

Solving overparametrized systems of random equations: I. Model and algorithms for approximate solutions

Andrea Montanari ^{*} Eliran Subag[†]

Abstract

Consider the problem of solving a system of equations $\mathbf{F}(\mathbf{x}) = \mathbf{0}$, subject to $\|\mathbf{x}\|_2 = 1$, whereby $\mathbf{F} : \mathbb{R}^d \rightarrow \mathbb{R}^n$ is a random nonlinear map. More precisely, $\mathbf{F}(\mathbf{x}) = (F_1(\mathbf{x}), \dots, F_n(\mathbf{x}))$ where the $F_i(\cdot)$'s are i.i.d. rotationally invariant Gaussian processes. We study this problem under the proportional asymptotics $n, d \rightarrow \infty$, $n/d \rightarrow \alpha \in [0, 1)$ and establish results about the existence of solutions and polynomial-time algorithms to find them.

First, we establish upper and lower bounds α_{UB} , α_{LB} on the threshold for existence of solutions. Namely, if the number of equations per variable satisfies $\alpha < \alpha_{\text{LB}}$, then the system admits exact solutions with high probability, while for $\alpha > \alpha_{\text{UB}}$, no solutions exist, even in an approximate sense.

We then analyze several algorithms to find solutions: gradient descent, Hessian descent, and a two-phase algorithm. In particular, for Hessian descent and the two-phase algorithm, we characterize their thresholds α_{HD} , α_{TP} . Namely, for $\alpha < \alpha_{\text{HD}}$ (or $\alpha < \alpha_{\text{TP}}$) the algorithm finds an approximate solution with high probability, while for $\alpha > \alpha_{\text{HD}}$ (respectively $\alpha > \alpha_{\text{TP}}$), it does not.

Finally, we compare the theoretical predictions within this model to empirical results obtained with structured systems of nonlinear equations.

1 Introduction

Given data $(\mathbf{z}_i, y_i) \in \mathbb{R}^D \times \mathbb{R}$, $i \leq n$, and a function class $\mathcal{F} \subseteq \{f : \mathbb{R}^D \rightarrow \mathbb{R}\}$, an *interpolator* of the data in \mathcal{F} is any function $f \in \mathcal{F}$ such that $f(\mathbf{z}_i) = y_i$ for all $i \leq n$. Existence and construction of interpolators are fundamental problems in pure and applied mathematics. For instance, Kirszbraun (Lipschitz) extension theorem completely characterizes the existence of interpolators when \mathcal{F} is the class of L -Lipschitz functions.

Recently, existence and algorithmic construction of interpolators has attracted new attention because of applications to machine learning [BHMM19, BLLT20, ZBH⁺21, BMR21]. In this setting, the data is normally assumed to be random $\{(\mathbf{z}_i, y_i)\}_{i \leq n} \sim_{i.i.d.} \mathbb{P}$, and the function class is parametric. Namely, there exists a parametric form $f : \mathbb{R}^D \times \mathbb{R}^d \rightarrow \mathbb{R}$ such that $\mathcal{F} = \{f(\cdot; \mathbf{x}) : \mathbf{x} \in \mathbb{S}^{d-1}\}$. Here \mathbb{S}^{d-1} is the unit sphere in d -dimensions and \mathbf{x} is a vector of parameters that parametrizes the function class. The condition $\mathbf{x} \in \mathbb{S}^{d-1}$ is a stylized version of typical constraints that are imposed (either implicitly or explicitly) on the parameters' vectors in these applications. The interpolation problem then reduces to

$$\text{Find } \mathbf{x} \in \mathbb{S}^{d-1} \text{ such that } f(\mathbf{z}_i; \mathbf{x}) = y_i \text{ for all } i \leq n. \quad (1)$$

^{*}Department of Electrical Engineering and Department of Statistics, Stanford University

[†]Department of Mathematics, Weizmann Institute of Science

Here a parametric form f is fixed, and the pairs (\mathbf{z}_i, y_i) are i.i.d. with common law \mathbb{P} .

We can abstract away the parametric form and data, thus leading to the formulation adopted in the present paper.

$$\text{Find } \mathbf{x} \in \mathbb{S}^{d-1} \text{ such that } F_i(\mathbf{x}) = 0 \text{ for all } i \leq n. \quad (2)$$

The original formulation is recovered by choosing the special case $F_i(\cdot) = y_i - f(\mathbf{z}_i; \cdot)$, i.e. by selecting a special distribution of the functions F_i . Since above we assumed the data (\mathbf{z}_i, y_i) to be i.i.d., the F_i 's will be i.i.d. random functions $F_i : \mathbb{S}^{d-1} \rightarrow \mathbb{R}$. It is convenient to define the mapping $\mathbf{F} : \mathbb{R}^d \rightarrow \mathbb{R}^n$ via $\mathbf{F}(\mathbf{x}) := (F_1(\mathbf{x}), \dots, F_n(\mathbf{x}))$.

In this paper, we study this problem under a simple distribution for the random functions F_i . Namely, we assume the F_i 's to be i.i.d. centered Gaussian processes with

$$\mathbb{E}[F_i(\mathbf{x}^1)F_j(\mathbf{x}^2)] = \delta_{ij}\xi(\langle \mathbf{x}^1, \mathbf{x}^2 \rangle). \quad (3)$$

Assuming this specific covariance structure is equivalent to the following two assumptions: (i) F_i and F_j are independent and identically distributed for $i \neq j$; (ii) F_i is invariant (in distribution) under rotations.

The fact that the covariance is positive semidefinite implies $\xi(t) = \sum_{k \geq 0} \xi_k t^k$, with $\xi_k \geq 0$ for all k [Sch42]. We will impose the random function $\mathbf{F} : \mathbf{x} \mapsto \mathbf{F}(\mathbf{x})$ to be smooth, by imposing $\xi(1 + \varepsilon) < \infty$ for some $\varepsilon > 0$ ($\sum_{k \geq 0} \xi_k k^{2\ell + \varepsilon} < \infty$ is sufficient for $\mathbf{F} \in C^\ell(\mathbb{S}^{d-1})$ by the Kolmogorov-Centsov theorem).

This model is closely related to the mixed p -spin model from spin glass theory, and indeed, each coordinate F_i is an independent copy of the so-called ‘mixed p -spin Hamiltonian’ [CS92, CS95, Sub17b]. The case of multiple equations $n \asymp d$, with energy function $\|\mathbf{F}(\mathbf{x})\|_2^2/2$ was recently introduced by Urbani [Urb23, KU23] as a model for confluent tissues. In particular, [Urb23] derives the phase diagram of this model using the replica method from spin glass physics. Our focus (and approach) will be different from the one of these works. In [Sub23], one of the authors used the second moment method to prove that, when $\xi(0) = 0$, the Hausdorff volume (or number) of solutions concentrates provided $n \leq d - 1$.

We will refer to the model defined above as the ‘Nonlinear random Equations Model’ (NEM). The random function \mathbf{F} is naturally extended to the unit ball $\mathbf{B}^d(1)$, by positing the same covariance of Eq. (3). In particular, for a fixed $\mathbf{x} \in \mathbb{R}^d$, $\mathbb{E}[\|\mathbf{F}(\mathbf{x})\|^2] = n\xi(\|\mathbf{x}\|)$. We define the set of ε -approximate solutions by

$$\text{Sol}_{n,d}(\varepsilon) := \left\{ \mathbf{x} \in \mathbb{S}^{d-1} : \|\mathbf{F}(\mathbf{x})\|_2^2 \leq n\xi(1) \cdot \varepsilon \right\}. \quad (4)$$

A number of natural questions arise within this model:

- Q1 Do exact solutions exist with high probability, i.e. is $\text{Sol}_{n,d}(0)$ non-empty?
- Q2 Do approximate solution exist, i.e. is $\text{Sol}_{n,d}(\varepsilon)$ non-empty for some $\varepsilon > 0$?
- Q3 Can we find these solutions (either exact or approximate ones) in polynomial time?

Questions Q1 and Q2 can be addressed (albeit non-rigorously) using the replica method [Urb23].

Here we will focus on the problem of efficiently finding approximate solutions in the high-dimensional limit $n, d \rightarrow \infty$, with $n/d \rightarrow \alpha \in [0, 1)$. Hence α is the asymptotic number of equations per unknown. Since $\alpha < 1$ the system of non-linear equations $\mathbf{F}(\mathbf{x})$ is overparametrized.

We present the following results, illustrated in Figure 1 which refers to the special case $\xi(q) = \xi_0 + q^p$, $p \in \{3, 7, 11, 15\}$.

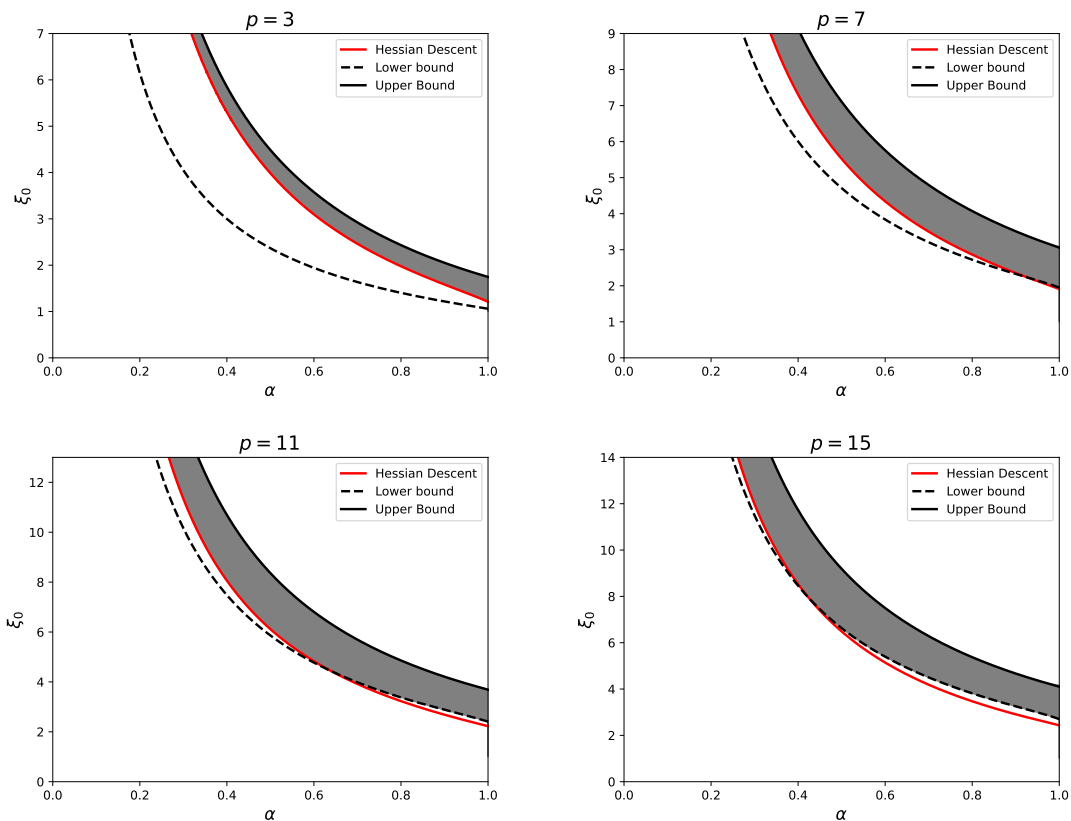


Figure 1: Phase diagram for $\xi(q) = \xi_0 + q^p$, $p \in \{3, 7, 11, 15\}$. Upper (solid) and lower (dashed) black lines are the upper and lower bounds on the threshold for existence of solutions α_{UB} and α_{LB} obtained respectively by Gaussian comparison and second moment method. Red line: The threshold α_{HD} below which the Hessian descent algorithm finds solutions with high probability.

Lower bound. We determine a lower bound threshold $\alpha_{\text{LB}} = \alpha_{\text{LB}}(\xi)$ such that, if $\alpha < \alpha_{\text{LB}}$ then with high probability the system admits exact solutions. Namely, for $\alpha < \alpha_{\text{LB}}$, $\text{Sol}_{n,d}(0) \neq \emptyset$ with probability converging to one as $n, d \rightarrow \infty$, with $n/d \rightarrow \alpha$. The lower bound is based on a second-moment calculation.

In the examples of Figure 1, this lower bound is traced as the dashed black line.

Upper bound. We determine an upper bound threshold $\alpha_{\text{UB}} = \alpha_{\text{UB}}(\xi)$ such that, if $\alpha > \alpha_{\text{UB}}$ then with high probability the system does not admit approximate solutions. Namely, for any $\alpha > \alpha_{\text{UB}}$ there exists $\varepsilon > 0$ such that $\text{Sol}_{n,d}(\varepsilon) = \emptyset$ with probability converging to one as $n, d \rightarrow \infty$, with $n/d \rightarrow \alpha$. Our upper bound is proved using Gaussian comparison inequalities.

In Figure 1, this corresponds to the upper continuous black line.

Gradient descent. We formulate the problem of finding a solution as an optimization problem, with cost function $H(\mathbf{x}) := \|\mathbf{F}(\mathbf{x})\|_2^2/2$. We analyze gradient descent with respect to this cost, and establish a lower bound $\alpha_{\text{GD}} = \alpha_{\text{GD}}(\xi)$ such that, for $\alpha < \alpha_{\text{GD}}$, gradient flow converges to an exact solution.

Our analysis of gradient descent uses techniques recently developed to analyze overparametrized neural networks [DZPS18, COB19, BMR21]. Because of its ubiquity, this approach provides a useful benchmark for the more precise analysis outlined below (which yields better guarantees in terms of overparametrization ratio α .)

Hessian descent. We introduce an Hessian descent algorithm, that instead of optimizing the first order approximation to the cost function $H(\mathbf{x})$, follows the best direction according to its Hessian. For this case, we characterize precisely its threshold $\alpha_{\text{HD}}(\xi)$. For $\alpha < \alpha_{\text{HD}}$ the algorithm converges to an ε -approximate solution, with ε arbitrarily small, with high probability (with respect to the random realization of \mathbf{F}). For $\alpha > \alpha_{\text{HD}}$ it only converges to an ε_0 -approximate solution for some $\varepsilon_0(\alpha)$ bounded away from 0. The threshold α_{HD} is reported as a red line in Figure 1.

For a broad set of random Gaussian functions \mathbf{F} , the guarantee we obtain for Hessian descent are significantly superior to the more standard ones for gradient descent.

Two-phase algorithm. In general the Hessian descent algorithm is not optimal because it disregards almost entirely gradient information. In particular, it is not optimal when $\xi'(0) \neq 0$ (and hence $\nabla H(\mathbf{0})$ is bounded away from $\mathbf{0}$). This is related to the fact that the set of solution is not centered around $\mathbf{0}$ in this case. We describe a two-phase algorithm that is well suited to these cases. In a first phase we use approximate message passing (AMP) to construct a point $\mathbf{m}_* \in \mathcal{B}^d(1)$, with $\nabla H(\mathbf{m}_*) \approx c\mathbf{m}_*$. In the second phase we run Hessian descent in the hyperplane orthogonal to \mathbf{m}_* .

We characterize the threshold α_{TP} below which the two-phase algorithm finds an approximate solution.

By comparing the second moment lower bound on the existence of solutions, and the analysis of various efficient algorithms, we identify regimes in which solutions exist with high probability, but we do not now of any polynomial-time algorithm to find them. Further, we expect the two-phase algorithm to be essentially optimal among polynomial-time algorithms. We will provide evidence towards this expectation in a forthcoming publication [MS23a].

The five set of results outlined above are presented, respectively, in Sections 2.1, 2.2, 3, 4, 5. Finally, Section 6 discusses, on the basis of numerical simulations, possible generalizations of the present work to other probabilistic models for the equations F_i .

Definitions and notations

The Gaussian random function $\mathbf{F} : \mathbb{R}^d \rightarrow \mathbb{R}^n$ defined in the previous section can be constructed explicitly by setting

$$\begin{aligned} F_i(\mathbf{x}) &:= \sum_{k \geq 0} \sqrt{\xi_k} \sum_{j_1, \dots, j_k=1}^d G_{i, j_1, \dots, j_k}^{(k)} x_{j_1} \cdots x_{j_k} \\ &= \sqrt{\xi_0} G_i^{(0)} + \sqrt{\xi_1} \sum_{j=1}^d G_{i, j}^{(1)} x_j + \dots, \end{aligned} \quad (5)$$

where the random coefficients $(\mathbf{G}^{(k)})_{k \geq 0} := (G_{i, j_1, \dots, j_k}^{(k)})_{k \geq 0, i \leq n, j_1, \dots, j_k \leq d} \sim i.i.d. \mathcal{N}(0, 1)$ are i.i.d. Gaussian random variables. The relation with the previous definition is given by the coefficients ξ_k .

Given a symmetric matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$, we denote by $\lambda_1(\mathbf{A}) \geq \lambda_2(\mathbf{A}) \geq \dots \geq \lambda_n(\mathbf{A})$ its eigenvalues in decreasing order. For a general matrix $\mathbf{M} \in \mathbb{R}^{m \times n}$, $m \leq n$, $\sigma_1(\mathbf{M}) \geq \sigma_2(\mathbf{M}) \geq \dots \geq \sigma_m(\mathbf{M}) \geq 0$ denote its singular values and $\|\mathbf{M}\|_{\text{op}}$ its operator norm.

We let $\mathbf{B}^d(r) := \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\| \leq r\}$ denote the ball of radius r in \mathbb{R}^d , with $\mathbf{B}^d := \mathbf{B}^d(1)$.

Throughout, we will write $\mathbf{W} \sim \text{GOE}(N)$ if $\mathbf{W} = \mathbf{W}^\top$ and $(W_{ij})_{i \leq j \leq N}$ are independent with $W_{ii} \sim \mathcal{N}(0, 2)$, $W_{ij} \sim \mathcal{N}(0, 1)$ for $i < j$. We write $\mathbf{Z} \sim \text{GOE}(M, N)$ if $(Z_{ij})_{i \leq M, j \leq N}$ are independent with $Z_{ij} \sim \mathcal{N}(0, 1)$.

2 Existence of solutions

In this section we derive upper and lower bounds for the maximum number of equations per variable α , below which solutions do exist. The lower bound α_{LB} (Section 2.1) is based on the second moment method and guarantees the existence of exact solutions for $\alpha < \alpha_{\text{LB}}$. The upper bound α_{UB} (Section 2.2) is based on Gaussian comparison inequalities and implies instead that even approximate solutions do not exist for $\alpha > \alpha_{\text{UB}}$.

It will be useful to keep in mind the following elementary fact that often allows us to focus on $\mathbb{E} \min_{\mathbf{x} \in \mathbb{S}^{d-1}} \|\mathbf{F}(\mathbf{x})\|_2$.

Remark 2.1. For each $i \leq n$, and each $\mathbf{x} \in \mathbb{S}^{d-1}$ the function $(G_{i, j_1, \dots, j_k}^{(k)}) \mapsto F_i(\mathbf{x})$ is Lipschitz continuous (in ℓ_2), with Lipschitz norm bounded by $\sqrt{\xi(1)}$ (recall that $\xi(1) = \sum_k \xi_k$). As a consequence $(G_{i, j_1, \dots, j_k}^{(k)}) \mapsto \min_{\mathbf{x} \in \mathbb{S}^{d-1}} \|\mathbf{F}(\mathbf{x})\|_2 / \sqrt{n}$ is also Lipschitz hence it concentrates:

$$\mathbb{P}\left(\left| \min_{\mathbf{x} \in \mathbb{S}^{d-1}} \|\mathbf{F}(\mathbf{x})\|_2 - \mathbb{E} \min_{\mathbf{x} \in \mathbb{S}^{d-1}} \|\mathbf{F}(\mathbf{x})\|_2 \right| \geq t\right) \leq 2e^{-t^2/(cn)}. \quad (6)$$

2.1 Lower bound

Theorem 1. Assume $\xi'(0) = \xi''(0) = 0$ and define $\Psi(\cdot; \alpha, \xi) : [0, 1] \rightarrow \mathbb{R}$ via

$$\Psi(r; \alpha, \xi) := \frac{1}{2} \log(1 - r^2) - \frac{\alpha}{2} \log\left(1 - \left(\frac{\xi(r) - \xi(0)}{\xi(1) - \xi(0)}\right)^2\right) - \frac{\alpha \xi(0)}{\xi(1) + \xi(r) - 2\xi(0)} + \frac{\alpha \xi(0)}{\xi(1) - \xi(0)}. \quad (7)$$

Further define

$$\alpha_{\text{LB}}(\xi) := \inf \left\{ \alpha \geq 0 : \sup_{r \in [0, 1]} \Psi(r; \alpha, \xi) > 0 \right\}. \quad (8)$$

If $\alpha < \alpha_{\text{LB}}(\xi)$ then $\text{Sol}_{n,d}(0) \neq \emptyset$ with probability converging to one as $n, d \rightarrow \infty$ with $n/d \rightarrow \alpha \in (0, 1)$.

Proof technique. The proof of this result is based on the second moment method. Namely, let $\mathbf{F}_{>0}(\mathbf{x}) := \mathbf{F}(\mathbf{x}) - \mathbf{F}(\mathbf{0})$, and note that $\mathbf{F}_{>0}(\cdot)$ is independent of $\mathbf{F}(\mathbf{0}) \sim \mathbf{N}(\mathbf{0}, \xi(0)\mathbf{I}_n)$ and distributed as the original process, with $\xi(q)$ replaced by $\xi_{>0}(q) = \xi(q) - \xi(0)$.

Consider the modified set of solutions

$$\text{Sol}_{n,d}(\mathbf{u}; \varepsilon) := \left\{ \mathbf{x} \in \mathbb{S}^{d-1} : \|\mathbf{F}_{>0}(\mathbf{x}) - \mathbf{u}\|_2^2 \leq n\xi(1) \cdot \varepsilon \right\}. \quad (9)$$

We prove that, for $\alpha < \alpha_{\text{LB}}$ and $\delta_d \geq 0$ any deterministic sequence such that $\delta_d \rightarrow 0$, we have

$$\lim_{n,d \rightarrow \infty} \sup_{\|\mathbf{u}\|_2/\sqrt{n} \in [\xi_0^{1/2} - \delta_d, \xi_0^{1/2} + \delta_d]} \mathbb{P}(\text{Sol}_{n,d}(\mathbf{u}; 0) = \emptyset) = 0. \quad (10)$$

This of course implies the claim of the theorem since $\|\mathbf{F}(\mathbf{0})\|_2/\sqrt{n}$ concentrates around $\xi_0^{1/2}$.

In order to upper bound the probability of $\text{Sol}_{n,d}(\mathbf{u}; 0) = \emptyset$, as per Eq. (10), we introduce the random variable

$$\mathbf{V}(\mathbf{u}) := \text{Vol}_{d-n-1}(\text{Sol}_{n,d}(\mathbf{u}; 0)), \quad (11)$$

where Vol_i is the Hausdorff measure of dimension i , or the counting measure when $i = 0$.

We use Kac-Rice formula [Kac43, Ric45] to compute the first two moments of $\mathbf{V}(\mathbf{u})$ and obtain, by a second moment argument, a sufficient condition for it to concentrate around the mean. For $\mathbf{u} = \mathbf{0}$ this calculation was carried out in [Sub23]. Here we generalize the argument to the case $\mathbf{u} \neq \mathbf{0}$. We present this calculation in Appendix B. \square

In the case $\xi(t) = \xi_0 + t^p$ (writing, with an abuse of notation, $\alpha_{\text{LB}}(\xi_0, p)$ instead of $\alpha_{\text{LB}}(\xi)$), we have the large- p asymptotics (cf. Appendix C)

$$\alpha_{\text{LB}}(\xi_0, p) = \frac{\log p}{\xi_0} \cdot (1 + o_p(1)). \quad (12)$$

(This should be interpreted as $\alpha_{\text{LB}}(\xi_0 = \gamma_0 \log p, p) = \gamma_0^{-1}(1 + o_p(1))$ for any $\gamma_0 > 1$.)

2.2 Upper bound

We will prove two upper bounds $\alpha_{\text{UB}}^{(1)}$ and $\alpha_{\text{UB}}^{(2)}$ on the satisfiability threshold. The first one holds for general ξ and has the advantage of being quite simple, while we prove the second only for the ‘pure’ model $\xi(q) = \xi_0 + q^p$.

Define

$$E_\star(\xi) := \lim_{d \rightarrow \infty} \frac{1}{\sqrt{d}} \mathbb{E} \max_{\mathbf{x} \in \mathbb{S}^{d-1}} F_1(\mathbf{x}). \quad (13)$$

The limit is known to exist and is given by the Parisi formula [AC17, JT17], which we recall below.

Given $\gamma : [0, 1] \rightarrow \mathbb{R}_{\geq 0}$ and $L \geq \int_0^1 \gamma(s) ds$, we define

$$\mathbb{P}(\gamma, L) = \frac{1}{2} \int_0^1 \left(\xi''(t) \Gamma(t) + \frac{1}{\Gamma(t)} \right) dt, \quad (14)$$

$$\Gamma(t) := L - \int_0^t \gamma(s) ds. \quad (15)$$

Equivalently, we can view \mathbb{P} as a function of $\Gamma : [0, 1] \rightarrow \mathbb{R}_{\geq 0}$ which is continuous and non-increasing. We then have

$$E_*(\xi) := \inf \left\{ \mathbb{P}(\gamma, L) : \gamma \text{ non-decreasing } L \geq \int_0^1 \gamma(t) dt \right\}. \quad (16)$$

Proposition 2.1. *Assume that $\xi(0) > 0$ and define $\xi_{>0}(t) := \xi(t) - \xi(0)$ and*

$$\alpha_{\text{UB}}^{(1)}(\xi) := \frac{E_*(\xi_{>0})^2}{\xi(0)}. \quad (17)$$

If $\alpha > \alpha_{\text{UB}}^{(1)}$ then there exists $\varepsilon > 0$ such that, for $n/d \rightarrow \alpha$:

$$\lim_{n,d \rightarrow \infty} \mathbb{P}(\text{Sol}_{n,d}(\varepsilon) = \emptyset) = 1.$$

In fact this holds for any $\varepsilon < \varepsilon_0 := (\xi(0)/\xi(1))(1 - \sqrt{\alpha_{\text{UB}}^{(1)}/\alpha})_+$.

Proof. As in the previous section, we write $\mathbf{F}(\mathbf{x}) = \mathbf{F}(\mathbf{0}) + \mathbf{F}_{>0}(\mathbf{x})$ and note that the two summands are independent. By Remark 2.1 (applied to $\mathbf{F}_{>0}$) and Eq. (13)

$$\text{p-lim}_{n,d \rightarrow \infty} \frac{1}{\sqrt{d}} \max_{\mathbf{x} \in \mathbb{S}^{d-1}} \left| \frac{\langle \mathbf{F}(\mathbf{0}), \mathbf{F}_{>0}(\mathbf{x}) \rangle}{\|\mathbf{F}(\mathbf{0})\|_2} \right| = E_*(\xi). \quad (18)$$

Therefore,

$$\frac{1}{\sqrt{n}} \min_{\mathbf{x} \in \mathbb{S}^{d-1}} \|\mathbf{F}(\mathbf{x})\|_2 \geq \frac{1}{\sqrt{n}} \|\mathbf{F}(\mathbf{0})\|_2 - \frac{1}{\sqrt{n}} \max_{\mathbf{x} \in \mathbb{S}^{d-1}} \left| \frac{\langle \mathbf{F}(\mathbf{0}), \mathbf{F}_{>0}(\mathbf{x}) \rangle}{\|\mathbf{F}(\mathbf{0})\|_2} \right| \quad (19)$$

$$\geq \sqrt{\xi(0)} - \frac{1}{\sqrt{\alpha}} E_*(\xi) - o_P(1), \quad (20)$$

where $o_P(1)$ denotes a term which converges in probability to 0 as $n, d \rightarrow \infty$, which implies the claim. \square

We next obtain a bound for pure models $\xi(t) = \xi_0 + t^p$. For $E \geq 2\sqrt{(p-1)/p}$, define

$$\Theta_p(E) := \frac{1}{2} \log(p-1) - \frac{p-2}{p-1} \frac{E^2}{4} - \sqrt{\frac{p}{p-1}} \frac{E}{4} \sqrt{\frac{p}{p-1} E^2 - 4} + \log \left(\sqrt{\frac{p}{p-1}} \frac{E^2}{4} - 1 + \sqrt{\frac{p}{p-1}} \frac{E}{2} \right).$$

The meaning of this function is established in [CS95] (at a heuristic level) and in [ABA13, Auf13, Sub17a, SZ21] (rigorously), and it is worth reminding it here. Consider the process $F_{>0,1}(\cdot)$ (the first coordinate of $\mathbf{F}_{>0}(\mathbf{x}) = \mathbf{F}(\mathbf{x}) - \mathbf{F}(\mathbf{0})$), which is a Gaussian process with covariance $\mathbb{E}[F_{>0,1}(\mathbf{x}^1)F_{>0,1}(\mathbf{x}^2)] = \langle \mathbf{x}^1, \mathbf{x}^2 \rangle^p$. Then, the number of local maxima \mathbf{x} of this process with $F_{>0,1}(\mathbf{x}) \approx E\sqrt{d}$ concentrates around $\exp(d\Theta_p(E))$. In particular, the function $E \mapsto \Theta_p(E)$ is monotone decreasing for $E > 2\sqrt{(p-1)/p}$ and vanishes at $E = E_*(\xi(t) = t^p)$. With an abuse of notation, we will write $E_*(p) = E_*(\xi(t) = t^p)$.

For $c \geq 0$, define

$$\varphi_1(c, p) = \sup_{t,s \geq 0} \left\{ c(E_*(p) + t + s) - \frac{1}{2} \frac{s^2}{\xi(0)} + \Theta_p(E_*(p) + t) \right\},$$

$$\varphi_2(c, \alpha) = \sup_{0 < t \leq 1} \left\{ -c\sqrt{\alpha\xi(1)}t - \alpha \frac{t^2 - 1}{2} + \alpha \log t \right\}$$

$$= -c\sqrt{\alpha\xi(1)}t_* - \alpha\frac{t_*^2 - 1}{2} + \alpha\log t_*,$$

where $t_* = \frac{1}{2}(\sqrt{c^2\xi(1)/\alpha + 4} - c\sqrt{\xi(1)}/\sqrt{\alpha})$. Note that $\varphi_2(c, \alpha)$ is decreasing in $\alpha > 0$.

Theorem 2. *Define*

$$\alpha_{\text{UB}}^{(2)}(\xi_0, p) := \inf \left\{ \alpha \geq 0 : \inf_{c>0} [\varphi_1(c, p) + \varphi_2(c, \alpha) - \frac{1}{2}c^2(1 + \xi_0)] < 0 \right\}.$$

For the pure model $\xi(q) = \xi_0 + q^p$ with $p \geq 2$, if $\alpha > \alpha_{\text{UB}}^{(2)}(\xi_0, p)$, then there exists $\varepsilon > 0$ such that

$$\lim_{n,d \rightarrow \infty} \mathbb{P}(\text{Sol}_{n,d}(\varepsilon) = \emptyset) = 1.$$

The proof of this result is presented in Appendix D.

3 Gradient descent

In this section we study the classical projected gradient descent algorithm on the sphere \mathbb{S}^{d-1} . We use the cost function $H(\mathbf{x}) = \|\mathbf{F}(\mathbf{x})\|_2^2/2$. This algorithm can be thought of as a discretization of the gradient flow dynamics

$$\dot{\mathbf{x}}(t) = -\text{P}_{\mathbb{T}, \mathbf{x}(t)} \nabla H(\mathbf{x}(t)), \quad (21)$$

where $\text{P}_{\mathbb{T}, \mathbf{x}}$ is the projection onto the tangent space to the sphere of radius $\|\mathbf{x}\|_2$ at \mathbf{x} , namely¹

$$\text{P}_{\mathbb{T}, \mathbf{x}} := \mathbf{I}_d - \frac{\mathbf{x}\mathbf{x}^\top}{\|\mathbf{x}\|_2^2}. \quad (22)$$

Appendix E states results about gradient flow as well.

In gradient descent, at iteration k we take a step along the direction $-\text{P}_{\mathbb{T}, \mathbf{x}^k} \nabla H(\mathbf{x}^k)$, with stepsize η , and then project back on the sphere \mathbb{S}^{d-1} . The algorithm, is defined by the pseudocode of Algorithm 1.

Algorithm 1: Projected Gradient Descent

Data: Couplings $\{\mathbf{G}^{(k)}\}_{k \geq 0}$, stepsize η , number of iterations K

Result: approximate optimizer $\mathbf{x} \in \mathbb{S}^{d-1}$

Initialize $\mathbf{x}^0 \sim \text{Unif}(\mathbb{S}^{d-1})$;

for $k \in \{0, \dots, K-1\}$ **do**

$\mathbf{z}^{k+1} = \mathbf{x}^k - \eta \text{P}_{\mathbb{T}, \mathbf{x}^k} \nabla H(\mathbf{x}^k)$;

$\mathbf{x}^{k+1} = \mathbf{z}^{k+1} / \|\mathbf{z}^{k+1}\|_2$;

end

return \mathbf{x}^K

Our analysis of this algorithm is based on a technique that became recently popular to analyze overparametrized neural networks in the so called ‘neural tangent’ or ‘lazy’ regime. A few pointers to this literature include [DZPS18, COB19, OS20, ADH⁺19, AZLL19, BMR21]. We deliberately

¹In Eq. (21) we have $\|\mathbf{x}\|_2 = 1$ but we define $\text{P}_{\mathbb{T}, \mathbf{x}}$ more generally for future reference. By convention, we set $\mathbf{v}/\|\mathbf{v}\|_2 = \mathbf{0}$ if $\mathbf{v} = \mathbf{0}$.

follow this type of analysis as it provides a useful comparison point for the sharper techniques in the next sections.

Below, for $\mathbf{x} \in \mathbb{S}^{d-1}$, $\mathbf{U}_{\mathbf{x}}$ is an orthonormal basis for $\mathbb{T}_{\mathbf{x}}$ the orthogonal space to \mathbf{x} . Further, for $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d$, define $\mathbf{U}_{\mathbf{x}_1, \mathbf{x}_2} := \mathbf{R}_{\mathbf{x}_1, \mathbf{x}_2} \mathbf{U}_{\mathbf{x}_1}$, where $\mathbf{R}_{\mathbf{x}_1, \mathbf{x}_2}$ is the rotation that keeps unchanged the space orthogonal to $\mathbf{x}_1, \mathbf{x}_2$, and maps \mathbf{x}_1 to \mathbf{x}_2 . Finally, let $\mathbf{DF}(\mathbf{x}) \in \mathbb{R}^{n \times d}$ be the Jacobian of \mathbf{F} at \mathbf{x} .

Lemma 3.1. *Consider the Gradient Descent Algorithm 1, and fix $\varepsilon_0 \in (0, 1/2)$. Let $\lambda_0 := \sigma_{\min}(\mathbf{DF}(\mathbf{x}^0)|_{\mathbb{T}, \mathbf{x}^0})$, and J_n, L_n, M_n be given by*

$$J_n := \sup_{\mathbf{x}_1 \neq \mathbf{x}_2 \in \mathbb{B}^d(1+\varepsilon_0)} \frac{\|\mathbf{DF}(\mathbf{x}_1) - \mathbf{DF}(\mathbf{x}_2)\|_{\text{op}}}{\|\mathbf{x}_1 - \mathbf{x}_2\|_2}, \quad (23)$$

$$L_n := \sup_{\mathbf{x}_1 \neq \mathbf{x}_2 \in \Omega} \frac{\|\mathbf{DF}(\mathbf{x}_1)\mathbf{U}_{\mathbf{x}_1} - \mathbf{DF}(\mathbf{x}_2)\mathbf{U}_{\mathbf{x}_1, \mathbf{x}_2}\|_{\text{op}}}{\|\mathbf{x}_1 - \mathbf{x}_2\|_2}, \quad (24)$$

$$M_n := \sup_{\mathbf{x} \in \mathbb{B}^d(1+\varepsilon_0)} \|\mathbf{DF}(\mathbf{x})\|_{\text{op}}, \quad (25)$$

and assume they satisfy

$$L_n \|\mathbf{F}(\mathbf{x}^0)\|_2 < \frac{\lambda_0^2}{16}, \quad (26)$$

Further assume the step size η to be such that

$$\eta \leq \min \left\{ \frac{\varepsilon_0}{\max_{\mathbf{x} \in \mathbb{S}^{d-1}} \|\mathbf{P}_{\mathbb{T}} \mathbf{DF}(\mathbf{x})^{\top} \mathbf{F}(\mathbf{x})\|_2}; \frac{1}{M_n^2}; \frac{1}{10\sqrt{n}(M_n + J_n) \max_{\mathbf{x} \in \mathbb{S}^{d-1}} \|\mathbf{F}(\mathbf{x})\|_2} \right\}. \quad (27)$$

Then, for all $k \geq 0$

$$\|\mathbf{F}(\mathbf{x}^k)\|_2^2 \leq \|\mathbf{F}(\mathbf{x}^0)\|_2^2 e^{-\lambda_0^2 \eta k / 8}. \quad (28)$$

The proof of this statement is presented in Appendix E.

By evaluating the conditions in this statement, we obtain the following.

Theorem 3. *Consider the Gradient Descent Algorithm 1, with \mathbf{F} the Gaussian process defined in Section 1, and initialization \mathbf{x}^0 independent of \mathbf{F} . Assume $n, d \rightarrow \infty$ with $n/d \rightarrow \alpha \in [0, 1)$. Define, for c_0 a sufficiently small absolute constant and*

$$\underline{\alpha}_{\text{GD}}(\xi) := \frac{c_0 \xi'(1)^2}{\xi''(1) \xi(1) (\log(\xi'''(1)/\xi''(1)) \vee 1)}. \quad (29)$$

If $\alpha < \underline{\alpha}_{\text{GD}}(\xi)$, and $\eta < 1/(C_1 d)$ with C_1 a suitable constant depending on ξ , then the following happens with high probability. For all $k \geq 1$,

$$\|\mathbf{F}(\mathbf{x}^k)\|_2^2 \leq 2n\xi(1) \exp\left(-\frac{\xi'(1)\eta}{16} (\sqrt{d} - \sqrt{n})^2 \cdot k\right). \quad (30)$$

Again, we refer to Appendix E for a proof.

Considering our running example $\xi(t) = \xi_0 + t^p$, $\underline{\alpha}_{\text{GD}}(\xi)$ is equivalent, up to constants, to

$$\underline{\alpha}_{\text{GD}}(\xi_0, p) := \frac{c_1}{\xi_0 \log p}. \quad (31)$$

It is instructive to compare this result with the lower bound on the threshold for existence of solutions, see Eq. (12). Roughly speaking, for large p , and $1/\log p \alpha \xi_0 \ll \log p$ we know that solutions exist with high probability, but the approach developed here does not guarantee that we can find them. We will see that this gap shrinks using the methods in next sections.

4 Hessian descent

We next consider the Hessian descent algorithm first introduced in [Sub21] to optimize the spherical spin glass Hamiltonian. We do not expect this algorithm to be optimal in general, and in particular not so unless $\xi'(0) = 0$ (i.e. $\mathbf{F}(\mathbf{x})$ does not contain terms linear in \mathbf{x}). In the next section we describe an extension that covers this case as well, but we think it is useful to first introduce Hessian descent in a simpler setting.

The Hessian descent algorithm presents two important differences with respect to gradient descent. First, we extend the objective function $H(\mathbf{x})$ to the unit ball $\mathbf{B}^d(1)$: iterates are initialized at $\mathbf{x}^0 = \mathbf{0}$, and after k iterations we have $\|\mathbf{x}^k\|_2^2 = k\delta$. Second, we disregard gradient information at \mathbf{x}^k and instead optimize the Hessian contribution to the cost function $H(\mathbf{x})$.

The pseudocode for Hessian descent is specified in Algorithm 2 below.

Algorithm 2: Hessian Descent

Data: Couplings $\{\mathbf{G}^{(k)}\}_{k \geq 0}$, stepsize δ , with $1/\delta \in \mathbb{N}$

Result: Approximate optimizer $\mathbf{x}^{\text{HD}} \in \mathbb{S}^{d-1}$

Initialize $\mathbf{x}^0 = \mathbf{0}$, $\mathbf{x}^1 \sim \sqrt{\delta} \cdot \text{Unif}(\mathbb{S}^{d-1})$;

for $k \in \{1, \dots, K := 1/\delta - 1\}$ **do**

 Compute $\mathbf{v} = \mathbf{v}(\mathbf{x}^k) \in \mathbb{T}_{\mathbf{x}^k}$ such that $\|\mathbf{v}\|_2 = 1$ and

$$\langle \mathbf{v}, \nabla^2 H(\mathbf{x}^k) \mathbf{v} \rangle \leq \lambda_{\min}(\nabla^2 H(\mathbf{x}^k)|_{\mathbb{T}_{\mathbf{x}^k}}) + d\delta; \quad (32)$$

$s_k := \text{sign}(\langle \mathbf{v}(\mathbf{x}^k), \nabla H(\mathbf{x}^k) \rangle)$;

$\mathbf{x}^{k+1} = \mathbf{x}^k - s_k \sqrt{\delta} \mathbf{v}(\mathbf{x}^k)$;

end

return $\mathbf{x}^{\text{HD}} = \mathbf{x}^K$;

The next theorem bounds the value achieved by the Hessian descent algorithm. (We formally state an upper bound on the value achieved, but we expect the bound to be tight.)

Theorem 4. For $\alpha \in (0, 1)$, $a, b \in \mathbb{R}_{\geq 0}$, define

$$Q(m; \alpha, a, b) := -\frac{1}{m} + \frac{\alpha b}{1 + bm} - a^2 m, \quad (33)$$

$$z_*(\alpha, a, b) := -\sup_{m > 0} Q(m; \alpha, a, b). \quad (34)$$

Let $u(\cdot; \alpha, \xi) : [0, 1] \rightarrow \mathbb{R}$ be the unique solution of the ordinary differential equation

$$\frac{du}{dt}(t) = -\frac{1}{2\alpha} z_*(\alpha; \sqrt{2\alpha u(t)\xi''(t)}, \xi'(t)), \quad u(0) = \frac{1}{2}\xi(0). \quad (35)$$

Then there exists constants $C_0 = C_0(\alpha, \xi)$, $\delta_0 = \delta_0(\alpha, \xi) > 0$ depending uniquely on α, ξ such that the Hessian descent algorithm, with stepsize parameter $\delta \leq \delta_0$, outputs $\mathbf{x}^{\text{HD}} \in \mathbb{S}^{d-1}$ such that, with probability converging to one as $n, d \rightarrow \infty$ ($n/d \rightarrow \alpha$)

$$\frac{1}{2n} \|\mathbf{F}(\mathbf{x}^{\text{HD}})\|_2^2 \leq u(1; \alpha, \xi) + C_0 \delta. \quad (36)$$

Further, the algorithm has complexity at most $(C_0 \chi_{n,d} / \delta) \log(1/\delta)$, where $\chi_{n,d}$ is the complexity of a single matrix vector multiplication by $\nabla^2 H(\mathbf{x})$ at a query point $\mathbf{x} \in \mathbf{B}^d(1)$.

The proof of this theorem is presented in Appendix F.

Remark 4.1. Theorem 4 implies that the Hessian descent algorithm achieves an approximate solution (in the sense that $\text{p-lim}_{n,d \rightarrow \infty} \|\mathbf{F}(\mathbf{x}^{\text{HD}})\|_2^2/n = 0$) provided $\alpha < \alpha_{\text{HD}}(\xi)$, where

$$\alpha_{\text{HD}}(\xi) := \inf \{ \alpha \geq 0 : u(1; \alpha, \xi) = 0 \}. \quad (37)$$

Remark 4.2. The complexity $\chi_{n,d}$ of matrix-vector multiplication by $\nabla^2 H(\mathbf{x})$ at a query point \mathbf{x} depends on the details of the computation model in use. In a model in which sums and products in \mathbb{R} can be carried out in $O(1)$ time, if $\xi_k = 0$ for all $k > k_{\text{max}}$ (i.e. \mathbf{F} is a polynomial), then $\chi_{n,d} = O(d^{k_{\text{max}}})$.

If $\xi_k \neq 0$ for infinitely many k (i.e. \mathbf{F} is not a polynomial), we can truncate it at a large level k_{max} , and hence approximate matrix-vector multiplication by using the truncated Hessian. Similarly, if only integer operations are allowed, we can use finite-precision approximations of the entries of $\nabla^2 H(\mathbf{x})$. It is easy to show that these modifications do not change the claim of Theorem 4.

It is clear that establishing Theorem 4 requires to analyze the eigen-structure of the Hessian $\nabla^2 H(\mathbf{x})$ at a point $\mathbf{x} \in \mathbb{R}^d$. In particular, the following simple lemma characterizes the joint distribution of $H(\mathbf{x})$, $\nabla H(\mathbf{x})$, $\nabla^2 H(\mathbf{x})$. In synthesis, the restriction of the Hessian on the tangent space is a linear combination of a Wigner and an independent Wishart matrix, with coefficients that depend on the energy level.

Lemma 4.1. For a fixed $\mathbf{x} \in \mathbb{R}^d$ with $\|\mathbf{x}\|_2^2 = q$, we have $\mathbf{F}(\mathbf{x}) = \sqrt{\xi(q)} \mathbf{g}$, $\mathbf{D}\mathbf{F}(\mathbf{x})\mathbf{U}_{\mathbf{x}} = \sqrt{\xi'(q)} \mathbf{Z}$, $\mathbf{U}_{\mathbf{x}}^{\top} \nabla^2 F_{\ell}(\mathbf{x}) \mathbf{U}_{\mathbf{x}} = \sqrt{\xi''(q)} \mathbf{W}_{\ell}$, where \mathbf{g} , $(\mathbf{W}_{\ell})_{\ell \leq n}$, \mathbf{Z} are mutually independent with

$$\mathbf{g} \sim \mathbf{N}(0, \mathbf{I}_n), \quad \mathbf{W}_{\ell} \sim \text{GOE}(d-1), \quad \mathbf{Z} \sim \text{GOE}(n, d-1). \quad (38)$$

As a consequence, letting $\mathcal{H}(\mathbf{x}) := \mathbf{U}_{\mathbf{x}}^{\top} \nabla^2 H(\mathbf{x}) \mathbf{U}_{\mathbf{x}}$ be the restriction of the Hessian to the tangent space, we have

$$\mathcal{H}(\mathbf{x}) = \sqrt{\xi(q)\xi''(q)} \|\mathbf{g}\|_2 \mathbf{W} + \xi'(q) \mathbf{Z}^{\top} \mathbf{Z}, \quad (39)$$

$$H(\mathbf{x}) = \frac{1}{2} \xi(q) \|\mathbf{g}\|_2^2, \quad (40)$$

where $(\mathbf{g}, \mathbf{W}, \mathbf{Z}) \sim \mathbf{N}(0, \mathbf{I}_n) \otimes \text{GOE}(d-1) \otimes \text{GOE}(n, d-1)$.

This lemma suggests to estimate the energy decrease at step k of Hessian descent, by computing the minimum eigenvalue of $\lambda_{\min}(\mathcal{H}(\mathbf{x}))$ at a point \mathbf{x} with $\|\mathbf{x}\|_2^2 = k\delta$ and $H(\mathbf{x})/n = u$. If \mathbf{x} is a point independent of the Gaussian process $\mathbf{F}(\cdot)$, $\lambda_{\min}(\mathcal{H}(\mathbf{x}))$ turns out to concentrate around $-z_{\#}(t = k\delta)d$ where $z_{\#}(t) := z_*(\alpha; \sqrt{2\alpha u \xi''(t)}, \xi'(t))$. By summing this energy decrement over $k \in \{1, \dots, \lfloor 1/\delta \rfloor - 1\}$ and letting δ be small, this calculation yields the value $u(1; \alpha, \xi)$ of Theorem 4.

At first sight, such a derivation might seem incorrect because \mathbf{x}^k is not independent of $\mathcal{H}(\mathbf{x}^k)$. However, the fast decay of the probability of upper deviations of the minimum eigenvalue allows to establish the claim nevertheless.

Remark 4.3 (Algorithm complexity). Recall that $\lambda_1(\mathbf{A}) \geq \lambda_2(\mathbf{A}) \geq \dots$ denote the eigenvalues of \mathbf{A} in decreasing order. Then the proof of Theorem 4 outlined above (together with the fact that the matrix $\mathcal{H}(\mathbf{x})$ of Eq. (39) has with high probability $c_0 \varepsilon^{3/2} \cdot d$ eigenvalues in $[-z_{\#}d, (-z_{\#} + \varepsilon)d]$) implies that condition (32) can be replaced by

$$\langle \mathbf{v}, \nabla^2 H(\mathbf{x}^k) \mathbf{v} \rangle \leq (-z_{\#}(k\delta) + c_* \delta) d, \quad (41)$$

which is essentially equivalent. Such a vector \mathbf{v} can be computed with $O(\log(1/\delta))$ matrix vector multiplications by $\nabla^2 H(\mathbf{x}^k)|_{\mathbb{T}, \mathbf{x}^k}$, via Chebyshev approximation [Saa11]. Also notice that computing $\lambda_{\min}(\nabla^2 H(\mathbf{x}^k)|_{\mathbb{T}, \mathbf{x}^k})$ is not needed.

Summing up, the total complexity of one step Hessian descent algorithm is $O(\log(1/\delta))$ matrix vector multiplications by the Hessian.

The formula $u(1; \alpha, \xi)$ for the energy achieved by Hessian descent, cf. Theorem 4, is somewhat implicit. The next corollary provides user-friendly upper and lower bounds. Its proof follows immediately from Theorem 4 using the bounds on z_* given in Lemma F.1, part 3.

Corollary 4.2. *Define*

$$u_{\text{LB}}(\alpha, \xi) := \frac{1}{2} \left(\sqrt{\xi(0)} - \sqrt{\frac{1}{\alpha}} \int_0^1 \sqrt{\xi''(s)} \right)_+^2, \quad (42)$$

$$u_{\text{UB}}(\alpha, \xi) := \frac{1}{2} \left(\sqrt{\xi(0)} - \sqrt{\frac{1-\alpha}{\alpha}} \int_0^1 \sqrt{\xi''(s)} \right)_+^2. \quad (43)$$

Then the energy achieved by Hessian descent satisfies

$$u_{\text{LB}}(\alpha, \xi) \leq \text{p-lim}_{d, n \rightarrow \infty} \frac{1}{2n} \|\mathbf{F}(\mathbf{x}^{\text{HD}})\|_2^2 \leq u_{\text{UB}}(\alpha, \xi). \quad (44)$$

In particular, the critical point of the algorithm satisfies

$$\frac{A(\xi)}{1 + A(\xi)} \leq \alpha_{\text{HD}}(\xi) \leq A(\xi), \quad A(\xi) := \left(\int_0^1 \sqrt{\frac{\xi''(t)}{\xi(0)}} dt \right)^2. \quad (45)$$

For our running example of a pure model $\xi(t) = \xi_0 + t^p$ (denoting the corresponding threshold by $\alpha^{\text{HD}}(\xi_0, p)$), the last bounds yields

$$\frac{4(p-1)}{p\xi_0 + 4(p-1)} \leq \alpha_{\text{HD}}(\xi_0, p) \leq \frac{4(p-1)}{p\xi_0}. \quad (46)$$

In particular, for large p , we obtain $\alpha_{\text{HD}}(\xi_0, p) \geq (4/(4 + \xi_0))(1 + o_p(1))$. This is substantially better than the guarantee $\underline{\alpha}_{\text{GD}}(\xi_0, p) \asymp 1/(\xi_0 \log p)$ that we obtained in the previous section for gradient descent, but still far from the maximum value of α in which we know that solutions exist, $\alpha < \alpha_{\text{LB}}(\xi_0, p) \asymp (\log p)/\xi_0$, cf. Eq. (12).

5 Two-phase algorithm

By analogy with the case of the spherical p -spin glass model [Sub21], we expect Hessian descent to be optimal (among algorithms with comparable complexity) if $\xi'(0) = 0$. On the other hand, it is easy to construct examples in which it is suboptimal if $\xi'(0) \neq 0$, as discussed in the next remark.

Remark 5.1. Consider the case $\xi(q) = \xi_0 + q$. Namely \mathbf{F} is a linear function $\mathbf{F}(\mathbf{x}) = \sqrt{\xi_0} \mathbf{h} + \mathbf{G}\mathbf{x}$ for $\mathbf{h} \in \mathbb{R}^n$, $\mathbf{G} \in \mathbb{R}^{n \times d}$ with i.i.d. standard normal entries. In this case the Hessian of $H(\mathbf{x})$ is always positive semidefinite and Hessian descent does not yield any improvement over random guessing, namely $H(\mathbf{x}^{\text{HD}})/n = \xi_0/2 + o_P(1)$.

On the other hand, it is simple to find a solution by linear algebra. More precisely, we can find the minimum norm solution of $\mathbf{G}\mathbf{x} = -\sqrt{\xi_0}\mathbf{h}$, namely $\mathbf{x}_0 = -\mathbf{G}^\top(\mathbf{G}\mathbf{G}^\top)^{-1}\sqrt{\xi_0}\mathbf{h}$ and add a vector in the null space of \mathbf{G} of suitable length to get a solution of unit norm. A simple random matrix calculation yields

$$\|\mathbf{x}_0\|_2^2 = \frac{\xi_0\alpha}{1-\alpha} + o_P(1). \quad (47)$$

We thus can efficiently construct a solution $\mathbf{x}_* \in \mathbb{S}^{d-1}$ provided $\alpha < 1/(\xi_0 + 1)$.

From the point of view of optimization, the reason for suboptimality is that Hessian descent does not exploit gradient information in a neighborhood of $\mathbf{x} = \mathbf{0}$. Spin-glass theory offers a more refined picture. If $\xi'(0) \neq 0$, the set of solutions of the system $\mathbf{F}(\mathbf{x}) = \mathbf{0}$ is centered at \mathbf{m}_* , with $\|\mathbf{m}_*\|_2$ bounded away from 0.

In order to find solutions in this situation, we proceed in two phases: (i) We first find a proxy \mathbf{m}^L for \mathbf{m}_* using an approximate message passing (AMP) algorithm. (ii) We use the vector \mathbf{m}^L as initialization for the next phase. We restrict to the hyperplane orthogonal to this vector and run Hessian descent to optimize the energy $H(\mathbf{x})$.

The AMP iteration in the first phase is defined by letting, for $\ell \geq 0$

$$\mathbf{h}^{\ell+1} = \frac{1}{\sqrt{n}}\mathbf{F}(\mathbf{m}^\ell) - \gamma B_\ell \mathbf{h}^{\ell-1}, \quad (48)$$

$$\mathbf{m}^{\ell+1} = \frac{\gamma}{\sqrt{d}}\mathbf{D}\mathbf{F}(\mathbf{m}^\ell)^\top \mathbf{h}^\ell - \gamma C_\ell \mathbf{m}^{\ell-1} - \gamma^2 D_\ell \mathbf{m}^{\ell-1}. \quad (49)$$

with initialization

$$\mathbf{m}^0 = \mathbf{m}^{-1} = \mathbf{h}^0 = \mathbf{0}. \quad (50)$$

Further, γ is a (non-random) constant to be fixed in the course of the proof and B_ℓ, C_ℓ, D_ℓ are given by

$$B_\ell := \frac{1}{\sqrt{\alpha}}\xi'(\langle \mathbf{m}^\ell, \mathbf{m}^{\ell-1} \rangle), \quad (51)$$

$$C_\ell := \sqrt{\alpha}\xi'(\langle \mathbf{m}^\ell, \mathbf{m}^{\ell-1} \rangle), \quad (52)$$

$$D_\ell := \xi'(\langle \mathbf{m}^\ell, \mathbf{m}^{\ell-1} \rangle) \langle \mathbf{h}^\ell, \mathbf{h}^{\ell-1} \rangle. \quad (53)$$

The pseudocode the algorithm is given in Algorithm 3.

Algorithm 3: Two-Phase Algorithm

Data: Couplings $\{\mathbf{G}^{(k)}\}_{0 \leq k \leq k_{\max}}$, iteration number L , stepsize δ , AMP parameter γ

Result: Approximate optimizer $\mathbf{x}^{\text{TP}} \in \mathbb{S}^{d-1}$

Initialize $\mathbf{m}^0 = \mathbf{m}^{-1} = \mathbf{h}^{-1} = \mathbf{0}$;

for $\ell \in \{0, \dots, L-1\}$ **do**

 | Compute $\mathbf{h}^\ell, \mathbf{m}^{\ell+1}$ via Eqs. (48), (49);

end

Set $\mathbf{x}^0 = \mathbf{m}^L, V_L := \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{x}, \mathbf{m}^L \rangle = 0\}$;

Set $K = (1 - \|\mathbf{m}^L\|^2)/\delta$;

for $k \in \{0, \dots, K-1\}$ **do**

 | Compute $\mathbf{v} = \mathbf{v}(\mathbf{x}^k) \in \mathbb{T}_{\mathbf{x}^k} \cap V_L$ such that $\|\mathbf{v}\|_2 = 1$ and

$$\langle \mathbf{v}, \nabla^2 H(\mathbf{x}^k) \mathbf{v} \rangle \leq \lambda_{\min}(\nabla^2 H(\mathbf{x}^k)|_{\mathbb{T}_{\mathbf{x}^k} \cap V_L}) + d\delta;$$

 | Set $s_k := \text{sign}(\langle \mathbf{v}(\mathbf{x}^k), \nabla H(\mathbf{x}^k) \rangle)$;

 | $\mathbf{x}^{k+1} = \mathbf{x}^k - s_k \sqrt{\delta} \mathbf{v}(\mathbf{x}^k)$;

end

return $\mathbf{x}^{\text{TP}} = \mathbf{x}^K / \|\mathbf{x}^K\|_2$;

We begin by stating a theorem that characterizes the first phase of the algorithm.

Theorem 5. Assume $\xi(0), \xi'(0) > 0$, $\alpha \in (0, 1)$. Let $q_{\text{RS}}(\xi) := \text{argmin}_{q>0} [\xi(q)\xi'(q)/q]$ and $q_0(\alpha) = q_0(\alpha, \xi)$ to be the unique positive solution of $\alpha = q\xi'(q)/\xi(q)$. For $q \in (0, q_{\text{RS}}(\xi))$, define $\gamma_*(q, \alpha, \xi)$ via

$$\gamma_*(q, \alpha, \xi) = -\sqrt{\frac{q}{\xi(q)\xi'(q)}}, \quad \underline{q} := q \wedge q_0(\alpha). \quad (54)$$

(If $\xi(t) = \xi_0 + \xi_1 t$ we set $q_{\text{RS}}(\xi) = \infty$.)

Then the first phase of Algorithm 3, with input parameter $\gamma = \gamma_*(q, \alpha, \xi)$, outputs \mathbf{m}^L such that

$$\left| \text{p-lim}_{n,d \rightarrow \infty} \frac{1}{d} \|\mathbf{m}^L\|^2 - q \wedge q_0(\alpha) \right| \leq C e^{-L/C}, \quad (55)$$

$$\left| \text{p-lim}_{n,d \rightarrow \infty} \frac{1}{2n} \|\mathbf{F}(\mathbf{m}^L)\|^2 - u_{\text{RS}}(q, \alpha, \xi) \right| \leq C e^{-L/C}, \quad (56)$$

$$u_{\text{RS}}(q, \alpha, \xi) := \frac{1}{2} \left(\sqrt{\xi(q)} - \sqrt{\frac{1}{\alpha} q \xi'(q)} \right)_+^2. \quad (57)$$

Note that for our problem statement, we should choose $q \in (0, q_{\text{RS}}(\xi) \wedge 1)$. However, the statement above makes sense and holds also for $q \geq 1$.

Remark 5.2. Note that $q \mapsto V(q) := \xi(q)\xi'(q)/q$ is strictly convex on $(0, q_{\max})$ with q_{\max} the maximum radius of convergence of ξ (see Appendix G), hence $q_{\text{RS}}(\xi)$ is well defined, with $q_{\text{RS}}(\xi) = \infty$ if and only if ξ is linear. Further $q \mapsto g(q) := q\xi'(q)/\xi(q)$ is strictly increasing with $g(0) = 0$, $\lim_{q \rightarrow \infty} g(q) = \infty$. Hence $q_0(\alpha, \xi)$ is well defined as well. In general neither $q_{\text{RS}}(\xi)$ nor $q_0(\alpha, \xi)$ need to take values in $[0, 1]$.

We use this result together with a generalization of the proof of Theorem 4 to prove the next statement.

Theorem 6. For $\alpha \in (0, 1)$, $a, b \in \mathbb{R}_{\geq 0}$, define $z_*(\alpha, a, b)$ as in Theorem 4. Further, define $q_{\text{RS}}(\xi)$, $q_0(\alpha, \xi)$, $\gamma_*(q, \alpha, \xi)$ as in Theorem 5 and set $q_* := q_{\text{RS}}(\xi) \wedge q_0(\alpha, \xi)$. Assume $q_* < 1$ (see Remark 5.3 for the case $q_* \geq 1$.)

Let $u(\cdot; \alpha, \xi, q) : [0, 1 - q] \rightarrow \mathbb{R}$ be the unique solution of the ordinary differential equation

$$\frac{du}{dt}(t) = -\frac{1}{2\alpha} z_*(\alpha; \sqrt{2\alpha u(t)\xi''(q+t)}, \xi'(q+t)), \quad u(0) = u_{\text{RS}}(q, \alpha, \xi). \quad (58)$$

Then for any $\varepsilon > 0$, there exists a constant $\delta_0 = \delta_0(\varepsilon, \alpha, \xi) > 0$ depending uniquely on α, ξ such that the two phase algorithm, with stepsize parameter $\delta \leq \delta_0$ and $\gamma = \gamma_*(q_* - \delta_0, \alpha, \xi)$ outputs $\mathbf{x}^{\text{TP}} \in \mathbb{S}^{d-1}$ such that

$$\text{p-lim}_{d, n \rightarrow \infty} \frac{1}{2n} \|\mathbf{F}(\mathbf{x}^{\text{TP}})\|_2^2 \leq u(1 - q_*; \alpha, \xi, q_*) + \varepsilon. \quad (59)$$

Further the algorithm has the complexity of $C(\delta)$ matrix vector multiplications by $\nabla^2 H$ and $D\mathbf{F}$.

The proofs of this theorem and of Theorem 5 are presented in Appendix G.

Remark 5.3. Theorem 6 only covers the case $q_* < 1$.

If $q_* \geq 1$, the first phase of the algorithm (with input parameter $\gamma = \gamma_*(q = 1, \alpha, \xi)$) achieves, by Theorem 5

$$\text{p-lim}_{n, d \rightarrow \infty} \frac{1}{2n} \|\mathbf{F}(\mathbf{m}^L)\|_2^2 = \frac{1}{2} \left(\sqrt{\xi(1)} - \sqrt{\frac{1}{\alpha} 2\xi'(1)}_+ \right)^2 + O(e^{-L/C}). \quad (60)$$

We expect this value cannot be improved cannot be improved by other choices of the algorithm parameters.

6 Interpolating random data: An empirical comparison

Throughout the paper, we studied the problem of solving the system of nonlinear equations $\mathbf{F}(\mathbf{W}) = \mathbf{0}$ with respect to unknowns $\mathbf{W} \in \mathbb{S}^{d-1}$, when \mathbf{F} is a Gaussian process. (Throughout this section, the unknowns will be denoted by \mathbf{W} to match the applied literature.)

It is natural to wonder whether the theory developed in this setting can provide any guidance towards understanding more complex cases in which the functions F_1, \dots, F_n are random but non-Gaussian.

In this section we consider the interpolation problem introduced in Section 1 and describe a simulation study comparing the Gaussian theory of the previous section to empirical results. More precisely, we consider the problem of interpolating random data $(y_i, \mathbf{z}_i) \in \{+1, -1\} \times \mathbb{R}^D$ using a two-layer neural network with weights \mathbf{W} .

The setup for these simulations is defined in Section 6.1. Section 6.2 investigates the robustness of our results with respect to various choices in the simulations. Finally, in Section 6.3 we compare empirical results in the interpolation problem with the Gaussian theory of the previous sections. Namely, we compare empirical results with theoretical ones within a Gaussian model that approximately matches the covariance structure of the interpolation model.

In general, we observe a gap between the Gaussian theory and the neural network model. However, the two appear to be in rough qualitative agreement and, in certain cases, surprisingly close to each other.

6.1 Setup and definitions

We are interested in the ability of large neural networks to interpolate completely unstructured (pure noise) data, a phenomenon that has attracted considerable attention over the last few years [ZBH⁺21, BHMM19, BMR21]. In order to capture the essence of this problem in a simple setting, we assume to be given i.i.d. data $\{(y_i, \mathbf{z}_i)\}_{i \leq n}$ with $y_i \sim \text{Unif}(\{+1, -1\})$ independent of $\mathbf{z}_i \sim \mathbf{N}(\mathbf{0}, \mathbf{I}_D)$. We consider a two-layer neural network with D inputs and m hidden units:

$$f(\mathbf{z}; \mathbf{W}) = \frac{a}{\sqrt{m}} \sum_{j=1}^m s_j \sigma(\langle \mathbf{w}_j, \mathbf{z} \rangle). \quad (61)$$

Here $a \in \mathbb{R}_{>0}$ is a scale parameter that will be fixed independently of the data. The signs $s_1, \dots, s_m \in \{+1, -1\}$ are also fixed (independent of the data) and uniformly random subject to $\#\{i \in [m] : s_i = +1\} = \#\{i \in [m] : s_i = -1\} = m/2$ (for simplicity we assume m even). The weights $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_m)^\top \in \mathbb{R}^{m \times D}$, are instead fit to the data as to minimize the empirical risk

$$\hat{R}_n(\mathbf{W}) := \frac{1}{2n} \sum_{i=1}^n (y_i - f(\mathbf{z}_i; \mathbf{W}))^2, \quad (62)$$

subject to the norm constraints

$$\|\mathbf{W}\|_F^2 = \sum_{i=1}^m \|\mathbf{w}_i\|_2^2 \leq m. \quad (63)$$

Note that $n \cdot \hat{R}_n(\mathbf{W})$ is nothing but the Hamiltonian $H(\mathbf{x}) = \|\mathbf{F}(\mathbf{x})\|_2^2/2$ of the previous pages, with the replacement $\mathbf{x} = \mathbf{W}/\sqrt{m}$ and $F_i(\mathbf{x}) = y_i - f(\mathbf{z}_i; \mathbf{W})$. We thus identify the dimension of the optimization problem as $d = mD$.

We attempt to minimize the cost $\hat{R}_n(\mathbf{W})$ using stochastic gradient descent (SGD). Below are some specifics of our experiments.

Activation function. We use the standard ELU activation:

$$\sigma(x) = \begin{cases} x & \text{if } x \geq 0, \\ e^x - 1 & \text{if } x < 0. \end{cases} \quad (64)$$

Optimization algorithm. We use SGD with batch size $|B(k)| = 4$:

$$\tilde{\mathbf{W}}^{k+1} = \mathbf{W}^k - \frac{\eta_k}{2} \sum_{i \in B(k)} \nabla_{\mathbf{w}} (y_i - f(\mathbf{z}_i; \mathbf{W}^k))^2, \quad (65)$$

$$\mathbf{W}^{k+1} = \left(1 \wedge \frac{\sqrt{m}}{\|\tilde{\mathbf{W}}^{k+1}\|_F}\right) \tilde{\mathbf{W}}^{k+1}, \quad (66)$$

and stepsize $\eta_k = \text{lr}/(1 + E(k))^{1/2}$, where $E(k)$ is the epoch index. We will vary the learning rate lr and number of epochs.

In order to satisfy the ℓ_2 constraint (63), we project back onto this set at each iteration, as per Eq. (66).

Initialization. We initialize the weights to $\mathbf{W} \sim \mathbf{N}(0, \varepsilon^2 \mathbf{I}_{mD}/D)$, and use $\varepsilon = 0.03$ in simulations. This corresponds to initializing close to the center of the ball (63), since $\|\mathbf{W}\|_F/\sqrt{m} = \varepsilon + o_P(1)$.

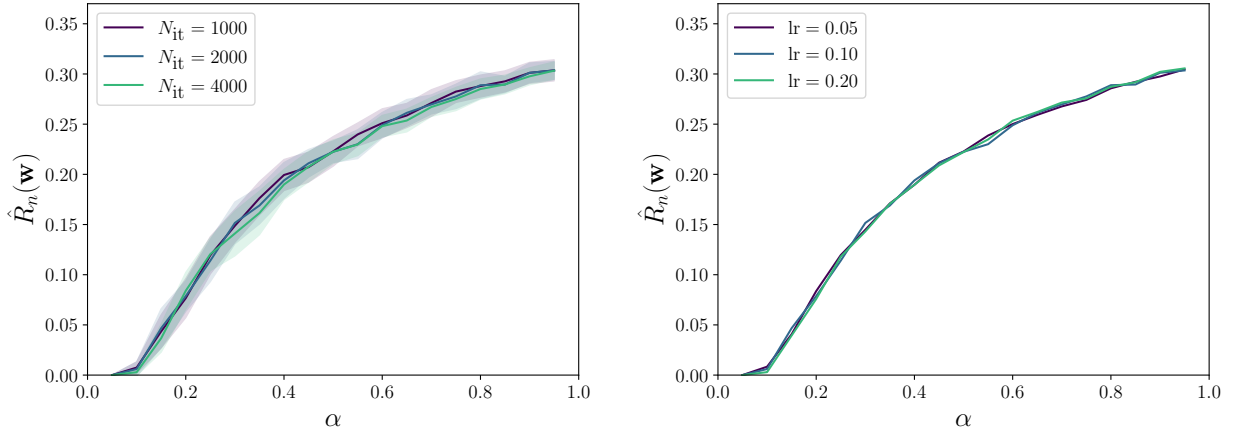


Figure 2: Fitting random binary labels using a 2-layer ELU network (61), using SGD and weights bounded as per Eq. (63), as a function of the number of samples per parameter $n/(mD)$. Left: Dependence of training error on the number of epochs N_{it} . Right: Dependence on the learning rate lr .

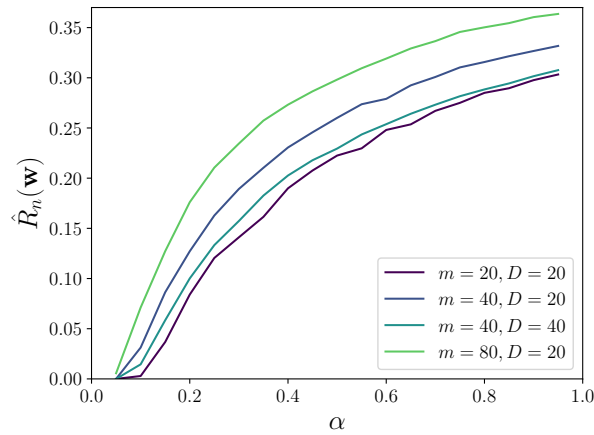


Figure 3: Same as in Fig. 2. Training error as a function of the number of samples per parameter $n/(mD)$. Here we vary the input dimension D and number of neurons m .

6.2 Dependence on the simulations parameters

In Figures 2, 3 we investigated the dependence of training error (energy) achieved over various parameters of the model and training algorithm. All of these figures are obtained by setting $a = 1$, and report the average training error over 20 realizations of the data, and independent runs of the algorithm. When present, bands represent an interval of plus/minus one standard deviation over these 20 realizations.

In more detail, we consider the following settings:

- Figure 2, left frame: we use $m = D = 20$, $\text{lr} = 0.1$, and increase the number of epochs N_{it} from 1000 to 4000.
- Figure 2, right frame: we use $m = D = 20$, $N_{\text{it}} = 2000$, and vary the learning rate lr between 0.05 and 0.20.

Our theory does not capture dependence on learning rate or number of iterations. In this respect, the empirical observation that this dependence is mild in the current setting is reassuring.

Next, the Gaussian theory only depends on the ratio $\alpha = n/d$, where d is the number of decision variables. In the current setting $d = mD$, and is therefore interesting to investigate the dependence on m, D while keeping $\alpha = n/(mD)$ fixed:

- In Figure 3 we use $\text{lr} = 0.1$, $N_{\text{it}} = 4000$, and vary the number of neurons m and input dimension D .

We observe that the training error achieved *does* depend on the architecture, and most noticeably on the ‘aspect ratio’ m/D . On the other hand, this dependence is not as dramatic as one might expect. A fourfold increase in aspect ratio only changes the training error by 20%.

6.3 Comparison with Gaussian theory

We next compare simulation results with the analytical prediction for a Gaussian model with matching covariance (see Section 6.4 for a description of this prediction). More precisely we compute the prediction for the two-phase algorithms of Theorem 6. We expect this to correspond to the optimal energy achieved (asymptotically as $n, d \rightarrow \infty$) by a broad class of efficient algorithms. Hence should provide an upper bound on the energy achieved by SGD (and an upper bound on the threshold.)

In Figure 4 we consider number of neurons $m = 20$, and input dimension $D = 20$, and vary N_{it}, lr as indicated in legends. We change $a \in \{1, 2, 5\}$, which of course impacts the Gaussian prediction. The agreement is surprisingly good for $a = 1$, and worsens for $a \in \{2, 5\}$. Despite the worse agreement for $a \gtrsim 2$, the Gaussian theory still has the correct qualitative behavior, and appears to be a good starting point to analyze the more complex model (61).

6.4 Covariance matching

We conclude this section by discussing the ‘‘covariance matching’’ used in to compute the theoretical predictions in Fig. 4. We begin by observing that, *in general*, we cannot expect the Gaussian model to capture the behavior of the system of nonlinear equations $y_i - f(\mathbf{z}_i; \mathbf{W}) = 0$, with $f(\cdot; \mathbf{W})$ given by Eq. (61).

To see this, consider the case of a linear activation $\sigma(x) = x$. Then the system of equations reads

$$\left(\frac{a}{\sqrt{m}} \sum_{j=1}^m s_j \mathbf{w}_j \right)^\top \mathbf{z}_i = y_i \quad \forall i \leq n. \quad (67)$$

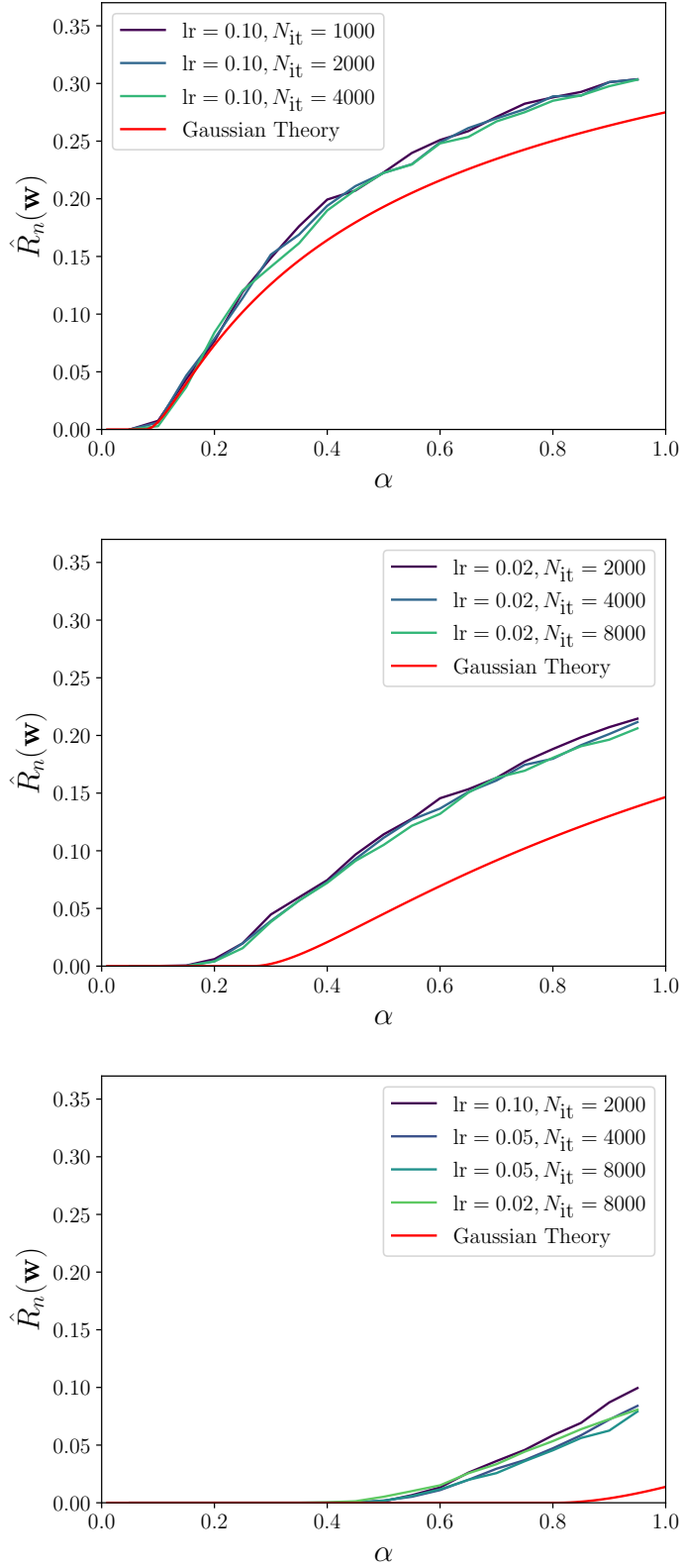


Figure 4: Training error of SGD as a function of the number of samples per parameter $\alpha = n/(mD)$. From top to bottom: $a \in \{1, 2, 5\}$. Red lines report the theoretical prediction of Theorem 6 for the optimal ‘two-phase’ algorithm.

In other words, the model depends on the weight vectors only through the D -dimensional projection $\mathbf{u} := am^{-1/2} \sum_{j=1}^m s_j \mathbf{w}_j$. In particular, a solution exists if and only if $D \geq n$ (for generic data $(\mathbf{z}_i)_{i \leq n}$), and the total number of parameters is irrelevant.

In order to take into account the fact that the linear component of σ has a much smaller number of degrees of freedom than in the corresponding Gaussian model, we project out the linear component of activations before matching covariances. Namely, we define $\sigma_{\#}(x) := \sigma(x) - \sigma_1 \cdot x$, where $\sigma_1 := \mathbb{E}_{G \sim \mathcal{N}(0,1)}[G\sigma(G)]$ is the linear coefficient in the Hermite expansion of σ , and

$$f_{\#}(\mathbf{z}; \mathbf{W}) = \frac{a}{\sqrt{m}} \sum_{j=1}^m s_j \sigma_{\#}(\langle \mathbf{w}_j, \mathbf{z} \rangle), \quad (68)$$

$$F_i(\mathbf{W}) := y_i - f_{\#}(\mathbf{z}; \mathbf{W}). \quad (69)$$

We now compute the covariance of the process $F_i(\mathbf{w})$ with respect to the random data \mathbf{z} . Letting $C_F(\mathbf{W}; \mathbf{W}') := \mathbb{E}_{\mathbf{z}_i, y_i}[F_i(\mathbf{W})F_i(\mathbf{W}')]$, we have

$$C_F(\mathbf{W}; \mathbf{W}') = 1 + \mathbb{E}_{\mathbf{z}}\{f_{\#}(\mathbf{z}; \mathbf{W})f_{\#}(\mathbf{z}; \mathbf{W}')\} \quad (70)$$

$$= 1 + \frac{a^2}{m} \sum_{j,l=1}^m s_j s_l \mathbb{E}_{\mathbf{z}} \sigma_{\#}(\langle \mathbf{w}_j, \mathbf{z} \rangle) \sigma_{\#}(\langle \mathbf{w}'_l, \mathbf{z} \rangle) \quad (71)$$

$$= 1 + \frac{a^2}{m} \sum_{j,k=1}^m s_j s_l \widehat{K}(\langle \mathbf{w}_j, \mathbf{w}'_l \rangle, \|\mathbf{w}_j\|_2^2, \|\mathbf{w}'_l\|_2^2), \quad (72)$$

where we introduced the kernel \widehat{K} defined by

$$\widehat{K}(r_{12}, r_{11}, r_{22}) := \mathbb{E}[\sigma_{\#}(G_1)\sigma_{\#}(G_2)], \quad (73)$$

$$(G_1, G_2) \sim \mathcal{N}(\mathbf{0}, \mathbf{R}) \quad \mathbf{R} := \begin{pmatrix} r_{11} & r_{12} \\ r_{12} & r_{22} \end{pmatrix}. \quad (74)$$

We observe that the covariance (72) is not invariant under the orthogonal group $\mathcal{O}(mD)$ but only under $\mathfrak{S}_{m/2} \otimes \mathfrak{S}_{m/2} \otimes \mathcal{O}(D)$ (with $\mathfrak{S}_{m/2}$ denoting the group of permutations in the space of neurons, for neurons with same s_j , and $\mathcal{O}(D)$ denoting rotations in the space of weights). In order to replace it by an orthogonally invariant covariance (depending only on $\langle \mathbf{W}, \mathbf{W}' \rangle$), we make the following approximations:

1. Since the signs $(s_j)_{j \leq m}$ in Eq. (72) are random, we keep only diagonal terms in the sum, and thus replace $C_F(\mathbf{W}; \mathbf{W}')$ by $C_F^{(1)}(\mathbf{W}; \mathbf{W}')$ whereby

$$C_F^{(1)}(\mathbf{W}; \mathbf{W}') = 1 + \frac{a^2}{m} \sum_{j=1}^m \widehat{K}(\langle \mathbf{w}_j, \mathbf{w}'_j \rangle, \|\mathbf{w}_j\|_2^2, \|\mathbf{w}'_j\|_2^2).$$

Notice that this is accurate at a fixed \mathbf{W}, \mathbf{W}' but not uniformly over \mathbf{W}, \mathbf{W}' .

2. We approximate the value of $C_F^{(1)}(\mathbf{W}; \mathbf{W}')$ by the one taken by $C_F^{(1)}(\mathbf{W}; \mathbf{W}')$ for a ‘typical’ pair \mathbf{W}, \mathbf{W}' with a given inner product $\langle \mathbf{W}, \mathbf{W}' \rangle$. More formally define the set

$$\mathcal{W}_{m,n}(q) := \left\{ \mathbf{W}, \mathbf{W}' \in \mathbb{R}^{m \times D} : \|\mathbf{W}\|_F^2 = m, \|\mathbf{W}'\|_F^2 = m, \langle \mathbf{W}, \mathbf{W}' \rangle = mq \right\}. \quad (75)$$

If we draw $(\mathbf{W}, \mathbf{W}') \sim \text{Unif}(\mathcal{W}_{m,n}(q))$, then for a fixed $j \in [m]$, $\|\mathbf{w}_j\|_2 = 1 + o_P(1)$, $\|\mathbf{w}'_j\|_2 = 1 + o_P(1)$ and $\langle \mathbf{w}_j, \mathbf{w}'_j \rangle = q + o_P(1)$. As a consequence, we will have

$$C_F^{(1)}(\mathbf{W}; \mathbf{W}') = 1 + a^2 \widehat{K}(q, 1, 1) + o_P(1). \quad (76)$$

Summarizing, and simplifying notations, the above approximation suggests to use a Gaussian process \mathbf{F}^G in $d = mD$ dimensions with

$$\mathbb{E}[F_i^G(\mathbf{W})F_j^G(\mathbf{W}')] = \xi(\langle \mathbf{W}, \mathbf{W}' \rangle / m), \quad (77)$$

$$\xi(q) := 1 + a^2 K_0(q) := 1 + a^2[K(q) - \sigma_1^2 q], \quad (78)$$

$$K(q) := \mathbb{E}[\sigma(G_1)\sigma(G_2)], \quad (G_1, G_2) \sim \mathbf{N}\left(\mathbf{0}, \begin{pmatrix} 1 & q \\ q & 1 \end{pmatrix}\right). \quad (79)$$

This is exactly the function $\xi(\cdot)$ used to compare simulations and theory in the previous sections.

We conclude this section with a warning. *We do not expect the theoretical predictions based on the covariance matching here to be asymptotically exact as $m, D, n \rightarrow \infty$.* In particular, as emphasized above, the actual covariance structure is not invariant under $\mathcal{O}(mD)$.

Nevertheless, the rough agreement between theory and empirical results suggests that the Gaussian model (possibly with a more complex covariance structure) might be a useful starting point.

7 Discussion

In this paper, presented some simple upper and lower bounds on the threshold for existence of solutions of a set of overparametrized nonlinear equations $\mathbf{F}(\mathbf{x}) = \mathbf{0}$ when $\mathbf{F} : \mathbb{S}^{d-1} \rightarrow \mathbb{R}^n$ is a smooth Gaussian process.

Our main focus was to analyze polynomial-time algorithms to construct approximate solutions, i.e. $\mathbf{x} \in \mathbb{S}^{d-1}$ such that $\|\mathbf{F}(\mathbf{x})\|_2^2 \leq \varepsilon \mathbb{E}\{\|\mathbf{F}(\mathbf{x}_0)\|_2^2\}$ (here $\mathbf{x}_0 \in \mathbb{S}^{d-1}$ is uniformly random). In particular, we presented a two-phase algorithm which we expect to be near-optimal, and characterized the value it achieves in Section 5. For cases in which $D\mathbf{F}(\mathbf{0}) = \mathbf{0}$ (i.e. $\xi'(0) = 0$), this algorithm reduces to the Hessian descent algorithm of [Sub21].

As shown in Figure 1, there exist cases in which we know that solutions exist but we are not able to find any in polynomial time.

These results naturally suggest a number of interesting questions. Among others: (i) Can we provide evidence of hardness for cases in which the algorithms presented here fail? (ii) Can we construct exact solutions? (iii) What happens when the system of equations is only modestly overparametrized? (iv) Can we characterize sharply the conditions under which gradient flow finds solutions?

We plan to report on questions (i)-(iii) in the near future [MS23a, MS23b].

Acknowledgements

AM was supported by the NSF through award DMS-2031883, the Simons Foundation through Award 814639 for the Collaboration on the Theoretical Foundations of Deep Learning, the NSF grant CCF-2006489 and the ONR grant N00014-18-1-2729. Part of this work was carried out while Andrea Montanari was on partial leave from Stanford and a Chief Scientist at Ndata Inc dba Project N. The present research is unrelated to AM's activity while on leave. ES was supported by the Israel Science Foundation (Grant Agreement No. 2055/21) and a research grant from the Center for Scientific Excellence at the Weizmann Institute of Science. ES is the incumbent of the Skirball Chair in New Scientists.

A Basic estimates

This appendix contains some basic estimates of the process $(\mathbf{F}(\mathbf{x}))_{\mathbf{x} \in \mathbb{R}^d}$, which will be useful in the following. Throughout $(F_\ell)_{\ell \geq 0}$ are i.i.d. Gaussian processes with $\mathbb{E}[F_\ell(\mathbf{x})F_\ell(\mathbf{y})] = \xi(\langle \mathbf{x}, \mathbf{y} \rangle)$.

We begin by bounding the expected maximum. Notice that the limiting value of this quantity is exactly given by Parisi's formula, cf. Eq. (16). However, we find it convenient to derive some explicit estimates.

Proposition A.1. *There exist absolute constants C_1, C_2 such that (writing $\log_+(t) := \max(1, \log t)$)*

$$C_1 \sqrt{\frac{\xi'(1)^2}{\xi(1)\xi''(1)} \log_+ \frac{\xi''(1)}{\xi'(1)}} \leq \frac{1}{\sqrt{d}\xi(1)} \mathbb{E} \max_{\mathbf{x} \in \mathbb{S}^{d-1}} F_1(\mathbf{x}) \leq C_2 \sqrt{\log_+ \frac{\xi'(1)}{\xi(1)}}. \quad (80)$$

The same bounds hold if the maximum is taken over $\mathbf{x} \in \mathbb{B}^d(1)$.

Proof. Given $\mathbf{x}, \mathbf{y} \in \mathbb{S}^{d-1}$, we denote their Euclidean and rescaled canonical distances by r and u respectively

$$r := \|\mathbf{x} - \mathbf{y}\|_2 = \sqrt{2(1 - \langle \mathbf{x}, \mathbf{y} \rangle)}, \quad (81)$$

$$u := \sqrt{\frac{1}{\xi(1)} \mathbb{E}\{[F_1(\mathbf{x}) - F_1(\mathbf{y})]^2\}} = \sqrt{2\left(1 - \frac{\xi(\langle \mathbf{x}, \mathbf{y} \rangle)}{\xi(1)}\right)}. \quad (82)$$

We denote by $r(u)$ the (strictly increasing) function that maps u to r , namely the unique solution of the equation

$$1 - \frac{u^2}{2} = \frac{\xi(1 - r^2/2)}{\xi(1)}. \quad (83)$$

Since $\xi'(t) \leq \xi'(1)$ on $[0, 1]$,

$$\frac{\xi(1 - r^2/2)}{\xi(1)} \geq 1 - \frac{\xi'(1)}{\xi(1)} \cdot \frac{r^2}{2}, \quad (84)$$

whence

$$r(u) \geq \sqrt{\frac{\xi(1)}{\xi'(1)}} u. \quad (85)$$

Let $N_d(r) \leq [(10/r) \vee 1]^d$ (respectively, $N_d^{F_1}(r)$) be the covering number of \mathbb{S}^{d-1} w.r.t. the Euclidean distance (respectively, the canonical distance of F_1). Noting the $1/\xi(1)$ factor in (82), we have that

$$N_d^{F_1}(\sqrt{\xi(1)}u) = N_d(r(u)) \leq N_d\left(\sqrt{\frac{\xi(1)}{\xi'(1)}}u\right) \leq \left(\sqrt{\frac{\xi'(1)}{\xi(1)}} \frac{10}{u} \vee 1\right)^d.$$

Finally notice that, under the canonical distance, $\text{diam}_{F_1}(\mathbb{S}^{d-1}) \leq 2\sqrt{\xi(1)}$. By Dudley's inequality, we have

$$\begin{aligned} \mathbb{E} \max_{\mathbf{x} \in \mathbb{S}^{d-1}} F_1(\mathbf{x}) &\leq 24 \int_0^{2\sqrt{\xi(1)}} \sqrt{\log N_d^{F_1}(u)} du \\ &\leq 24\sqrt{\xi(1)} \int_0^2 \sqrt{\log N_d^{F_1}(\sqrt{\xi(1)}u)} du \end{aligned}$$

$$\begin{aligned}
&\leq 24\sqrt{d\xi(1)} \int_0^2 \sqrt{\log_+ \left(\sqrt{\frac{\xi'(1)}{\xi(1)}} \cdot \frac{10}{u} \right)} du \\
&\leq C\sqrt{d\xi(1) \log_+ \frac{\xi'(1)}{\xi(1)}}.
\end{aligned}$$

This proves the desired upper bound.

In order to prove the lower bound, we use Sudakov's inequality. Recall that the r -packing number is lower bounded by the $2r$ -covering number. For any $0 < r < 1$, letting $t = r^2/2$, we have

$$\begin{aligned}
\mathbb{E} \max_{\mathbf{x} \in \mathbb{S}^{d-1}} F_1(\mathbf{x}) &\geq C' \sqrt{(\xi(1) - \xi(1 - (r^2/2))) \log N_d(2r)} \\
&\geq C' \sqrt{d(\xi(1) - \xi(1 - t)) \log_+ \frac{1}{t}} \\
&\geq C' \sqrt{d(\xi'(1)t - \xi''(1)t^2/2) \log_+ \frac{1}{t}},
\end{aligned}$$

where $C' > 0$ is a constant that may change from line to line and where in the last step we used the intermediate value theorem, and the fact that $\xi''(s) \leq \xi''(1)$ for $s \leq 1$. The desired inequality follows by choosing $t = \xi'(1)/(2\xi''(1))$. \square

Remark A.1. Let L be an integer-valued random variable with probability distribution $\mathbb{P}(L = \ell) = \xi_\ell/\xi(1)$. Then the upper and lower bounds of the last proposition are within a constant provided $\mathbb{E}[L^2] \leq C\mathbb{E}[L]^2$.

We will next establish upper bounds on the derivatives of the process $F_\ell(\cdot)$. Before doing it, we state a useful lemma.

Lemma A.2. Let $(Z(\mathbf{x}, \mathbf{v}) : \mathbf{x} \in \mathbf{B}^d(1), \mathbf{v} \in \mathbb{S}^{d-1})$ be a continuous centered Gaussian process, and assume that, for some $r_0 > 0$, the following holds for all sets $A \subseteq \mathbf{B}^d(1)$:

$$\text{diam}(A) \leq r_0 \Rightarrow \mathbb{E} \left[\sup_{\mathbf{x} \in A, \mathbf{v} \in \mathbb{S}^{d-1}} Z(\mathbf{x}, \mathbf{v}) \right] \leq M(r_0). \quad (86)$$

Then there exists a universal constant C such that

$$\mathbb{E} \left[\sup_{\mathbf{x} \in \mathbf{B}^d(1), \mathbf{v} \in \mathbb{S}^{d-1}} Z(\mathbf{x}, \mathbf{v}) \right] \leq M(r_0) + C\sqrt{dV_* \log_+(1/r_0)}, \quad (87)$$

$$V_* := \sup_{\mathbf{x} \in \mathbf{B}^d(1), \mathbf{v} \in \mathbb{S}^{d-1}} \text{Var}(Z(\mathbf{x}, \mathbf{v})). \quad (88)$$

Proof. Let $(A_\ell : \ell \leq N)$ be a covering of $\mathbf{B}^d(1)$ with sets of diameter $\text{diam}(A_\ell) \leq r_0$, such that $\log N \leq Cd \log_+(1/r_0)$. Define

$$Z_\ell := \sup_{\mathbf{x} \in A_\ell, \mathbf{v} \in \mathbb{S}^{d-1}} Z(\mathbf{x}, \mathbf{v}). \quad (89)$$

We then have

$$\mathbb{P} \left(\sup_{\mathbf{x} \in \mathbf{B}^d(1), \mathbf{v} \in \mathbb{S}^{d-1}} Z(\mathbf{x}, \mathbf{v}) \geq M(r_0) + t \right) = \mathbb{P} \left(\max_{\ell \leq N} Z_\ell \geq M(r_0) + t \right)$$

$$\begin{aligned}
&\leq N \cdot \max_{\ell \leq N} \mathbb{P}\left(Z_\ell \geq \mathbb{E}Z_\ell + t\right) \\
&\leq \exp\left\{Cd \log_+(1/r_0) - \frac{t^2}{2V_*}\right\},
\end{aligned}$$

where the last step follows from the Borell-TIS inequality. This proves the claim. \square

Lemma A.3. *Recall that $P_{\mathbb{T}, \mathbf{x}}$ denotes the projector onto the tangent space to the sphere of radius $\|\mathbf{x}\|_2$ at \mathbf{x} . Write, as above, $\log_+(t) = \max(\log t; 1)$. Then there exist an absolute constant C such that, letting $\hat{\mathbf{x}} := \mathbf{x}/\|\mathbf{x}\|_2$, we have*

$$\mathbb{E} \max_{\mathbf{x} \in \mathbb{B}^d(1)} \|P_{\mathbb{T}, \mathbf{x}} \nabla F_1(\mathbf{x})\|_2 \leq C \sqrt{d \xi'(1) \log_+ \frac{\xi''(1)}{\xi'(1)}}, \quad (90)$$

$$\mathbb{E} \max_{\mathbf{x} \in \mathbb{B}^d(1)} |\langle \hat{\mathbf{x}}, \nabla F_1(\mathbf{x}) \rangle| \leq C \sqrt{d \xi''(1) \log_+ \frac{\xi'''(1)}{\xi''(1)}}, \quad (91)$$

$$\mathbb{E} \max_{\mathbf{x} \in \mathbb{B}^d(1)} \|P_{\mathbb{T}, \mathbf{x}} \nabla^2 F_1(\mathbf{x}) P_{\mathbb{T}, \mathbf{x}}\|_{\text{op}} \leq C \sqrt{d \xi''(1) \log_+ \frac{\xi^{(3)}(1)}{\xi''(1)}}, \quad (92)$$

$$\mathbb{E} \max_{\mathbf{x} \in \mathbb{B}^d(1)} \|P_{\mathbb{T}, \mathbf{x}} \nabla^2 F_1(\mathbf{x}) \hat{\mathbf{x}}\|_{\text{op}} \leq C \sqrt{d \xi^{(3)}(1) \log_+ \frac{\xi^{(4)}(1)}{\xi^{(3)}(1)}}, \quad (93)$$

$$\mathbb{E} \max_{\mathbf{x} \in \mathbb{B}^d(1)} \langle \hat{\mathbf{x}}, \nabla^2 F_1(\mathbf{x}) \hat{\mathbf{x}} \rangle \leq C \sqrt{d \xi^{(4)}(1) \log_+ \frac{\xi^{(5)}(1)}{\xi^{(4)}(1)}}. \quad (94)$$

(Here $\xi^{(\ell)}$ denotes the ℓ -th derivative of ξ .)

Proof. All these inequalities are analogous to the upper bound in Proposition A.1. We will repeatedly use Lemma A.2 to restrict ourselves to sets of some small diameter r_0 .

Consider Eq. (90). Define the following Gaussian process indexed by $\mathbf{x} \in \mathbb{B}^d(1)$, $\mathbf{v} \in \mathbb{T}_{\mathbf{x}}$, $\|\mathbf{v}\|_2 = 1$

$$G_{\perp}(\mathbf{v}, \mathbf{x}) := \langle \mathbf{v}, \nabla F_1(\mathbf{x}) \rangle. \quad (95)$$

We have

$$\mathbb{E}\{G_{\perp}(\mathbf{v}_1, \mathbf{x}_1) G_{\perp}(\mathbf{v}_2, \mathbf{x}_2)\} = \langle \mathbf{v}_1, \mathbf{x}_2 \rangle \langle \mathbf{x}_1, \mathbf{v}_2 \rangle \xi''(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle) + \langle \mathbf{v}_1, \mathbf{v}_2 \rangle \xi'(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle). \quad (96)$$

This, and other covariance calculations below, can be checked either by writing the derivatives of F_1 explicitly or using that the order of taking derivatives and expectations can be interchanged. Letting $d_{G, \perp}(\mathbf{v}_1, \mathbf{x}_1; \mathbf{v}_2, \mathbf{x}_2)$ denote the canonical distance associated to this process, we have

$$\begin{aligned}
d_{G, \perp}(\mathbf{v}_1, \mathbf{x}_1; \mathbf{v}_2, \mathbf{x}_2)^2 &= \xi'(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle) \|\mathbf{v}_1 - \mathbf{v}_2\|_2^2 + \xi'(\|\mathbf{x}_1\|_2^2) - 2\xi'(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle) + \xi'(\|\mathbf{x}_2\|_2^2) \\
&\quad + (\langle \mathbf{v}_1, \mathbf{x}_2 \rangle - \langle \mathbf{v}_2, \mathbf{x}_1 \rangle)^2 \xi''(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle) \\
&\quad - \langle \mathbf{v}_1, \mathbf{x}_2 \rangle^2 \xi''(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle) - \langle \mathbf{v}_1, \mathbf{x}_2 \rangle^2 \xi''(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle) \\
&\leq \xi'(1) \|\mathbf{v}_1 - \mathbf{v}_2\|_2^2 + \xi'(\|\mathbf{x}_1\|_2^2) - 2\xi'(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle) + \xi'(\|\mathbf{x}_2\|_2^2) \\
&\quad + \xi''(1) \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 \|\mathbf{v}_1 - \mathbf{v}_2\|_2^2 \\
&\leq 2\xi'(1) \|\mathbf{v}_1 - \mathbf{v}_2\|_2^2 + \xi'(\|\mathbf{x}_1\|_2^2) - 2\xi'(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle) + \xi'(\|\mathbf{x}_2\|_2^2),
\end{aligned}$$

where the last inequality follows for $\|\mathbf{x}_1 - \mathbf{x}_2\|_2 \leq r_0 := \xi'(1)/\xi''(1)$. Therefore, for any set $A \subseteq \mathbb{B}^d(1)$, $\text{diam}(A) \leq r_0$, using the Sudakov-Fernique inequality we have

$$\mathbb{E} \max_{\mathbf{x} \in A} \max_{\mathbf{v} \in \mathbb{T}_{\mathbf{x}}, \|\mathbf{v}\|_2=1} G_{\perp}(\mathbf{v}, \mathbf{x}) \leq C \mathbb{E} \max_{\|\mathbf{v}\|_2=1} G_1(\mathbf{v}) + \mathbb{E} \max_{\|\mathbf{x}\|_2 \leq 1} G_2(\mathbf{x}), \quad (97)$$

where G_1 is a process with canonical distance $\xi'(1)\|\mathbf{v}_1 - \mathbf{v}_2\|_2^2$ and G_2 is a process with canonical distance $\xi'(\|\mathbf{x}_1\|_2^2) - 2\xi'(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle) + \xi'(\|\mathbf{x}_2\|_2^2)$. By applying the previous proposition to G_2 (with ξ replaced by ξ'), we obtain

$$\mathbb{E} \max_{\mathbf{x} \in A} \max_{\mathbf{v} \in \mathbb{T}_{\mathbf{x}}, \|\mathbf{v}\|_2=1} G_{\perp}(\mathbf{v}, \mathbf{x}) \leq C \sqrt{s\xi'(1)} + C \sqrt{d\xi'(1) \log_+ \frac{\xi''(1)}{\xi'(1)}}. \quad (98)$$

Finally, using Lemma A.2, we see that

$$\mathbb{E} \max_{\mathbf{x} \in \mathbb{B}^d(1)} \max_{\mathbf{v} \in \mathbb{T}_{\mathbf{x}}, \|\mathbf{v}\|_2=1} G_{\perp}(\mathbf{v}, \mathbf{x}) \leq C \sqrt{d\xi'(1) \log_+ \frac{\xi''(1)}{\xi'(1)}} + C \sqrt{d\xi'(1) \log_+(1/r_0)}, \quad (99)$$

which yields the bound of Eq. (90).

Next consider Eq. (91). Recalling the representation (5), we have

$$\langle \hat{\mathbf{x}}, \nabla F_1(\mathbf{x}) \rangle = \frac{1}{\|\mathbf{x}\|_2} \sum_{k \geq 0} k \sqrt{\xi_k} \sum_{j_1, \dots, j_k=1}^d G_{1, j_1 \dots j_k}^{(k)} x_{j_1} \cdots x_{j_k}, \quad (100)$$

and therefore

$$\mathbb{E} \{ \langle \hat{\mathbf{x}}_1, \nabla F_1(\mathbf{x}_1) \rangle \langle \hat{\mathbf{x}}_2, \nabla F_1(\mathbf{x}_2) \rangle \} = \frac{1}{\|\mathbf{x}_1\|_2 \|\mathbf{x}_2\|_2} [\langle \mathbf{x}_1, \mathbf{x}_2 \rangle^2 \xi''(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle) + \langle \mathbf{x}_1, \mathbf{x}_2 \rangle \xi'(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle)]. \quad (101)$$

In other words, $\langle \hat{\mathbf{x}}, \nabla F_1(\mathbf{x}) \rangle$ has the same structure as $F_1(\mathbf{x})$ provided we replace ξ by $\xi'' + \xi'$. Therefore the claim follows by argument in Proposition A.1.

The bounds (92), (93), (94) follow from similar arguments, and we limit ourselves to defining the relevant Gaussian process and computing its covariance:

- For Eq. (92), we define $H_{\perp}(\mathbf{v}, \mathbf{x})$ indexed by $\mathbf{x} \in \mathbb{B}^d(1)$ and $\mathbf{v} \in \mathbb{T}_{\mathbf{x}}$ via

$$H_{\perp}(\mathbf{v}, \mathbf{x}) := \langle \mathbf{v}, \nabla^2 F_1(\mathbf{x}) \mathbf{v} \rangle. \quad (102)$$

Its covariance is (denoting by $\xi^{(\ell)}$ the ℓ -th derivative of ξ and letting $q := \langle \mathbf{x}_1, \mathbf{x}_2 \rangle$)

$$\begin{aligned} \mathbb{E} \{ H_{\perp}(\mathbf{v}_1, \mathbf{x}_1) H_{\perp}(\mathbf{v}_2, \mathbf{x}_2) \} &= \xi^{(4)}(q) \langle \mathbf{x}_1, \mathbf{v}_2 \rangle^2 \langle \mathbf{v}_1, \mathbf{x}_2 \rangle^2 + 4 \xi^{(3)}(q) \langle \mathbf{x}_1, \mathbf{v}_2 \rangle \langle \mathbf{v}_1, \mathbf{x}_2 \rangle \langle \mathbf{v}_1, \mathbf{v}_2 \rangle \\ &\quad + 2 \xi^{(2)}(q) q^2 \langle \mathbf{v}_1, \mathbf{v}_2 \rangle^2. \end{aligned}$$

The associated canonical distance is given by

$$d_{H, \perp}(\mathbf{v}_1, \mathbf{x}_1; \mathbf{v}_2, \mathbf{x}_2)^2 = \Delta_{1,2} + 2\xi^{(2)}(\|\mathbf{x}_1\|_2^2) - 4\xi^{(2)}(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle) + 2\xi^{(2)}(\|\mathbf{x}_2\|_2^2),$$

where $\Delta_{1,2}$ can be bounded as follows for $\|\mathbf{x}_1 - \mathbf{x}_2\|_2 \leq c_0$ with c_0 a small absolute constant (recall that $\mathbb{P}_{\mathbb{T}, \mathbf{x}_i}$ is the projector orthogonal to \mathbf{x}_i)

$$\Delta_{1,2} = -2\xi^{(4)}(q) \langle \mathbf{x}_1, \mathbf{v}_2 \rangle^2 \langle \mathbf{v}_1, \mathbf{x}_2 \rangle^2 - 8 \xi^{(3)}(q) \langle \mathbf{x}_1, \mathbf{v}_2 \rangle \langle \mathbf{v}_1, \mathbf{x}_2 \rangle \langle \mathbf{v}_1, \mathbf{v}_2 \rangle$$

$$\begin{aligned}
& + 4\xi^{(2)}(q)q^2(1 - \langle \mathbf{v}_1, \mathbf{v}_2 \rangle)^2 \\
& \leq 8\xi^{(3)}(q) |\langle \mathbf{P}_{\mathbb{T}, \mathbf{x}_2} \mathbf{x}_1, \mathbf{v}_1 - \mathbf{v}_2 \rangle \langle \mathbf{P}_{\mathbb{T}, \mathbf{x}_1} \mathbf{x}_2, \mathbf{v}_1 - \mathbf{v}_2 \rangle| \\
& \quad + 8\xi^{(2)}(q) \|\mathbf{v}_1 - \mathbf{v}_2\|_2^2 \\
& \leq 8(\xi^{(3)}(q) \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 + \xi^{(2)}(q)) \|\mathbf{v}_1 - \mathbf{v}_2\|_2^2 \\
& \leq 8(\xi^{(3)}(q)(1 - q) + \xi^{(2)}(q)) \|\mathbf{v}_1 - \mathbf{v}_2\|_2^2 \\
& \leq 8\xi^{(2)}(1) \|\mathbf{v}_1 - \mathbf{v}_2\|_2^2.
\end{aligned}$$

Hence the proof of Eq. (92) follows by applying Proposition A.1 and Lemma A.2.

- For Eq. (93), we define $H_2(\mathbf{v}, \mathbf{x})$ indexed by $\mathbf{x} \in \mathbb{B}^d(1)$ and $\mathbf{v} \in \mathbb{T}_{\mathbf{x}}$ via

$$H_2(\mathbf{v}, \mathbf{x}) := \langle \hat{\mathbf{x}}, \nabla^2 F_1(\mathbf{x}) \mathbf{v} \rangle. \quad (103)$$

Its covariance is (here $q := \langle \mathbf{x}_1, \mathbf{x}_2 \rangle$)

$$\begin{aligned}
\mathbb{E}\{H_2(\mathbf{v}_1, \mathbf{x}_1)H_2(\mathbf{v}_2, \mathbf{x}_2)\} &= \frac{\langle \mathbf{v}_1, \mathbf{v}_2 \rangle}{\|\mathbf{x}_1\|_2 \|\mathbf{x}_2\|_2} [\xi^{(3)}(q)q^2 + \xi^{(2)}(q)q] \\
& \quad + \frac{\langle \mathbf{v}_1, \mathbf{x}_2 \rangle \langle \mathbf{x}_1, \mathbf{v}_2 \rangle}{\|\mathbf{x}_1\|_2 \|\mathbf{x}_2\|_2} [\xi^{(4)}(q)q^2 + 3\xi^{(3)}(q)q + \xi^{(2)}(q)].
\end{aligned}$$

- Finally, for (94) we define

$$H_3(\mathbf{x}) := \langle \hat{\mathbf{x}}, \nabla^2 F_1(\mathbf{x}) \hat{\mathbf{x}} \rangle,$$

and note that (again $q := \langle \mathbf{x}_1, \mathbf{x}_2 \rangle$)

$$\mathbb{E}\{H_3(\mathbf{x}_1)H_3(\mathbf{x}_2)\} = \frac{1}{\|\mathbf{x}_1\|_2^2 \|\mathbf{x}_2\|_2^2} [\xi^{(4)}(q)q^4 + 4\xi^{(3)}(q)q^3 + 2\xi^{(2)}(q)q^2].$$

□

Proposition A.4. *Let $\mathbf{DF}(\mathbf{x}) \in \mathbb{R}^{n \times d}$ be the Jacobian of \mathbf{F} at \mathbf{x} and $\mathbf{D}^2\mathbf{F}(\mathbf{x}) \in \mathbb{R}^{n \times d \times d}$ the tensor of second derivatives. We equivalently view the latter as a linear operator $\mathbf{D}^2\mathbf{F}(\mathbf{x}) : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}^n$, and write $\mathbf{D}^2\mathbf{F}(\mathbf{x})|_{\mathbb{T} \otimes \mathbb{T}}$ for its restriction to $\mathbb{T}_{\mathbf{x}} \otimes \mathbb{T}_{\mathbf{x}}$ and similarly for $\mathbf{D}^2\mathbf{F}(\mathbf{x})|_{\mathbb{T} \otimes \hat{\mathbf{x}}}$ and $\mathbf{D}^2\mathbf{F}(\mathbf{x})|_{\hat{\mathbf{x}} \otimes \hat{\mathbf{x}}}$.*

If $n \leq d$, there exists an absolute constant C such that (using the notation $\log_+(t) := \max(\log(t); 1)$ introduced above)

$$\mathbb{E}\left[\max_{\mathbf{x} \in \mathbb{B}^d(1)} \|\mathbf{DF}(\mathbf{x})|_{\mathbb{T}}\|_{\text{op}}\right] \leq C \sqrt{d \xi'(1) \log_+ \frac{\xi''(1)}{\xi'(1)}}, \quad (104)$$

$$\mathbb{E}\left[\max_{\mathbf{x} \in \mathbb{B}^d(1)} \|\mathbf{DF}(\mathbf{x})|_{\hat{\mathbf{x}}}\|_{\text{op}}\right] \leq C \sqrt{d \xi''(1) \log_+ \frac{\xi^{(3)}(1)}{\xi''(1)}}, \quad (105)$$

$$\mathbb{E}\left[\max_{\mathbf{x} \in \mathbb{B}^d(1)} \|\mathbf{D}^2\mathbf{F}(\mathbf{x})|_{\mathbb{T} \otimes \mathbb{T}}\|_{\text{op}}\right] \leq C \sqrt{d \xi''(1) \log_+ \frac{\xi^{(3)}(1)}{\xi''(1)}}, \quad (106)$$

$$\mathbb{E}\left[\max_{\mathbf{x} \in \mathbb{B}^d(1)} \|\mathbf{D}^2\mathbf{F}(\mathbf{x})|_{\mathbb{T} \otimes \hat{\mathbf{x}}}\|_{\text{op}}\right] \leq C \sqrt{d \xi^{(3)}(1) \log_+ \frac{\xi^{(4)}(1)}{\xi^{(3)}(1)}}, \quad (107)$$

$$\mathbb{E}\left[\max_{\mathbf{x} \in \mathbb{B}^d(1)} \|\mathbf{D}^2\mathbf{F}(\mathbf{x})|_{\hat{\mathbf{x}} \otimes \hat{\mathbf{x}}}\|_{\text{op}}\right] \leq C \sqrt{d \xi^{(4)}(1) \log_+ \frac{\xi^{(5)}(1)}{\xi^{(4)}(1)}}. \quad (108)$$

Proof. All of these bounds follow from Lemma A.3 via the same argument. We will focus to be definite on the bound (104). We define the following Gaussian process indexed by $\mathbf{u} \in \mathbb{R}^n$, $\|\mathbf{u}\|_2 = 1$, $\mathbf{x} \in \mathbb{B}(1)$, $\mathbf{v} \in \mathbb{T}_{\mathbf{x}}$, $\|\mathbf{v}\|_2 = 1$:

$$Z(\mathbf{u}, \mathbf{v}, \mathbf{x}) := \langle \mathbf{u}, \mathbf{DF}(\mathbf{x})\mathbf{v} \rangle, \quad (109)$$

and notice that of course

$$\max_{\mathbf{x} \in \mathbb{B}^d(1)} \|\mathbf{DF}(\mathbf{x})|_{\mathbb{T}}\|_{\text{op}} = \max_{\mathbf{u} \in \mathbb{R}^n, \|\mathbf{u}\|_2=1} \max_{\mathbf{v} \in \mathbb{T}_{\mathbf{x}}, \|\mathbf{v}\|_2=1} \max_{\mathbf{x} \in \mathbb{B}^d(1)} Z(\mathbf{u}, \mathbf{v}, \mathbf{x}). \quad (110)$$

We next compute the canonical distance of Z , to get

$$d_Z(\mathbf{u}_1, \mathbf{v}_1, \mathbf{x}_1; \mathbf{u}_2, \mathbf{v}_2, \mathbf{x}_2)^2 = \frac{1}{2} \|\mathbf{u}_1 - \mathbf{u}_2\|^2 (\xi'(\|\mathbf{x}_1\|^2) + \xi'(\|\mathbf{x}_2\|^2)) + \langle \mathbf{u}_1, \mathbf{u}_2 \rangle d_{G, \perp}(\mathbf{v}_1, \mathbf{x}_1; \mathbf{v}_2, \mathbf{x}_2)^2, \quad (111)$$

where $d_{G, \perp}(\mathbf{v}_1, \mathbf{x}_1; \mathbf{v}_2, \mathbf{x}_2)^2$ is the canonical distance of the process $G_{\perp}(\mathbf{v}, \mathbf{x}) := \langle \mathbf{v}, \nabla F_1(\mathbf{x}) \rangle$ (this was derived in the proof of Lemma A.3 but will not be needed here). We bound this distance as

$$d_Z(\mathbf{u}_1, \mathbf{v}_1, \mathbf{x}_1; \mathbf{u}_2, \mathbf{v}_2, \mathbf{x}_2)^2 \leq \xi'(1) \|\mathbf{u}_1 - \mathbf{u}_2\|^2 + d_{G, \perp}(\mathbf{v}_1, \mathbf{x}_1; \mathbf{v}_2, \mathbf{x}_2)^2, \quad (112)$$

Therefore by Sudakov-Fernique, letting $\mathbf{g} \sim \mathcal{N}(0, \mathbf{I}_n)$,

$$\mathbb{E} \left[\max_{\mathbf{x} \in \mathbb{B}^d(1)} \|\mathbf{DF}(\mathbf{x})|_{\mathbb{T}}\|_{\text{op}} \right] \leq \mathbb{E} \left[\max_{\mathbf{u} \in \mathbb{S}^{n-1}} \sqrt{\xi'(1)} \langle \mathbf{g}, \mathbf{u} \rangle \right] + \mathbb{E} \max_{\mathbf{x} \in \mathbb{B}^d(1)} \|\mathbf{P}_{\mathbb{T}, \mathbf{x}} \nabla F_1(\mathbf{x})\|_2. \quad (113)$$

The bound (104) then follows from Lemma A.3. The other bounds are proved analogously. \square

We next use the above estimates to control the Lipschitz constant of F and its derivatives. Recall that, for $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d$, $\mathbf{x}_i \neq \mathbf{0}$, we let $\mathbf{R}_{\mathbf{x}_1, \mathbf{x}_2} \in \mathbb{R}^{d \times d}$ be the rotation in the plane $\mathbf{x}_1, \mathbf{x}_2$ such that $\mathbf{R}_{\mathbf{x}_1, \mathbf{x}_2} \mathbf{x}_1 / \|\mathbf{x}_1\| = \mathbf{x}_2 / \|\mathbf{x}_2\|$. Among the two rotations that leave the orthogonal complement of $\mathbf{x}_1, \mathbf{x}_2$ unchanged, choose the one with smallest angle, and break ties arbitrarily. (In other words, $\mathbf{R}_{\mathbf{x}_1, \mathbf{x}_2}$ is the parallel transport on \mathbb{S}^{d-1} .)

Remark A.2. By Borell-TIS inequality, the bounds in Proposition A.1, Lemma A.3 and Proposition A.4 with modified constants also hold with probability at least $1 - \exp(-C(\xi)d)$ for some $C(\xi) > 0$.

Definition A.5. For $\mathbf{x} \in \mathbb{R}^d$, we choose $\mathbf{U}_{\mathbf{x}} \in \mathbb{R}^{d \times (d-1)}$ such that $\mathbf{U}_{\mathbf{x}}^{\top} \mathbf{U}_{\mathbf{x}} = \mathbf{I}_{d-1}$, $\mathbf{U}_{\mathbf{x}}^{\top} \mathbf{x} = \mathbf{0}$ (i.e. $\mathbf{U}_{\mathbf{x}}$ is an orthonormal basis for the tangent space $\mathbb{T}_{\mathbf{x}}$.) Further, for $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d$, define $\mathbf{U}_{\mathbf{x}_1, \mathbf{x}_2} := \mathbf{R}_{\mathbf{x}_1, \mathbf{x}_2} \mathbf{U}_{\mathbf{x}_1}$.

For $\Omega \subseteq \mathbb{R}^d$, we define the following Lipschitz constants

$$\text{Lip}(F; \Omega) := \sup_{\mathbf{x}_1 \neq \mathbf{x}_2 \in \Omega} \frac{\|\mathbf{F}(\mathbf{x}_1) - \mathbf{F}(\mathbf{x}_2)\|}{\|\mathbf{x}_1 - \mathbf{x}_2\|_2}, \quad (114)$$

$$\text{Lip}_{\perp}(\mathbf{DF}; \Omega) := \sup_{\mathbf{x}_1 \neq \mathbf{x}_2 \in \Omega} \frac{\|\mathbf{DF}(\mathbf{x}_1) \mathbf{U}_{\mathbf{x}_1} - \mathbf{DF}(\mathbf{x}_2) \mathbf{U}_{\mathbf{x}_1, \mathbf{x}_2}\|_{\text{op}}}{\|\mathbf{x}_1 - \mathbf{x}_2\|_2}, \quad (115)$$

$$\text{Lip}_{\perp}(\nabla^2 F; \Omega) := \sup_{\mathbf{x}_1 \neq \mathbf{x}_2 \in \Omega} \max_{\ell \leq n} \frac{\|\mathbf{U}_{\mathbf{x}_1}^{\top} \nabla^2 F_{\ell}(\mathbf{x}_1) \mathbf{U}_{\mathbf{x}_1} - \mathbf{U}_{\mathbf{x}_1, \mathbf{x}_2}^{\top} \nabla^2 F_{\ell}(\mathbf{x}_2) \mathbf{U}_{\mathbf{x}_1, \mathbf{x}_2}\|_{\text{op}}}{\|\mathbf{x}_1 - \mathbf{x}_2\|_2}. \quad (116)$$

Lemma A.6. Assume $n \leq d$, and recall the notation $\log_+(t) := \max(\log(t); 1)$. Then there exists a universal constant C_1 , and a ξ -dependent constant $C_*(\xi)$ such that the following hold with probability at least $1 - \exp(-d/C_*(\xi))$:

$$\text{Lip}(\mathbf{F}; \mathbb{B}^d(1)) \leq C_1 \sqrt{d \xi''(1) \log_+ \frac{\xi^{(3)}(1)}{\xi''(1)}}, \quad (117)$$

$$\text{Lip}_\perp(\mathbf{DF}; \mathbb{B}^d(1)) \leq C_1 \sqrt{d \xi^{(4)}(1) \log_+ \frac{\xi^{(5)}(1)}{\xi^{(4)}(1)}}, \quad (118)$$

$$\text{Lip}_\perp(\nabla^2 \mathbf{F}; \mathbb{B}^d(1)) \leq C_*(\xi) \sqrt{d}. \quad (119)$$

Further, the following tighter Lipschitz constant holds on \mathbb{S}^{d-1} :

$$\text{Lip}(\mathbf{F}; \mathbb{S}^{d-1}) \leq C_1 \sqrt{d \xi'(1) \log_+ \frac{\xi''(1)}{\xi'(1)}}, \quad (120)$$

$$\text{Lip}_\perp(\mathbf{DF}; \mathbb{S}^{d-1}) \leq C_1 \sqrt{d \xi''(1) \log_+ \frac{\xi^{(3)}(1)}{\xi''(1)}}. \quad (121)$$

Proof. Throughout, we make use of Remark A.2.

- The bound (117) follows from Eqs. (104), (105).
- The bound (118) follows from Eqs. (106), (107), (108).
- The proof of Eq. (119) follows by bounding $\max_{\mathbf{x}} \|\nabla^3 F_\ell(\mathbf{x})\|$, which in turn can be done by the same technique as in the proof of Lemma A.3. Since we are not seeking a precise characterization of the constant $C_*(\xi)$ it suffices to notice that the canonical distance of the process $\langle \nabla^3 F_\ell(\mathbf{x}), \mathbf{v}^{\otimes 3} \rangle$ is bounded by a smooth function of the distances $\|\mathbf{v}_1 - \mathbf{v}_2\|$, $\|\mathbf{x}_1 - \mathbf{x}_2\|$.
- The bound (120) follows from Eqs. (104) by noting that, for $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{S}^{d-1}$, letting $\mathbf{x}(t)$ denote a geodesic on the sphere parametrized by arclength, we have

$$\mathbf{F}(\mathbf{x}_2) - \mathbf{F}(\mathbf{x}_1) = \int_0^{t_*} \mathbf{DF}(\mathbf{x}(t)) \dot{\mathbf{x}}(t) dt. \quad (122)$$

Here $t_* := \arccos(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle)$ and we notice that $\dot{\mathbf{x}}(t) \in \mathbb{T}_{\mathbf{x}(t)}$.

- Finally, for the bound (121) we consider a geodesic path between \mathbf{x}_1 and \mathbf{x}_2 in the sphere. For any two vectors $\mathbf{v}_1 \in \mathbb{T}_{\mathbf{x}_1}$, $\mathbf{v}_2 = \mathbf{R}_{\mathbf{x}_1, \mathbf{x}_2} \mathbf{v}_1 \in \mathbb{T}_{\mathbf{x}_2}$, let $\mathbf{v}(t)$ be the parallel transport of \mathbf{v}_1 along the geodesic connecting \mathbf{x}_1 to \mathbf{x}_2 . We then have

$$\mathbf{DF}(\mathbf{x}_2) \mathbf{v}_2 - \mathbf{DF}(\mathbf{x}_1) \mathbf{v}_1 = \int_0^{t_*} \{ \mathbf{D}^2 \mathbf{F}(\mathbf{x}(t)) \{ \dot{\mathbf{x}}(t), \mathbf{v}(t) \} + \mathbf{DF}(\mathbf{x}(t)) \dot{\mathbf{v}}(t) \} dt. \quad (123)$$

Hence the claim follows from Eqs. (105) and (106)

□

B Second moment method: Proof of Theorem 1

In this section we study the volume $V(\mathbf{u})$ of the set of solutions for $\mathbf{F}(\mathbf{x}) = \mathbf{u}$ for a given vector $\mathbf{u} \in \mathbb{R}^n$, as defined in Eqs. (9) and (11). We will use the Kac-Rice formula [Kac43, Ric45, Sub23] to prove that, for any deterministic sequence $\delta_d > 0$ with $\lim_{d \rightarrow \infty} \delta_d = 0$,

$$\lim_{n, d \rightarrow \infty} \sup_{\|\mathbf{u}\|_2 / \sqrt{n} \in [\xi_0^{1/2} - \delta_d, \xi_0^{1/2} + \delta_d]} \frac{\mathbb{E}\{V(\mathbf{u})^2\}}{\mathbb{E}\{V(\mathbf{u})\}^2} = 1. \quad (124)$$

This will complete the proof sketch given in Section 2.1.

We begin by computing the first moment. Recall that Vol_i is the Hausdorff measure of dimension i , or the counting measure when $i = 0$.

Proposition B.1. *For any $\mathbf{u} \in \mathbb{R}^n$,*

$$\mathbb{E}V(\mathbf{u}) = \mathcal{V}_{d-n-1} \left(\frac{\xi'(1)}{\xi(1)} \right)^{\frac{n}{2}} e^{-\frac{\|\mathbf{u}\|^2}{2\xi(1)}}, \quad (125)$$

where $\mathcal{V}_i := \text{Vol}_i(\{\mathbf{x} \in \mathbb{R}^{i+1} : \|\mathbf{x}\| = 1\}) = 2\pi^{\frac{i+1}{2}} / \Gamma(\frac{i+1}{2})$ is the volume of the sphere.

Proof. The proof here is identical to the case $\mathbf{u} = 0$ treated in [Sub23]. We repeat it for completeness, using the notation of the current paper.

Recall that, for $\mathbf{x} \in \mathbb{R}^d$ on the sphere $\mathbf{U}_{\mathbf{x}} \in \mathbb{R}^{d \times (d-1)}$ denotes a matrix whose columns form a basis of the tangent space $\mathbf{T}_{\mathbf{x}}$ to the sphere of radius $\|\mathbf{x}\|_2$ at \mathbf{x} (equivalently, to the orthogonal complement of \mathbf{x} in \mathbb{R}^d). We write $\mathbf{D}_{\perp} \mathbf{F}(\mathbf{x}) := \mathbf{D}\mathbf{F}(\mathbf{x})\mathbf{U}_{\mathbf{x}} \in \mathbb{R}^{n \times d-1}$.

By a variant of the Kac-Rice formula in [AW09, Theorem 6.8],

$$\mathbb{E}V(\mathbf{u}) = \int_{S^{d-1}} \mathbb{E} \left[J(\mathbf{D}_{\perp} \mathbf{F}(\mathbf{x})) \mid \mathbf{F}(\mathbf{x}) = \mathbf{u} \right] p_{\mathbf{F}(\mathbf{x})}(\mathbf{u}) d\text{Vol}_{d-1}(\mathbf{x}), \quad (126)$$

where $J(\mathbf{A}) = \sqrt{\det \mathbf{A}\mathbf{A}^{\top}}$, $d\text{Vol}_{d-1}(\mathbf{x})$ is the $(d-1)$ -dimensional volume element on the sphere and

$$p_{\mathbf{F}(\mathbf{x})}(\mathbf{u}) = (2\pi\xi(1))^{-\frac{n}{2}} e^{-\frac{\|\mathbf{u}\|^2}{2\xi(1)}}$$

is the density of $\mathbf{F}(\mathbf{x})$ at \mathbf{u} . Using the Kac-Rice formula requires checking that certain conditions are satisfied, this was done in Remark 3 of [Sub23].

Recall that by Lemma 4.1, $\mathbf{D}_{\perp} \mathbf{F}(\mathbf{x}) \stackrel{d}{=} \sqrt{\xi'(q)} \mathbf{Z}$ where $q = \|\mathbf{x}\|^2$, $\mathbf{Z} \sim \text{GOE}(n, d-1)$ and $\mathbf{D}_{\perp} \mathbf{F}(\mathbf{x})$ is independent of $\mathbf{F}(\mathbf{x})$. Thus,

$$\mathbb{E}V(\mathbf{u}) = \mathcal{V}_{d-1} \left(\frac{\xi'(1)}{2\pi\xi(1)} \right)^{n/2} e^{-\frac{\|\mathbf{u}\|^2}{2\xi(1)}} \mathbb{E}J(\mathbf{Z}). \quad (127)$$

In the case $\xi(t) = t$ and $\mathbf{u} = 0$, the above is the volume of the set of points on the sphere orthogonal to n independent Gaussian vectors, hence equal to \mathcal{V}_{d-n-1} . Hence,

$$\mathcal{V}_{d-n-1} = \mathcal{V}_{d-1} \left(\frac{1}{2\pi} \right)^{n/2} \mathbb{E}J(\mathbf{Z}) \quad (128)$$

and (125) follows. \square

The next theorem establishes the desired second moment bound. To verify this note that the function $\Phi(r)$ defined here matches $\Psi(r; \alpha, \xi)$ of Eq. (7), if we replace γ^2 by $\alpha\xi(0)$ and $\xi(t)$ by $\xi_{>0}(t) = \xi(t) - \xi(0)$, and that condition (129) is implied by $\alpha < \alpha_{\text{LB}}(\xi)$.

Theorem 7. Assume that $\xi(0) = \xi'(0) = \xi''(0) = 0$. Let $\gamma \geq 0$ and $\alpha \in (0, 1]$ and define

$$\Phi(r) := \frac{1}{2} \log(1 - r^2) - \frac{\alpha}{2} \log\left(1 - \frac{\xi(r)^2}{\xi(1)^2}\right) + \frac{\gamma^2}{\xi(1)} - \frac{\gamma^2}{\xi(1) + \xi(r)}.$$

If for any $\varepsilon > 0$,

$$\Phi(0) > \sup_{r \in (\varepsilon, 1)} \Phi(r), \quad (129)$$

then for any deterministic sequence $\delta_d > 0$ such that $\lim_{d \rightarrow \infty} \delta_d = 0$,

$$\lim_{d, n \rightarrow \infty} \frac{\mathbb{E}(\mathbf{V}(\mathbf{u})^2)}{(\mathbb{E}\mathbf{V}(\mathbf{u}))^2} = 1 \quad (130)$$

uniformly in $n = n(d) < d - 1$ and $\mathbf{u} = \mathbf{u}^{(n, d)} \in \mathbb{R}^n$ such that $\lim_{d \rightarrow \infty} n(d) \rightarrow \alpha$ and $\sqrt{d}(\gamma - \delta_d) \leq \|\mathbf{u}\| \leq \sqrt{d}(\gamma + \delta_d)$.

Proof. To simplify the proof, throughout we will assume that $n \leq d - 2$. The case in which $n = d - 1$ infinitely often (in particular $\alpha = 1$) was treated for $\mathbf{u} = \mathbf{0}$ in [Sub23] and the argument can be adapted to our situation with $\mathbf{u} \neq \mathbf{0}$, but to save space we will omit the proof. We refer the interested reader to Lemmas 8 and 9 in [Sub23].

For any union of intervals $I \subset [-1, 1]$, define

$$\mathbf{V}^{(2)}(\mathbf{u}, I) := \text{Vol}_{2(d-n-1)}\left(\left\{(\mathbf{x}, \mathbf{y}) : \langle \mathbf{x}, \mathbf{y} \rangle \in I, \mathbf{F}(\mathbf{x}) = \mathbf{F}(\mathbf{y}) = \mathbf{u}, \mathbf{x}, \mathbf{y} \in \mathbb{S}^{d-1}\right\}\right).$$

As in the proof of Proposition B.1, we may apply the Kac-Rice formula to the random field $(\mathbf{x}, \mathbf{y}) \rightarrow (\mathbf{F}(\mathbf{x}), \mathbf{F}(\mathbf{y}))$ defined on the domain

$$\mathcal{D}(I) := \{(\mathbf{x}, \mathbf{y}) : \langle \mathbf{x}, \mathbf{y} \rangle \in I, \|\mathbf{x}\| = \|\mathbf{y}\| = 1\} \quad (131)$$

to compute the expectation of $\mathbf{V}^{(2)}(\mathbf{u}, I)$, see Section 3 of [Sub23] for the details. This yields that

$$\begin{aligned} \mathbb{E} \mathbf{V}^{(2)}(\mathbf{u}, I) &= (\xi'(1))^n \int_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}(I)} p_{\mathbf{F}(\mathbf{x}), \mathbf{F}(\mathbf{y})}(\mathbf{u}, \mathbf{u}) T(\mathbf{u}, \langle \mathbf{x}, \mathbf{y} \rangle) d\text{Vol}_{d-1}(\mathbf{x}) d\text{Vol}_{d-1}(\mathbf{y}) \\ &= \mathcal{V}_{d-1} \cdot (\xi'(1))^n \int_{\{\mathbf{y} \in \mathbb{S}^{d-1} : \langle \mathbf{x}, \mathbf{y} \rangle \in I\}} p_{\mathbf{F}(\mathbf{x}), \mathbf{F}(\mathbf{y})}(\mathbf{u}, \mathbf{u}) T(\mathbf{u}, \langle \mathbf{x}, \mathbf{y} \rangle) d\text{Vol}_{d-1}(\mathbf{y}), \end{aligned} \quad (132)$$

where in the second line \mathbf{x} is an arbitrary point on the sphere, and writing $\langle \mathbf{x}, \mathbf{y} \rangle = r$, we have

$$p_{\mathbf{F}(\mathbf{x}), \mathbf{F}(\mathbf{y})}(\mathbf{u}, \mathbf{u}) = (2\pi)^{-n} (\xi(1)^2 - \xi(r)^2)^{-\frac{n}{2}} e^{-\frac{\|\mathbf{u}\|^2}{\xi(1) + \xi(r)}}$$

is the density of $(\mathbf{F}(\mathbf{x}), \mathbf{F}(\mathbf{y}))$ at (\mathbf{u}, \mathbf{u}) since the covariance matrix of $(F_i(\mathbf{x}), F_i(\mathbf{y}))$ is

$$\begin{pmatrix} \xi(1) & \xi(r) \\ \xi(r) & \xi(1) \end{pmatrix},$$

and we define (always for $\langle \mathbf{x}, \mathbf{y} \rangle = r$)

$$T(\mathbf{u}, r) := \mathbb{E} \left[J\left(\frac{\mathbf{D}_\perp \mathbf{F}(\mathbf{x})}{\sqrt{\xi'(1)}}\right) J\left(\frac{\mathbf{D}_\perp \mathbf{F}(\mathbf{y})}{\sqrt{\xi'(1)}}\right) \middle| \mathbf{F}(\mathbf{x}) = \mathbf{F}(\mathbf{y}) = \mathbf{u} \right]. \quad (133)$$

By rotational invariance $T(\mathbf{u}, r)$ does not depend on the choice of \mathbf{U}_x and \mathbf{U}_y and of \mathbf{x} and \mathbf{y} as long as $\langle \mathbf{x}, \mathbf{y} \rangle = r$.

Using the co-area formula with the function $\rho(\mathbf{y}) = \langle \mathbf{x}, \mathbf{y} \rangle$, we may express the integral w.r.t. \mathbf{y} in Eq. (132) as a one-dimensional integral over a parameter r (the volume of the inverse-image $\rho^{-1}(r)$ and the inverse of the Jacobian are given by $\mathcal{V}_{d-2}(1-r^2)^{(d-2)/2}$ and $(1-r^2)^{-1/2}$, respectively). This yields

$$\mathbb{E} V^{(2)}(\mathbf{u}, I) = \mathcal{V}_{d-1} \mathcal{V}_{d-2} (2\pi)^{-n} \int_I \left(\frac{\xi'(1)^2}{\xi(1)^2 - \xi(r)^2} \right)^{n/2} (1-r^2)^{\frac{d-3}{2}} T(\mathbf{u}, r) e^{-\frac{\|\mathbf{u}\|^2}{\xi(1) + \xi(r)}} dr. \quad (134)$$

From Eqs. (127) and (134), we get

$$\frac{\mathbb{E} V^{(2)}(\mathbf{u}, I)}{(\mathbb{E} V(\mathbf{u}))^2} = \frac{\mathcal{V}_{d-2}}{\mathcal{V}_{d-1}} \int_I \frac{T(\mathbf{u}, r)}{(\mathbb{E} J(\mathbf{Z}))^2} \left(1 - \frac{\xi(r)^2}{\xi(1)^2} \right)^{-n/2} (1-r^2)^{(d-3)/2} e^{\frac{\|\mathbf{u}\|^2}{\xi(1)} - \frac{\|\mathbf{u}\|^2}{\xi(1) + \xi(r)}} dr. \quad (135)$$

The next lemma provide a useful bound on the integrand.

Lemma B.2. *Assume that $\xi(0) = \xi'(0) = 0$. Define*

$$\kappa = \kappa(r) = \frac{\xi'(r)}{\xi(1) + \xi(r)} \sqrt{\frac{1-r^2}{\xi'(1)}}. \quad (136)$$

Then for some universal constant c and sequence $\tau_d \rightarrow 0$, for any $r \in (-1, 1)$,

$$\frac{T(\mathbf{u}, r)}{(\mathbb{E} J(\mathbf{Z}))^2} \leq (1 + \tau_d) (1 + 2r^2 \sqrt{d}) (1 + c\kappa(r) \|\mathbf{u}\| + c\kappa(r)^2 \|\mathbf{u}\|^2).$$

Before proving this lemma, let us complete the proof of the theorem. Note that $\frac{\mathcal{V}_{d-2}}{\mathcal{V}_{d-1}} = \sqrt{\frac{d}{2\pi}} (1 + o_d(1))$. Hence, we have that

$$\frac{1}{d} \log \frac{\mathbb{E} V^{(2)}(\mathbf{u}, I)}{(\mathbb{E} V(\mathbf{u}))^2} \leq \frac{1}{d} \log \int_I \nu(\mathbf{u}, r) dr + o_d(1),$$

where the $o_d(1)$ term is independent of \mathbf{u} and n , and

$$\nu(\mathbf{u}, r) := \left(1 - \frac{\xi(r)^2}{\xi(1)^2} \right)^{-n/2} (1-r^2)^{(d-3)/2} (1 + c\kappa(r) \|\mathbf{u}\| + c\kappa(r)^2 \|\mathbf{u}\|^2) e^{\frac{\|\mathbf{u}\|^2}{\xi(1)} - \frac{\|\mathbf{u}\|^2}{\xi(1) + \xi(r)}}.$$

Note that $\nu(\mathbf{u}, r) \leq (1-r^2)^{-1/2} e^{d\Phi(r) + o(d)}$, uniformly in $n \leq d-2$, $\sqrt{d}(\gamma - \delta_d) \leq \|\mathbf{u}\| \leq \sqrt{d}(\gamma + \delta_d)$ and $r \in (-1, 1)$. Assuming that $n = n(d) \rightarrow \alpha$, since $(1-r^2)^{-1/2}$ is integrable on $(-1, 1)$, we have that for any $\varepsilon > 0$, uniformly in $\sqrt{d}(\gamma - \delta_d) \leq \|\mathbf{u}\| \leq \sqrt{d}(\gamma + \delta_d)$,

$$\limsup_{d \rightarrow \infty} \frac{1}{d} \log \frac{\mathbb{E} V^{(2)}(\mathbf{u}, [-1, 1] \setminus (-\varepsilon, \varepsilon))}{(\mathbb{E} V(\mathbf{u}))^2} \leq \sup_{r \in [\varepsilon, 1]} \Phi(r) < \Phi(0) = 0,$$

where the second inequality follows by assumption.

Of course, the same bound holds with some sequence $\varepsilon_d \rightarrow 0$ instead of ε . Since the ratio of moments in (130) is lower bounded by 1, to prove the same equation it remains to show that, uniformly,

$$\limsup_{d \rightarrow \infty} \frac{\mathbb{E} V^{(2)}(\mathbf{u}, (-\varepsilon_d, \varepsilon_d))}{(\mathbb{E} V(\mathbf{u}))^2} \leq 1. \quad (137)$$

Note that $\frac{\xi(r)}{\xi(1)} = O(r^3)$ and $\kappa(r) = O(r^2)$ for small r , since we assume that $\xi(0) = \xi'(0) = \xi''(0) = 0$, and recall that $\frac{\nu_{d-2}}{\nu_{d-1}} = \sqrt{\frac{d}{2\pi}}(1 + o_d(1))$. From (135) we obtain that the ratio in (137) is equal to

$$\begin{aligned} & (1 + o_d(1)) \sqrt{\frac{d}{2\pi}} \int_{-\varepsilon_d}^{\varepsilon_d} \frac{T(\mathbf{u}, r)}{(\mathbb{E}J(\mathbf{Z}))^2} (1 - O(r^6))^{-n/2} (1 - r^2)^{d/2} e^{\frac{\|\mathbf{u}\|^2}{\xi(1)} O(r^3)} dr \\ &= (1 + o_d(1)) \sqrt{\frac{d}{2\pi}} \int_{-\varepsilon_d}^{\varepsilon_d} \exp\left(-\frac{d}{2}r^2 + 2\sqrt{d}r^2 + \frac{\gamma^2 d}{\xi(1)} O(r^3) + dO(r^6) + c\gamma\sqrt{d}O(r^2) + c\gamma^2 dO(r^4)\right) dr \end{aligned}$$

where we used Lemma B.2 and that $\log(1+t) \leq t$. This proves (137) and completes the proof. \square

Proof of Lemma B.2. Fix some $r \in (-1, 1)$. Let \mathbf{M}^1 and \mathbf{M}^2 be jointly Gaussian, centered random matrices such that

$$\text{cov}(\mathbf{M}_{ij}^t, \mathbf{M}_{kl}^s) = \delta_{ik}\delta_{jl} \cdot \begin{cases} 1 & t = s, j < d-1 \\ 1 - \frac{\xi(1)}{\xi'(1)} \frac{\xi'(r)^2(1-r^2)}{\xi(1)^2 - \xi(r)^2} & t = s, j = d-1 \\ \frac{\xi'(r)}{\xi'(1)} & t \neq s, j < d-1 \\ r \frac{\xi'(r)}{\xi'(1)} - \frac{\xi''(r)}{\xi'(1)}(1-r^2) - \frac{\xi(r)}{\xi'(1)} \frac{\xi'(r)^2(1-r^2)}{\xi(1)^2 - \xi(r)^2} & t \neq s, j = d-1 \\ 0 & \text{otherwise.} \end{cases} \quad (138)$$

Suppose \mathbf{x} and \mathbf{y} are such that $\langle \mathbf{x}, \mathbf{y} \rangle = r \in (-1, 1)$. In Lemma 6 of [Sub23] it was shown that for an appropriate choice of $\mathbf{U}_{\mathbf{x}}$ and $\mathbf{U}_{\mathbf{y}}$ (bases of the tangent spaces $\mathbb{T}_{\mathbf{x}}, \mathbb{T}_{\mathbf{y}}$),

$$\left(\frac{\mathbf{D}_{\perp} \mathbf{F}(\mathbf{x})}{\sqrt{\xi'(1)}}, \frac{\mathbf{D}_{\perp} \mathbf{F}(\mathbf{y})}{\sqrt{\xi'(1)}} \right) \stackrel{d}{=} \left(\mathbf{M}^1 - \kappa \mathbf{u} \cdot \mathbf{e}_{d-1}^{\top}, \mathbf{M}^2 + \kappa \mathbf{u} \cdot \mathbf{e}_{d-1}^{\top} \right),$$

where $\mathbf{e}_{d-1} \in \mathbb{R}^{d-1}$ is the last (column) vector in the standard basis and κ is defined in (136). (To be precise, only the case $\mathbf{u} = 0$ was treated there, but only the conditional expectation depends on \mathbf{u} and is easily computed in the same way.)

Let \mathbf{O} be some orthogonal matrix such that $\mathbf{O}\mathbf{u} = (0, \dots, 0, \|\mathbf{u}\|)^{\top}$. Since $J(\mathbf{A}) = J(\mathbf{O}\mathbf{A})$ and $(\mathbf{O}\mathbf{M}^1, \mathbf{O}\mathbf{M}^2) \stackrel{d}{=} (\mathbf{M}^1, \mathbf{M}^2)$,

$$\left(J\left(\frac{\mathbf{D}_{\perp} \mathbf{F}(\mathbf{x})}{\sqrt{\xi'(1)}}\right), J\left(\frac{\mathbf{D}_{\perp} \mathbf{F}(\mathbf{y})}{\sqrt{\xi'(1)}}\right) \right) \stackrel{d}{=} \left(J(\hat{\mathbf{M}}^1), J(\hat{\mathbf{M}}^2) \right),$$

where we denote $\hat{\mathbf{M}}^1 := \mathbf{M}^1 - \kappa\|\mathbf{u}\|\mathbf{e}_n \mathbf{e}_{d-1}^{\top}$ and $\hat{\mathbf{M}}^2 := \mathbf{M}^2 + \kappa\|\mathbf{u}\|\mathbf{e}_n \mathbf{e}_{d-1}^{\top}$. Recall that for any matrix \mathbf{A} , $J(\mathbf{A}) = \prod_{i=1}^n \Theta_i(\mathbf{A})$, where we define $\Theta_i(\mathbf{A})$ as the norm of the projection of the i -th row of \mathbf{A} to the orthogonal space to its first $i-1$ rows. Note that $\Theta_i(\hat{\mathbf{M}}^j) = \Theta_i(\mathbf{M}^j)$ for $i < n$ and $\Theta_n(\hat{\mathbf{M}}^j) \leq \Theta_n(\mathbf{M}^j) + \kappa\|\mathbf{u}\|$. Therefore, using independence of the rows and the orthogonal invariance of Gaussians,

$$\begin{aligned} T(\mathbf{u}, r) &\leq \mathbb{E} \left[\prod_{i=1}^{n-1} \Theta_i(\mathbf{M}^1) (\Theta_n(\mathbf{M}^1) + \kappa\|\mathbf{u}\|) \prod_{i=1}^{n-1} \Theta_i(\mathbf{M}^2) (\Theta_n(\mathbf{M}^2) + \kappa\|\mathbf{u}\|) \right] \\ &= \mathbb{E} [J(\mathbf{M}^1) J(\mathbf{M}^2)] \left(1 + \frac{2\kappa\|\mathbf{u}\| \mathbb{E} \Theta_n(\mathbf{M}^1)}{\mathbb{E} [\Theta_n(\mathbf{M}^1) \Theta_n(\mathbf{M}^2)]} + \frac{\kappa^2 \|\mathbf{u}\|^2}{\mathbb{E} [\Theta_n(\mathbf{M}^1) \Theta_n(\mathbf{M}^2)]} \right). \end{aligned}$$

Using that $\Theta_n(\mathbf{M}^i) \sim \chi_{n-d}$, we have that, for some universal constant c ,

$$T(\mathbf{u}, r) \leq \mathbb{E}[J(\mathbf{M}^1)J(\mathbf{M}^2)](1 + c\kappa(r)\|\mathbf{u}\| + c\kappa(r)^2\|\mathbf{u}\|^2).$$

Lemma 10 of [Sub23] bounds $T(\mathbf{u}, r)$ for $\mathbf{u} = 0$, which is exactly equal to $\mathbb{E}[J(\mathbf{M}^1)J(\mathbf{M}^2)]$. Precisely, it states that

$$T(\mathbf{0}, r) = \mathbb{E}[J(\mathbf{M}^1)J(\mathbf{M}^2)] \leq (1 + \tau_d)(1 + 2r^2\sqrt{d}),$$

for some universal $\tau_d \rightarrow 0$. This completes the proof. \square

C Large- p Asymptotics of the lower bound: Eq. (12)

For $\xi(t) = \xi_0 + t^p$, the function $\Psi(r) = \Psi(r; \alpha, \xi_0, p)$ of Eq. (7) reads

$$\Psi(r; \alpha, \xi_0, p) := \frac{1}{2} \log(1 - r^2) - \frac{\alpha}{2} \log(1 - r^{2p}) + \alpha \xi_0 \frac{r^p}{1 + r^p}, \quad (139)$$

and since $\Psi''(0; \alpha, \xi_0, p) = -1$ for all α , we have

$$\alpha_{\text{LB}}(\xi_0, p) := \inf \left\{ \alpha \geq 0 : \sup_{r \in [0, 1]} \Psi(r; \alpha, \xi_0, p) > 0 \right\}. \quad (140)$$

To get an upper bound on $\alpha_{\text{LB}}(\xi_0, p)$, we choose $r = 1 - M/p$ for some constant $M > 0$. Recalling that $\xi_0 = \gamma_0 \log p$ for some constant $\gamma_0 > 1$, it is easy to compute that, as $p \rightarrow \infty$,

$$\Psi\left(r = 1 - \frac{M}{p}\right) = -\frac{1}{2} \log p + \alpha \xi_0 \frac{e^{-M}}{1 + e^{-M}} + O(1) \quad (141)$$

$$= \left(-\frac{1}{2} + \frac{\alpha \gamma_0}{1 + e^M} \right) \log p + O(1). \quad (142)$$

Therefore, for all p large enough

$$\alpha_{\text{LB}}(\xi_0 = \gamma_0 \log p, p) \leq \frac{e^M + 1}{2\gamma_0}. \quad (143)$$

Since $M > 0$ is arbitrary, we get $\limsup_{p \rightarrow \infty} \alpha_{\text{LB}}(\xi_0 = \gamma_0 \log p, p) \leq 1/\gamma_0$.

Next we prove the lower bound. For that purpose, we will fix $\gamma_0 > 1$, assume that $\xi_0 = \gamma_0 \log p$, $\alpha \leq (1 - \varepsilon)\gamma_0^{-1}$ and prove that $\sup_{r \in [0, 1]} \Psi(r; \alpha, \xi_0, p) \leq 0$ for all p large enough. Note that, under these definitions, $\Psi(r; \alpha, \xi_0, p) \leq \Psi_*(r)$, where

$$\Psi_*(r) := \frac{1}{2} \log(1 - r^2) - \frac{1}{2} \log(1 - r^{2p}) + (1 - \varepsilon) \frac{r^p}{1 + r^p} \log p. \quad (144)$$

We split the maximization over r in two regions (depending on a constant Δ to be fixed below):

1. $0 \leq r \leq 1 - p^{-1+\Delta}$. For an absolute constant C , we have

$$\Psi_*(r) \leq -\frac{1}{2} r^2 + C(\log p) r^p \quad (145)$$

$$\leq \left(-\frac{1}{2} + C'(\log p) e^{-p^\Delta} \right) r^2. \quad (146)$$

2. $1 - p^{-1+\Delta} < r < 1$. Setting $r = 1 - x/p$, $x \in (0, p^\Delta)$, we have

$$\begin{aligned}\Psi_*(r) &\leq \frac{1}{2} \log\left(\frac{2x}{p}\right) - \frac{1}{2} \log(1 - e^{-2x}) + \frac{1 - \varepsilon}{1 + e^x} \log p \\ &\leq -\frac{\varepsilon}{2} \log p + \frac{1}{2} \log(2x) - \frac{1}{2} \log(1 - e^{-2x}) \\ &\leq -\frac{\varepsilon}{2} \log p + \frac{1}{2} \log(2p^\Delta) + 1 \\ &\leq \left(\Delta - \frac{\varepsilon}{2}\right) \log p + 2.\end{aligned}$$

Therefore the claim follows by choosing $\Delta = \varepsilon/4$.

D Upper bounds on the existence threshold: Proof of Theorem 2

Note that

$$\text{Sol}_{n,d}(\varepsilon) \neq \emptyset \iff \min_{\mathbf{x} \in \mathbb{S}^{d-1}} \max_{\mathbf{v} \in \mathbb{S}^{n-1}} \langle \mathbf{F}(\mathbf{x}), \mathbf{v} \rangle \leq \sqrt{n\xi(1)\varepsilon}. \quad (147)$$

We define the two Gaussian processes on $\mathbb{S}^{d-1} \times \mathbb{S}^{n-1}$,

$$f_1(\mathbf{x}, \mathbf{v}) := \langle \mathbf{F}(\mathbf{x}), \mathbf{v} \rangle + \sqrt{\xi(1)} z, \quad f_2(\mathbf{x}, \mathbf{v}) := F_1(\mathbf{x}) + \sqrt{\xi(1)} \langle \mathbf{v}, \mathbf{g} \rangle,$$

where $\mathbf{g} \sim \mathbf{N}(0, \mathbf{I}_n)$ and $z \sim \mathbf{N}(0, 1)$, independent of each other and $\mathbf{F}(\mathbf{x})$. Clearly,

$$\begin{aligned}\mathbb{E} f_1(\mathbf{x}^1, \mathbf{v}^1) f_1(\mathbf{x}^2, \mathbf{v}^2) &= \xi(\langle \mathbf{x}^1, \mathbf{x}^2 \rangle) \langle \mathbf{v}^1, \mathbf{v}^2 \rangle + \xi(1), \\ \mathbb{E} f_2(\mathbf{x}^1, \mathbf{v}^1) f_2(\mathbf{x}^2, \mathbf{v}^2) &= \xi(\langle \mathbf{x}^1, \mathbf{x}^2 \rangle) + \xi(1) \langle \mathbf{v}^1, \mathbf{v}^2 \rangle.\end{aligned}$$

For $\mathbf{x}^1 = \mathbf{x}^2$ the two covariance functions coincide and for general $\mathbf{x}^1, \mathbf{x}^2 \in \mathbb{S}^{d-1}$,

$$\mathbb{E} f_1(\mathbf{x}^1, \mathbf{v}^1) f_1(\mathbf{x}^2, \mathbf{v}^2) \geq \mathbb{E} f_2(\mathbf{x}^1, \mathbf{v}^1) f_2(\mathbf{x}^2, \mathbf{v}^2).$$

Hence, by Gordon's Gaussian comparison inequality [Gor85], we have

$$\min_{\mathbf{x} \in \mathbb{S}^{d-1}} \max_{\mathbf{v} \in \mathbb{S}^{n-1}} f_1(\mathbf{x}^2, \mathbf{v}^2) \succeq \min_{\mathbf{x} \in \mathbb{S}^{d-1}} \max_{\mathbf{v} \in \mathbb{S}^{n-1}} f_2(\mathbf{x}^2, \mathbf{v}^2),$$

where \succeq denotes stochastic domination. In particular, for any $c > 0$,

$$\begin{aligned}e^{\frac{1}{2}c^2 d \xi(1)} \cdot \mathbb{E} \exp \left\{ -c\sqrt{d} \min_{\mathbf{x} \in \mathbb{S}^{d-1}} \max_{\mathbf{v} \in \mathbb{S}^{n-1}} \langle \mathbf{F}(\mathbf{x}), \mathbf{v} \rangle \right\} &= \mathbb{E} \exp \left\{ -c\sqrt{d} \min_{\mathbf{x} \in \mathbb{S}^{d-1}} \max_{\mathbf{v} \in \mathbb{S}^{n-1}} f_1(\mathbf{x}, \mathbf{v}) \right\} \\ &\leq \mathbb{E} \exp \left\{ -c\sqrt{d} \left(\min_{\mathbf{x} \in \mathbb{S}^{d-1}} F_1(\mathbf{x}) + \sqrt{\xi(1)} \|\mathbf{g}\| \right) \right\}.\end{aligned}$$

Using Eq. (147) and Markov's inequality,

$$\begin{aligned}\mathbb{P}\left(\text{Sol}_{n,d}(\varepsilon) \neq \emptyset\right) &\leq e^{c\sqrt{dn\xi(1)\varepsilon}} \mathbb{E} \exp \left\{ -c\sqrt{d} \min_{\mathbf{x} \in \mathbb{S}^{d-1}} \max_{\mathbf{v} \in \mathbb{S}^{n-1}} \langle \mathbf{F}(\mathbf{x}), \mathbf{v} \rangle \right\} \\ &\leq e^{d\varepsilon' - \frac{1}{2}c^2 d \xi(1)} \mathbb{E} e^{-c\sqrt{d\xi(1)}\|\mathbf{g}\|} \mathbb{E} \exp \left\{ c\sqrt{d} \max_{\mathbf{x} \in \mathbb{S}^{d-1}} F_1(\mathbf{x}) \right\},\end{aligned} \quad (148)$$

where in the second line we defined $\varepsilon' := c\sqrt{\alpha\xi(1)\varepsilon}$.

By Cramer's theorem, for $t \in (0, 1]$,

$$\mathbb{P}(\|\mathbf{g}\| < t\sqrt{n}) = \mathbb{P}(\|\mathbf{g}\|^2 < t^2n) = \exp\left(-n\frac{t^2-1}{2} + n\log t + o(n)\right)$$

and therefore

$$\lim_{n,d \rightarrow \infty} \frac{1}{d} \log \mathbb{E} e^{-c\sqrt{d\xi(1)}\|\mathbf{g}\|} = \varphi_2(c, \alpha). \quad (149)$$

Let $F_1^p(\mathbf{x})$ be the random process as defined in (3) with $\xi(t) = t^p$ and let $z \sim \mathbf{N}(0, 1)$ be independent of $F_1^p(\mathbf{x})$. Then, $F_1(\mathbf{x}) \stackrel{d}{=} F_1^p(\mathbf{x}) + \xi_0 z$ as processes. By Markov's inequality and [Auf13, Theorem 2.8], for $t \geq 0$,

$$\mathbb{P}\left(\max_{\mathbf{x} \in \mathbb{S}^{d-1}} F_1^p(\mathbf{x}) \geq \sqrt{d}(E+t)\right) \leq e^{d\Theta_p(E+t)+o(d)}.$$

For $s \geq 0$, of course, $\mathbb{P}(\xi_0 z \geq \sqrt{dt}) \leq e^{-\frac{1}{2}\frac{dt^2}{\xi_0^2} + o(d)}$. Combining these facts, one easily sees that

$$\lim_{n,d \rightarrow \infty} \frac{1}{d} \log \mathbb{E} \exp\left\{c\sqrt{d} \max_{\mathbf{x} \in \mathbb{S}^{d-1}} F_1(\mathbf{x})\right\} = \varphi_1(c, p). \quad (150)$$

The lemma follows from (148), (149) and (150).

E Analysis of gradient descent: Proof of Theorem 3

Since calculations are more transparent in the case of gradient flow, we will first treat this case, and then outline the modifications that arise for discrete time. Both arguments are standard in the machine learning literature [DZPS18, COB19, OS20, ADH⁺19, AZLL19, BMR21] and we therefore present them succinctly.

E.1 Gradient flow

For gradient flow the time $t \in \mathbb{R}_{\geq 0}$ is continuous and the state is updated according to

$$\dot{\mathbf{x}}(t) = -\mathbf{P}_{\top, \mathbf{x}(t)} \nabla H(\mathbf{x}(t)). \quad (151)$$

(Recall the definition of cost function $H(\mathbf{x}) = \|\mathbf{F}(\mathbf{x})\|_2^2/2$.)

Lemma E.1. *Let $\lambda_0 := \sigma_{\min}(\mathbf{DF}(\mathbf{x}(0))|_{\top, \mathbf{x}(0)})$ and $L_n := \text{Lip}_{\perp}(\mathbf{DF}; \mathbb{S}^{d-1})$, with the Lipschitz constant of Definition A.5. If*

$$L_n \|\mathbf{F}(\mathbf{x}(0))\|_2 < \frac{\lambda_0^2}{4}, \quad (152)$$

then for all $t \geq 0$,

$$\|\mathbf{F}(\mathbf{x}(t))\|_2^2 \leq \|\mathbf{F}(\mathbf{x}(0))\|_2^2 e^{-\lambda_0^2 t/4}. \quad (153)$$

Proof. We define $\mathbf{K}(\mathbf{x}) := \mathbf{DF}(\mathbf{x})\mathbf{P}_{\top, \mathbf{x}}\mathbf{DF}(\mathbf{x})^{\top}$. Then the gradient flow equation implies

$$\frac{d}{dt} \mathbf{F}(\mathbf{x}(t)) = -\mathbf{K}(\mathbf{x}(t)) \cdot \mathbf{F}(\mathbf{x}(t)). \quad (154)$$

Let

$$B := \left\{ \mathbf{x} \in \mathbb{S}^{d-1} : \sigma_{\min}(\mathbf{DF}(\mathbf{x})|_{\mathbb{T}, \mathbf{x}}) \leq \frac{\lambda_0}{2} \right\}. \quad (155)$$

Obviously $\mathbf{x}(0) \in B^c$. Further, letting $d_{\mathbb{S}^{d-1}}(\mathbf{x}_1, \mathbf{x}_2) := \arccos(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle)$ denote the geodesic distance on the unit sphere, we have

$$d_{\mathbb{S}^{d-1}}(\mathbf{x}(0), B) := \inf \left\{ d_{\mathbb{S}^{d-1}}(\mathbf{x}(0), \mathbf{x}) : \mathbf{x} \in B \right\} \geq \frac{\lambda_0}{2L_n}. \quad (156)$$

Define $t_* := \inf\{t : \mathbf{x}(t) \in B\}$. For all $t \leq t_*$, we have

$$\frac{d}{dt} \|\mathbf{F}(\mathbf{x}(t))\|_2^2 \leq -\frac{\lambda_0^2}{4} \|\mathbf{F}(\mathbf{x}(t))\|_2^2, \quad (157)$$

which implies Eq. (153) for all $t \leq t_*$.

We next prove that $t_* = \infty$. Indeed, note that for $t \leq t_*$,

$$\frac{d}{dt} \|\mathbf{F}(\mathbf{x}(t))\|_2 = -\frac{1}{\|\mathbf{F}(\mathbf{x}(t))\|_2} \left\| \mathbf{P}_{\mathbb{T}, \mathbf{x}(t)} \mathbf{DF}(\mathbf{x}(t))^\top \mathbf{F}(\mathbf{x}(t)) \right\|_2^2 \quad (158)$$

$$\leq -\sigma_{\min}(\mathbf{DF}(\mathbf{x}(t)) \mathbf{P}_{\mathbb{T}, \mathbf{x}(t)}) \left\| \mathbf{P}_{\mathbb{T}, \mathbf{x}(t)} \mathbf{DF}(\mathbf{x}(t))^\top \mathbf{F}(\mathbf{x}(t)) \right\|_2 \quad (159)$$

$$\leq -\frac{\lambda_0}{2} \left\| \mathbf{P}_{\mathbb{T}, \mathbf{x}(t)} \mathbf{DF}(\mathbf{x}(t))^\top \mathbf{F}(\mathbf{x}(t)) \right\|_2. \quad (160)$$

Further, denoting by $\mathbf{u}(t) \in \mathbb{T}_{\mathbf{x}(t)}$ the unit vector that is tangent to the geodesic between $\mathbf{x}(0)$ and $\mathbf{x}(t)$ at $\mathbf{x}(t)$, we have

$$\frac{d}{dt} d_{\mathbb{S}^{d-1}}(\mathbf{x}(t), \mathbf{x}(0)) = -\langle \mathbf{P}_{\mathbb{T}, \mathbf{x}(t)} \mathbf{DF}(\mathbf{x}(t))^\top \mathbf{F}(\mathbf{x}(t)), \mathbf{u}(t) \rangle. \quad (161)$$

Therefore,

$$\frac{d}{dt} \left\{ d_{\mathbb{S}^{d-1}}(\mathbf{x}(t), \mathbf{x}(0)) + \frac{2}{\lambda_0} \|\mathbf{F}(\mathbf{x}(t))\|_2 \right\} \leq 0. \quad (162)$$

whence, for all $t \leq t_*$, we have $d_{\mathbb{S}^{d-1}}(\mathbf{x}(t), \mathbf{x}(0)) \leq 2\|\mathbf{F}(\mathbf{x}(0))\|_2/\lambda_0$. Recalling Eq. (152), we get $t_* = \infty$ as claimed. \square

We next prove a version of Theorem 3 for gradient flow.

Theorem 8. *Consider gradient flow, as defined in Eq. (151), with respect to the energy function $H(\mathbf{x}) = \|\mathbf{F}(\mathbf{x})\|_2^2/2$ with \mathbf{F} the Gaussian process defined in Section 1, and initialization $\mathbf{x}(0)$ independent of \mathbf{F} . Assume $n, d \rightarrow \infty$ with $n/d \rightarrow \alpha \in [0, 1)$. Define, for c_0 a sufficiently small absolute constant,*

$$\underline{\alpha}_{\text{GF}}(\xi) := \frac{c_0 \xi'(1)^2}{\xi''(1) \xi(1) (\log(\xi^{(3)}(1)/\xi''(1)) \vee 1)}. \quad (163)$$

If $\alpha < \underline{\alpha}_{\text{GF}}(\xi)$, then the following happens with high probability. For all $t \geq 0$

$$\|\mathbf{F}(\mathbf{x}(t))\|_2^2 \leq 2n\xi(1) \exp\left(-\frac{3\xi'(1)}{16} (\sqrt{d} - \sqrt{n})^2 \cdot t\right). \quad (164)$$

Proof. This is an immediate application of Lemma E.1, whereby we note that

$$\|\mathbf{F}(\mathbf{x}_0)\|_2^2 = n\xi(1) + o_P(n), \quad (165)$$

$$\sigma_{\min}(\mathbf{DF}(\mathbf{x}(0))|_{\mathbb{T}, \mathbf{x}(0)}) = \sqrt{\xi'(1)}(\sqrt{d} - \sqrt{n}) + o_P(\sqrt{n}), \quad (166)$$

$$\text{Lip}_{\perp}(\mathbf{DF}; \mathbb{S}^{d-1}) \leq C \sqrt{d \xi''(1) \log \frac{\xi^{(3)}(1)}{\xi''(1)}}. \quad (167)$$

where the last inequality holds with high probability for a universal constant C . The first estimate is by the law of large numbers, the second by Lemma 4.1, which implies $\mathbf{DF}(\mathbf{x})\mathbf{U}_{\mathbf{x}} = \sqrt{\xi'(q)}\mathbf{Z}$ with $\mathbf{Z} \sim \text{GOE}(n, d-1)$ and $q = \|\mathbf{x}\|_2^2$ and the Bai-Yin law, and for the last one, we refer to Lemma A.6. \square

E.2 Gradient descent

The analogue of Lemma E.1 for gradient descent is stated in the main text as Lemma 3.1.

Proof of Lemma 3.1. The key step is to prove that an inequality analogous to the upper bound (157) holds for all $k \leq k_{\star} := \min\{k : \mathbf{x}^k \in B\}$ with B defined as per Eq. (155), which we copy here for the reader's convenience

$$B := \left\{ \mathbf{x} \in \mathbb{S}^{d-1} : \sigma_{\min}(\mathbf{DF}(\mathbf{x})|_{\mathbb{T}, \mathbf{x}}) \leq \frac{\lambda_0}{2} \right\}. \quad (168)$$

To this end first note that

$$\|\mathbf{z}^{k+1} - \mathbf{x}^k\| = \eta \|\mathbf{P}_{\mathbb{T}} \mathbf{DF}(\mathbf{x}^k)^{\top} \mathbf{F}(\mathbf{x}^k)\|_2 \leq \varepsilon_0. \quad (169)$$

(In what follows we omit the reference to the point on the sphere from the projector $\mathbf{P}_{\mathbb{T}}$.)

Further there exist $\xi_{(i)}^k \in [\mathbf{x}^k, \mathbf{z}^{k+1}]$ (here $[\mathbf{u}, \mathbf{v}]$ denotes the segment between \mathbf{u} and $\mathbf{v} \in \mathbb{R}^d$) such that

$$\begin{aligned} F_i(\mathbf{z}^{k+1}) &= F_i(\mathbf{x}^k) - \eta \nabla F_i(\mathbf{x}^k)^{\top} \mathbf{P}_{\mathbb{T}} \mathbf{DF}(\mathbf{x}^k)^{\top} \mathbf{F}(\mathbf{x}^k) - \eta [\nabla F_i(\xi_{(i)}^k) - \nabla F_i(\mathbf{x}^k)]^{\top} \mathbf{P}_{\mathbb{T}} \mathbf{DF}(\mathbf{x}^k)^{\top} \mathbf{F}(\mathbf{x}^k) \\ &=: F_i(\mathbf{x}^k) - \eta \nabla F_i(\mathbf{x}^k)^{\top} \mathbf{P}_{\mathbb{T}} \mathbf{DF}(\mathbf{x}^k)^{\top} \mathbf{F}(\mathbf{x}^k) + \overline{\Delta}_i^k. \end{aligned}$$

We have the estimate

$$|\overline{\Delta}_i^k| \leq \eta J_n \cdot \|\mathbf{z}^{k+1} - \mathbf{x}^k\|_2 \cdot \|\mathbf{P}_{\mathbb{T}} \mathbf{DF}(\mathbf{x}^k)^{\top} \mathbf{F}(\mathbf{x}^k)\|_2 \quad (170)$$

$$\leq J_n \eta^2 \|\mathbf{P}_{\mathbb{T}} \mathbf{DF}(\mathbf{x}^k)^{\top} \mathbf{F}(\mathbf{x}^k)\|_2^2. \quad (171)$$

Further by Pythagoras' theorem $\|\mathbf{z}^{k+1}\|_2^2 = 1 + \|\mathbf{z}^{k+1} - \mathbf{x}^k\|_2^2$, whence, for some $\zeta_{(i)}^k \in [\mathbf{z}^{k+1}, \mathbf{x}^k]$,

$$\begin{aligned} |F_i(\mathbf{x}^{k+1}) - F_i(\mathbf{z}^{k+1})| &= |\langle \nabla F_i(\zeta_{(i)}^k), \mathbf{x}^{k+1} \rangle| \cdot \|\mathbf{z}^{k+1}\|_2 - 1 \\ &\leq \sup_{\mathbf{x} \in \mathbb{B}^d(1+\varepsilon_0)} \|\nabla F_i(\mathbf{x})\|_2 \cdot \|\mathbf{z}^{k+1} - \mathbf{x}^k\|_2^2 \\ &\leq M_n \eta^2 \|\mathbf{P}_{\mathbb{T}} \mathbf{DF}(\mathbf{x}^k)^{\top} \mathbf{F}(\mathbf{x}^k)\|_2^2. \end{aligned}$$

We therefore obtain that

$$\mathbf{F}(\mathbf{x}^{k+1}) = \mathbf{F}(\mathbf{x}^k) - \eta \mathbf{DF}(\mathbf{x}^k) \mathbf{P}_{\mathbb{T}} \mathbf{DF}(\mathbf{x}^k)^{\top} \mathbf{F}(\mathbf{x}^k) + \mathbf{\Delta}^k, \quad (172)$$

$$\begin{aligned}\|\Delta^k\|_2 &\leq \sqrt{n}(J_n + M_n)\eta^2 \|\mathbf{P}_\top \mathbf{D}\mathbf{F}(\mathbf{x}^k)^\top \mathbf{F}(\mathbf{x}^k)\|_2^2 \\ &\leq \frac{\eta}{10 \max_{\mathbf{x} \in \mathbb{S}^{d-1}} \|\mathbf{F}(\mathbf{x})\|_2} \|\mathbf{P}_\top \mathbf{D}\mathbf{F}(\mathbf{x}^k)^\top \mathbf{F}(\mathbf{x}^k)\|_2^2,\end{aligned}\quad (173)$$

where the last inequality follows from Eq. (27). Further

$$\begin{aligned}\|\mathbf{F}(\mathbf{x}^k) - \eta \mathbf{D}\mathbf{F}(\mathbf{x}^k) \mathbf{P}_\top \mathbf{D}\mathbf{F}(\mathbf{x}^k)^\top \mathbf{F}(\mathbf{x}^k)\|_2 &\leq \|\mathbf{F}(\mathbf{x}^k)\|_2 + \eta M_n \|\mathbf{P}_\top \mathbf{D}\mathbf{F}(\mathbf{x}^k)^\top \mathbf{F}(\mathbf{x}^k)\|_2 \\ &\leq (1 + \eta M_n^2) \|\mathbf{F}(\mathbf{x}^k)\|_2 \\ &\leq 2 \|\mathbf{F}(\mathbf{x}^k)\|_2.\end{aligned}\quad (174)$$

Again, in the last step, we used condition (27). Also, note that by Eq. (173) and condition (27), we have

$$\|\Delta^k\|_2 \leq \frac{\eta}{8 \max_{\mathbf{x} \in \mathbb{S}^{d-1}} \|\mathbf{F}(\mathbf{x})\|_2} M_n^2 \|\mathbf{F}(\mathbf{x}^k)\|_2^2 \leq \frac{1}{10} \|\mathbf{F}(\mathbf{x}^k)\|_2. \quad (175)$$

Using Eqs. (172) and (174), we get

$$\|\mathbf{F}(\mathbf{x}^{k+1})\|_2^2 \leq \|\mathbf{F}(\mathbf{x}^k) - \eta \mathbf{D}\mathbf{F}(\mathbf{x}^k) \mathbf{P}_\top \mathbf{D}\mathbf{F}(\mathbf{x}^k)^\top \mathbf{F}(\mathbf{x}^k)\|_2^2 + 4 \|\Delta^k\|_2 \|\mathbf{F}(\mathbf{x}^k)\|_2 + \|\Delta^k\|_2^2.$$

Using Eqs. (172) to (175), we get

$$\begin{aligned}\|\mathbf{F}(\mathbf{x}^{k+1})\|_2^2 &\leq \|\mathbf{F}(\mathbf{x}^k) - \eta \mathbf{D}\mathbf{F}(\mathbf{x}^k) \mathbf{P}_\top \mathbf{D}\mathbf{F}(\mathbf{x}^k)^\top \mathbf{F}(\mathbf{x}^k)\|_2^2 + 5 \|\Delta^k\|_2 \|\mathbf{F}(\mathbf{x}^k)\|_2 \\ &\leq \|\mathbf{F}(\mathbf{x}^k) - \eta \mathbf{D}\mathbf{F}(\mathbf{x}^k) \mathbf{P}_\top \mathbf{D}\mathbf{F}(\mathbf{x}^k)^\top \mathbf{F}(\mathbf{x}^k)\|_2^2 + \frac{\eta}{2} \|\mathbf{P}_\top \mathbf{D}\mathbf{F}(\mathbf{x}^k)^\top \mathbf{F}(\mathbf{x}^k)\|_2^2 \\ &\leq \langle \mathbf{F}(\mathbf{x}^k), \left(\mathbf{I} - \frac{3}{2} \eta \mathbf{K}(\mathbf{x}^k) + \eta^2 \mathbf{K}(\mathbf{x}^k)^2 \right) \mathbf{F}(\mathbf{x}^k) \rangle.\end{aligned}\quad (176)$$

Since $\eta \lambda_{\max}(\mathbf{K}(\mathbf{x}_k)) \leq \eta M_n^2 \leq 1$ by Eq. (27), we conclude that

$$\|\mathbf{F}(\mathbf{x}^{k+1})\|_2^2 \leq \left(1 - \frac{1}{2} \eta \lambda_{\min}(\mathbf{K}(\mathbf{x}^k))\right) \|\mathbf{F}(\mathbf{x}^k)\|_2^2,$$

and therefore, for any $k \leq k_*$,

$$\|\mathbf{F}(\mathbf{x}^k)\|_2^2 \leq e^{-\eta \lambda_0^2 k / 8} \|\mathbf{F}(\mathbf{x}^0)\|_2^2.$$

We are left with the task of proving that $k_* = \infty$. To this end we proceed as in the case of gradient flow. Namely, we note that, for $k \leq k_*$ Eq. (176) implies (for $\eta \lambda_{\max}(\mathbf{K}(\mathbf{x}_k)) \leq 1$)

$$\|\mathbf{F}(\mathbf{x}^{k+1})\|_2^2 \leq \|\mathbf{F}(\mathbf{x}^k)\|_2^2 - \frac{\eta}{2} \|\mathbf{P}_\top \mathbf{D}\mathbf{F}(\mathbf{x}^k)^\top \mathbf{F}(\mathbf{x}^k)\|_2^2 \quad (177)$$

$$\leq \|\mathbf{F}(\mathbf{x}^k)\|_2^2 - \frac{\lambda_0 \eta}{4} \|\mathbf{P}_\top \mathbf{D}\mathbf{F}(\mathbf{x}^k)^\top \mathbf{F}(\mathbf{x}^k)\|_2 \|\mathbf{F}(\mathbf{x}^k)\|_2, \quad (178)$$

and therefore

$$\|\mathbf{F}(\mathbf{x}^{k+1})\|_2 \leq \|\mathbf{F}(\mathbf{x}^k)\|_2 - \frac{\lambda_0 \eta}{8} \|\mathbf{P}_\top \mathbf{D}\mathbf{F}(\mathbf{x}^k)^\top \mathbf{F}(\mathbf{x}^k)\|_2. \quad (179)$$

Further

$$d_{\mathbb{S}^{d-1}}(\mathbf{x}^{k+1}, \mathbf{x}^0) \leq d_{\mathbb{S}^{d-1}}(\mathbf{x}^{k+1}, \mathbf{x}^0) + \eta \|\mathbf{P}_\top \mathbf{D}\mathbf{F}(\mathbf{x}^k)^\top \mathbf{F}(\mathbf{x}^k)\|_2. \quad (180)$$

Therefore, defining $\Psi(\mathbf{x}) := d_{\mathbb{S}^{d-1}}(\mathbf{x}, \mathbf{x}^0) + (8/\lambda_0) \|\mathbf{F}(\mathbf{x})\|_2$, we have $\Psi(\mathbf{x}^{k+1}) \leq \Psi(\mathbf{x}^k)$, whence the proof follows. \square

We are now in position to prove Theorem 3.

Proof of Theorem 3. The proof consists in using the estimates of Section A to check the assumptions of Lemma 3.1.

More precisely, the assumption in (26) holds with high probability under condition (29), using the estimates (165) to (167) that we obtained in the case of gradient flow. The convergence rate (30) follows from Eqs. (28) and (166).

Finally, the assumption in (27) on the stepsize, holds with high probability when $\eta < 1/(C_1 d)$ with C_1 large enough, because the following inequalities hold with high probability for $C = C(\xi)$ a sufficiently large constant:

$$\max_{\mathbf{x} \in \mathbb{S}^{d-1}} \|\mathbf{F}(\mathbf{x})\|_2 \leq C\sqrt{d}, \quad (181)$$

$$\max_{\mathbf{x} \in \mathbb{S}^{d-1}} \|\mathbf{P}_\top \mathbf{D}\mathbf{F}(\mathbf{x})^\top \mathbf{F}(\mathbf{x})\|_2 \leq C d, \quad (182)$$

$$J_n \vee M_n \leq C\sqrt{d}. \quad (183)$$

These inequalities follow from Propositions A.1 and A.4, thus completing the proof. \square

F Analysis of Hessian descent: Proof of Theorem 4

We will collect some random matrix theory results (mainly about the distribution of the Hessian) in Section F.1, and use them to prove Theorem 4 in Section F.2.

F.1 Random matrix theory

Throughout, for $M \leq N$, we let $\mathbf{W} \sim \text{GOE}(N)$ independent of $\mathbf{Z} \sim \text{GOE}(M, N)$ and define

$$\mathbf{A} = \mathbf{A}_{M,N} := a\sqrt{N}\mathbf{W} + b\mathbf{Z}^\top \mathbf{Z}. \quad (184)$$

Lemma F.1. *Assume $a, b, \alpha \in \mathbb{R}_{\geq 0}$, and recall the definition of Q, z_* from Theorem 4, namely*

$$Q(m; \alpha, a, b) := -\frac{1}{m} + \frac{\alpha b}{1 + bm} - a^2 m, \quad (185)$$

$$z_*(\alpha, a, b) := -\sup_{m>0} Q(m; \alpha, a, b). \quad (186)$$

Further, let $S(\cdot; \alpha, a, b) : \mathbb{H} \rightarrow \mathbb{C}$ be the only analytic function on the upper half plane, such that:

(i) $S(\cdot; \alpha, a, b) = -1/z + o(1/z)$ as $z \rightarrow i\infty$; (ii) $S(z; \alpha, a, b)$ solves $Q(S; \alpha, a, b) = z$.

Then the following hold almost surely in the limit $M, N \rightarrow \infty$ with $M/N \rightarrow \alpha \in (0, 1]$.

1.

$$\lim_{M, N \rightarrow \infty} \frac{1}{N} \lambda_{\min}(\mathbf{A}_{M,N}) = -z_*(\alpha, a, b). \quad (187)$$

2. For $z \in \mathbb{H}$ (the upper half complex plane)

$$\lim_{M, N \rightarrow \infty} \frac{1}{N} \text{Tr}((\mathbf{A}_{M,N}/N - z\mathbf{I})^{-1}) = S(z; \alpha, a, b). \quad (188)$$

3. If $a > 0$, $\alpha < 1$ then,

$$2a\sqrt{1-\alpha} < z_*(\alpha, a, b) < 2a. \quad (189)$$

4. Letting $\hat{\nu}_{M,N}(\cdot; a, b)$ denote the empirical spectral distribution of $\mathbf{A}_{M,N}/N$, there exists a non-decreasing, deterministic function $\nu_0(\cdot; a, b) : \mathbb{R} \rightarrow \mathbb{R}_{\geq 0}$ such that $\nu_0(t; a, b) > 0$ for all $a, b \geq 0$ and $z_* = z_*(\alpha, a, b)$, almost surely

$$\lim_{M,N \rightarrow \infty} \hat{\nu}_{M,N}(-z_* + t; a, b) \geq \nu_0(t; a, b). \quad (190)$$

5. For any $a_0, b_0 > 0$, $t > t' > 0$, we have

$$\lim_{M,N \rightarrow \infty} \inf_{a \in [0, a_0], b \in [0, b_0]} \hat{\nu}_{M,N}(-z_* + t; a, b) \geq \nu_{\min}(t'; a_0, b_0), \quad (191)$$

$$\nu_{\min}(t'; a_0, b_0) := \inf_{a \in [0, a_0], b \in [0, b_0]} \nu_0(t'; a, b) > 0. \quad (192)$$

Proof. The asymptotic R-transforms of the random matrices $\mathbf{X} := \mathbf{W}/\sqrt{N}$, $\mathbf{Y} := \mathbf{Z}^\top \mathbf{Z}/N$ are [MS17]

$$R_{\mathbf{X}}(z) = z, \quad R_{\mathbf{Y}}(z) = \frac{\alpha}{1-z}. \quad (193)$$

Since \mathbf{X}, \mathbf{Y} are asymptotically free,

$$R_{\mathbf{A}/N}(z) = aR_{\mathbf{X}}(az) + bR_{\mathbf{Y}}(bz) \quad (194)$$

$$= a^2z + \frac{ab}{1-bz}. \quad (195)$$

Therefore, we have $Q(m; \alpha, a, b) = R_{\mathbf{A}/N}(-m) - m^{-1}$. Point 2 follows then from the standard connection between Stieltjes transform and S-transform.

Point 1 follows from [CDMFF11].

For part 3, Eq. (189) follows from the inequalities

$$-\frac{1}{m} - a^2m < Q(m; \alpha, a, b) < -\frac{1}{m} + \frac{\alpha}{m} - a^2m. \quad (196)$$

For part 4, note that $Q(S; \alpha, a, b) = z$ is a third order algebraic equation for the Stieltjes transform S , with coefficients that depend continuously on a, b, z . The equation reduces to a second order one if $a = 0$ or $b = 0$. By the definition of z_* , there exists $r > 0$ such that, for $z \in (-z_*, -z_* + r)$ this equation has three solutions, of which two are complex conjugates. The imaginary part of these solutions gives (up to a constant factor) the asymptotic density of empirical spectral distribution, which is strictly positive, hence implying the claim (see e.g. [AGZ09, Theorem 2.4.3]).

Finally, we consider part 5. On the favorable event $\mathcal{G} := \{\|\mathbf{W}\|_F^2 \leq 2N^2, \|\mathbf{Z}^\top \mathbf{Z}\|_F^2 \leq 2N^3\}$ (which holds with probability at least $1 - \exp(-cN)$), we have

$$\left\| \frac{1}{N} \mathbf{A}_{M,N}(a_1, b_1) - \frac{1}{N} \mathbf{A}_{M,N}(a_2, b_2) \right\|_F \leq 4\sqrt{N} \|(a_1, b_1) - (a_2, b_2)\|_2, \quad (197)$$

where we noted explicitly the dependence of $\mathbf{A}_{M,N}$ on parameters a, b . Hence, by the Wielandt-Hoffman inequality, for any Lipschitz function $f : \mathbb{R} \rightarrow \mathbb{R}$, the following holds on \mathcal{F} for all $a_1, a_2, b_1, b_2 \geq 0$:

$$\left| \int f(x) \hat{\nu}_{M,N}(dx; a_1, b_1) - \int f(x) \hat{\nu}_{M,N}(dx; a_2, b_2) \right| \leq 4\|f\|_{\text{Lip}} \cdot \|(a_1, b_1) - (a_2, b_2)\|_2. \quad (198)$$

Using the function

$$f_{t,\varepsilon}(x) = \begin{cases} 1 & \text{if } x \leq -z_* + t - \varepsilon, \\ (-z_* + t - x)/\varepsilon & \text{if } -z_* + t - \varepsilon < x < -z_* + t, \\ 0 & \text{if } -z_* + t \leq x, \end{cases} \quad (199)$$

we get, for all $a_1, a_2, b_1, b_2 \geq 0$:

$$\begin{aligned} \hat{\nu}_{M,N}(-z_* + t; a_1, b_1) &\geq \int f_{t,\varepsilon}(x) \hat{\nu}_{M,N}(\mathrm{d}x; a_1, b_1) \\ &\geq \int f_{t,\varepsilon}(x) \hat{\nu}_{M,N}(\mathrm{d}x; a_2, b_2) - \frac{4}{\varepsilon} \|(a_1, b_1) - (a_2, b_2)\|_2 \\ &\geq \hat{\nu}_{M,N}(-z_* + t - \varepsilon; a_2, b_2) - \frac{4}{\varepsilon} \|(a_1, b_1) - (a_2, b_2)\|_2. \end{aligned}$$

Let S_δ be a finite δ -net in $[0, a_0] \times [0, b_0]$. Using the last inequality and the result at point 4 on S_δ , alongside Borel-Cantelli (which implies that \mathcal{G} holds eventually almost surely), we get

$$\lim_{M,N \rightarrow \infty} \inf_{a \in [0, a_0], b \in [0, b_0]} \hat{\nu}_{M,N}(-z_* + t; a, b) \geq \min_{(a,b) \in S_\delta} \nu_0(t - \varepsilon; a, b) - \frac{4\delta}{\varepsilon}.$$

The lower bound (191) follows by taking $\delta \rightarrow 0$.

Finally to prove Eq. (192), i.e. $\nu_{\min}(t'; a, b) > 0$ strictly, we note that

$$\nu_0(t'; a, b) \geq \nu_0(f_{t',\varepsilon}; a, b) := \int f_{t',\varepsilon}(x) \nu_0(\mathrm{d}x; a, b). \quad (200)$$

Further, $(a, b) \mapsto \nu_0(f_{t',\varepsilon}; a, b)$ is Lipschitz continuous by point 4 and the argument given above (with Lipschitz modulus $4/\varepsilon$), and $\nu_0(f_{t',\varepsilon}; a, b) \geq \nu_0(t' - \varepsilon; a, b) > 0$ for any a, b . Therefore, by taking $\varepsilon \in (0, t')$, we get $\nu_{\min}(t'; a, b) > \inf_{a,b \in [0, a_0] \times [0, b_0]} \nu_0(f_{t',\varepsilon}; a, b) > 0$. \square

Lemma F.2. *Let $a_0, b_0 \in \mathbb{R}_{\geq 0}$, $\alpha \in (0, 1]$ and, for $(a, b) \in [0, a_0] \times [0, b_0]$ let $z_* = z_*(\alpha; a, b)$ be defined as in Lemma F.1. Let $M = M(N)$ be a sequence such that $M/N \rightarrow \alpha$ as $N \rightarrow \infty$. Then for any fixed $a_0, b_0, t, c > 0$, there exists $C(t) = C(t, c, a_0, b_0, \alpha)$ such that, if $(a, b) \in [0, a_0] \times [0, b_0]$ then for large enough N ,*

$$\mathbb{P}(\lambda_{\min}(\mathbf{A}_{M,N}) \geq N(-z_* + t)) \leq C(t) e^{-N^2/C(t)}. \quad (201)$$

Proof. Throughout the proof, we denote by C constants that depend on a_0, b_0, α, t, c , and possibly other quantities as indicated, but not on a, b .

Define $\mathbf{W}/\sqrt{N} := \frac{1}{\sqrt{2}}(\mathbf{G} + \mathbf{G}^\top)$ and $\mathbf{S} := \mathbf{Z}/\sqrt{N}$ so that $N^{-1}\mathbf{A}_{M,N} = \frac{a}{\sqrt{2}}(\mathbf{G} + \mathbf{G}^\top) + b\mathbf{S}^\top\mathbf{S}$, and \mathbf{G}, \mathbf{S} are matrices containing i.i.d. random variables $\mathbf{N}(0, 1/N)$. Further, for $\varepsilon > 0$, define

$$\mathbf{B}^\varepsilon := \frac{a}{\sqrt{2}}(\mathbf{G} + \mathbf{G}^\top) + b\mathbf{M}^\varepsilon(\mathbf{S}), \quad (202)$$

$$\mathbf{M}^\varepsilon(\mathbf{S}) := \mathbf{S}^\top(\mathbf{I} + \varepsilon\mathbf{S}\mathbf{S}^\top)^{-1}\mathbf{S}. \quad (203)$$

Since $N^{-1}\mathbf{A}_{M,N} \succeq \mathbf{B}^\varepsilon$, it is sufficient to prove the following claim:

Claim. For any $a_0, b_0, \alpha, t, c > 0$, there exists $\varepsilon, C \geq 0$ such that $\mathbb{P}(\lambda_{\min}(\mathbf{B}^\varepsilon) \geq -z_* + t) \leq C \exp(-N^2/C)$ for all $(a, b) \in [0, a_0] \times [0, b_0]$ and large N .

In order to prove this claim, we view \mathbf{B}^ε as a function of the $N^2 + NM$ random variables (G_{ij}, S_{ij}) . Note that

$$\|\mathbf{B}^\varepsilon(\mathbf{G}_1, \mathbf{S}_1) - \mathbf{B}^\varepsilon(\mathbf{G}_2, \mathbf{S}_2)\|_F \leq 2a\|\mathbf{G}_1 - \mathbf{G}_2\|_F + b\|\mathbf{M}^\varepsilon(\mathbf{S}_1) - \mathbf{M}^\varepsilon(\mathbf{S}_2)\|_F. \quad (204)$$

Further

$$\begin{aligned} \|\mathbf{M}^\varepsilon(\mathbf{S}_1) - \mathbf{M}^\varepsilon(\mathbf{S}_2)\|_F &\stackrel{(a)}{=} \frac{1}{\varepsilon} \left\| (\mathbf{I} + \varepsilon \mathbf{S}_1^\top \mathbf{S}_1)^{-1} - (\mathbf{I} + \varepsilon \mathbf{S}_2^\top \mathbf{S}_2)^{-1} \right\|_F \\ &\stackrel{(b)}{=} \left\| (\mathbf{I} + \varepsilon \mathbf{S}_1^\top \mathbf{S}_1)^{-1} (\mathbf{S}_1^\top \mathbf{S}_1 - \mathbf{S}_2^\top \mathbf{S}_2) (\mathbf{I} + \varepsilon \mathbf{S}_2^\top \mathbf{S}_2)^{-1} \right\|_F \\ &\stackrel{(c)}{\leq} \left\| (\mathbf{I} + \varepsilon \mathbf{S}_1^\top \mathbf{S}_1)^{-1} \mathbf{S}_1^\top (\mathbf{S}_1 - \mathbf{S}_2) \right\|_F + \left\| (\mathbf{I} + \varepsilon \mathbf{S}_2^\top \mathbf{S}_2)^{-1} \mathbf{S}_2^\top (\mathbf{S}_1 - \mathbf{S}_2) \right\|_F \\ &\leq \left(\left\| (\mathbf{I} + \varepsilon \mathbf{S}_1^\top \mathbf{S}_1)^{-1} \mathbf{S}_1^\top \right\|_{\text{op}} + \left\| (\mathbf{I} + \varepsilon \mathbf{S}_2^\top \mathbf{S}_2)^{-1} \mathbf{S}_2^\top \right\|_{\text{op}} \right) \|\mathbf{S}_1 - \mathbf{S}_2\|_F \\ &\stackrel{(d)}{\leq} \frac{1}{\sqrt{\varepsilon}} \|\mathbf{S}_1 - \mathbf{S}_2\|_F. \end{aligned}$$

The equality (a) can be proved by using the singular value decomposition of \mathbf{S}_i to show that the two matrices whose Frobenius norm is computed define the same bilinear form. To verify (b) add and subtract \mathbf{I} from the middle term. For (c) add and subtract $\mathbf{S}_1^\top \mathbf{S}_2$ from the middle term in the left-hand side of the inequality and use the fact that $\mathbf{I} + \varepsilon \mathbf{S}_i^\top \mathbf{S}_i \succeq \mathbf{I}$. Finally, for (d) we use the elementary inequality $x/(1 + \varepsilon x^2) \leq \varepsilon^{-1/2}/2$.

Recall that, by [AGZ09, Lemma 2.3.1], for any L -Lipschitz function $f : \mathbb{R} \rightarrow \mathbb{R}$, $N^{-1} \sum_{i=1}^N f(\lambda_i(\mathbf{B}^\varepsilon))$ is $\sqrt{2/N} \cdot L$ Lipschitz function of \mathbf{B}^ε . Therefore, by the above and Gaussian concentration, for any such function, any $u \geq 0$, and any $(a, b) \in [0, a_0] \times [0, b_0]$

$$\mathbb{P} \left(\left| \mathcal{F}(\mathbf{B}^\varepsilon) - \mathbb{E} \left\{ \mathcal{F}(\mathbf{B}^\varepsilon) \right\} \right| \geq u \right) \leq 2 \exp \left(- \frac{\varepsilon N^2 u^2}{C(a_0, b_0) L^2} \right), \quad (205)$$

$$\mathcal{F}(\mathbf{B}^\varepsilon) := \frac{1}{N} \sum_{i=1}^N f(\lambda_i(\mathbf{B}^\varepsilon)), \quad (206)$$

where we can take $C(a_0, b_0) = C_0(a_0^2 + b_0^2)$ for a suitable numerical constant $C_0 > 0$.

We next take

$$f(x) = \begin{cases} 0 & \text{if } -z_* + t \leq x, \\ (2/t)(-x - z_* + t) & \text{if } -z_* + t/2 < x < -z_* + t, \\ 1 & \text{if } x \leq -z_* + t/2. \end{cases} \quad (207)$$

We have

$$\frac{1}{N} \sum_{i=1}^N f(\lambda_i(\mathbf{B}^\varepsilon)) \geq \frac{1}{N} \sum_{i=1}^N f(\lambda_i(\mathbf{B}^0)) - \frac{2}{t} \cdot \frac{1}{N} \sum_{i=1}^N |\lambda_i(\mathbf{B}^0) - \lambda_i(\mathbf{B}^\varepsilon)| \quad (208)$$

$$\stackrel{(a)}{\geq} \frac{1}{N} \sum_{i=1}^N \mathbf{1} \{ \lambda_i(\mathbf{B}^0) \leq -z_* + t/2 \} - \frac{2b_0}{t\sqrt{N}} \|\mathbf{B}^0 - \mathbf{B}^\varepsilon\|_F \quad (209)$$

$$\stackrel{(b)}{\geq} \nu_{\min}(t/3; a_0, b_0) - \frac{2b_0}{t\sqrt{N}} \left\| \varepsilon (\mathbf{S}^\top \mathbf{S})^2 (\mathbf{I} + \varepsilon \mathbf{S}^\top \mathbf{S})^{-1} \right\|_F \quad (210)$$

$$\geq \nu_{\min}(t/3; a_0, b_0) - \frac{2b_0}{t} \varepsilon \|\mathbf{S}\|_{\text{op}}^4. \quad (211)$$

Here in (a) we used Wielandt-Hoffman and (b) holds eventually almost surely, by part 5 of Lemma F.1. Recall that $\|\mathbf{S}\|_{\text{op}} \leq 2 + \sqrt{\alpha} \leq 3$ with probability $1 - \exp(-cN)$. Hence, taking $\varepsilon = t\nu_{\min}(t/3; a_0, b_0)/(4 \cdot 3^4 b_0)$, we get that, with high probability

$$\frac{1}{N} \sum_{i=1}^N f(\lambda_i(\mathbf{B}^\varepsilon)) \geq \frac{1}{2} \nu_{\min}(t/3; a_0, b_0). \quad (212)$$

Since f is $2/t$ Lipschitz, Eq. (205) implies

$$\mathbb{P} \left(\frac{1}{N} \sum_{i=1}^N f(\lambda_i(\mathbf{B}^\varepsilon)) \leq \frac{1}{2} \nu_{\min}(t/3; a_0, b_0) - u \right) \leq 2 \exp \left(- \frac{\varepsilon N^2 t^2 u^2}{C(a_0, b_0)} \right) \quad (213)$$

$$\leq 2 \exp \left(- \frac{N^2 t^4 \nu_{\min}(t/3; a_0, b_0) u^2}{C'(a_0, b_0)} \right), \quad (214)$$

and the proof of the claim follows by noting that

$$\mathbb{P}(\lambda_{\min}(\mathbf{B}^\varepsilon) \geq -z_* + t) \leq \mathbb{P} \left(\frac{1}{N} \sum_{i=1}^N f(\lambda_i(\mathbf{B}^\varepsilon)) \leq 0 \right). \quad (215)$$

□

Lemma F.3. *Assume we let $n, d \rightarrow \infty$ with $n/d \rightarrow \alpha \in (0, 1]$. Recall the definition of Hamiltonian $H(\mathbf{x}) := \|\mathbf{F}(\mathbf{x})\|_2^2/2$, and that $\mathbb{T}_{\mathbf{x}}$ denotes the tangent space to the sphere of radius $\|\mathbf{x}\|_2$ at \mathbf{x} .*

Let $\lambda(\mathbf{x}) := \lambda_{\min}(\nabla^2 H(\mathbf{x})|_{\mathbb{T}_{\mathbf{x}}})$ and define

$$z_0(\mathbf{x}) := -z_*(\alpha; \sqrt{2\alpha\xi''(\|\mathbf{x}\|^2)H(\mathbf{x})/n}, \xi'(\|\mathbf{x}\|^2)). \quad (216)$$

Then, for any $\varepsilon > 0$ there exists a constant $C(\varepsilon) > 0$ such that for large enough n and d ,

$$\mathbb{P}(\forall \mathbf{x} \in \mathbf{B}^d(1) : \lambda(\mathbf{x}) \leq d(z_0(\mathbf{x}) + \varepsilon)) \geq 1 - C(\varepsilon) e^{-d/C(\varepsilon)}. \quad (217)$$

Proof. Lemmas 4.1 and F.2 imply that for any $\varepsilon, C_0 > 0$, there exists a constant $C(C_0, \varepsilon) > 0$ (depending on ξ as well) such that, for large n, d for any $\mathbf{x} \in \mathbf{B}^d(1)$

$$\mathbb{P}(\lambda(\mathbf{x}) \geq d(z_0(\mathbf{x}) + \varepsilon); H(\mathbf{x}) \leq C_0 d) \leq C(C_0, \varepsilon) e^{-2d^2/C(C_0, \varepsilon)}. \quad (218)$$

Applying Proposition A.1, Lemma A.6, and Weyl's inequality, we can work on the event

$$\mathcal{G} := \left\{ \max_{\mathbf{x} \in \mathbf{B}^d(1)} H(\mathbf{x}) \leq C_0 d, \max_{\mathbf{x} \in \mathbf{B}^d(1)} \|\mathbf{D}\mathbf{F}|_{\mathbb{T}_{\mathbf{x}}}\|_{\text{op}} \leq C_0 \sqrt{d}, \text{Lip}(H; \mathbf{B}^d(\mathbf{0}; 1)) \leq C_0 d, \right. \\ \left. \text{Lip}(\mathbf{D}\mathbf{F}; \mathbf{B}^d(1)) \leq C_0 d, \text{Lip}(\lambda; \mathbf{B}^d(1)) \leq C_0 d \right\},$$

for C_0 a sufficiently large constant (dependent on ξ). Indeed $\mathbb{P}(\mathcal{G}) \geq 1 - C_1 e^{-d/C_1}$ for a certain constant C_1 and large n, d . Further, on the same event $\text{Lip}(z_0; \mathbf{B}^d(1)) \leq C_0$ (eventually enlarging C_0).

Let $N^d(\delta_d)$ be a δ_d -net in $\mathbf{B}^d(1)$, with $\delta_d = 1/d$. Then for large n, d ,

$$\mathbb{P}(\exists \mathbf{x} \in \mathbf{B}^d(1) : \lambda(\mathbf{x}) \geq d(z_0(\mathbf{x}) + \varepsilon); \mathcal{G}) \leq \mathbb{P}(\exists \mathbf{x} \in N^d(\delta_d) : \lambda(\mathbf{x}) \geq d(z_0(\mathbf{x}) + \varepsilon) - C'_0 d \delta_d; \mathcal{G}) \\ \leq |N^d(\delta_d)| \max_{\mathbf{x} \in N^d(\delta_d)} \mathbb{P}(\lambda(\mathbf{x}) \geq d(z_0(\mathbf{x}) + \varepsilon/2); H(\mathbf{x}) \leq C_0 d) \\ \leq \exp \left(C_1 d \log d - \frac{d^2}{C(C_0, \varepsilon/2)} \right),$$

which proves the desired claim. □

F.2 Proof of Theorem 4

Recall the definition of $z_0(\mathbf{x})$ in Lemma F.3 and define the event $\mathcal{E}(\varepsilon)$ by

$$\begin{aligned} \mathcal{E}(\varepsilon) := & \left\{ \forall \mathbf{x} \in \mathbf{B}^d(1) : \lambda_{\min}(\nabla^2 H(\mathbf{x})|_{\mathbb{T}, \mathbf{x}}) \leq d(z_0(\mathbf{x}) + \varepsilon); \right. \\ & \left. \text{Lip}(H; \mathbf{B}^d(1)) \leq C_0 d; \text{Lip}(\nabla^2 H; \mathbf{B}^d(1)) \leq C_0 d \right\}. \end{aligned} \quad (219)$$

By Lemma A.6 and Lemma F.3 we can choose C sufficiently large, so that, for any $\varepsilon > 0$, $\mathbb{P}(\mathcal{E}(\varepsilon)) \geq 1 - C(\varepsilon) \exp(-d/C(\varepsilon))$ for large n, d .

Consider the sequence \mathbf{x}^k produced by the Hessian descent algorithm. By Taylor's theorem, on the event $\text{Lip}(\nabla^2 H; \mathbf{B}^d(1)) \leq C_0 d$, we have

$$\begin{aligned} H(\mathbf{x}^{k+1}) & \leq H(\mathbf{x}^k) - s_k \sqrt{\delta} \langle \mathbf{v}(\mathbf{x}^k), \nabla H(\mathbf{x}^k) \rangle + \frac{1}{2} \delta \langle \mathbf{v}(\mathbf{x}^k), \nabla^2 H(\mathbf{x}^k) \mathbf{v}(\mathbf{x}^k) \rangle + C_0 d \delta^{3/2} \\ & \leq H(\mathbf{x}^k) + \frac{1}{2} \delta \lambda_{\min}(\nabla^2 H(\mathbf{x}^k)|_{\mathbb{T}, \mathbf{x}^k}) + 2C_0 d \delta^{3/2}, \end{aligned} \quad (220)$$

where in the second line we used the definition of s_k and $\mathbf{v}(\mathbf{x}^k)$. We will hereafter work on the event $\mathcal{E}(\varepsilon)$, with $\varepsilon = \delta^{3/2}$. We thus have that

$$H(\mathbf{x}^{k+1}) \leq H(\mathbf{x}^k) - \frac{1}{2} \delta d z_0(\mathbf{x}^k) + 3C_0 d \delta^{3/2}. \quad (221)$$

For $t = k\delta \in [0, 1]$ with integer k define $U(t) := H(\mathbf{x}^{t/\delta})/n$, and on $(k\delta, (k+1)\delta)$ define $U(t)$ by linearly interpolating the two end values. Defining $[t]_\delta := \delta \lfloor t/\delta \rfloor$ and

$$\Phi(u, t) := -\frac{1}{2\alpha} z_*(\alpha; \sqrt{2\alpha \xi''(t)u}, \xi'(t)), \quad (222)$$

we can rewrite Eq. (220) as (excluding the points $t = k\delta$)

$$\frac{dU}{dt}(t) \leq \Phi(U([t]_\delta), [t]_\delta) + C' \delta^{3/2}, \quad (223)$$

$$U(0) \leq \frac{1}{2} \xi(0) + \delta. \quad (224)$$

where the bound in the last line holds with high probability since $H(\mathbf{0})/n$ concentrates around its expectation. Since $u \mapsto \Phi(u, t)$ is a Lipschitz function of u (a fact that follows from Weyl's inequality), for all δ small enough, $U((k+1)\delta)$ is a monotone increasing function of $U(k\delta)$, and hence it follows that $U(t) \leq U_*(t)$ where

$$\frac{dU_*}{dt}(t) = \Phi(U_*([t]_\delta), [t]_\delta) + C' \delta^{3/2}, \quad (225)$$

$$U_*(0) = \frac{1}{2} \xi(0) + \delta. \quad (226)$$

Finally, letting $u(t) = u(t; \alpha, \xi)$ denote the solution of the ODE (35), and using the fact that $t \mapsto \Phi(u, t)$ is a Lipschitz function of t , we obtain $|U_*(t) - u(t)| \leq C'' \sqrt{\delta}$ by standard discretization arguments for ODEs.

G Analysis of the two-phase algorithm

G.1 Proof of Theorem 5

The proof is based on the following state evolution characterization of the AMP iteration of Eqs. (48), (49).

Proposition G.1. *Consider the AMP iteration of Eqs. (48), (49) with initialization given by Eq. (50), and Onsager coefficients in Eqs. (51) to (53). Let $(M_\ell : \ell \geq 0)$, $(H_\ell : \ell \geq 0)$ be two independent centered Gaussian processes with covariances $\mathbf{Q}^M = (Q_{k,l}^M)_{k,l \geq 0}$, $\mathbf{Q}^H = (Q_{k,l}^H)_{k,l \geq 0}$ defined recursively via*

$$Q_{k+1,l+1}^H = \xi(Q_{k,l}^M), \quad Q_{k+1,l+1}^M = \gamma^2 \xi'(Q_{k,l}^M) Q_{k,l}^H. \quad (227)$$

with initialization $Q_{0,\ell}^M = Q_{0,\ell}^H = 0$ for all $\ell \geq 0$.

Then, for any L , and any locally Lipschitz function $\psi : \mathbb{R}^{L+1} \rightarrow \mathbb{R}$, with $|\psi(\mathbf{x})| \leq C(1 + \|\mathbf{x}\|_2^2)$, we have

$$\text{p-lim}_{n,d \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \psi(\sqrt{nh_i^0}, \dots, \sqrt{nh_i^L}) = \mathbb{E}\{\psi(H_0, \dots, H_L)\}, \quad (228)$$

$$\text{p-lim}_{n,d \rightarrow \infty} \frac{1}{d} \sum_{i=1}^d \psi(\sqrt{dm_i^0}, \dots, \sqrt{dm_i^L}) = \mathbb{E}\{\psi(M_0, \dots, M_L)\}. \quad (229)$$

We next use this result to prove the theorem.

Proof of Theorem 5. By construction, for any $k \geq 0$, we have that $Q_{\ell,\ell+k}^M = Q_{\ell,\ell+k}^H = q_\ell$, where the sequence $(q_\ell)_{\ell \geq 0}$ satisfies

$$q_{\ell+1} = \gamma^2 \xi(q_\ell) \xi'(q_\ell), \quad q_0 = 0. \quad (230)$$

Fixed points of this map are solutions of the equation

$$\frac{1}{\gamma^2} = V(q), \quad V(q) := \frac{\xi(q) \xi'(q)}{q}. \quad (231)$$

Notice that $q \mapsto V(q)$ is a strictly convex function on $(0, q_{\max})$, where $q_{\max} \in [1, \infty]$ is the maximum radius of convergence of ξ (because $\xi(q)$ is real analytic with non-negative coefficients, and hence so is $(\xi(q)\xi'(q) - \xi(0)\xi'(0))/q$). Further, since $\xi(0), \xi'(0) > 0$, we have $V(q) \uparrow \infty$ as $q \downarrow 0$. Therefore $V(q)$ has a unique minimum at q_{RS} . Further $\xi(q) \uparrow \infty$ as $q \uparrow \infty$ unless ξ is linear, in which case $V(q)$ is monotone decreasing. Therefore $q_{\text{RS}} < \infty$ if and only if ξ is nonlinear. In the nonlinear case, this implies that the fixed point equation (231) has two solutions $q_1(\gamma) < q_{\text{RS}} < q_2(\gamma)$ if $\gamma^2 < 1/V_* := 1/\min_{q \in (0, q_{\max})} V(q)$ and no solution for $\gamma^2 > 1/V_*$. In the linear case, we have only one solution $q_1(\gamma) < q_{\text{RS}} = \infty$.

Since $q_1(\gamma)$ decreases continuously from q_{RS} to 0 as γ decreases from $1/\sqrt{V_*}$ to 0, for any $q \in (0, q_{\text{RS}})$ we can select $\gamma = \gamma_0(q, \xi)$, where

$$\gamma_0(q, \xi) := -\sqrt{\frac{q}{\xi(q)\xi'(q)}} \quad (232)$$

and this yields that the smallest fixed point of (230) coincides with q . Further, $V'(q) < 0$ (again by strict convexity). Letting $f(r) := \gamma^2 \xi(r)\xi'(r)$, this implies $f'(q) < 1$ strictly. Since f is convex

non-decreasing, we conclude that there exists a constant $C > 0$ (depending on ξ, q) such that, if $\gamma = \gamma_0(q, \xi)$, then

$$q - C e^{-\ell/C} \leq q_\ell \leq q \quad \forall \ell. \quad (233)$$

Notice that this implies that, for all $\ell, k \geq 0$ (the constant C will change from line to line):

$$|Q_{\ell, \ell+k}^M - q| \leq C e^{-\ell/C}, \quad |Q_{\ell, \ell+k}^H - \xi(q)| \leq C e^{-\ell/C}. \quad (234)$$

By Proposition G.1, this implies that, for all $\ell, k \geq 0$

$$\text{p-lim}_{n, d \rightarrow \infty} \left| \|\mathbf{m}^\ell\|_2^2 - q \right| \leq C e^{-\ell/C} \quad \text{p-lim}_{n, d \rightarrow \infty} \left\| \mathbf{m}^\ell - \mathbf{m}^{\ell+k} \right\|_2^2 \leq C e^{-\ell/C}, \quad (235)$$

$$\text{p-lim}_{n, d \rightarrow \infty} \left| \|\mathbf{h}^\ell\|_2^2 - \xi(q) \right| \leq C e^{-\ell/C} \quad \text{p-lim}_{n, d \rightarrow \infty} \left\| \mathbf{h}^\ell - \mathbf{h}^{\ell+k} \right\|_2^2 \leq C e^{-\ell/C}, \quad (236)$$

and therefore,

$$\text{p-lim}_{n, d \rightarrow \infty} \left| B_\ell - \frac{1}{\sqrt{\alpha}} \xi'(q) \right| \leq C e^{-\ell/C}, \quad (237)$$

$$\text{p-lim}_{n, d \rightarrow \infty} \left| C_\ell - \sqrt{\alpha} \xi'(q) \right| \leq C e^{-\ell/C}, \quad (238)$$

$$\text{p-lim}_{n, d \rightarrow \infty} \left| D_\ell - \xi(q) \xi'(q) \right| \leq C e^{-\ell/C}. \quad (239)$$

Therefore, using Eq. (48), we get

$$\begin{aligned} \text{p-lim}_{n, d \rightarrow \infty} \frac{1}{n} \|\mathbf{F}(\mathbf{m}^\ell)\|^2 &= \text{p-lim}_{n, d \rightarrow \infty} \|\mathbf{h}^{\ell+1} + \gamma B_\ell \mathbf{h}^{\ell-1}\|^2 \\ &= \text{p-lim}_{n, d \rightarrow \infty} (1 + \gamma_0(q, \xi) B_\ell)^2 \cdot \|\mathbf{h}^\ell\|_2^2 + O(e^{-\ell/C}) \\ &= \left(1 + \gamma_0(q, \xi) \frac{1}{\sqrt{\alpha}} \xi'(q) \right)^2 \xi(q) + O(e^{-\ell/C}) \\ &= \left(\sqrt{\xi(q)} - \sqrt{\frac{1}{\alpha} q \xi'(q)} \right)^2 + O(e^{-\ell/C}). \end{aligned}$$

The claim of the theorem now follows by fixing

$$\gamma = \gamma_*(q, \alpha, \xi) = \gamma_0(q \wedge q_0(\alpha), \xi). \quad (240)$$

While not needed for the proof of the theorem we also note that our analysis implies that

$$\lim_{n, d \rightarrow \infty} \frac{1}{n} \left\| \mathbf{P}_{\mathbf{T}, \mathbf{m}^L} \nabla H(\mathbf{m}^L) \right\|_2 \leq C e^{-L/C}. \quad (241)$$

(Recall that the typical scale of $\|\nabla H(\mathbf{m})\|_2$ at a non-random point \mathbf{m} is $\Theta(n)$.) To prove this claim, recall that $\nabla H(\mathbf{m}^\ell) = \mathbf{D}\mathbf{F}(\mathbf{m}^\ell)^\top \mathbf{F}(\mathbf{m}^\ell)$ and therefore, using Eqs. (48) and (49)

$$\begin{aligned} \frac{1}{n} \nabla H(\mathbf{m}^\ell) &= \frac{1}{\sqrt{n}} \mathbf{D}\mathbf{F}(\mathbf{m}^\ell)^\top (\mathbf{h}^{\ell+1} + \gamma B_\ell \mathbf{h}^{\ell-1}) \\ &= \frac{1}{\sqrt{n}} (1 + \gamma B_\ell) \mathbf{D}\mathbf{F}(\mathbf{m}^\ell)^\top \mathbf{h}^\ell + \mathbf{e}_1^\ell \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{\gamma\sqrt{\alpha}}(1 + \gamma B_\ell)(\mathbf{m}^{\ell+1} + \gamma C_\ell \mathbf{m}^{\ell-1} + \gamma^2 D_\ell \mathbf{m}^{\ell-1}) + \mathbf{e}_1^\ell \\
&= \frac{1}{\gamma\sqrt{\alpha}}(1 + \gamma B_\ell)(1 + \gamma C_\ell + \gamma^2 D_\ell)\mathbf{m}^\ell + \mathbf{e}_1^\ell + \mathbf{e}_2^\ell,
\end{aligned}$$

where, by Eqs. (235), (236), we have

$$\text{p-lim}_{n,d \rightarrow \infty} \left\{ \|\mathbf{e}_1^\ell\|_2 + \|\mathbf{e}_2^\ell\|_2 \right\} \leq C \exp(-\ell/C).$$

This immediately implies Eq. (241). \square

G.2 Proof of Theorem 6

Theorem 6 follows from Theorem 5 using the same exactly the same argument as for Theorem 4, with the following adaptations:

1. The initial value of the energy for the second phase is given by $H(\mathbf{m}^L)/n = u_{\text{RS}}(q, \alpha, \xi) + o_P(1)$.
2. The Hessian descent algorithm takes place in the $(d-1)$ -dimensional space $V_L := \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{x}, \mathbf{m}^L \rangle = 0\}$.
3. Within this subspace, we attempt to solve an optimization problem that has the same form as the original one, namely

$$\text{minimize} \quad \|\tilde{\mathbf{F}}(\mathbf{x}; \mathbf{m}^L)\|_2^2, \quad \tilde{\mathbf{F}}(\mathbf{x}; \mathbf{m}^L) := \mathbf{F}(\mathbf{m}^L + \mathbf{x}), \quad (242)$$

$$\text{subj. to} \quad \mathbf{x} \in V_L, \quad \|\mathbf{x}\|_2^2 = 1 - \|\mathbf{m}^L\|_2^2 \quad (243)$$

Of course, the reduction from dimension d to $d-1$ is immaterial, and the change in the radius constraint does not affect the argument given for Theorem 5.

4. For a non-random $\mathbf{v} \in \mathbb{B}^d(1)$, the random function $\tilde{\mathbf{F}}(\cdot; \mathbf{v})$ restricted to \mathbf{v}^\perp is again a centered Gaussian process with covariance

$$\mathbb{E}\{\tilde{F}_i(\mathbf{x}_1)\tilde{F}_j(\mathbf{x}_2; \mathbf{v})\} = \delta_{ij} \xi(\|\mathbf{v}\|_2^2 + \langle \mathbf{x}_1, \mathbf{x}_2 \rangle). \quad (244)$$

Hence, by Theorem 4, Hessian descent on the subspace \mathbf{v}^\perp with respect to $\tilde{\mathbf{F}}$ outputs $\mathbf{x}^* \in \mathbb{S}^{d-1}$ with energy $\|\mathbf{F}(\mathbf{x}^*)\|_2^2/(2n) \leq u(t_*; \mathbf{v}) + C\delta + o_P(1)$, where $t_* := 1 - \|\mathbf{v}\|_2^2$, where $u(\cdot; \mathbf{v})$ solves

$$\frac{du}{dt}(t; \mathbf{v}) = -\frac{1}{2\alpha} z_*(\alpha; \sqrt{2\alpha u(t)\xi''(\|\mathbf{v}\|^2 + t)}, \xi'(\|\mathbf{v}\|^2 + t)), \quad u(0) = \frac{1}{2n} \|\mathbf{F}(\mathbf{v})\|_2^2. \quad (245)$$

5. The proof of Theorem 4 only uses a uniform upper bound over the minimum eigenvalue of Hessian of the energy function, and therefore the last claim actually applies uniformly over $\mathbf{v} \in \mathbb{B}^d(1)$. As a consequence it also applies to the random point \mathbf{m}^L .

G.3 Proof of Proposition G.1

This statement is a direct consequence of [EAMS21, Theorem 6]. We begin by restating the latter in a form that is adapted to our use here. For each $k \geq 1$, let $\mathbf{W}^{(k)} \in (\mathbb{R}^N)^{\otimes k}$ be an independent,

Gaussian symmetric tensor. More precisely, if $(\mathbf{G}^{(k)})_{k \geq 0}$ is a collection of tensors with i.i.d. $\mathbf{N}(0, 1)$ entries, then

$$\mathbf{W}^{(k)} := \frac{1}{k!} \sqrt{\frac{k}{N^{k-1}}} \sum_{\pi \in \mathfrak{S}_k} \mathbf{G}_{\pi}^{(k)}, \quad (246)$$

where the sum is over permutations of k objects, and $\mathbf{G}_{\pi}^{(k)}$ is obtained by permuting the indices of $\mathbf{G}^{(k)}$. Let $\nu(t) = \sum_{k \geq 1} c_k t^k$ for some coefficients c_k , $D \geq 1$ be an integer, $\Phi_{\ell} : \mathbb{R}^D \times \mathbb{R} \rightarrow \mathbb{R}^D$ be Lipschitz functions, and denote by $\mathbf{D}\Phi_{\ell}(z; v) \in \mathbb{R}^{D \times D}$ the Jacobian of Φ_{ℓ} with respect to its first argument (which exists in weak sense).

Define the sequence of random variables $V, (Z_{\ell})_{\ell \geq 0}$ by letting $(Z_0, V) \sim p_0$ be an arbitrary random vector taking values in $\mathbb{R}^D \times \mathbb{R}$, and $(Z_{\ell})_{\ell \geq 1}$ be a centered Gaussian process, again with Z_{ℓ} again taking values in \mathbb{R}^D and covariance $\mathbf{Q}_{jk} := \mathbb{E}[Z_j Z_k^{\top}]$ determined recursively via

$$\mathbf{Q}_{j+1, k+1} = \nu \left(\mathbb{E} \left\{ \Phi_j(Z_j; V) \Phi_k(Z_k; V)^{\top} \right\} \right). \quad (247)$$

(Throughout this section we will not use boldface for D -dimensional vectors.)

For any dimension $N \in \mathbb{N}$, consider iterates $\mathbf{z}^{\ell} \in \mathbb{R}^{N \times D}$ given by

$$\mathbf{z}^{\ell+1} = \sum_{k=1}^{\infty} c_k \mathbf{W}^{(k)} \{ \Phi_{\ell}(\mathbf{z}^{\ell}; \mathbf{v}) \} - \Phi_{\ell-1}(\mathbf{z}^{\ell-1}; \mathbf{v}) \mathbf{B}_{\ell}, \quad (248)$$

$$\mathbf{B}_{\ell} := \nu' \left(\mathbb{E} [\Phi_{\ell}(Z_{\ell}; V) \Phi_{\ell-1}(Z_{\ell-1}; V)^{\top}] \right) \odot \mathbb{E} [\mathbf{D}\Phi_{\ell}(Z_{\ell}; V)], \quad (249)$$

where $\mathbf{v} \in \mathbb{R}^N$ and $\mathbf{z}^0 \in \mathbb{R}^{N \times D}$ are non-random, and ν' is applied entrywise, and \odot denotes entrywise multiplication. Further $\Phi_{\ell}(\mathbf{z}^{\ell}; \mathbf{v}) \in \mathbb{R}^{N \times D}$ denotes the matrix whose i -th row is $\Phi_{\ell}(\mathbf{z}^{\ell}; \mathbf{v})_i = \Phi_{\ell}(\mathbf{z}_i^{\ell}; v_i)$. Further, for $\mathbf{x} \in \mathbb{R}^{N \times D}$, $\mathbf{x} = (x_1, \dots, x_N)^{\top}$, $x_i \in \mathbb{R}^D$, we denote by $\mathbf{W}^{(k)}\{\mathbf{x}\} \in \mathbb{R}^{N \times D}$ the matrix whose i -th row is

$$\mathbf{W}^{(k)}\{\mathbf{x}\}_i = \sum_{j_1, \dots, j_{k-1}=1}^N W_{i j_1 \dots j_{k-1}}^{(k)} \mathbf{x}_{j_1} \odot \dots \odot \mathbf{x}_{j_{k-1}}. \quad (250)$$

With these definitions, we have the following.

Proposition G.2 (Theorem 6 in [EAMS21]). *Assume that the empirical distribution of \mathbf{v}, \mathbf{z}^0 , namely $\hat{p}_{\mathbf{v}, \mathbf{z}^0} := N^{-1} \sum_{i \leq N} \delta_{v_i, \mathbf{z}_i^0}$ converges in W_2 distance to p_0 . Then, for any $L \geq 0$ and any function $\psi : \mathbb{R}^{D(L+1)+1} \rightarrow \mathbb{R}$ that is locally Lipschitz and at most quadratic growth (i.e. $|\psi(\mathbf{x})| \leq C(1 + \|\mathbf{x}\|_2^2)$), we have*

$$\text{p-lim}_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \psi(\mathbf{z}_i^0, \dots, \mathbf{z}_i^L; v_i) = \mathbb{E} \left\{ \psi(Z_0, \dots, Z_L; V) \right\} \quad (251)$$

The same holds if \mathbf{B}_{ℓ} is replaced by its empirical version

$$\hat{\mathbf{B}}_{\ell} := \nu' \left(\frac{1}{N} \sum_{i=1}^N \Phi_{\ell}(\mathbf{z}_i^{\ell}; v_i) \Phi_{\ell-1}(\mathbf{z}_i^{\ell-1}; v_i)^{\top} \right) \odot \frac{1}{N} \sum_{i=1}^N \mathbf{D}\Phi_{\ell}(\mathbf{z}_i^{\ell}; v_i). \quad (252)$$

Proof. As mentioned, this is an adaptation from [EAMS21, Theorem 6], with three points of difference. However, each of these points can be reduced to the version stated in [EAMS21].

1. In [EAMS21] it is assumed that $D = 1$, but memory across iterations is allowed. Namely, iterates $\mathbf{x}^\ell \in \mathbb{R}^N$ are updated via

$$\mathbf{x}^{\ell+1} = \sum_{k=1}^{\infty} c_k \mathbf{W}^{(k)} \{f_\ell(\mathbf{x}^0, \dots, \mathbf{x}^\ell)\} - f_{\ell-1}(\mathbf{x}^0, \dots, \mathbf{x}^\ell) \tilde{\mathbf{B}}_\ell, \quad (253)$$

It is clear that this case also covers the seemingly more general one with iterates $\mathbf{z}^\ell \in \mathbb{R}^{N \times D}$ by grouping the $\mathbf{z}^\ell = (\mathbf{x}^{\ell D}, \mathbf{x}^{\ell D+1}, \dots, \mathbf{x}^{\ell D+D-1})$. Indeed, it covers a more general form of the present statement in which Φ_ℓ depends on $\mathbf{z}^0, \dots, \mathbf{z}^\ell$ rather than just on \mathbf{z}^ℓ .

2. The statement in [EAMS21] does not allow for dependency on the vector \mathbf{v} . However, the case with \mathbf{v} can be reduced to the one without \mathbf{v} simply by encoding \mathbf{v} as an additional column of \mathbf{z}^0 , and increasing D to $D + 1$. Since, as mentioned in the previous point, the result of [EAMS21] implies the case in which both D is arbitrary, and Φ_ℓ depends on $\mathbf{z}^0, \dots, \mathbf{z}^\ell$, dependence on \mathbf{v} can be captured.
3. Finally, [EAMS21] uses the deterministic version of coefficient \mathbf{B}_ℓ . One can show by induction over ℓ that: (i) $\tilde{\mathbf{B}}_\ell = \mathbf{B}_\ell + o_P(1)$; (ii) Denoting by $\hat{\mathbf{z}}^\ell$ the iterates that result from using the empirical version, we have $\|b\mathbf{z}^\ell - \hat{\mathbf{z}}^\ell\|_2 = o_P(1)$.

Normalizations are different but equivalent to the ones of [EAMS21]. \square

We next set $N = n + d$ and

$$\mathbf{z}^\ell = \begin{pmatrix} \mathbf{x}^\ell \\ \mathbf{y}^\ell \end{pmatrix}, \quad \mathbf{x}^\ell \in \mathbb{R}^{n \times D}, \quad \mathbf{y}^\ell \in \mathbb{R}^{d \times D}. \quad (254)$$

Further, we set $v_i = 0$ for $i \leq n$, $v_i = 1$ for $i > n$, and write $J_\ell(x) = \Phi_\ell(x; 0)$, $H_\ell(x) = \Phi_\ell(x; 1)$, and write $\overline{W}_{i; i_1, \dots, i_q, j_1, \dots, j_{k-q}}^{(k,q)} = W_{i, i_1, \dots, i_q, n-1+j_1, \dots, n-1+j_{k-q}}^{(k)}$, for $q \leq k-1$. In what follows indices denoted by i, i_1, i_2, \dots run over $[n]$, and indices denoted by j, j_1, j_2, \dots run over $[d]$. We can then rewrite the iteration as

$$\begin{aligned} x_i^{\ell+1} &= \sum_{k=1}^{\infty} c_k \sum_{q=0}^{k-1} \binom{k-1}{q} \sum_{i_1 \dots i_q \leq n} \sum_{j_{q+1} \dots j_{k-1} \leq n} \overline{W}_{i; i_1, \dots, i_q, j_1, \dots, j_{k-1-q}}^{(k,q)} \\ &\quad \cdot J_\ell(x_{i_1}^\ell) \odot \dots \odot J_\ell(x_{i_q}^\ell) \odot H_\ell(y_{j_1}^\ell) \odot \dots \odot H_\ell(y_{j_{k-q}}^\ell) - \mathbf{B}_\ell^\top F_{\ell-1}(x_i^{\ell-1}), \end{aligned} \quad (255)$$

$$\begin{aligned} y_j^{\ell+1} &= \sum_{k=1}^{\infty} c_k \sum_{q=0}^{k-1} \binom{k-1}{q} \sum_{i_1 \dots i_q \leq n} \sum_{j_{q+1} \dots j_{k-1} \leq n} \overline{W}_{i_1; i_2, \dots, i_q, j, j_1, \dots, j_{k-1-q}}^{(k,q)} \\ &\quad \cdot J_\ell(x_{i_1}^\ell) \odot \dots \odot J_\ell(x_{i_q}^\ell) \odot H_\ell(y_{j_1}^\ell) \odot \dots \odot H_\ell(y_{j_{k-q}}^\ell) - \mathbf{B}_\ell^\top H_{\ell-1}(y_j^{\ell-1}). \end{aligned} \quad (256)$$

We next set $D = 3$ and

$$J_\ell(x = (x_1, x_2, x_3)) = \left(0, \bar{f}_\ell(x_1), 0\right), \quad (257)$$

$$H_\ell(y = (y_1, y_2, y_3)) = \left(\bar{h}_\ell(y_2 - y_3), \bar{g}_\ell(y_2 - y_3), \bar{g}_\ell(y_2 - y_3)\right), \quad (258)$$

Writing $\hat{\mathbf{x}}^\ell$ for the first column of \mathbf{x}^ℓ and $\hat{\mathbf{y}}^\ell, \hat{\mathbf{y}}_0^\ell$ for the second and third columns of \mathbf{y}^ℓ . In terms of these variables, the AMP iteration reads

$$\hat{x}_i^{\ell+1} = \sum_{k=1}^{\infty} c_k \sum_{j_1, \dots, j_{k-1} > n} \overline{W}_{i; i_1 \dots i_{k-1}}^{(k,0)} \bar{h}_\ell(\hat{y}_{j_1}^\ell - \hat{y}_{0, j_1}^\ell) \dots \bar{h}_\ell(\hat{y}_{j_{k-1}}^\ell - \hat{y}_{0, j_{k-1}}^\ell) - b_{1,\ell} \bar{f}_{\ell-1}(\hat{x}_i^{\ell-1}), \quad (259)$$

$$\hat{y}_j^{\ell+1} = \sum_{k=1}^{\infty} c_k \sum_{q=0}^{k-1} \binom{k-1}{q} \sum_{i_1 \dots i_q \leq n} \sum_{j_{q+1} \dots j_{k-1} \leq n} \overline{W}_{i_1; i_2, \dots, i_q, j_1 \dots j_{k-1-q}}^{(k,q)} \cdot \overline{f}_\ell(\hat{x}_{i_1}^\ell) \cdots \overline{f}_\ell(\hat{x}_{i_q}^\ell) \overline{g}_\ell(\hat{y}_{j_1}^\ell - \hat{y}_{0,j_1}^\ell) \cdots \overline{g}_\ell(\hat{y}_{j_{k-q}}^\ell - \hat{y}_{0,j_{k-q}}^\ell) - b_{2,\ell} \overline{g}_{\ell-1}(\hat{y}_j^{\ell-1} - \hat{y}_{0,j}^{\ell-1}) - b_{3,\ell} \overline{h}_{\ell-1}(\hat{y}_j^{\ell-1} - \hat{y}_{0,j}^{\ell-1}), \quad (260)$$

$$\hat{y}_{0,j}^{\ell+1} = \sum_{k=1}^{\infty} c_k \sum_{j_1, \dots, j_{k-1} \leq d} \overline{W}_{j_1; j_2, \dots, j_{k-1}}^{(k,0)} \overline{g}_\ell(\hat{y}_{j_1}^\ell - \hat{y}_{0,j_1}^\ell) \cdots \overline{g}_\ell(\hat{y}_{j_{k-1}}^\ell - \hat{y}_{0,j_{k-1}}^\ell) - b_{4,\ell} \overline{g}_{\ell-1}(\hat{y}_j^{\ell-1} - \hat{y}_{0,j}^{\ell-1}) - b_{5,\ell} \overline{h}_{\ell-1}(\hat{y}_j^{\ell-1} - \hat{y}_{0,j}^{\ell-1}), \quad (261)$$

for suitable sequence of constants $b_{1,\ell}, \dots, b_{5,\ell}$. It is straightforward (albeit tedious) to write expressions for these constants, as well as the state evolution characterization that follows from Proposition G.2.

We next introduce a small parameter ε , and set $\overline{f}_\ell(x) = \varepsilon f_\ell(x/\varepsilon)$, $\overline{g}_\ell(x) = g_\ell(x/\varepsilon)$, $\overline{h}_\ell(x) = h_\ell(x/\varepsilon)$. By defining $x_i^\ell := (\hat{y}_i^\ell - \hat{y}_{0,i}^\ell)/\varepsilon$, and taking the difference of Eqs. (260) and (261), we get

$$\hat{x}_i^{\ell+1} = \sum_{k=1}^{\infty} c_k \sum_{j_1, \dots, j_{k-1} \leq d} \overline{W}_{i; j_1 \dots j_{k-1}}^{(k,0)} h_\ell(x_{j_1}^\ell) \cdots h_\ell(x_{j_{k-1}}^\ell) - b_{*,\ell} f_{\ell-1}(\hat{x}_i^{\ell-1}) + O(\varepsilon), \quad (262)$$

$$x_j^{\ell+1} = \sum_{k=1}^{\infty} c_k (k-1) \sum_{i_1 \leq n} \sum_{j_1 \dots j_{k-2} \leq d} \overline{W}_{i_1; j, j_1 \dots j_{k-2}}^{(k,0)} f_\ell(\hat{x}_{i_1}^\ell) g_\ell(x_{j_1}^\ell) \cdots g_\ell(x_{j_{k-2}}^\ell) - a_{1,\ell} g_{\ell-1}(x_j^{\ell-1}) - a_{2,\ell} h_{\ell-1}(x_j^{\ell-1}) + O(\varepsilon). \quad (263)$$

The order ε term can be controlled uniformly over n, d on a high probability event (controlling the operator norms of tensors $\mathbf{W}^{(k)}$), hence allowing to derive a state evolution characterization of the recursion in which $O(\varepsilon)$ terms are dropped.

We finally define

$$F_i(\mathbf{x}) := \sum_{k=1}^{\infty} c_k \sum_{j_1, \dots, j_{k-1} \leq d} \overline{W}_{i; j_1 \dots j_{k-1}}^{(k,0)} x_{j_1}^\ell \cdots x_{j_{k-1}}^\ell, \quad (264)$$

and note that the above recursion (dropping $O(\varepsilon)$ terms) can be rewritten as

$$\hat{\mathbf{x}}^{\ell+1} = \mathbf{F}(h_\ell(\mathbf{x}^\ell)) - b_{*,\ell} f_{\ell-1}(\hat{\mathbf{x}}^{\ell-1}), \quad (265)$$

$$\mathbf{x}^{\ell+1} = \mathbf{D}\mathbf{F}(\mathbf{x}^\ell)^\top f_\ell(\hat{\mathbf{x}}^\ell) - a_{1,\ell} g_{\ell-1}(\mathbf{x}^{\ell-1}) - a_{2,\ell} h_{\ell-1}(\mathbf{x}^{\ell-1}). \quad (266)$$

This is exactly the form of the algorithm introduced in the main text, (48), (49), where functions f_ℓ, h_ℓ, g_ℓ are all linear.

The proof of Proposition G.1 follows from keeping track of the coefficients in the Onsager terms (the memory terms in the AMP recursions), as well as of the state evolution recursion, along the chain of reductions just defined.

References

- [ABA13] Antonio Auffinger and Gérard. Ben Arous, *Complexity of random smooth functions on the high-dimensional sphere*, Ann. Probab. **41** (2013), no. 6, 4214–4247. MR 3161473
- [AC17] Antonio Auffinger and Wei-Kuo Chen, *Parisi formula for the ground state energy in the mixed p -spin model*, The Annals of Probability **45** (2017), no. 6b, 4617–4631.
- [ADH⁺19] Sanjeev Arora, Simon Du, Wei Hu, Zhiyuan Li, and Ruosong Wang, *Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks*, International Conference on Machine Learning, PMLR, 2019, pp. 322–332.
- [AGZ09] Greg W. Anderson, Alice Guionnet, and Ofer Zeitouni, *An introduction to random matrices*, Cambridge University Press, 2009.
- [Auf13] Auffinger, Antonio and Ben Arous, Gérard and Cerný, Jirí, *Random matrices and complexity of spin glasses*, Communications on Pure and Applied Mathematics **66** (2013), no. 2, 165–201.
- [AW09] Jean-Marc Azaïs and Mario Wschebor, *Level sets and extrema of random processes and fields*, John Wiley & Sons, 2009. MR 2478201
- [AZLL19] Zeyuan Allen-Zhu, Yuanzhi Li, and Yingyu Liang, *Learning and generalization in overparameterized neural networks, going beyond two layers*, Advances in Neural Information Processing Systems **32** (2019), 6158–6169.
- [BHMM19] Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal, *Reconciling modern machine-learning practice and the classical bias–variance trade-off*, Proceedings of the National Academy of Sciences **116** (2019), no. 32, 15849–15854.
- [BLLT20] Peter L Bartlett, Philip M Long, Gábor Lugosi, and Alexander Tsigler, *Benign overfitting in linear regression*, Proceedings of the National Academy of Sciences (2020).
- [BMR21] Peter L. Bartlett, Andrea Montanari, and Alexander Rakhlin, *Deep learning: a statistical viewpoint*, Acta Numerica **30** (2021), 87–201.
- [CDMFF11] Mireille Capitaine, Catherine Donati-Martin, Delphine Féral, and Maxime Février, *Free convolution with a semicircular distribution and eigenvalues of spiked deformations of wigner matrices*, Electron. J. Probab **16** (2011), no. 64, 1750–1792.
- [COB19] Lenaic Chizat, Edouard Oyallon, and Francis Bach, *On lazy training in differentiable programming*, Advances in Neural Information Processing Systems, 2019, pp. 2937–2947.
- [CS92] Andrea Crisanti and H-J Sommers, *The spherical p -spin interaction spin glass model: the statics*, Zeitschrift für Physik B Condensed Matter **87** (1992), no. 3, 341–354.
- [CS95] Andrea Crisanti and Hans-J. Sommers, *Thouless-Anderson-Palmer Approach to the Spherical p -Spin Spin Glass Model*, J. Phys. I France **5** (1995), no. 7, 805–813.
- [DZPS18] Simon S Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh, *Gradient descent provably optimizes over-parameterized neural networks*, International Conference on Learning Representations, 2018.

- [EAMS21] Ahmed El Alaoui, Andrea Montanari, and Mark Sellke, *Optimization of mean-field spin glasses*, *The Annals of Probability* **49** (2021), no. 6, 2922–2960.
- [Gor85] Yehoram Gordon, *Some inequalities for gaussian processes and applications*, *Israel Journal of Mathematics* **50** (1985), no. 4, 265–289.
- [JT17] Aukosh Jagannath and Ian Tobasco, *Low temperature asymptotics of spherical mean field spin glasses*, *Communications in Mathematical Physics* **352** (2017), no. 3, 979–1017.
- [Kac43] Mark Kac, *On the average number of real roots of a random algebraic equation*, *Bull. Amer. Math. Soc.* **49** (1943), 314–320. MR 7812
- [KU23] Persia Jana Kamali and Pierfrancesco Urbani, *Dynamical mean field theory for models of confluent tissues and beyond*, arXiv:2306.06420 (2023).
- [MS17] James A Mingo and Roland Speicher, *Free probability and random matrices*, vol. 35, Springer, 2017.
- [MS23a] Andrea Montanari and Eliran Subag, *Solving overparametrized systems of random equations II: On Smale’s 17th problem over the reals*, 2023, In preparation.
- [MS23b] ———, *Solving overparametrized systems of random equations III: Lipschitz hardness*, 2023, In preparation.
- [OS20] Samet Oymak and Mahdi Soltanolkotabi, *Towards moderate overparameterization: global convergence guarantees for training shallow neural networks*, *IEEE Journal on Selected Areas in Information Theory* (2020).
- [Ric45] Stephen O. Rice, *Mathematical analysis of random noise*, *Bell System Tech. J.* **24** (1945), 46–156. MR 11918
- [Saa11] Yousef Saad, *Numerical methods for large eigenvalue problems: revised edition*, SIAM, 2011.
- [Sch42] Isaac J. Schoenberg, *Positive definite functions on spheres*, *Duke Math. J.* **9** (1942), 96–108. MR 0005922 (3,232c)
- [Sub17a] Eliran Subag, *The complexity of spherical p -spin models—a second moment approach*, *Ann. Probab.* **45** (2017), no. 5, 3385–3450. MR 3706746
- [Sub17b] ———, *The geometry of the Gibbs measure of pure spherical spin glasses*, *Invent. Math.* **210** (2017), no. 1, 135–209. MR 3698341
- [Sub21] ———, *Following the ground states of full-rsb spherical spin glasses*, *Communications on Pure and Applied Mathematics* **74** (2021), no. 5, 1021–1044.
- [Sub23] ———, *Concentration for the zero set of random polynomial systems*, arXiv preprint arXiv:2303.11924 (2023).
- [SZ21] Eliran Subag and Ofer Zeitouni, *Concentration of the complexity of spherical pure p -spin models at arbitrary energies*, *J. Math. Phys.* **62** (2021), no. 12, Paper No. 123301, 15. MR 4346481

- [Urb23] Pierfrancesco Urbani, *A continuous constraint satisfaction problem for the rigidity transition in confluent tissues*, *Journal of Physics A: Mathematical and Theoretical* **56** (2023), no. 11, 115003.
- [ZBH⁺21] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals, *Understanding deep learning (still) requires rethinking generalization*, *Communications of the ACM* **64** (2021), no. 3, 107–115.