research, Vol. 107, edited by J. Lu and R. Ward (PMLR, 2020) pp. 165–192.
- [29] B. Hanin, Random fully connected neural networks as perturbatively solvable hierarchies (2023), arXiv:2204.01058 [math.PR].
- [30] J. Zavatone-Veth and C. Pehlevan, Exact marginal prior distributions of finite bayesian neural networks, in *Advances in Neural Information Processing Systems*, Vol. 34, edited by M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan (Curran Associates, Inc., 2021) pp. 3364–3375.
- [31] Y. Bengio and O. Delalleau, On the expressive power of deep architectures, in *International conference on algorithmic learning theory* (Springer, 2011) pp. 18–36.
- [32] P. L. Bartlett, N. Harvey, C. Liaw, and A. Mehrabian, Nearly-tight vc-dimension and pseudodimension bounds for piecewise linear neural networks, *The Journal of Machine Learning Research* **20**, 2285 (2019).
- [33] P. Rotondo, M. C. Lagomarsino, and M. Gherardi, Counting the learnable functions of geometrically struc-

- tured data, *Phys. Rev. Research* **2**, 023169 (2020).
- [34] P. Rotondo, M. Pastore, and M. Gherardi, Beyond the storage capacity: Data-driven satisfiability transition, *Phys. Rev. Lett.* **125**, 120601 (2020).
- [35] M. Pastore, P. Rotondo, V. Erba, and M. Gherardi, Statistical learning theory of structured data, *Phys. Rev. E* **102**, 032119 (2020).
- [36] M. Gherardi, Solvable model for the linear separability of structured data, *Entropy* **23**, 10.3390/e23030305 (2021).
- [37] M. Pastore, Critical properties of the SAT/UNSAT transitions in the classification problem of structured data, *Journal of Statistical Mechanics: Theory and Experiment* **2021**, 113301 (2021).
- [38] F. Aguirre-López, M. Pastore, and S. Franz, Satisfiability transition in asymmetric neural networks, *Journal of Physics A: Mathematical and Theoretical* **55**, 305001 (2022).
- [39] A. M. Saxe, J. L. McClelland, and S. Ganguli, Exact solutions to the nonlinear dynamics of learning in deep linear neural networks, in *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, edited by Y. Bengio and Y. LeCun (2014).
- [40] A. M. Saxe, J. L. McClelland, and S. Ganguli, A mathematical theory of semantic development in deep neural networks, *Proceedings of the National Academy of Sciences* **116**, 11537 (2019), <https://www.pnas.org/doi/pdf/10.1073/pnas.1820226116>.
- [41] J. Zavatone-Veth, A. Canatar, B. Ruben, and C. Pehlevan, Asymptotics of representation learning in finite bayesian neural networks, in *Advances in Neural Information Processing Systems*, Vol. 34, edited by M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan (Curran Associates, Inc., 2021) pp. 24765–24777.
- [42] G. Naveh and Z. Ringel, A self consistent theory of gaussian processes captures feature learning effects in finite cnns, in *Advances in Neural Information Processing Systems*, Vol. 34, edited by M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan (Curran Associates, Inc., 2021) pp. 21352–21364.
- [43] J. A. Zavatone-Veth, W. L. Tong, and C. Pehlevan, Contrasting random and learned features in deep bayesian linear regression, *Phys. Rev. E* **105**, 064118 (2022).
- [44] P. Breuer and P. Major, Central limit theorems for non-linear functionals of gaussian fields, *Journal of Multivariate Analysis* **13**, 425 (1983).
- [45] I. Nourdin, G. Peccati, and M. Podolskij, *Quantitative Breuer-Major theorems* (2010).
- [46] J.-M. Bardet and D. Surgailis, Moment bounds and central limit theorems for gaussian subordinated arrays, *Journal of Multivariate Analysis* **114**, 457 (2013).
- [47] F. Gerace, B. Loureiro, F. Krzakala, M. Mezard, and L. Zdeborova, Generalisation error in learning with random features and the hidden manifold model, in *Proceedings of the 37th International Conference on Machine Learning*, Proceedings of Machine Learning Research, Vol. 119, edited by H. D. III and A. Singh (PMLR, 2020) pp. 3452–3462.
- [48] B. Loureiro, C. Gerbelot, H. Cui, S. Goldt, F. Krzakala, M. Mezard, and L. Zdeborová, Learning curves of generic features maps for realistic datasets with a teacher-student model, *Advances in Neural Information Processing Systems* **34** (2021).
- [49] S. Goldt, B. Loureiro, G. Reeves, F. Krzakala, M. Mézard, and L. Zdeborová, The gaussian equivalence of generative models for learning with shallow neural networks, arXiv preprint arXiv:2006.14709 (2020).
- [50] E. Dobriban and S. Wager, High-dimensional asymptotics of prediction: ridge regression and classification, *The Annals of Statistics* **46**, 247 (2018).
- [51] S. Mei and A. Montanari, The generalization error of random features regression: Precise asymptotics and the double descent curve, *Communications on Pure and Applied Mathematics* (2019).
- [52] B. Ghorbani, S. Mei, T. Misiakiewicz, and A. Montanari, Linearized two-layers neural networks in high dimension, *The Annals of Statistics* **49**, 1029 (2021).
- [53] S. Ariosto, R. Pacelli, F. Ginelli, M. Gherardi, and P. Rotondo, Universal mean-field upper bound for the generalization gap of deep neural networks, *Phys. Rev. E* **105**, 064309 (2022).
- [54] A. Shah, A. Wilson, and Z. Ghahramani, Student-t Processes as Alternatives to Gaussian Processes, in *Proceedings of the Seventeenth International Conference on Artificial Intelligence and Statistics*, Proceedings of Machine Learning Research, Vol. 33, edited by S. Kaski and J. Corander (PMLR, Reykjavik, Iceland, 2014) pp. 877–885.
- [55] J. A. Zavatone-Veth, W. L. Tong, and C. Pehlevan, Contrasting random and learned features in deep bayesian linear regression, *Phys. Rev. E* **105**, 064118 (2022).
- [56] J. Zavatone-Veth, A. Canatar, B. Ruben, and C. Pehlevan, Asymptotics of representation learning in finite bayesian neural networks, in *Advances in Neural Information Processing Systems*, Vol. 34, edited by M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan (Curran Associates, Inc., 2021) pp. 24765–24777.
- [57] J. A. Zavatone-Veth, A. Canatar, B. S. Ruben, and C. Pehlevan, Asymptotics of representation learning in finite bayesian neural networks*, *Journal of Statistical Mechanics: Theory and Experiment* **2022**, 114008 (2022).
- [58] S. Ariosto, R. Pacelli, M. Pastore, F. Ginelli, M. Gherardi, and P. Rotondo, *Supplemental Material for "A statistical mechanics framework for Bayesian deep neural networks beyond the infinite-width limit"* (2023).
- [59] B. Hanin and A. Zlokapa, Bayesian interpolation with deep linear networks (2023), [arXiv:2212.14457 \[stat.ML\]](https://arxiv.org/abs/2212.14457).
- [60] A. C. C. Coolen, M. Sheikh, A. Mozeika, F. Aguirre-López, and F. Antenucci, Replica analysis of overfitting in generalized linear regression models, *Journal of Physics A: Mathematical and Theoretical* **53**, 365001 (2020).
- [61] A. Mozeika, M. Sheikh, F. Aguirre-López, F. Antenucci, and A. C. C. Coolen, Exact results on high-dimensional linear regression via statistical physics, *Phys. Rev. E* **103**, 042142 (2021).
- [62] Y. Uchiyama, H. Oka, and A. Nono, Student’s t-process regression on the space of probability density functions, *Proceedings of the ISCTE International Symposium on Stochastic Systems Theory and its Applications* **2021**, 1 (2021).
- [63] H. Lee, E. Yun, H. Yang, and J. Lee, Scale mixtures of neural network gaussian processes, in *International Conference on Learning Representations* (2022).
- [64] L. Aitchison, Why bigger is not always better: on finite and infinite neural networks, in *Proceedings of the 37th*

- International Conference on Machine Learning*, Proceedings of Machine Learning Research, Vol. 119, edited by H. D. III and A. Singh (PMLR, 2020) pp. 156–164.
- [65] J. A. Zavattone-Veth and C. Pehlevan, Depth induces scale-averaging in overparameterized linear bayesian neural networks, in *2021 55th Asilomar Conference on Signals, Systems, and Computers* (2021) pp. 600–607.
 - [66] A. X. Yang, M. Robeyns, E. Milsom, N. Schoots, and L. Aitchison, A theory of representation learning in deep neural networks gives a deep generalisation of kernel methods (2023), [arXiv:2108.13097 \[stat.ML\]](https://arxiv.org/abs/2108.13097).
 - [67] Y. Cho and L. Saul, Kernel methods for deep learning, in *Advances in Neural Information Processing Systems*, Vol. 22, edited by Y. Bengio, D. Schuurmans, J. Lafferty, C. Williams, and A. Culotta (Curran Associates, Inc., 2009).
 - [68] B. Poole, S. Lahiri, M. Raghu, J. Sohl-Dickstein, and S. Ganguli, Exponential expressivity in deep neural networks through transient chaos, in *Advances in Neural Information Processing Systems*, Vol. 29, edited by D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett (Curran Associates, Inc., 2016).
 - [69] G. Yang and S. Schoenholz, Mean field residual networks: On the edge of chaos, in *Advances in Neural Information Processing Systems*, Vol. 30, edited by I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Curran Associates, Inc., 2017).
 - [70] B. D. Tracey and D. Wolpert, Upgrading from gaussian processes to student's-t processes, in *2018 AIAA Non-Deterministic Approaches Conference*, <https://arc.aiaa.org/doi/pdf/10.2514/6.2018-1659>.
 - [71] D. A. Roberts, S. Yaida, and B. Hanin, *The Principles of Deep Learning Theory* (Cambridge University Press, 2022) <https://deeplearningtheory.com>, [arXiv:2106.10165 \[cs.LG\]](https://arxiv.org/abs/2106.10165).
 - [72] F. Gerace, F. Krzakala, B. Loureiro, L. Stephan, and L. Zdeborová, Gaussian universality of linear classifiers with random labels in high-dimension, [arXiv preprint arXiv:2205.13303](https://arxiv.org/abs/2205.13303) (2022).
 - [73] H. Cui, F. Krzakala, and L. Zdeborová, Optimal learning of deep random networks of extensive-width (2023), [arXiv:2302.00375 \[stat.ML\]](https://arxiv.org/abs/2302.00375).
 - [74] J. Lee, S. Schoenholz, J. Pennington, B. Adlam, L. Xiao, R. Novak, and J. Sohl-Dickstein, Finite versus infinite neural networks: an empirical study, in *Advances in Neural Information Processing Systems*, Vol. 33, edited by H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin (Curran Associates, Inc., 2020) pp. 15156–15172.
 - [75] G. Pang, L. Yang, and G. E. Karniadakis, Neural-net-induced gaussian process regression for function approximation and pde solution, *Journal of Computational Physics* **384**, 270 (2019).
 - [76] R. Pacelli, [rpacelli/FC_deep_bayesian_networks](https://github.com/rpacelli/FC_deep_bayesian_networks): [FC_deep_bayesian_networks](https://github.com/rpacelli/FC_deep_bayesian_networks) (2023).

ACKNOWLEDGEMENTS

M. P. has been supported by a grant from the Simons Foundation (grant No. 454941, S. Franz). P. R. acknowledges funding from the Fellini program under the H2020-MSCA-COFUND action, Grant Agreement No. 754496, INFN (IT). The authors would like to thank Silvio Franz, Luca Molinari, Fabian Aguirre-López, Raffaella Burioni, Alessandro Vezzani, Riccardo Aiudi, Federico Bassetti, Bruno Bassetti and the Computing Sciences group at Bocconi University in Milan for discussions and suggestions.

AUTHOR CONTRIBUTION

P.R., S.A. and M.P. performed the analytical calculations, supported by F.G, M.G. and R.P. Numerical experiments, data analysis and data visualization were carried out by R.P. All the authors contributed to discussing and interpreting the results and to writing and editing the manuscript. S.A and R.P. contributed equally to the work.

COMPETING INTERESTS

The authors declare no competing interests.