# A Mean Field View of the Landscape of Two-Layer Neural Networks

Song Mei\*

Andrea Montanari<sup>†</sup>

Phan-Minh Nguyen<sup>‡</sup>

August 29, 2018

#### Abstract

Multi-layer neural networks are among the most powerful models in machine learning, yet the fundamental reasons for this success defy mathematical understanding. Learning a neural network requires to optimize a non-convex high-dimensional objective (risk function), a problem which is usually attacked using stochastic gradient descent (SGD). Does SGD converge to a global optimum of the risk or only to a local optimum? In the first case, does this happen because local minima are absent, or because SGD somehow avoids them? In the second, why do local minima reached by SGD have good generalization properties?

In this paper we consider a simple case, namely two-layers neural networks, and prove that –in a suitable scaling limit– SGD dynamics is captured by a certain non-linear partial differential equation (PDE) that we call *distributional dynamics* (DD). We then consider several specific examples, and show how DD can be used to prove convergence of SGD to networks with nearly-ideal generalization error. This description allows to 'average-out' some of the complexities of the landscape of neural networks, and can be used to prove a general convergence result for noisy SGD.

# Contents

_	Introduction			
	1.1	An informal overview		
2	Exa	amples		
	2.1	Centered isotropic Gaussians		
	2.2	Centered anisotropic Gaussians		
	2.3	A better activation function		
	2.4	Predicting failure		
3	Ger	neral results		
	3.1	Convergence: noisy SGD		
	3.2	Convergence: noiseless SGD		

<sup>\*</sup>Institute for Computational and Mathematical Engineering, Stanford University

<sup>&</sup>lt;sup>†</sup>Department of Electrical Engineering and Department of Statistics, Stanford University

<sup>&</sup>lt;sup>‡</sup>Department of Electrical Engineering, Stanford University

4	Discussion and future directions	15	
5	Notations	19	
6	General results: Statics	20	
	6.1 Proof of Proposition 1		
7	General results: Dynamics	23	
	7.1 Proof of Theorem 3: Convergence to the PDE		
	7.2 Proof of Theorem 3: Generalization to $\beta < \infty$		
	7.3 Proof of Proposition 2: Monotonicity of the risk		
	7.4 A general continuity result		
	7.5 Some properties of the solution of the PDE (7.1)		
	7.6 Proof of Theorems 6: Stability conditions		
	7.7 Proof of Theorem 7: Instability conditions	40	
8	Centered isotropic Gaussians	47	
	8.1 Statics		
	8.2 Dynamics: Fixed points		
	8.3 Dynamics: Convergence to global minimum for $d = \infty$		
	8.4 Proof of Theorem 1		
	8.5 Checking conditions S0–S4 for the running example	62	
9	Centered anisotropic Gaussians		
	9.1 Statics		
	9.2 Dynamics: Fixed points		
	9.3 Dynamics: Convergence to global minimum for $d=\infty$		
	9.4 Dynamics: Proof of Theorem 2	73	
10	7 Finite temperature	76	
	10.1 Statics	77	
	10.2 Dynamics	80	
	10.3 Proof of Proposition 3, Theorem 4, and Theorem 5		
	10.4 Dependence of convergence time on $D$ and $\eta$	90	
11	1 Numerical Experiments	90	
	11.1 Isotropic Gaussians	90	
	11.1.1 Empirical validation of distributional dynamics	91	
	11.1.2 Empirical validation of the statics	93	
	11.1.3 Checking the condition of Lemma 1 in the main text	96	
	11.2 Centered anisotropic Gaussians with ReLU Activation	98	
	11.3 Isotropic Gaussians: Predictable Failure of SGD	98	
$\mathbf{A}$	Concentration inequalities	102	
В	On the generalization to other loss functions	103	

# 1 Introduction

Multi-layer neural networks are one of the oldest approaches to statistical machine learning, dating back at least to the 1960's [Ros62]. Over the last ten years, under the impulse of increasing computer power and larger data availability, they have emerged as a powerful tool for a wide variety of learning tasks [KSH12, GBCB16].

In this paper we focus on the classical setting of supervised learning, whereby we are given data points  $(\boldsymbol{x}_i, y_i) \in \mathbb{R}^d \times \mathbb{R}$ , indexed by  $i \in \mathbb{N}$ , which are assumed to be independent and identically distributed from an unknown distribution  $\mathbb{P}$  on  $\mathbb{R}^d \times \mathbb{R}$ . Here  $\boldsymbol{x}_i \in \mathbb{R}^d$  is a feature vector (e.g. a set of descriptors of an image), and  $y_i \in \mathbb{R}$  is a label (e.g. labeling the object in the image). Our objective is to model the dependence of the label  $y_i$  on the feature vector  $\boldsymbol{x}_i$  in order to assign labels to previously unlabeled examples. In a two-layers neural network, this dependence is modeled as

$$\hat{y}(\boldsymbol{x};\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^{N} \sigma_*(\boldsymbol{x};\boldsymbol{\theta}_i).$$
 (1)

Here N is the number of hidden units (neurons),  $\sigma_* : \mathbb{R}^d \times \mathbb{R}^D \to \mathbb{R}$  is an activation function, and  $\boldsymbol{\theta}_i \in \mathbb{R}^D$  are parameters, which we collectively denote by  $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N)$ . The factor (1/N) is introduced for convenience and can be eliminated by redefining the activation. Often  $\boldsymbol{\theta}_i = (a_i, b_i, \boldsymbol{w}_i)$  and

$$\sigma_*(\boldsymbol{x};\boldsymbol{\theta}_i) = a_i \, \sigma(\langle \boldsymbol{w}_i, \boldsymbol{x} \rangle + b_i) \,, \tag{2}$$

for some  $\sigma : \mathbb{R} \to \mathbb{R}$ . Ideally, the parameters  $\boldsymbol{\theta} = (\boldsymbol{\theta}_i)_{i \leq N}$  should be chosen as to minimize the risk (generalization error)  $R_N(\boldsymbol{\theta}) = \mathbb{E}\{\ell(y, \hat{y}(\boldsymbol{x}; \boldsymbol{\theta}))\}$  where  $\ell : \mathbb{R} \times \mathbb{R} \to \mathbb{R}$  is a certain loss function. For the sake of simplicity, we will focus on the square loss  $\ell(y, \hat{y}) = (y - \hat{y})^2$  but more general choices can be treated along the same lines.

In practice, the parameters of neural networks are learned by stochastic gradient descent [RM51] (SGD) or its variants. In the present case, this amounts to the iteration

$$\boldsymbol{\theta}_i^{k+1} = \boldsymbol{\theta}_i^k + 2s_k \left( y_k - \hat{y}(\boldsymbol{x}_k; \boldsymbol{\theta}^k) \right) \nabla_{\boldsymbol{\theta}_i} \sigma_*(\boldsymbol{x}_k; \boldsymbol{\theta}_i^k). \tag{3}$$

Here  $\boldsymbol{\theta}^k = (\boldsymbol{\theta}_i^k)_{i \leq N}$  denotes the parameters after k iterations,  $s_k$  is a step size, and  $(\boldsymbol{x}_k, y_k)$  is the k-th example. Throughout the paper, we make the following assumption:

One-pass assumption. Training examples are never revisited. Equivalently,  $\{(\boldsymbol{x}_k,y_k)\}_{k\geq 1}$  are i.i.d.  $(\boldsymbol{x}_k,y_k)\sim \mathbb{P}$ .

In large scale applications, this is not far from truth: the data is so large that each example is visited at most a few times [Bot10]. Further, theoretical guarantees suggest that there is limited advantage to be gained from multiple passes [SSBD14]. For recent work deriving scaling limits under such assumption (in different problems) see [WML17].

Understanding the optimization landscape of two-layers neural networks is largely an open problem even when we have access to an infinite number of examples, i.e. to the population risk  $R_N(\theta)$ . Several studies have focused on special choices of the activation function  $\sigma_*$  and of the data distribution  $\mathbb{P}$ , proving that the population risk has no bad local minima [SJL17, GLM17, BG17]. This type of analysis requires delicate calculations that are somewhat sensitive to the specific choice of the model. Another line of work proposes new algorithms with theoretical guarantees [ABGM14, SA15, JSA15, ZLJ16, Tia17, ZSJ+17], which use initializations based on tensor factorization.

In this paper, we prove that -in a suitable scaling limit—the SGD dynamics admits an asymptotic description in terms of a certain non-linear partial differential equation (PDE). This PDE has a remarkable mathematical structure, in that it corresponds to a gradient flow in the metric space  $(\mathcal{P}(\mathbb{R}^D), W_2)$ : the space of probability measures on  $\mathbb{R}^D$ , endowed with the Wasserstein metric. This gradient flow minimizes an asymptotic version of the population risk which is defined for  $\rho \in \mathcal{P}(\mathbb{R}^D)$  and will be denoted by  $R(\rho)$ . This description simplifies the analysis of the landscape of two-layers neural networks, for instance by exploiting underlying symmetries. We illustrate this by obtaining new results on several concrete examples, as well as a general convergence result for 'noisy SGD.' In the next section, we provide an informal outline, focusing on basic intuitions rather than on formal results. We then present the consequences of these ideas on a few specific examples, and subsequently state our general results.

#### 1.1 An informal overview

A good starting point is to rewrite the population risk  $R_N(\theta) = \mathbb{E}\{[y - \hat{y}(x;\theta)]^2\}$  as

$$R_N(\boldsymbol{\theta}) = R_\# + \frac{2}{N} \sum_{i=1}^N V(\boldsymbol{\theta}_i) + \frac{1}{N^2} \sum_{i,j=1}^N U(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j),$$
 (4)

where we defined the potentials  $V(\boldsymbol{\theta}) = -\mathbb{E}\{y\,\sigma_*(\boldsymbol{x};\boldsymbol{\theta})\}$ ,  $U(\boldsymbol{\theta}_1,\boldsymbol{\theta}_2) = \mathbb{E}\{\sigma_*(\boldsymbol{x};\boldsymbol{\theta}_1)\sigma_*(\boldsymbol{x};\boldsymbol{\theta}_2)\}$ . In particular  $U(\cdot,\cdot)$  is a symmetric positive semidefinite kernel. The constant  $R_{\#} = \mathbb{E}\{y^2\}$  is the risk of the trivial predictor  $\hat{y} = 0$ .

Notice that  $R_N(\boldsymbol{\theta})$  only depends on  $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N$  through their empirical distribution  $\hat{\rho}^{(N)} = N^{-1} \sum_{i=1}^N \delta_{\boldsymbol{\theta}_i}$ . This suggests to consider a risk function defined for  $\rho \in \mathscr{P}(\mathbb{R}^D)$  (we denote by  $\mathscr{P}(\Omega)$  the space of probability distributions on  $\Omega$ ):

$$R(\rho) = R_{\#} + 2 \int V(\boldsymbol{\theta}) \ \rho(\mathrm{d}\boldsymbol{\theta}) + \int U(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) \ \rho(\mathrm{d}\boldsymbol{\theta}_1) \ \rho(\mathrm{d}\boldsymbol{\theta}_2) \ . \tag{5}$$

Formal relationships can be established between  $R_N(\theta)$  and  $R(\rho)$ . For instance, under mild assumptions,  $\inf_{\theta} R_N(\theta) = \inf_{\rho} R(\rho) + O(1/N)$ . We refer to the next sections for mathematical statements of this type.

Roughly speaking,  $R(\rho)$  corresponds to the population risk when the number of hidden units goes to infinity, and the empirical distribution of parameters  $\hat{\rho}^{(N)}$  converges to  $\rho$ . Since  $U(\cdot, \cdot)$  is positive semidefinite, we obtain that the risk becomes convex in this limit. The fact that learning can be viewed as convex optimization in an infinite-dimensional space was indeed pointed out in the past [LBW96, BRV<sup>+</sup>06]. Does this mean that the landscape of the population risk simplifies for large N and descent algorithms will converge to a unique (or nearly unique) global optimum?

The answer to the last question is generally negative, and a physics analogy can explain why. Think of  $\theta_1, \ldots, \theta_N$  as the positions of N particles in a D-dimensional space. When N is large, the behavior of such a 'gas' of particles is effectively described by a density  $\rho_t(\theta)$  (with t indexing time). However, not all 'small' changes of this density profile can be realized in the actual physical dynamics: the dynamics conserves mass locally because particles cannot move discontinuously. For instance, if  $\sup(\rho_t) = S_1 \cup S_2$  for two disjoint compact sets  $S_1, S_2 \subseteq \mathbb{R}^D$ , and all  $t \in [t_1, t_2]$ , then the total mass in each of these regions cannot change over time, i.e.  $\rho_t(S_1) = 1 - \rho_t(S_2)$  does not depend on  $t \in [t_1, t_2]$ .

We will prove that stochastic gradient descent is well approximated (in a precise quantitative sense described below) by a continuum dynamics that enforces this local mass conservation principle. Namely, assume that the step size in SGD given by  $s_k = \varepsilon \, \xi(k\varepsilon)$ , for  $\xi : \mathbb{R}_{\geq 0} \to \mathbb{R}_{\geq 0}$  a sufficiently regular function. Denoting by  $\hat{\rho}_k^{(N)} = N^{-1} \sum_{i=1}^N \delta_{\boldsymbol{\theta}_i^k}$  the empirical distribution of parameters after k SGD steps, we prove that

$$\hat{\rho}_{t/\varepsilon}^{(N)} \Rightarrow \rho_t \tag{6}$$

when  $N \to \infty$ ,  $\varepsilon \to 0$  (here  $\Rightarrow$  denotes weak convergence). The asymptotic dynamics of  $\rho_t$  is defined by the following PDE, which we shall refer to as distributional dynamics (DD)

$$\partial_t \rho_t = 2\xi(t) \, \nabla_{\boldsymbol{\theta}} \cdot \left( \rho_t \nabla_{\boldsymbol{\theta}} \Psi(\boldsymbol{\theta}; \rho_t) \right), \tag{7}$$

$$\Psi(\boldsymbol{\theta}; \rho) \equiv V(\boldsymbol{\theta}) + \int U(\boldsymbol{\theta}, \boldsymbol{\theta}') \, \rho(\mathrm{d}\boldsymbol{\theta}') \,. \tag{8}$$

(Here  $\nabla_{\boldsymbol{\theta}} \cdot \boldsymbol{v}(\boldsymbol{\theta})$  denotes the divergence of the vector field  $\boldsymbol{v}(\boldsymbol{\theta})$ .) This should be interpreted as an evolution equation in  $\mathscr{P}(\mathbb{R}^D)$ . While we described the convergence to this dynamics in asymptotic terms, the results in the next sections provide explicit non-asymptotic bounds. In particular,  $\rho_t$  is a good approximation of  $\hat{\rho}_k^{(N)}$ ,  $k = t/\varepsilon$ , as soon as  $\varepsilon \ll 1/D$  and  $N \gg D$ .

Using these results, analyzing learning in two-layer neural networks reduces to analyzing the PDE (7). While this is far from being an easy task, the PDE formulation leads to several simplifications and insights. First of all, it factors out the invariance of the risk (4) (and of the SGD dynamics (3)), with respect to permutations of the units  $\{1, \ldots, N\}$ .

Second, it allows to exploit symmetries in the data distribution  $\mathbb{P}$ . If  $\mathbb{P}$  is left invariant under a group of transformations (e.g. rotations), we can look for a solution  $\rho_t$  of the DD (7) that enjoys the same symmetry, hence reducing the dimensionality of the problem. This is impossible for the finite-N dynamics (3), since no arrangement of the points  $\{\theta_1, \ldots, \theta_N\} \subseteq \mathbb{R}^D$  is left invariant –say–under rotations. We will provide examples of this approach in the next sections.

Third, there is rich mathematical literature on the PDE (7) which was motivated by the study of interacting particle systems in mathematical physics. As mentioned above, a key structure exploited in this line of work is that (7) can be viewed as a gradient flow for the cost function  $R(\rho)$  in the space ( $\mathcal{P}(\mathbb{R}^D)$ ,  $W_2$ ), of probability measures on  $\mathbb{R}^D$  endowed with the Wasserstein metric [JKO98, AGS08, CMV<sup>+</sup>03]. Roughly speaking, this means that the trajectory  $t \mapsto \rho_t$  attempts to minimize the risk  $R(\rho)$  while maintaining the 'local mass conservation' constraint. Recall that Wasserstein distance is defined as

$$W_2(\rho_1, \rho_2) = \left(\inf_{\gamma \in \mathcal{C}(\rho_1, \rho_2)} \int \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_2^2 \gamma(\mathrm{d}\boldsymbol{\theta}_1, \mathrm{d}\boldsymbol{\theta}_2)\right)^{1/2},\tag{9}$$

where the infimum is taken over all couplings of  $\rho_1$  and  $\rho_2$ . Informally, the fact that  $\rho_t$  is a gradient flow means that (7) is equivalent, for small  $\tau$ , to

$$\rho_{t+\tau} \approx \arg\min_{\rho \in \mathcal{P}(\mathbb{R}^D)} \left\{ R(\rho) + \frac{1}{2\xi(t)\tau} W_2(\rho, \rho_t)^2 \right\}. \tag{10}$$

Powerful tools from the mathematical literature on gradient flows in measure spaces [AGS08] can be exploited to study the behavior of (7).

Most importantly, the scaling limit elucidates the dependence of the landscape of two-layer neural networks on the number of hidden units N.

A remarkable feature of neural networks is the observation that, while they might be dramatically over parametrized, this does not lead to performance degradation. In the case of bounded activation functions, this phenomenon was clarified in the nineties for empirical risk minimization algorithms, see e.g. [Bar98]. The present work provides analogous insight for the SGD dynamics: roughly speaking, our results imply that the landscape remains essentially unchanged as N grows, provided  $N \gg D$ . In particular, assume that the PDE (7) converges close to an optimum in time  $t_*(D)$ . This might depend on D, but does not depend on the number of hidden units N (which does not appear in the DD PDE (7)). If  $t_*(D) = O_D(1)$ , we can then take N arbitrarily (as long as  $N \gg D$ ) and will achieve a population risk which is independent of N (and corresponds to the optimum), using  $k = t_*/\varepsilon = O(D)$  samples.

Our analysis can accommodate some important variants of SGD, a particularly interesting one being noisy SGD:

$$\boldsymbol{\theta}_{i}^{k+1} = (1 - 2\lambda s_{k})\boldsymbol{\theta}_{i}^{k} + 2s_{k} (y_{k} - \hat{y}_{k}) \nabla_{\boldsymbol{\theta}_{i}} \sigma_{*}(\boldsymbol{x}_{k}; \boldsymbol{\theta}_{i}^{k}) + \sqrt{2s_{k}/\beta} \boldsymbol{g}_{i}^{k},$$

$$(11)$$

where  $\boldsymbol{g}_i^k \sim \mathsf{N}(0, \boldsymbol{I}_D)$  and  $\hat{y}_k = \hat{y}(\boldsymbol{x}_k; \boldsymbol{\theta}^k)$ . (The term  $-2\lambda s_k \boldsymbol{\theta}_i^k$  corresponds to an  $\ell_2$  regularization and will be useful for our analysis below.) The resulting scaling limit differ from (7) by the addition of a diffusion term:

$$\partial_t \rho_t = 2\xi(t) \, \nabla_{\boldsymbol{\theta}} \cdot \left( \rho_t \nabla_{\boldsymbol{\theta}} \Psi_{\lambda}(\boldsymbol{\theta}; \rho_t) \right) + 2\xi(t) \beta^{-1} \Delta_{\boldsymbol{\theta}} \rho_t \,, \tag{12}$$

where  $\Psi_{\lambda}(\boldsymbol{\theta}; \rho) = \Psi(\boldsymbol{\theta}; \rho) + (\lambda/2) \|\boldsymbol{\theta}\|_{2}^{2}$ , and  $\Delta_{\boldsymbol{\theta}} f(\boldsymbol{\theta}) = \sum_{i=1}^{d} \partial_{\theta_{i}}^{2} f(\boldsymbol{\theta})$  denotes the usual Laplacian. This can be viewed as a gradient flow for the free energy  $F_{\beta,\lambda}(\rho) = (1/2)R(\rho) + (\lambda/2) \int \|\boldsymbol{\theta}\|_{2}^{2} \rho(\mathrm{d}\boldsymbol{\theta}) - \beta^{-1}\mathrm{Ent}(\rho)$ , where  $\mathrm{Ent}(\rho) = -\int \rho(\boldsymbol{\theta}) \log \rho(\boldsymbol{\theta}) \,\mathrm{d}\boldsymbol{\theta}$  is the entropy of  $\rho$  (by definition  $\mathrm{Ent}(\rho) = -\infty$  if  $\rho$  is singular).  $F_{\beta,\lambda}(\rho)$  is an entropy-regularized risk, which penalizes strongly non-uniform  $\rho$ .

We will prove below that, for  $\beta < \infty$ , the evolution (12) generically converges to the minimizer of  $F_{\beta,\lambda}(\rho)$ , hence implying global convergence of noisy SGD in a number of steps *independent of* N.

# 2 Examples

In this section, we discuss some simple applications of the general approach outlined above. Let us emphasize that these examples are not realistic. First, the data distribution  $\mathbb{P}$  is extremely simple: we made this choice in order to be able to carry out explicit calculations. Second, the activation function  $\sigma_*(x;\theta)$  is not necessarily optimal: we made this choice in order to illustrate some interesting phenomena.

#### 2.1 Centered isotropic Gaussians

One-neuron neural networks perform well with (nearly) linearly separable data. The simplest classification problem which requires multilayer networks is –arguably– the one of distinguishing two Gaussians with the same mean. Assume the joint law  $\mathbb{P}$  of (y, x) to be as follows:

With probability 1/2: 
$$y = +1$$
,  $\boldsymbol{x} \sim N(0, (1 + \Delta)^2 \boldsymbol{I}_d)$ 

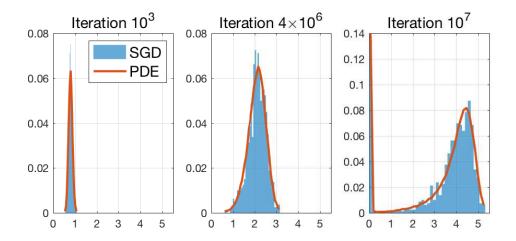


Figure 1: Evolution of the radial distribution  $\bar{\rho}_t$  for the isotropic Gaussian model, with  $\Delta = 0.8$ . Histograms are obtained from SGD experiments with d = 40, N = 800, initial weights distribution  $\rho_0 = N(\mathbf{0}, 0.8^2/d \cdot \mathbf{I}_d)$ , step size  $\epsilon = 10^{-6}$  and  $\xi(t) = 1$ . Continuous lines correspond to a numerical solution of the DD (13).

With probability 1/2: y = -1,  $\boldsymbol{x} \sim N(0, (1 - \Delta)^2 \boldsymbol{I}_d)$ .

(This example will be generalized later.) Of course, optimal classification in this model becomes entirely trivial if we compute the feature  $h(\mathbf{x}) = ||\mathbf{x}||_2$ . However, it is non-trivial that a SGD-trained neural network will succeed.

We choose an activation function without offset or output weights, namely  $\sigma_*(\boldsymbol{x}; \boldsymbol{\theta}_i) = \sigma(\langle \boldsymbol{w}_i, \boldsymbol{x} \rangle)$ . While qualitatively similar results are obtained for other choices of  $\sigma$ , we will use a simple piecewise linear function as a running example:  $\sigma(t) = s_1$  if  $t \leq t_1$ ,  $\sigma(t) = s_2$  if  $t \geq t_2$ , and  $\sigma(t)$  interpolated linearly for  $t \in (t_1, t_2)$ . In simulations we use  $t_1 = 0.5$ ,  $t_2 = 1.5$ ,  $t_3 = -2.5$ ,  $t_4 = 0.5$ .

We run SGD with initial weights  $(\boldsymbol{w}_i^0)_{i\leq N} \sim_{iid} \rho_0$ , where  $\rho_0$  is spherically symmetric. Figure 1 reports the result of such an experiment. Due to the symmetry of the distribution  $\mathbb{P}$ , the distribution  $\rho_t$  remains spherically symmetric for all t, and hence is completely determined by the distribution  $\overline{\rho}_t$  of the norm  $r = ||\boldsymbol{w}||_2$ . This distribution satisfies a one-dimensional reduced DD:

$$\partial_t \overline{\rho}_t = 2\xi(t) \,\partial_r \left( \overline{\rho}_t \partial_r \psi(r; \overline{\rho}_t) \right), \tag{13}$$

where the form of  $\psi(r;\rho)$  can be derived from  $\Psi(\theta;\rho)$ . This reduced PDE can be efficiently solved numerically, see Supplementary Information (SI) for technical details. As illustrated by Fig. 1, the empirical results match closely the predictions produced by this PDE.

In Fig. 2, we compare the asymptotic risk achieved by SGD with the prediction obtained by minimizing  $R(\rho)$ , cf. (5) over spherically symmetric distributions. It turns out that, for certain values of  $\Delta$ , the minimum is achieved by the uniform distribution over a sphere of radius  $\|\boldsymbol{w}\|_2 = r_*$ , to be denoted by  $\rho_{r_*}^{\text{unif}}$ . The value of  $r_*$  is computed by minimizing

$$\overline{R}_{J}^{(1)}(r) = 1 + 2v(r) + u_{J}(r, r), \qquad (14)$$

where expressions for v(r),  $u_d(r_1, r_2)$  can be readily derived from  $V(\boldsymbol{w})$ ,  $U(\boldsymbol{w}_1, \boldsymbol{w}_2)$  and are given in the SI.

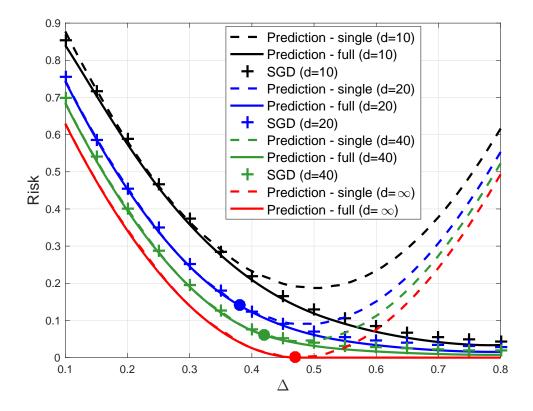


Figure 2: Population risk in the problem of separating two isotropic Gaussians, as a function of the separation parameter  $\Delta$ . We use a a two-layers network with piecewise linear activation, no offset, and output weights equal to one. Empirical results obtained by SGD (a single run per data point) are marked "+". Continuous lines are theoretical predictions obtained by numerically minimizing  $R(\rho)$  (see SI for details). Dashed lines are theoretical predictions from the single-delta ansatz of (14). Notice that this ansatz is incorrect for  $\Delta > \Delta_d^{\rm h}$ , which is marked as a solid round dot. Here N=800.

**Lemma 1.** Let  $r_*$  be a global minimizer of  $r \mapsto R_d^{(1)}(r)$ . Then  $\rho_{r_*}^{\text{unif}}$  is a global minimizer of  $\rho \mapsto R(\rho)$  if and only if  $v(r) + u_d(r, r_*) \ge v(r_*) + u_d(r_*, r_*)$  for all  $r \ge 0$ .

Checking numerically this condition yields that  $\rho_{r_*}^{\text{unif}}$  is a global minimizer for  $\Delta$  in an interval  $[\Delta_d^{\text{l}}, \Delta_d^{\text{h}}]$ , where  $\lim_{d \to \infty} \Delta_d^{\text{l}} = 0$  and  $\lim_{d \to \infty} \Delta_d^{\text{h}} = \Delta_{\infty} \approx 0.47$ .

Figure 2 shows good quantitative agreement between empirical results and theoretical predictions, and suggests that SGD achieves a value of the risk which is close to optimum. Can we prove that this is indeed the case, and that the SGD dynamics does not get stuck in local minima? It turns out that we can use our general theory (see next section) to prove that this is the case for large d. In order to state this result, we need to introduce a class of good uninformative initializations  $\mathscr{P}_{\text{good}} \subseteq \mathscr{P}(\mathbb{R}_{\geq 0})$  for which convergence to the optimum takes place. For  $\overline{\rho} \in \mathscr{P}(\mathbb{R}_{\geq 0})$ , we let  $\overline{R}_d(\overline{\rho}) \equiv R(\overline{\rho} \times \text{Unif}(\mathbb{S}^{d-1}))$ . This risk has a well defined limit as  $d \to \infty$ . We say that  $\overline{\rho} \in \mathscr{P}_{\text{good}}$  if:  $(i) \overline{\rho}$  is absolutely continuous with respect to Lebesgue measure, with bounded density; (ii)

 $\overline{R}_{\infty}(\overline{\rho}) < 1.$ 

**Theorem 1.** For any  $\eta, \Delta, \delta > 0$ , and  $\overline{\rho}_0 \in \mathscr{P}_{good}$ , there exists  $d_0 = d_0(\eta, \overline{\rho}_0, \Delta)$ ,  $T = T(\eta, \overline{\rho}_0, \Delta)$ , and  $C_0 = C_0(\eta, \overline{\rho}_0, \Delta, \delta)$ , such that the following holds for the problem of classifying isotropic Gaussians. For any dimension  $d \geq d_0$ , number of neurons  $N \geq C_0 d$ , consider SGD initialized with  $(\boldsymbol{w}_i^0)_{i \leq N} \sim_{iid} \overline{\rho}_0 \times \text{Unif}(\mathbb{S}^{d-1})$  and step size  $\varepsilon \in [1/N^{10}, 1/(C_0 d)]$ . Then we have

$$R_N(\boldsymbol{\theta}^k) \le \inf_{\boldsymbol{\theta} \in \mathbb{R}^{N \times d}} R(\boldsymbol{\theta}) + \eta$$
 (15)

for any  $k \in [T/\varepsilon, 10T/\varepsilon]$  with probability at least  $1 - \delta$ .

In particular, if we set  $\varepsilon = 1/(C_0 d)$ , then the number of SGD steps is  $k \in [(C_0 T) d, (10C_0 T) d]$ : the number of samples used by SGD does not depend on the number of hidden units N, and is only linear in the dimension. Unfortunately the proof does not provide the dependence of T on  $\eta$ , but Theorem 6 below suggests exponential local convergence.

While we stated Theorem 1 for the piecewise linear sigmoids, the SI presents technical conditions under which it holds for a general monotone function  $\sigma : \mathbb{R} \to \mathbb{R}$ .

### 2.2 Centered anisotropic Gaussians

We can generalize the previous result to a problem in which the network needs to select a subset of relevant nonlinear features out of many a priori equivalent ones. We assume the joint law of (y, x) to be as follows:

With probability 1/2: y = +1,  $\boldsymbol{x} \sim N(0, \boldsymbol{\Sigma}_{+})$ , and

With probability 1/2: y = -1,  $\boldsymbol{x} \sim N(0, \boldsymbol{\Sigma}_{-})$ .

Given a linear subspace  $\mathcal{V} \subseteq \mathbb{R}^d$  of dimension  $s_0 \leq d$ , we assume that  $\Sigma_+$ ,  $\Sigma_-$  differ uniquely along  $\mathcal{V}$ :  $\Sigma_{\pm} = \mathbf{I}_d + (\tau_{\pm}^2 - 1)\mathbf{P}_{\mathcal{V}}$ , where  $\tau_{\pm} = (1 \pm \Delta)$  and  $\mathbf{P}_{\mathcal{V}}$  is the orthogonal projector onto  $\mathcal{V}$ . In other words, the projection of  $\boldsymbol{x}$  on the subspace  $\mathcal{V}$  is distributed according to a isotropic Gaussian with variance  $\tau_+^2$  (if y = +1) or  $\tau_-^2$  (if y = -1). The projection orthogonal to  $\mathcal{V}$  has instead the same variance in the two classes. A successful classifier must be able to learn the relevant subspace  $\mathcal{V}$ . We assume the same class of activations  $\sigma_*(\boldsymbol{x}; \boldsymbol{\theta}) = \sigma(\langle \boldsymbol{w}, \boldsymbol{x} \rangle)$  as for the isotropic case.

The distribution  $\mathbb{P}$  is invariant under a reduced symmetry group  $\mathcal{O}(s_0) \times \mathcal{O}(d - s_0)$ . As a consequence, letting  $r_1 = \|\mathbf{P}_{\mathcal{V}}\mathbf{w}\|_2$  and  $r_2 \equiv \|(\mathbf{I}_d - \mathbf{P}_{\mathcal{V}})\mathbf{w}\|_2$ , it is sufficient to consider distributions  $\rho$  that are uniform, conditional on the values of  $r_1$  and  $r_2$ . If we initialize  $\rho_0$  to be uniform conditional on  $(r_1, r_2)$ , this property is preserved by the evolution (7). As in the isotropic case, we can use our general theory to prove convergence to a near-optimum if d is large enough.

**Theorem 2.** For any  $\eta, \Delta, \delta > 0$ , and  $\overline{\rho}_0 \in \mathscr{P}_{good}$ , there exists  $d_0 = d_0(\eta, \overline{\rho}_0, \Delta, \gamma)$ ,  $T = T(\eta, \overline{\rho}_0, \Delta, \gamma)$ , and  $C_0 = C_0(\eta, \overline{\rho}_0, \Delta, \delta, \gamma)$ , such that the following holds for the problem of classifying anisotropic Gaussians with  $s_0 = \gamma d$ ,  $\gamma \in (0, 1)$  fixed. For any dimension parameters  $s_0 = \gamma d \geq d_0$ , number of neurons  $N \geq C_0 d$ , consider SGD initialized with initialization  $(\mathbf{w}_i^0)_{i \leq N} \sim_{iid} \overline{\rho}_0 \times \text{Unif}(\mathbb{S}^{d-1})$  and step size  $\varepsilon \in [1/N^{10}, 1/(C_0 d)]$ . Then we have  $R_N(\boldsymbol{\theta}^k) \leq \inf_{\boldsymbol{\theta} \in \mathbb{R}^{N \times d}} R_N(\boldsymbol{\theta}) + \eta$  for any  $k \in [T/\varepsilon, 10T/\varepsilon]$  with probability at least  $1 - \delta$ .

Even with a reduced degree of symmetry, SGD converges to a network with nearly-optimal risk, after using a number of samples k = O(d), which is independent of the number of hidden units N.

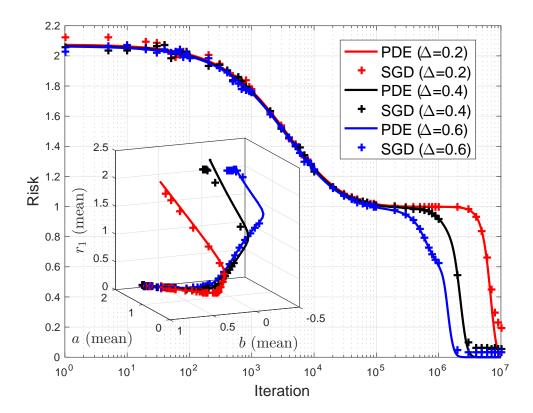


Figure 3: Evolution of the population risk for the variable selection problem using a two-layers neural network with ReLU activations. Here d=320,  $s_0=60$ , N=800, and we used  $\xi(t)=t^{-1/4}$  and  $\varepsilon=2\times 10^{-4}$  to set the step size. Numerical simulations using SGD (one run per data point) are marked "+", and curves are solutions of the reduced PDE with  $d=\infty$ . Inset: evolution of three parameters of the reduced distribution  $\overline{\rho}_t$  (average output weights a, average offsets b and average  $\ell_2$  norm in the relevant subspace  $r_1$ ) for the same setting.

#### 2.3 A better activation function

Our previous examples use activation functions  $\sigma_*(\boldsymbol{x};\boldsymbol{\theta}) = \sigma(\langle \boldsymbol{w}, \boldsymbol{x} \rangle)$  without output weights or offset, in order to simplify the analysis and illustrate some interesting phenomena. Here we consider instead a standard rectified linear unit (ReLU) activation, and fit both the output weight and the offset:  $\sigma_*(\boldsymbol{x};\boldsymbol{\theta}) = a \, \sigma_{\text{ReLU}}(\langle \boldsymbol{w}, \boldsymbol{x} \rangle + b)$  where  $\sigma_{\text{ReLU}}(x) = \max(x,0)$ . Hence  $\boldsymbol{\theta} = (\boldsymbol{w},a,b) \in \mathbb{R}^{d+2}$ .

We consider the same data distribution introduced in the last section (anisotropic Gaussians). Figure 3 reports the evolution of the risk  $R_N(\boldsymbol{\theta}^k)$  for three experiments with d=320,  $s_0=60$  and different values of  $\Delta$ . SGD is initialized by setting  $a_i=1$ ,  $b_i=1$  and  $\boldsymbol{w}_i^0 \sim_{iid} \mathsf{N}(\boldsymbol{0}, 0.8^2/d \cdot \boldsymbol{I}_d)$  for  $i \leq N$ . We observe that SGD converges to a network with very small risk, but this convergence has a nontrivial structure and presents long flat regions.

The empirical results are well captured by our predictions based on the continuum limit. In this case we obtain a reduced PDE for the joint distribution of the four quantities  $\mathbf{r} = (a, b, r_1 = \|\mathbf{P}_{\mathcal{V}} \mathbf{w}\|_2, r_2 = \|\mathbf{P}_{\mathcal{V}}^{\perp} \mathbf{w}\|_2)$ , denoted by  $\overline{\rho}_t$ . The reduced PDE is analogous to (13) albeit in 4 rather than 1 dimensions. In Figure 3 we consider the evolution of the risk, alongside three properties of

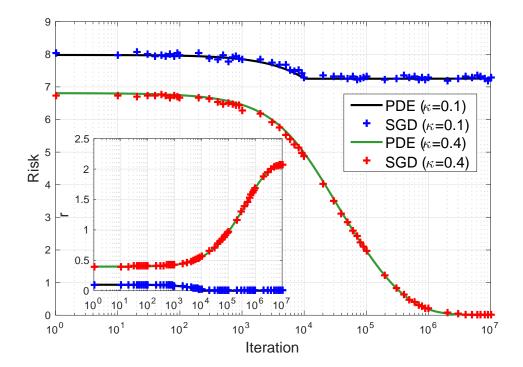


Figure 4: Separating two isotropic Gaussians, with a non-monotone activation function (see text for details). Here  $N=800,\ d=320,\ \Delta=0.5$ . The main frame presents the evolution of the population risk along the SGD trajectory, starting from two different initializations of  $(\boldsymbol{w}_i^0)_{i\leq N} \sim_{iid} N(\boldsymbol{0},\kappa^2/d\cdot\mathbf{I}_d)$  for either  $\kappa=0.1$  or  $\kappa=0.4$ . In the inset, we plot the evolution of the average of  $\|\boldsymbol{w}\|_2$  for the same conditions. Symbols are empirical results. Continuous lines are prediction obtained with the reduced PDE (13).

the distribution  $\overline{\rho}_t$  -the means of the output weight a, of the offset b, and of  $r_1$ .

#### 2.4 Predicting failure

SGD does not always converge to a near global optimum. Our analysis allows to construct examples in which SGD fails. For instance, Figure 4 reports results for isotropic Gaussians problem. We violate the assumptions of Theorem 1 by using non monotone activation function. Namely, we use  $\sigma_*(\boldsymbol{x};\boldsymbol{\theta}) = \sigma(\langle \boldsymbol{w}, \boldsymbol{x} \rangle)$ , where  $\sigma(t) = -2.5$  for  $t \leq 0$ ,  $\sigma(t) = 7.5$  for  $t \geq 1.5$ , and  $\sigma(t)$  linearly interpolates from (0, -2.5) to (0.5, -4), and from (0.5, -4) to (1.5, 7.5).

Depending on the initialization, SGD converges to two different limits, one with a small risk, and the second with high risk. Again this behavior is well tracked by solving a one-dimensional PDE for the distribution  $\overline{\rho}_t$  of  $r = \|\boldsymbol{w}\|_2$ .

# 3 General results

In this section we return to the general supervised learning problem described in the introduction and describe our general results. Proofs are deferred to the SI.

First, we note that the minimum of the asymptotic risk  $R(\rho)$  of (5) provides a good approximation of the minimum of the finite-N risk  $R_N(\theta)$ .

**Proposition 1.** Assume that either one of the following conditions hold: (a)  $\inf_{\rho} R(\rho)$  is achieved by a distribution  $\rho_*$  such that  $\int U(\boldsymbol{\theta}, \boldsymbol{\theta}) \, \rho_*(\mathrm{d}\boldsymbol{\theta}) \leq K$ ; (b) There exists  $\varepsilon_0 > 0$  such that, for any  $\rho \in \mathscr{P}(\mathbb{R}^D)$  such that  $R(\rho) \leq \inf_{\rho} R(\rho) + \varepsilon_0$  we have  $\int U(\boldsymbol{\theta}, \boldsymbol{\theta}) \, \rho(\mathrm{d}\boldsymbol{\theta}) \leq K$ . Then

$$\left|\inf_{\boldsymbol{\theta}} R_N(\boldsymbol{\theta}) - \inf_{\boldsymbol{\rho}} R(\boldsymbol{\rho})\right| \le K/N. \tag{16}$$

Further, assume that  $\boldsymbol{\theta} \mapsto V(\boldsymbol{\theta})$  and  $(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) \mapsto U(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)$  are continuous, with U bounded below. A probability measure  $\rho_*$  is a global minimum of R if  $\inf_{\boldsymbol{\theta} \in \mathbb{R}^D} \Psi(\boldsymbol{\theta}; \rho_*) > -\infty$  and

$$\operatorname{supp}(\rho_*) \subseteq \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^D} \Psi(\boldsymbol{\theta}; \rho_*). \tag{17}$$

We next consider the distributional dynamics (7) and (12). These should be interpreted to hold in weak sense, cf. SI. In order to establish that these PDEs indeed describe the limit of the SGD dynamics, we make the following assumptions.

- A1.  $t \mapsto \xi(t)$  is bounded Lipschitz:  $\|\xi\|_{\infty}$ ,  $\|\xi\|_{\text{Lip}} \leq K_1$ , with  $\int_0^{\infty} \xi(t) dt = \infty$ .
- A2. The activation function  $(\boldsymbol{x}, \boldsymbol{\theta}) \mapsto \sigma_*(\boldsymbol{x}; \boldsymbol{\theta})$  is bounded, with sub-Gaussian gradient:  $\|\sigma_*\|_{\infty} \leq K_2$ ,  $\|\nabla_{\boldsymbol{\theta}}\sigma_*(\boldsymbol{X}; \boldsymbol{\theta})\|_{\psi_2} \leq K_2$ . Labels are bounded  $|y_k| \leq K_2$ .
- A3. The gradients  $\boldsymbol{\theta} \mapsto \nabla V(\boldsymbol{\theta})$ ,  $(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) \mapsto \nabla_{\boldsymbol{\theta}_1} U(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)$  are bounded, Lipschitz continuous (namely  $\|\nabla_{\boldsymbol{\theta}} V(\boldsymbol{\theta})\|_2$ ,  $\|\nabla_{\boldsymbol{\theta}_1} U(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)\|_2 \leq K_3$ ,  $\|\nabla_{\boldsymbol{\theta}} V(\boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}} V(\boldsymbol{\theta}')\|_2 \leq K_3 \|\boldsymbol{\theta} \boldsymbol{\theta}'\|_2$ ,  $\|\nabla_{\boldsymbol{\theta}_1} U(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) \nabla_{\boldsymbol{\theta}_1} U(\boldsymbol{\theta}_1', \boldsymbol{\theta}_2')\|_2 \leq K_3 \|(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) (\boldsymbol{\theta}_1', \boldsymbol{\theta}_2')\|_2$ ).

We also introduce the following error term which quantifies in a non-asymptotic sense the accuracy of our PDE model

$$\operatorname{err}_{N,D}(z) \equiv \sqrt{1/N \vee \varepsilon} \cdot \left[ \sqrt{D + \log(N/\varepsilon)} + z \right].$$
 (18)

The convergence of the SGD process to the PDE model is an example of a phenomenon which is known in probability theory as *propagation of chaos* [Szn91].

**Theorem 3.** Assume that conditions A1, A2, A3 hold. For  $\rho_0 \in \mathscr{P}(\mathbb{R}^D)$ , consider SGD with initialization  $(\theta_i^0)_{i \leq N} \sim_{iid} \rho_0$  and step size  $s_k = \varepsilon \xi(k\varepsilon)$ . For  $t \geq 0$ , let  $\rho_t$  be the solution of PDE (7). Then, for any fixed  $t \geq 0$ ,  $\hat{\rho}_{\lfloor t/\varepsilon \rfloor}^{(N)} \Rightarrow \rho_t$  almost surely along any sequence  $(N, \varepsilon = \varepsilon_N)$  such that  $N \to \infty$ ,  $\varepsilon_N \to 0$ ,  $N/\log(N/\varepsilon_N) \to \infty$  and  $\varepsilon_N \log(N/\varepsilon_N) \to 0$ . Further, there exists a constant C (depending uniquely on the parameters  $K_i$  of conditions A1-A3) such that, for any  $f : \mathbb{R}^D \times \mathbb{R} \to \mathbb{R}$ , with  $\|f\|_{\infty}$ ,  $\|f\|_{\text{Lip}} \leq 1$ ,  $\varepsilon \leq 1$ ,

$$\sup_{k \in [0, T/\varepsilon] \cap \mathbb{N}} \left| \frac{1}{N} \sum_{i=1}^{N} f(\boldsymbol{\theta}_{i}^{k}) - \int f(\boldsymbol{\theta}) \rho_{k\varepsilon} (d\boldsymbol{\theta}) \right| \leq C e^{CT} \operatorname{err}_{N, D}(z),$$

$$\sup_{k \in [0, T/\varepsilon] \cap \mathbb{N}} \left| R_{N}(\boldsymbol{\theta}^{k}) - R(\rho_{k\varepsilon}) \right| \leq C e^{CT} \operatorname{err}_{N, D}(z),$$
(19)

with probability  $1 - e^{-z^2}$ . The same statements hold for noisy SGD (11), provided (7) is replaced by (12), and if  $\beta \geq 1$ ,  $\lambda \leq 1$ , and  $\rho_0$  is  $K_0$  sub-Gaussian for some  $K_0 > 0$ .

Notice that dependence of the error terms in N and D is rather benign. On the other hand, the error grows exponentially with the time horizon T, which limits its applicability to cases in which the DD converges rapidly to a good solution. We do not expect this behavior to be improvable within the general setting of 3, which a priori includes cases in which the dynamics is unstable.

We can regard  $J(\theta; \rho_t) = \rho_t(\theta) \nabla_{\theta} \Psi(\theta; \rho_t)$  as a current. The fixed points of the continuum dynamics are densities that correspond to zero current, as stated below.

**Proposition 2.** Assume  $V(\cdot), U(\cdot, \cdot)$  to be differentiable with bounded gradient. If  $\rho_t$  is a solution of the PDE (7), then  $R(\rho_t)$  is non-increasing. Further, probability distribution  $\rho$  is a fixed point of the PDE (7) if and only if

$$\operatorname{supp}(\rho) \subseteq \left\{ \boldsymbol{\theta} : \nabla_{\boldsymbol{\theta}} \Psi(\boldsymbol{\theta}; \rho) = \mathbf{0} \right\}. \tag{20}$$

Note that global optimizers of  $R(\rho)$ , defined by condition (17), are fixed points, but the set of fixed points is, in general, larger than the set of optimizers. Our next proposition provides an analogous characterization of the fixed points of diffusion DD (12) (see [CMV<sup>+</sup>03] for related results).

**Proposition 3.** Assume that conditions A1-A3 hold and that  $\rho_0$  is absolutely continuous with respect to Lebesgue measure, with  $F_{\beta,\lambda}(\rho_0) < \infty$ . If  $(\rho_t)_{t\geq 0}$  is a solution of the diffusion PDE (12), then  $\rho_t$  is absolutely continuous. Further, there is at most one fixed point  $\rho_* = \rho_*^{\beta,\lambda}$  of (12) satisfying  $F_{\beta,\lambda}(\rho_*) < \infty$ . This fixed point is absolutely continuous and its density satisfies

$$\rho_*(\boldsymbol{\theta}) = \frac{1}{Z(\beta)} \exp\left\{-\beta \Psi_{\lambda}(\boldsymbol{\theta}; \rho_*)\right\}. \tag{21}$$

In the next sections we state our results about convergence of the distributional dynamics to its fixed point. In the case of noisy SGD (and for the diffusion PDE (12)), a general convergence result can be established (although at the cost of an additional regularization). For noiseless SGD (and the continuity equation (12)), we do not have such general result. However, we obtain a stability condition for fixed point containing one point mass, which is useful to characterize possible limiting points (and is used in treating the examples in the previous section).

#### 3.1 Convergence: noisy SGD

Remarkably, the diffusion PDE (12) generically admits a unique fixed point, which is the global minimum of  $F_{\beta,\lambda}(\rho)$  and the evolution (12) converges to it, if initialized so that  $F_{\beta,\lambda}(\rho_0) < \infty$ . This statement requires some qualifications. First of all, we introduce sufficient regularity assumptions to guarantee the existence of sufficiently smooth solutions of (12).

A4 
$$V \in C^4(\mathbb{R}^D), U \in C^4(\mathbb{R}^D \times \mathbb{R}^D), \nabla^k_{\theta_1} U(\theta_1, \theta_2)$$
 is uniformly bounded for  $0 \le k \le 4$ .

Next notice that the right-hand side of the fixed point equation (21) is not necessarily normalizable (for instance, it is not when  $V(\cdot)$ ,  $U(\cdot, \cdot)$  are bounded). In order to ensure the existence of a fixed point, we need  $\lambda > 0$ .

**Theorem 4.** Assume that conditions A1-A4 hold, and  $1/K_0 \le \lambda \le K_0$  for some  $K_0 > 0$  Then  $F_{\beta,\lambda}(\rho)$  has a unique minimizer, denoted by  $\rho_*^{\beta,\lambda}$ , which satisfies

$$R(\rho_*^{\beta,\lambda}) \le \inf_{\boldsymbol{\theta} \in \mathbb{R}^{N \times D}} R_N(\boldsymbol{\theta}) + C D/\beta, \qquad (22)$$

where C is a constant depending on  $K_0, K_1, K_2, K_3$ . Further, letting  $\rho_t$  be a solution of the diffusion PDE (12) with initialization satisfying  $F_{\beta,\lambda}(\rho_0) < \infty$ , we have, as  $t \to \infty$ ,

$$\rho_t \Rightarrow \rho_*^{\beta,\lambda} \,. \tag{23}$$

The proof of this theorem is based on the following formula that describes the free energy decrease along the trajectories of the distributional dynamics (12):

$$\frac{\mathrm{d}F_{\beta,\lambda}(\rho_t)}{\mathrm{d}t} = -2\xi(t) \int_{\mathbb{R}^D} \|\nabla_{\boldsymbol{\theta}}(\Psi_{\lambda}(\boldsymbol{\theta}; \rho_t) + 1/\beta \cdot \log \rho_t(\boldsymbol{\theta}))\|_2^2 \rho_t(\boldsymbol{\theta}) \mathrm{d}\boldsymbol{\theta}. \tag{24}$$

(A key technical hurdle is of course proving that this expression makes sense, which we do by showing the existence of strong solutions.) It follows that the right-hand side must vanish as  $t \to \infty$ , from which we prove that (eventually taking subsequences)  $\rho_t \Rightarrow \rho_*$  where  $\rho_*$  must satisfy  $\beta \Psi_{\lambda}(\theta; \rho_*) + \log \rho_*(\theta) = \text{const.}$  This in turns mean  $\rho_*$  is a solution of the fixed point condition 21 and is in fact a global minimum of  $F_{\beta,\lambda}$  by convexity.

This result can be used in conjunction with Theorem 3, in order to analyze the regularized noisy SGD algorithm (11).

**Theorem 5.** Assume that conditions A1-A4 hold. Let  $\rho_0 \in \mathscr{P}(\mathbb{R}^D)$  be absolutely continuous with  $F_{\beta,\lambda}(\rho_0) < \infty$  and  $K_0$  sub-Gaussian. Consider regularized noisy SGD, cf. (11), at inverse temperature  $\beta < \infty$ , regularization  $1/K_0 \le \lambda \le K_0$  with initialization  $(\boldsymbol{\theta}_i^0)_{i \le N} \sim_{iid} \rho_0$ . Then for any  $\eta > 0$ , there exists  $K = K(\eta, \{K_i\})$  and setting  $\beta \ge KD$ , there exists  $T = T(\eta, V, U, \{K_i\}, D, \beta) < \infty$  and  $C_0 = C_0(\eta, \{K_i\}, \delta)$  (independent of the dimension D and temperature  $\beta$ ) such that the following happens for  $N, (1/\varepsilon) \ge C_0 e^{C_0 T} D$ ,  $\varepsilon \ge 1/N^{10}$ : for any  $k \in [T/\varepsilon, 10T/\varepsilon]$  we have, with probability  $1 - \delta$ ,

$$R_N(\boldsymbol{\theta}^k) \le \inf_{\rho \in \mathscr{P}(\mathbb{R}^D)} R_\lambda(\rho) + \eta.$$
 (25)

Let us emphasize that the convergence time T in the last theorem can depend on the dimension D and on the data distribution  $\mathbb{P}$ , but is independent of the number of hidden units N. As illustrated by the examples in the previous section, understanding the dependence of T on D requires further analysis, but examining the proof of this theorem suggests  $T = e^{O(D)}$  quite generally (examples in which T = O(1) or  $T = e^{\Theta(D)}$  can be constructed). We expect that our techniques could be pushed to investigate the dependence of T on  $\eta$  (see discussion in SI). In highly structured cases, the dimension D can be of constant order, and be much smaller than d.

# 3.2 Convergence: noiseless SGD

The next theorems provide necessary and sufficient conditions for distributions containing a single point mass to be a stable fixed point of the evolution. This result is useful in order to characterize the large time asymptotics of the dynamics (7). Here, we write  $\nabla_1 U(\theta_1, \theta_2)$  for the gradient of

U with respect to its first argument, and  $\nabla_{1,1}^2 U$  for the corresponding Hessian. Further, for a probability distribution  $\rho_*$ , we define

$$\boldsymbol{H}_0(\rho_*) = \nabla^2 V(\boldsymbol{\theta}_*) + \int \nabla_{1,1}^2 U(\boldsymbol{\theta}_*, \boldsymbol{\theta}) \, \rho_*(\mathrm{d}\boldsymbol{\theta}) \,. \tag{26}$$

Note that  $\mathbf{H}_0(\rho_*)$  is nothing but the Hessian of  $\boldsymbol{\theta} \mapsto \Psi(\boldsymbol{\theta}; \rho_*)$  at  $\boldsymbol{\theta}_*$ .

**Theorem 6.** Assume V, U to be twice differentiable with bounded gradient and bounded continuous Hessian. Let  $\boldsymbol{\theta}_* \in \mathbb{R}^D$  be given. Then  $\rho_* = \delta_{\boldsymbol{\theta}_*}$  is a fixed point of the evolution (7) if and only if  $\nabla V(\boldsymbol{\theta}_*) + \nabla_1 U(\boldsymbol{\theta}_*, \boldsymbol{\theta}_*) = \mathbf{0}$ .

Define  $\mathbf{H}_0(\delta_{\boldsymbol{\theta}_*}) \in \mathbb{R}^{D \times D}$  as per (26). If  $\lambda_{\min}(\mathbf{H}_0(\delta_{\boldsymbol{\theta}_*})) > 0$ , then there exists  $r_0 > 0$  such that, if  $\operatorname{supp}(\rho_{t_0}) \subseteq \mathsf{B}(\boldsymbol{\theta}_*; r_0) \equiv \{\boldsymbol{\theta} : \|\boldsymbol{\theta} - \boldsymbol{\theta}_*\|_2 \le r_0\}$ , then  $\rho_t \Rightarrow \rho_*$  as  $t \to \infty$ . In fact, convergence is exponentially fast, namely  $\int \|\boldsymbol{\theta} - \boldsymbol{\theta}_*\|_2^2 \rho_t(\mathrm{d}\boldsymbol{\theta}) \le e^{-\lambda(t-t_0)}$  for some  $\lambda > 0$ .

**Theorem 7.** Under the same assumptions of Theorem 6, let  $\rho_* = p_* \delta_{\boldsymbol{\theta}_*} + (1 - p_*) \tilde{\rho}_* \in \mathscr{P}(\mathbb{R}^D)$  be a fixed point of dynamics (7), with  $p_* \in (0,1]$  and  $\nabla \Psi(\boldsymbol{\theta}_*; \rho_*) = \mathbf{0}$  (which, in particular, is implied by the fixed point condition (20)). Define the level sets  $\mathcal{L}(\eta) \equiv \{\boldsymbol{\theta} : \Psi(\boldsymbol{\theta}; \rho_*) \leq \Psi(\boldsymbol{\theta}_*; \rho_*) - \eta\}$  and make the following assumptions: B1. The eigenvalues of  $\boldsymbol{H}_0 = \boldsymbol{H}_0(\rho_*)$  are all different from 0, with  $\lambda_{\min}(\boldsymbol{H}_0) < 0$ ; B2.  $\tilde{\rho}_*(\mathcal{L}(\eta)) \uparrow 1$  as  $\eta \downarrow 0$ ; B3. There exists  $\eta_0 > 0$  such that the sets  $\partial \mathcal{L}(\eta)$  are compact for all  $\eta \in (0, \eta_0)$ .

If  $\rho_0$  has a bounded density with respect to Lebesgue measure, then it cannot be that  $\rho_t$  converges weakly to  $\rho_*$  as  $t \to \infty$ .

# 4 Discussion and future directions

In this paper we developed a new approach to the analysis of two-layers neural networks. Using a propagation-of-chaos argument, we proved that –if the number of hidden units satisfies  $N\gg D$ –SGD dynamics is well approximated by the PDE in (7), while noisy SGD is well approximated by (12). Both of these asymptotic descriptions correspond to Wasserstein gradient flows for certain energy (or free energy) functionals. While empirical risk minimization is known to be insensitive to overparametrization [Bar98], the present work clarifies that the SGD behavior is also independent of the number of hidden units, as soon as this is large enough.

We illustrated our approach on several concrete examples, by proving convergence of SGD to a near-global optimum. This type of analysis provides a new mechanism for avoiding the perils of non-convexity. We do not prove that the finite-N risk  $R_N(\theta)$  has a unique local minimum, or that all local minima are close to each other. Such claims have often been the target of earlier work, but might be too strong for the case of neural networks. We prove instead that the PDE (7) converges to a near global optimum, when initialized with a bounded density. This effectively gets rid of some exceptional stationary points of  $R_N(\theta)$ , and merges multiple finite N stationary points that result into similar distributions  $\rho$ .

In the case of noisy SGD (11), we prove that it converges generically to a near-global minimum of the regularized risk, in time independent of the number of hidden units.

We emphasize that while we focused here on the case of square loss, our approach should be generalizable to other loss functions as well, cf. SI.

The present work opens the way to several interesting research directions. We will mention two of them. (i) The PDE (7) corresponds to gradient flow in the Wasserstein metric for the risk  $R(\rho)$ ,

see [AGS08]. Building on this remark, tools from optimal transportation theory can be used to prove convergence. (ii) Multiple finite-N local minima can correspond to the same minimizer  $\rho_*$  of  $R(\rho)$  in the limit  $N \to \infty$ . Ideas from glass theory [MP99] might be useful to investigate this structure.

Let us finally mention that, after a first version of this paper appeared as a preprint, several other groups obtained results that are closely related to Theorem 3 [RVE18, SS18, CB18].

# Acknowledgements

This work was partially supported by grants NSF DMS-1613091, NSF CCF-1714305 and NSF IIS-1741162. S. M. was partially supported by Office of Technology Licensing Stanford Graduate Fellowship. P.-M. N. was partially supported by William R. Hewlett Stanford Graduate Fellowship. The authors would like to thank Jiajun Tong for helpful discussions concerning strong solutions for parabolic PDEs.

### References

- [ABGM14] Sanjeev Arora, Aditya Bhaskara, Rong Ge, and Tengyu Ma. Provable bounds for learning some deep representations. In *Proc. of International Conference on Machine Learning (ICML)*, pages 584–592, 2014.
- [AGS08] Luigi Ambrosio, Nicola Gigli, and Giuseppe Savaré. Gradient flows: in metric spaces and in the space of probability measures. Springer Science & Business Media, 2008.
- [Bar98] Peter L Bartlett. The sample complexity of pattern classification with neural networks: the size of the weights is more important than the size of the network. *IEEE transactions on Information Theory*, 44(2):525–536, 1998.
- [BG17] Alon Brutzkus and Amir Globerson. Globally optimal gradient descent for a convnet with gaussian inputs. arXiv:1702.07966, 2017.
- [Bot10] Léon Bottou. Large-scale machine learning with stochastic gradient descent. In *Proceedings of COMPSTAT'2010*, pages 177–186. Springer, 2010.
- [BRV<sup>+</sup>06] Yoshua Bengio, Nicolas L Roux, Pascal Vincent, Olivier Delalleau, and Patrice Marcotte. Convex neural networks. In *Advances in neural information processing systems*, pages 123–130, 2006.
- [CB18] Lenaic Chizat and Francis Bach. On the global convergence of gradient descent for over-parameterized models using optimal transport. arXiv:1805.09545, 2018.
- [CDF<sup>+</sup>11] José A Carrillo, Marco DiFrancesco, Alessio Figalli, Thomas Laurent, Dejan Slepčev, et al. Global-in-time weak measure solutions and finite-time aggregation for nonlocal interaction equations. *Duke Mathematical Journal*, 156(2):229–271, 2011.
- [CMV<sup>+</sup>03] José A Carrillo, Robert J McCann, Cédric Villani, et al. Kinetic equilibration rates for granular media and related equations: entropy dissipation and mass transportation estimates. *Revista Matematica Iberoamericana*, 19(3):971–1018, 2003.

- [CMV06] José A Carrillo, Robert J McCann, and Cédric Villani. Contractions in the 2-wasserstein length space and thermalization of granular media. *Archive for Rational Mechanics and Analysis*, 179(2):217–263, 2006.
- [Eva09] Lawrence C. Evans. Partial Differential Equations. Springer, 2009.
- [GBCB16] Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1. MIT press Cambridge, 2016.
- [GLM17] Rong Ge, Jason D Lee, and Tengyu Ma. Learning one-hidden-layer neural networks with landscape design. arXiv:1711.00501, 2017.
- [GP10] Victor Guillemin and Alan Pollack. *Differential topology*, volume 370. American Mathematical Soc., 2010.
- [JKO98] Richard Jordan, David Kinderlehrer, and Felix Otto. The variational formulation of the fokker–planck equation. SIAM journal on mathematical analysis, 29(1):1–17, 1998.
- [JSA15] Majid Janzamin, Hanie Sedghi, and Anima Anandkumar. Beating the perils of non-convexity: Guaranteed training of neural networks using tensor methods. arXiv:1506.08473, 2015.
- [KSH12] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [Lan13] Serge Lang. Complex analysis, volume 103. Springer Science & Business Media, 2013.
- [LBW96] Wee Sun Lee, Peter L Bartlett, and Robert C Williamson. Efficient agnostic learning of neural networks with bounded fan-in. *IEEE Transactions on Information Theory*, 42(6):2118–2132, 1996.
- [LSU88] Olga Aleksandrovna Ladyzhenskaia, Vsevolod Alekseevich Solonnikov, and Nina N Ural'tseva. *Linear and quasi-linear equations of parabolic type*, volume 23. American Mathematical Soc., 1988.
- [MBM16] Song Mei, Yu Bai, and Andrea Montanari. The landscape of empirical risk for non-convex losses. arXiv:1607.06534, 2016.
- [Mit15] Boris Mityagin. The zero set of a real analytic function. arXiv:1512.07276, 2015.
- [MP99] Marc Mézard and Giorgio Parisi. Thermodynamics of glasses: A first principles computation. *Journal of Physics: Condensed Matter*, 11(10A):A157, 1999.
- [MV00] Peter A Markowich and Cédric Villani. On the trend to equilibrium for the fokkerplanck equation: an interplay between physics and functional analysis. *Mat. Contemp*, 19:1–29, 2000.
- [RM51] Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathematical statistics*, pages 400–407, 1951.

- [Ros62] Frank Rosenblatt. Principles of neurodynamics. Spartan Book, 1962.
- [RVE18] Grant M Rotskoff and Eric Vanden-Eijnden. Neural networks as interacting particle systems: Asymptotic convexity of the loss landscape and universal scaling of the approximation error. arXiv:1805.00915, 2018.
- [SA15] Hanie Sedghi and Anima Anandkumar. Provable methods for training neural networks with sparse connectivity. In *Proc. of International Conference on Learning Representation (ICLR)*, 2015.
- [San15] Filippo Santambrogio. Optimal Transport for Applied Mathematicians: Calculus of Variations, PDEs, and Modeling, volume 87. Birkhäuser, 2015.
- [SJL17] Mahdi Soltanolkotabi, Adel Javanmard, and Jason D. Lee. Theoretical insights into the optimization landscape of over-parameterized shallow neural networks. arXiv:1707.04926, 2017.
- [SS18] Justin Sirignano and Konstantinos Spiliopoulos. Mean field analysis of neural networks. arXiv:1805.01053, 2018.
- [SSBD14] Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From the-ory to algorithms*. Cambridge university press, 2014.
- [Szn91] Alain-Sol Sznitman. Topics in propagation of chaos. In *Ecole d'été de probabilités de Saint-Flour XIX—1989*, pages 165–251. Springer, 1991.
- [Tia17] Yuandong Tian. Symmetry-breaking convergence analysis of certain two-layered neural networks with ReLU nonlinearity. In Workshop at International Conference on Learning Representation (ICLR), 2017.
- [WML17] Chuang Wang, Jonathan Mattingly, and Yue M Lu. Scaling limit: Exact and tractable analysis of online learning algorithms with applications to regularized regression and pca. arXiv:1712.04332, 2017.
- [ZLJ16] Yuchen Zhang, Jason D. Lee, and Michael I. Jordan. L1-regularized neural networks are improperly learnable in polynomial time. In *Proc. of International Conference on Machine Learning (ICML)*, pages 993–1001, 2016.
- [ZSJ<sup>+</sup>17] Kai Zhong, Zhao Song, Prateek Jain, Peter L. Bartlett, and Inderjit S. Dhillon. Recovery guarantees for one-hidden-layer neural networks. arXiv:1706.03175, 2017.

# Supplementary information

We present here proofs and additional technical details for our mathematical results, as well as additional information concerning the numerical experiments.

### 5 Notations

We use lowercase bold for vectors (e.g. u, v, ...), uppercase bold for matrices (e.g. A, B, ...), and lowercase plain for scalar (x, y, ...).

- Given a measurable space  $\Omega$ , we denote by  $\mathscr{P}(\Omega)$  the set of probability measures on  $\Omega$ .
- $\mathsf{B}^d(x;r)$  denotes the Euclidean ball with center x and radius r in  $\mathbb{R}^d$ . We will drop the dimension superscript whenever clear from the context.
- Given a measurable function f, and a measure  $\mu$ , we denote by  $\langle f, \mu \rangle = \langle \mu, f \rangle = \int f \, d\mu$  the corresponding integral.
- For a univariate function  $f : \mathbb{R} \to \mathbb{R}$ , we denote by f'(x) its derivative at x. If the argument is time, we will also use  $\dot{f}(t)$ .
- $||f||_{\text{Lip}} \equiv \sup_{x \neq y} |f(x) f(y)| / ||x y||_2$  denotes the Lipshitz constant of a function f.
- $d_{\text{BL}}(\,\cdot\,,\,\cdot\,)$  is the bounded Lipschitz distance between probability measures

$$d_{\mathrm{BL}}(\mu,\nu) = \sup \left\{ \left| \int f(\boldsymbol{x}) \, \mu(\mathrm{d}\boldsymbol{x}) - \int f(\boldsymbol{x}) \, \nu(\mathrm{d}\boldsymbol{x}) \right| : \|f\|_{\infty} \le 1, \|f\|_{\mathrm{Lip}} \le 1 \right\}$$
 (5.1)

$$\leq 2 \inf_{\gamma \in \mathcal{C}(\mu,\nu)} \int (\|\boldsymbol{x} - \boldsymbol{y}\|_{2} \wedge 1) \gamma(d\boldsymbol{x}, d\boldsymbol{y}) \leq 4 d_{BL}(\mu,\nu).$$
 (5.2)

Here  $C(\mu, \nu)$  is the set of couplings of  $\mu$  and  $\nu$ .

•  $W_p(\cdot,\cdot)$  is the Wasserstein distance between probability measures

$$W_p(\mu, \nu) = \left(\inf_{\gamma \in \mathcal{C}(\mu, \nu)} \int \|\boldsymbol{x} - \boldsymbol{y}\|_2^p \gamma(\mathrm{d}\boldsymbol{x}, \mathrm{d}\boldsymbol{y})\right)^{1/p}.$$
 (5.3)

For p = 1, the Kantorovich-Rubinstein duality gives

$$W_1(\mu, \nu) = \sup \left\{ \left| \int f(\boldsymbol{x}) \, \mu(\mathrm{d}\boldsymbol{x}) - \int f(\boldsymbol{x}) \, \nu(\mathrm{d}\boldsymbol{x}) \right| : \|f\|_{\mathrm{Lip}} \le 1 \right\}. \tag{5.4}$$

- K is a generic constant depending on  $K_0, K_1, K_2, K_3$ , where  $K_i$ 's are constants which will be specified from the context.
- $\mathbb{N} = \{0, 1, 2, \ldots\}$  denote the set of natural numbers.

# 6 General results: Statics

In this section, we discuss some properties of the population risk,  $R_N(\theta)$ , and its continuum counterpart  $R(\rho)$ . For future reference, we copy the key definitions from the main text:

$$R_N(\boldsymbol{\theta}) \equiv R_\# + \frac{2}{N} \sum_{i=1}^N V(\boldsymbol{\theta}_i) + \frac{1}{N^2} \sum_{i,j=1}^N U(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j), \qquad (6.1)$$

$$R(\rho) \equiv R_{\#} + 2 \int V(\boldsymbol{\theta}) \, \rho(\mathrm{d}\boldsymbol{\theta}) + \int U(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) \, \rho(\mathrm{d}\boldsymbol{\theta}_1) \, \rho(\mathrm{d}\boldsymbol{\theta}_2) \,, \tag{6.2}$$

$$R_{\#} = \mathbb{E}\{y^2\}, \qquad V(\boldsymbol{\theta}) = -\mathbb{E}\{y\,\sigma_*(\boldsymbol{x};\boldsymbol{\theta})\}, \tag{6.3}$$

$$U(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) = \mathbb{E} \{ \sigma_*(\boldsymbol{x}; \boldsymbol{\theta}_1) \sigma_*(\boldsymbol{x}; \boldsymbol{\theta}_2) \}.$$
(6.4)

We further recall the notation

$$\Psi(\boldsymbol{\theta}; \rho) = V(\boldsymbol{\theta}) + \int U(\boldsymbol{\theta}, \boldsymbol{\theta}') \, \rho(\mathrm{d}\boldsymbol{\theta}') \,. \tag{6.5}$$

We will always assume that the expectations defining  $V(\boldsymbol{\theta}), U(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)$  exist finite for all  $\boldsymbol{\theta}, \boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \mathbb{R}^D$ . A necessary and sufficient condition for this is that  $\mathbb{E}\{\sigma_*(\boldsymbol{x};\boldsymbol{\theta})^2\} < \infty$  for all  $\boldsymbol{\theta}$ . Since in most cases of interest  $|\sigma_*(\boldsymbol{x};\boldsymbol{\theta})| \leq M(\boldsymbol{\theta}) ||\boldsymbol{x}||_2$ , for this to happen, it is sufficient that  $\boldsymbol{x}$  has a finite second moment.

Note that this  $\rho \mapsto R(\rho)$  is a convex function on the set of probability measures on  $\mathbb{R}^D$ . We will denote by  $\mathscr{P}_{V,U}$  the subset of probability measures  $\rho$  such that the expectations on the right-hand side are finite. We define  $R(\rho) = \infty$  if  $\rho \in \mathscr{P}(\mathbb{R}^D) \setminus \mathscr{P}_{V,U}$ .

#### 6.1 Proof of Proposition 1

The proof is divided in two parts:

1. We show that minimizing the population risk  $R_N(\boldsymbol{\theta})$  yields similar results to minimizing its continuum counterpart  $R(\rho)$ :

$$\left|\inf_{\boldsymbol{\theta}} R_N(\boldsymbol{\theta}) - \inf_{\rho} R(\rho)\right| \le \frac{K}{N}. \tag{6.6}$$

2. We establish the condition for  $\rho_*$  to be a minimizer:

$$\operatorname{supp}(\rho_*) \subseteq \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^D} \Psi(\boldsymbol{\theta}; \rho_*). \tag{6.7}$$

First notice that, for any  $\boldsymbol{\theta} = (\boldsymbol{\theta}_i)_{i \leq N}$ , we have

$$R_N(\boldsymbol{\theta}) \ge \inf_{\boldsymbol{\theta}} R(\boldsymbol{\rho}) \,. \tag{6.8}$$

Indeed,  $R_N(\boldsymbol{\theta}) = R(\rho)$  for  $\rho = (1/N) \sum_{i=1}^N \delta_{\boldsymbol{\theta}_i}$ .

In order to prove Eq. (6.6), let  $\rho_* \in \mathscr{P}(\mathbb{R}^D)$  be such that  $R(\rho_*) = R_*$  under assumption (a), or  $R(\rho_*) \leq R_* + \varepsilon$  under assumption (b). Let  $(\theta_i)_{i < N} \sim_{iid} \rho_*$ . A simple calculation shows that

$$\mathbb{E}_{\boldsymbol{\theta}}[R_N(\boldsymbol{\theta})] - R(\rho_*) = \frac{1}{N} \left\{ \int U(\boldsymbol{\theta}, \boldsymbol{\theta}) \, \rho_*(\mathrm{d}\boldsymbol{\theta}) - \int U(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) \, \rho_*(\mathrm{d}\boldsymbol{\theta}_1) \, \rho_*(\mathrm{d}\boldsymbol{\theta}_2) \right\}$$
(6.9)

$$\leq \frac{1}{N} \int U(\boldsymbol{\theta}, \boldsymbol{\theta}) \, \rho_*(\mathrm{d}\boldsymbol{\theta}) \leq \frac{K}{N},$$
 (6.10)

where the first inequality follows since  $\int U(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) \rho_*(\mathrm{d}\boldsymbol{\theta}_1) \rho_*(\mathrm{d}\boldsymbol{\theta}_2) = \mathbb{E}\{y(\boldsymbol{x})^2\} \geq 0$  for  $y(\boldsymbol{x}) = \int \sigma_*(\boldsymbol{x}; \boldsymbol{\theta}) \rho_*(\mathrm{d}\boldsymbol{\theta})$ , and the second inequality follows by assumption. It follows that

$$\inf_{\boldsymbol{\theta}} R_N(\boldsymbol{\theta}) \le R_* + \frac{K}{N} + \varepsilon, \tag{6.11}$$

whence the claim (6.6) follows since  $\varepsilon$  is arbitrary.

We next establish the minimum condition (6.7). Notice that since  $V(\cdot)$  is continuous, and  $U(\cdot, \cdot)$  is bounded below, it follows from Fatou's lemma that, for any  $\rho$ , the function  $\boldsymbol{\theta} \mapsto \Psi(\boldsymbol{\theta}; \rho)$  is lower semicontinuous and takes values in  $(-\infty, \infty]$ . In particular the set  $S_0(\rho) \equiv \arg \min_{\boldsymbol{\theta}} \Psi(\boldsymbol{\theta}; \rho)$  must be closed.

We first prove that any minimizer must satisfy (6.7). Let  $\rho_*$  be a minimizer and define  $\Psi_* = \inf_{\theta} \Psi(\theta; \rho_*)$ . By rearranging terms, for any probability measure  $\rho$ , we have

$$R(\rho) - R(\rho_*) = 2\langle \Psi(\cdot; \rho_*), (\rho - \rho_*) \rangle + \langle U, (\rho - \rho_*)^{\otimes 2} \rangle. \tag{6.12}$$

First we will assume  $\Psi_* > -\infty$  (whence, by lower semicontinuity,  $S_0(\rho_*)$  must be a non-empty closed set). Let  $\theta_1 \in S_0(\rho_*)$ , and assume by contradiction that there exist  $\theta_0 \in \operatorname{supp}(\rho_*)$ ,  $\theta_0 \notin S_0(\rho_*)$ . Let  $\mathsf{B}(\theta_0;\varepsilon)$  be a ball of radius  $\varepsilon$  around  $\theta_0$ . By lower semicontinuity, we can find  $\varepsilon_0, \Delta > 0$  such that  $\inf_{\theta \in \mathsf{B}(\theta_0;\varepsilon_0)} \Psi(\theta;\rho_*) = \Psi_* + \Delta > \Psi_*$ . Further  $t_0 \equiv \rho_*(\mathsf{B}(\theta_0;\varepsilon_0)) > 0$  because  $\theta_0 \in \operatorname{supp}(\rho_*)$ .

Let  $\nu \equiv \mathbf{1}_{\mathsf{B}(\boldsymbol{\theta}_0;\varepsilon_0)}\rho_*/t_0$  (i.e.  $\nu$  is the conditional distribution given  $\boldsymbol{\theta} \in \mathsf{B}(\boldsymbol{\theta}_0;\varepsilon_0)$ ). Define, for  $t \in [0,t_0]$ , the probability measure

$$\rho_t = \rho_* - t\nu + t\delta_{\theta_1} \,. \tag{6.13}$$

Using Eq. (6.12), we get

$$R(\rho_t) - R(\rho_*) = 2\langle \Psi(\cdot; \rho_*), (\delta_{\theta_1} - \nu) \rangle t + \langle U, (\delta_{\theta_1} - \nu)^{\otimes 2} \rangle t^2$$
(6.14)

$$\leq 2(\Psi_* - \Psi_* - \Delta)t + C_0t^2 = -2\Delta t + C_0t^2, \tag{6.15}$$

where the second inequality follows from the fact that U is continuous and  $\delta_{\theta_1}$ ,  $\nu$  have bounded support. By taking t small enough, we get  $R(\rho) < R(\rho_*)$  hence reaching a contradiction.

Next consider the case in which  $\Psi_* \equiv \inf_{\boldsymbol{\theta}} \Psi(\boldsymbol{\theta}; \rho_*) = -\infty$ . For  $M \in \mathbb{N}$ ,  $M \geq 1$ , let  $\boldsymbol{\theta}_M \in \mathbb{R}^D$  be such that  $\Psi(\boldsymbol{\theta}_M; \rho_*) \leq -M$ . For  $\boldsymbol{\theta}_0 \in \text{supp}(\boldsymbol{\theta}_*)$ , construct  $\nu$  as before. Note that, and call  $\inf_{\boldsymbol{\theta} \in \mathsf{B}(\boldsymbol{\theta}_0; \varepsilon_0)} \Psi(\boldsymbol{\theta}; \rho_*) = \Psi_0$ . Define, for  $t \in [0, t_0]$ 

$$\rho_{M,t} = \rho_* - t\nu + t\delta_{\boldsymbol{\theta}_M} \,. \tag{6.16}$$

By applying again Eq. (6.12), we get

$$R(\rho_{M,t}) - R(\rho_*) = 2\langle \Psi(\cdot; \rho_*), (\delta_{\boldsymbol{\theta}_M} - \nu) \rangle t + \langle U, (\delta_{\boldsymbol{\theta}_M} - \nu)^{\otimes 2} \rangle t^2$$
(6.17)

$$\leq -2(M + \Psi_0)t + C_0(M)t^2$$
. (6.18)

By selecting  $t = t_M = \min(t_0, (M + \Psi_0)/C_0(M))$  (which is positive for all M large enough), we obtain  $R(\rho_{M,t}) - R(\rho_*) < 0$  for all M large and hence reach a contradiction.

We finally prove that condition (6.7) is sufficient for  $\rho_*$  to be a minimizer. Indeed, for any non-negative measurable function  $\mu: \mathbb{R}^D \to \mathbb{R}$ , letting  $\Psi_* = \min_{\theta} \Psi(\theta; \rho_*)$ ,

$$R(\rho) \ge R_{\#} + 2\langle V, \rho \rangle + \langle U, \rho^{\otimes 2} \rangle - \langle \mu, \rho \rangle \tag{6.19}$$

$$= R(\rho_*) + 2\langle \Psi(\cdot; \rho_*), \rho - \rho_* \rangle + \langle U, (\rho - \rho_*)^{\otimes 2} \rangle - \langle \mu, \rho \rangle$$
(6.20)

$$= R(\rho_*) + 2\langle \Psi(\cdot; \rho_*) - \Psi_*, \rho - \rho_* \rangle + \langle U, (\rho - \rho_*)^{\otimes 2} \rangle - \langle \mu, \rho \rangle.$$
 (6.21)

Setting  $\mu = 2[\Psi(\cdot; \rho_*) - \Psi_*]$ , and noticing that condition (6.7) implies  $\langle \Psi(\cdot; \rho_*) - \Psi_*, \rho_* \rangle = 0$ , we get  $R(\rho) \ge R(\rho_*) + \langle U, (\rho - \rho_*)^{\otimes 2} \rangle \ge R(\rho_*).$ 

#### 6.2Some additional results

We often find empirically that the optimal density  $\rho_*$  is supported on a set of Lebesgue measure 0 (sometimes on a finite set of points). The following consequence of the previous results partially explains these findings.

Corollary 1. Assume  $\theta \mapsto V(\theta)$  to be an analytic function and  $(\theta_1, \theta_2) \mapsto U(\theta_1, \theta_2)$  to be analytic with respect to  $\theta_1$ , uniformly in  $\theta_2$ . Namely there exists a locally bounded function  $\theta \mapsto B(\theta)$  such that  $\|\nabla_{\boldsymbol{\theta}_1}^k U(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)\|_2 \leq k! B(\boldsymbol{\theta}_1)^k$  for all k,  $\boldsymbol{\theta}_1$ ,  $\boldsymbol{\theta}_2$ . If  $\rho_*$  is a minimizer of  $R(\rho)$ , then one of the following holds

- (a)  $\Psi(\boldsymbol{\theta}; \rho_*) = \Psi_*$  for some constant  $\Psi_*$  and all  $\boldsymbol{\theta} \in \mathbb{R}^D$ .
- (b) The support of  $\rho_*$  has zero Lebesque measure.

If D=1, then (b) can be replaced by: (b')  $\rho_*$  is a convex combination of countably many point masses with no accumulation point (finitely many if  $\Psi(\theta; \rho_*) \to \infty$  as  $|\theta| \to \infty$ ).

*Proof.* Note that, under the stated conditions  $f(\theta) \equiv \int U(\theta, \theta') \rho_*(d\theta')$  is analytic. by a standard dominated convergence argument, we have that  $\nabla^k f$  is given by the integral of  $\int \nabla^k U(\theta_1, \theta_2) \, \rho_*(\mathrm{d}\theta_2)$  for any  $k \geq 0$ . Further, by an application of the intermediate value theorem there exists  $t_{\theta_1,\theta_2,\delta} \in [0,1]$  such that

$$\left| f(\boldsymbol{\theta}_1 + \boldsymbol{\delta}) - \sum_{\ell=0}^{k-1} \frac{1}{\ell!} \left\langle \nabla^{\ell} f(\boldsymbol{\theta}_1), \boldsymbol{\delta}^{\otimes \ell} \right\rangle \right| \leq \frac{1}{k!} \left| \int \left\langle \nabla_{\boldsymbol{\theta}_1}^k U(\boldsymbol{\theta}_1 + t_{\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\delta}} \boldsymbol{\delta}, \boldsymbol{\theta}_2), \boldsymbol{\delta}^{\otimes k} \right\rangle \rho_*(\mathrm{d}\boldsymbol{\theta}_2) \right|$$
(6.22)

$$\leq \int B(\boldsymbol{\theta}_1 + t_{\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\delta}} \boldsymbol{\delta})^k \|\boldsymbol{\delta}\|_2^k \, \rho_*(\mathrm{d}\boldsymbol{\theta}_2) \tag{6.23}$$

$$\leq \int B(\boldsymbol{\theta}_1 + t_{\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\delta}} \boldsymbol{\delta})^k \|\boldsymbol{\delta}\|_2^k \, \rho_*(\mathrm{d}\boldsymbol{\theta}_2) \tag{6.23}$$

$$\leq \sup_{\boldsymbol{\theta} \in \mathsf{B}(\boldsymbol{\theta}_1; \|\boldsymbol{\delta}\|_2)} B(\boldsymbol{\theta})^k \|\boldsymbol{\delta}\|_2^k, \tag{6.24}$$

which vanishes as  $k \to \infty$  for uniformly over  $\|\boldsymbol{\delta}\|_2 \le \delta_0$  for  $\delta_0$  small enough.

Let  $\Psi_* = \min_{\boldsymbol{\theta} \in \mathbb{R}^D} \Psi(\boldsymbol{\theta}; \rho_*)$ . We thus have that  $\boldsymbol{\theta} \mapsto \Psi(\boldsymbol{\theta}; \rho_*)$  is also analytic and so is  $\boldsymbol{\theta} \mapsto$  $\Psi(\theta; \rho_*) - \Psi_*$ . Since  $\operatorname{supp}(\rho_*) \subseteq \{\theta : \Psi(\theta; \rho_*) = \Psi_*\}$ , the claim follows from the fact that the set of zeros of a non-trivial analytic function has vanishing Lebesgue measure [Mit15]. In the case D=1, the set of zeros of an analytic function cannot have any accumulation point [Lan13], which therefore allows to replace (b) with (b'). 

# 7 General results: Dynamics

In this section we consider the SGD dynamics with step size  $s_k = \varepsilon \xi(k\varepsilon)$ , under the assumptions A1, A2, A3 stated in the main text. For the readers convenience, we reproduce here the form of the limiting PDE

$$\partial_t \rho_t(\boldsymbol{\theta}) = 2\xi(t) \nabla \cdot \left[ \rho_t(\boldsymbol{\theta}) \nabla \Psi(\boldsymbol{\theta}; \rho_t) \right], \tag{7.1}$$

$$\Psi(\boldsymbol{\theta}; \rho) = V(\boldsymbol{\theta}) + \int U(\boldsymbol{\theta}, \boldsymbol{\theta}') \, \rho(\mathrm{d}\boldsymbol{\theta}') \,. \tag{7.2}$$

Recall that this is an evolution in the space of probability measures in  $\mathbb{R}^D$ , and is to be interpreted in weak sense. Namely  $\rho_t$  is a solution of Eq. (7.1), if, for any bounded differentiable function  $\varphi : \mathbb{R}^D \to \mathbb{R}$  with bounded gradient:

$$\frac{\mathrm{d}}{\mathrm{d}t}\langle \rho_t, \varphi \rangle = -2\xi(t) \int \langle \nabla \varphi(\boldsymbol{\theta}), \nabla \Psi(\boldsymbol{\theta}; \rho_t) \rangle \rho_t(\mathrm{d}\boldsymbol{\theta}). \tag{7.3}$$

For background on this and similar PDEs (and the analogous ones at finite temperature, cf. Section 10), we refer to [MV00, CMV<sup>+</sup>03, CMV06, AGS08, CDF<sup>+</sup>11]. Our treatment will be mostly self-contained because of some differences between our setting and the one in these papers.

**Remark 7.1.** Recall assumptions A1, A2, A3 in the main text. By [Szn91, Theorem 1.1], assumptions A1 and A3 are sufficient for the existence and uniqueness of solution of PDE (7.1).

A very useful tool for the analysis of the PDE (7.1) is provided by the following nonlinear dynamics. We introduce trajectories  $(\overline{\boldsymbol{\theta}}_i^t)_{1 \leq i \leq N, t \in \mathbb{R}_{\geq 0}}$  by letting  $\overline{\boldsymbol{\theta}}_i^0 = \boldsymbol{\theta}_i^0$  to be the same initialization as for SGD and, for  $t \geq 0$  (here  $P_X$  denotes the law of the random variable X):

$$\overline{\boldsymbol{\theta}}_{i}^{t} = \boldsymbol{\theta}_{i}^{0} - 2 \int_{0}^{t} \xi(s) \nabla \Psi(\overline{\boldsymbol{\theta}}_{i}^{s}; \rho_{s}) \, \mathrm{d}s, \qquad (7.4)$$

$$\rho_s = P_{\overline{\theta}_z^s} \,. \tag{7.5}$$

This should be regarded as an equation for the law of the trajectory  $(\overline{\theta}_i^t)_{t \in \mathbb{R}_{\geq 0}}$ , with boundary condition determined by  $\overline{\theta}_i^0 \sim \rho_0$ . As implied by [Szn91, Theorem 1.1], under the same assumptions A1 and A3, the nonlinear dynamics has a unique solution, with  $\rho_t$  satisfying Eq. (7.1).

**Lemma 7.1.** Assume conditions A1 and A3 hold. Let  $(\rho_t)_{t\geq 0}$  be the solution of the PDE (7.1). Let  $(\overline{\theta}_i^t)_{t\geq 0}$  be the solution of nonlinear dynamics (7.4). Then  $t\mapsto \overline{\theta}_i^t$  is  $K_1K_3$ -Lipschitz continuous, and  $t\mapsto \rho_t$  is  $K_1K_3$ -Lipschitz continuous in  $W_2$  Wasserstein distance, with  $K_1$  and  $K_3$  as per conditions A1 and A3. In particular,  $t\mapsto \rho_t$  is continuous in the topology of weak convergence.

*Proof.* Since  $\xi$  and  $\nabla \Psi$  are  $K_1$  and  $K_3$  bounded respectively,  $t \mapsto \overline{\theta}_i^t$  is  $K_1K_3$ -Lipschitz continuous. Further, Eq. (5.2) implies that  $t \mapsto \rho_t$  is Lipschitz continuous in  $W_2$  Wasserstein distance, namely

$$d_{\mathrm{BL}}(\rho_t, \rho_s) \le W_2(\rho_t, \rho_s) \le (\mathbb{E}[\|\overline{\boldsymbol{\theta}}_i^t - \overline{\boldsymbol{\theta}}_i^s\|_2^2])^{1/2} \le K_1 K_3 |t - s|. \tag{7.6}$$

We notice that, under the nonlinear dynamics, the trajectories  $(\overline{\boldsymbol{\theta}}_1^t)_{t \in \mathbb{R}_{\geq 0}}, \ldots, (\overline{\boldsymbol{\theta}}_N^t)_{t \in \mathbb{R}_{\geq 0}}$  are independent and identically distributed. In particular, this implies that, almost surely,

$$\frac{1}{N} \sum_{i=1}^{N} \delta_{\overline{\theta}_{i}^{t}} \stackrel{\mathrm{d}}{\Rightarrow} \rho_{t} . \tag{7.7}$$

### 7.1 Proof of Theorem 3: Convergence to the PDE

The proof follows a 'propagation of chaos' argument [Szn91]. Throughout this proof, we will use K to denote generic constant depending on the constants  $K_1, K_2, K_3$  in conditions A1, A2, A3.

It is convenient to introduce the notations  $z_k = (x_k, y_k)$  to denote the k-th example and define

$$\mathbf{F}_{i}(\boldsymbol{\theta}; \boldsymbol{z}_{k}) = (y_{k} - \hat{y}(\boldsymbol{x}_{k}; \boldsymbol{\theta})) \nabla_{\boldsymbol{\theta}_{i}} \sigma_{*}(\boldsymbol{x}_{k}; \boldsymbol{\theta}_{i}), \qquad \boldsymbol{\theta} = (\boldsymbol{\theta}_{i})_{i < N} \in \mathbb{R}^{D \times N},$$
(7.8)

$$G(\boldsymbol{\theta}; \rho) = -\nabla \Psi(\boldsymbol{\theta}; \rho) = -\nabla V(\boldsymbol{\theta}) - \int \nabla_{\boldsymbol{\theta}} U(\boldsymbol{\theta}, \boldsymbol{\theta}') \, \rho(\mathrm{d}\boldsymbol{\theta}') \,, \quad \boldsymbol{\theta} \in \mathbb{R}^{D}. \tag{7.9}$$

Note that the assumption of bounded Lipschitz  $\nabla V$ ,  $\nabla_1 U$  (here and below  $\nabla_1 U(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)$  denotes the gradient of U with respect to its first argument) implies  $\|\boldsymbol{G}(\boldsymbol{\theta}; \rho)\|_2 \leq K$  and  $\|\boldsymbol{G}(\boldsymbol{\theta}_1; \rho) - \boldsymbol{G}(\boldsymbol{\theta}_2; \rho)\|_2 \leq K \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_2$ . Further

$$\|\boldsymbol{G}(\boldsymbol{\theta}; \rho_1) - \boldsymbol{G}(\boldsymbol{\theta}; \rho_2)\|_2 = \left\| \int \nabla_{\boldsymbol{\theta}} U(\boldsymbol{\theta}; \boldsymbol{\theta}') (\rho_1 - \rho_2) (\mathrm{d}\boldsymbol{\theta}') \right\|_2 \le K d_{\mathrm{BL}}(\rho_1, \rho_2).$$
 (7.10)

With these notations, we can rewrite the SGD dynamics [3] in the main text as

$$\boldsymbol{\theta}_{i}^{k+1} = \boldsymbol{\theta}_{i}^{k} + 2\varepsilon \, \xi(k\varepsilon) \, \boldsymbol{F}_{i}(\boldsymbol{\theta}_{i}^{k}; \boldsymbol{z}_{k+1}) \,, \tag{7.11}$$

which yields

$$\boldsymbol{\theta}_{i}^{k} = \boldsymbol{\theta}_{i}^{0} + 2\varepsilon \sum_{\ell=0}^{k-1} \xi(\ell\varepsilon) \, \boldsymbol{F}_{i}(\boldsymbol{\theta}_{i}^{\ell}; \boldsymbol{z}_{\ell+1}) \,. \tag{7.12}$$

Recall  $(\boldsymbol{\theta}_i^0)_{i \le N} \sim \rho_0$  independently.

For  $t \in \mathbb{R}_{\geq 0}$  we will define  $[t] = \varepsilon \lfloor t/\varepsilon \rfloor$ . Eq. (7.12) should be compared with the nonlinear dynamics (7.4), which reads

$$\overline{\boldsymbol{\theta}}_{i}^{t} = \boldsymbol{\theta}_{i}^{0} + 2 \int_{0}^{t} \xi(s) \, \boldsymbol{G}(\overline{\boldsymbol{\theta}}_{i}^{s}; \rho_{s}) \, \mathrm{d}s \,. \tag{7.13}$$

We next state and prove the key estimate controlling the difference between the original dynamics and the nonlinear dynamics.

**Lemma 7.2.** Under the assumptions of Theorem 3, there exists a constant K depending uniquely on  $K_1, K_2, K_3$  in conditions A1, A2, and A3, such that for any  $T \ge 0$ , we have

$$\max_{i \le N} \sup_{k \in [0, T/\varepsilon] \cap \mathbb{N}} \left\| \boldsymbol{\theta}_i^k - \overline{\boldsymbol{\theta}}_i^{k\varepsilon} \right\|_2 \le K e^{KT} \cdot \sqrt{1/N \vee \varepsilon} \cdot \left[ \sqrt{D + \log(N(T/\varepsilon \vee 1))} + z \right]$$
(7.14)

with probability at least  $1 - e^{-z^2}$ .

*Proof.* Consider for simplicity of notation  $t \in \mathbb{N}\varepsilon \cap [0,T]$ . Taking the difference of Eqs. (7.12) and

(7.13), we get

$$\|\boldsymbol{\theta}_{i}^{t/\varepsilon} - \overline{\boldsymbol{\theta}}_{i}^{t}\|_{2} = 2 \| \int_{0}^{t} \xi(s) \, \boldsymbol{G}(\overline{\boldsymbol{\theta}}_{i}^{s}; \rho_{s}) \, \mathrm{d}s - \varepsilon \sum_{k=0}^{t/\varepsilon - 1} \xi(k\varepsilon) \, \boldsymbol{F}_{i}(\boldsymbol{\theta}^{k}; \boldsymbol{z}_{k+1}) \, \mathrm{d}s \|_{2}$$

$$\leq 2 \int_{0}^{t} \| \xi(s) \, \boldsymbol{G}(\overline{\boldsymbol{\theta}}_{i}^{s}; \rho_{s}) - \xi([s]) \, \boldsymbol{G}(\overline{\boldsymbol{\theta}}_{i}^{[s]}; \rho_{[s]}) \|_{2} \, \mathrm{d}s$$

$$+ 2 \int_{0}^{t} \| \xi([s]) \, \boldsymbol{G}(\overline{\boldsymbol{\theta}}_{i}^{[s]}; \rho_{[s]}) - \xi([s]) \, \boldsymbol{G}(\boldsymbol{\theta}_{i}^{\lfloor s/\varepsilon \rfloor}; \rho_{[s]}) \|_{2} \, \mathrm{d}s$$

$$+ 2 \| \varepsilon \sum_{k=0}^{t/\varepsilon - 1} \xi(k\varepsilon) \left\{ \boldsymbol{F}_{i}(\boldsymbol{\theta}^{k}; \boldsymbol{z}_{k+1}) - \boldsymbol{G}(\boldsymbol{\theta}_{i}^{k}; \rho_{k\varepsilon}) \right\} \|_{2}$$

$$\equiv 2E_{1}^{i}(t) + 2E_{2}^{i}(t) + 2E_{3}^{i}(t) .$$

$$(7.15)$$

We next consider the three terms above. Using the Lipschitz continuity of  $G(\theta; \rho)$  with respect to  $\theta$  and  $\rho$  (see Eq. (7.10)), and due to condition A1 and Lemma 7.1 (implying that  $\xi$ ,  $\overline{\theta}_i^t$ , and  $\rho_s$  are Lipschitz continuous), we get

$$E_{1}^{i}(t) \leq t \sup_{s \in [0,t]} \left\{ \left\| \xi(s) \, \boldsymbol{G}(\overline{\boldsymbol{\theta}}_{i}^{s}; \rho_{s}) - \xi([s]) \, \boldsymbol{G}(\overline{\boldsymbol{\theta}}_{i}^{s}; \rho_{s}) \right\|_{2} + \left\| \xi([s]) \, \boldsymbol{G}(\overline{\boldsymbol{\theta}}_{i}^{s}; \rho_{s}) - \xi([s]) \, \boldsymbol{G}(\overline{\boldsymbol{\theta}}_{i}^{[s]}; \rho_{s}) \right\|_{2} + \left\| \xi([s]) \, \boldsymbol{G}(\overline{\boldsymbol{\theta}}_{i}^{[s]}; \rho_{s}) - \xi([s]) \, \boldsymbol{G}(\overline{\boldsymbol{\theta}}_{i}^{[s]}; \rho_{s}) \right\|_{2} \right\}$$

$$\leq K \, t \, \varepsilon \, . \tag{7.16}$$

Bounding the second term yields (by using the Lipschitz continuity of G with respect to its first argument):

$$E_2^i(t) \le K \int_0^t \left\| \boldsymbol{G}(\overline{\boldsymbol{\theta}}_i^{[s]}; \rho_{[s]}) - \boldsymbol{G}(\boldsymbol{\theta}_i^{\lfloor s/\varepsilon \rfloor}; \rho_{[s]}) \right\|_2 \mathrm{d}s \le K^2 \int_0^t \left\| \overline{\boldsymbol{\theta}}_i^{[s]} - \boldsymbol{\theta}_i^{\lfloor s/\varepsilon \rfloor} \right\|_2 \mathrm{d}s. \tag{7.17}$$

In order to bound the last term we denote by  $\mathcal{F}_k$ , for  $k \in \mathbb{N}$ , the sigma-algebra generated by  $(\boldsymbol{\theta}_i^0)_{i \leq N}$  and  $\boldsymbol{z}_1, \dots, \boldsymbol{z}_k$ . Note that

$$\mathbb{E}\left\{\boldsymbol{F}_{i}(\boldsymbol{\theta}^{k};\boldsymbol{z}_{k+1})\big|\mathcal{F}_{k}\right\} = -\nabla V(\boldsymbol{\theta}_{i}^{k}) - \frac{1}{N}\sum_{j=1}^{N}\nabla_{1}U(\boldsymbol{\theta}_{i}^{k},\boldsymbol{\theta}_{j}^{k}) = \boldsymbol{G}(\boldsymbol{\theta}_{i}^{k};\hat{\rho}_{k}^{(N)}),$$
(7.18)

where  $\hat{\rho}_k^{(N)} \equiv (1/N) \sum_{i \leq N} \delta_{\theta_i^k}$ . Hence

$$E_{3}^{i}(t) \leq \left\| \varepsilon \sum_{k=0}^{t/\varepsilon - 1} \xi(k\varepsilon) \left\{ \boldsymbol{G}(\boldsymbol{\theta}_{i}^{k}; \hat{\rho}_{k}^{(N)}) - \boldsymbol{G}(\boldsymbol{\theta}_{i}^{k}; \rho_{k\varepsilon}) \right\} \right\|_{2} + \left\| \varepsilon \sum_{k=0}^{t/\varepsilon - 1} \xi(k\varepsilon) \boldsymbol{Z}_{k}^{i} \right\|_{2}$$

$$\equiv E_{3,0}^{i}(t) + Q_{1}^{i}(t),$$

$$(7.19)$$

where we introduced the martingale differences  $\mathbf{Z}_k^i \equiv \mathbf{F}_i(\boldsymbol{\theta}^k; \mathbf{z}_{k+1}) - \mathbb{E}\{\mathbf{F}_i(\boldsymbol{\theta}^k; \mathbf{z}_{k+1}) | \mathcal{F}_k\}$ . We can apply Azuma-Hoeffding inequality, cf. Lemma A.1. Indeed, condition (A.1) follows from the fact that  $\sigma_*(\mathbf{x}; \boldsymbol{\theta})$  is bounded and  $\nabla_{\boldsymbol{\theta}} \sigma_*(\mathbf{x}; \boldsymbol{\theta})$  is sub-Gaussian (the product of a sub-Gaussian random

vector and a bounded random variable is sub-Gaussian, cf. for instance Lemma 1.(d) in [MBM16]), hence each  $\xi(k\varepsilon)\mathbf{Z}_k^i$  are  $K^2$ -sub-Gaussian. We therefore get

$$\mathbb{P}\Big(\max_{k\in[0,t/\varepsilon]\cap\mathbb{N}}Q_1^i(k\varepsilon)\geq K\sqrt{t\varepsilon}(\sqrt{D}+u)\Big)\leq e^{-u^2},\qquad(7.21)$$

and taking union bound over  $i \leq N$ , we get

$$\mathbb{P}\Big(\max_{i\leq N}\max_{k\in[0,t/\varepsilon]\cap\mathbb{N}}Q_1^i(k\varepsilon)\leq K\sqrt{t\varepsilon}\left(\sqrt{D+\log N}+z\right)\Big)\geq 1-e^{-z^2}\ . \tag{7.22}$$

For the term  $E_{3,0}^{i}(t)$ , we use the Lipschitz continuity property (7.10), whence

$$\begin{aligned} & \left\| \boldsymbol{G}(\boldsymbol{\theta}_{i}^{k}; \hat{\rho}_{k}^{(N)}) - \boldsymbol{G}(\boldsymbol{\theta}_{i}^{k}; \rho_{k\varepsilon}) \right\|_{2} \\ \leq & \left\| \frac{1}{N} \sum_{j=1}^{N} \left[ \nabla_{1} U(\boldsymbol{\theta}_{i}^{k}, \boldsymbol{\theta}_{j}^{k}) - \nabla_{1} U(\boldsymbol{\theta}_{i}^{k}, \overline{\boldsymbol{\theta}}_{j}^{k\varepsilon}) \right] \right\|_{2} + \left\| \frac{1}{N} \sum_{j=1}^{N} \left[ \nabla_{1} U(\boldsymbol{\theta}_{i}^{k}, \overline{\boldsymbol{\theta}}_{j}^{k\varepsilon}) - \mathbb{E}_{\overline{\boldsymbol{\theta}}} \nabla_{1} U(\boldsymbol{\theta}_{i}^{k}, \overline{\boldsymbol{\theta}}_{j}^{k\varepsilon}) \right] \right\|_{2} \\ \leq & \frac{K}{N} \sum_{j=1}^{N} \left\| \boldsymbol{\theta}_{j}^{k} - \overline{\boldsymbol{\theta}}_{j}^{k\varepsilon} \right\|_{2} + Q_{2}^{i}(k\varepsilon) + \frac{K}{N}. \end{aligned}$$
(7.23)

Here  $Q_2^i(k\varepsilon)$  for  $k \in \mathbb{N}$  is defined as

$$Q_2^i(k\varepsilon) = \left\| \frac{1}{N} \sum_{j \le N, j \ne i} \left[ \nabla_1 U(\boldsymbol{\theta}_i^k, \overline{\boldsymbol{\theta}}_j^{k\varepsilon}) - \mathbb{E}_{\overline{\boldsymbol{\theta}}} \nabla_1 U(\boldsymbol{\theta}_i^k, \overline{\boldsymbol{\theta}}_j^{k\varepsilon}) \right] \right\|_2.$$

Since for any fixed k,  $(\overline{\boldsymbol{\theta}}_{j}^{k\varepsilon})_{j\leq N, j\neq i}$  are i.i.d. and independent of  $\boldsymbol{\theta}_{i}^{k}$ , and  $\nabla_{1}U$  is bounded, we get by another application of Azuma-Hoeffding inequality, cf. Lemma A.1,

$$\mathbb{P}\left(Q_2^i(k\varepsilon) \ge K\sqrt{1/N}(\sqrt{D} + u)\right) \le e^{-u^2}. \tag{7.24}$$

Therefore, the union bound for  $k \in [0, t/\varepsilon] \cap \mathbb{N}$ , and  $i \leq N$  gives

$$\mathbb{P}\left(\max_{i < N} \max_{k \in [0, t/\varepsilon] \cap \mathbb{N}} Q_2^i(k\varepsilon) \le K\sqrt{1/N} \cdot \left(\sqrt{D + \log(N(t/\varepsilon \vee 1))} + z\right)\right) \ge 1 - e^{-z^2}. \tag{7.25}$$

Conditional on the good events in Eq. (7.22) and (7.25), Eq. (7.20) thus yields

$$E_3^i(t) \le \frac{K}{N} \sum_{j=1}^N \int_0^t \|\boldsymbol{\theta}_j^{\lfloor s/\varepsilon \rfloor} - \overline{\boldsymbol{\theta}}_j^{[s]}\|_2 \, \mathrm{d}s + Q(t) + \frac{Kt}{N},$$
 (7.26)

where

$$\begin{split} Q(t) &\equiv \max_{i \leq N} Q_1^i(t) + t \cdot \max_{i \leq N} \max_{k \in [0, t/\varepsilon] \cap \mathbb{N}} Q_2^i(k\varepsilon) \\ &\leq K \sqrt{t\varepsilon} \left( z + \sqrt{D + \log N} \right) + tK \sqrt{1/N} \left( \sqrt{D + \log(N(t/\varepsilon \vee 1))} + z \right) \\ &\leq K (\sqrt{t} \vee t) \cdot \sqrt{1/N \vee \varepsilon} \cdot \left[ \sqrt{D + \log(N(t/\varepsilon \vee 1))} + z \right]. \end{split} \tag{7.27}$$

with probability at least  $1 - e^{-z^2}$ .

We finally define the random variable

$$\Delta(t; N, \varepsilon) \equiv \max_{i \le N} \sup_{k \in [0, t/\varepsilon] \cap \mathbb{N}} \|\boldsymbol{\theta}_i^k - \overline{\boldsymbol{\theta}}_i^{k\varepsilon}\|_2.$$
 (7.28)

Using the bounds (7.16), (7.17), (7.26) in Eq. (7.15), we get

$$\Delta(t; N, \varepsilon) \le K \int_0^t \Delta(s; N, \varepsilon) ds + K t \varepsilon + \frac{Kt}{N} + Q(t).$$
 (7.29)

By Gronwall's inequality, we have

$$\Delta(t; N, \varepsilon) \le K e^{Kt} \left\{ \varepsilon + \frac{1}{N} + Q(t) \right\}. \tag{7.30}$$

Using the bound (7.27), the claim follows.

**Lemma 7.3.** Under the assumptions of Theorem 3, we have

$$\max_{k \in [0, T/\varepsilon] \cap \mathbb{N}} \left| R_N(\overline{\boldsymbol{\theta}}^{k\varepsilon}) - R_N(\boldsymbol{\theta}^k) \right| \le K \cdot \max_{i \le N} \max_{k \in [0, T/\varepsilon] \cap \mathbb{N}} \left\| \boldsymbol{\theta}_i^k - \overline{\boldsymbol{\theta}}_i^{k\varepsilon} \right\|_2. \tag{7.31}$$

*Proof.* Let  $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_i, \dots, \boldsymbol{\theta}_n)$  and  $\boldsymbol{\theta}' = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_i', \dots, \boldsymbol{\theta}_n)$  be two configurations that differ only in position i. Then

$$\begin{aligned}
& \left| R_{N}(\boldsymbol{\theta}) - R_{N}(\boldsymbol{\theta}') \right| \\
& \leq \frac{1}{N} |V(\boldsymbol{\theta}_{i}) - V(\boldsymbol{\theta}'_{i})| + \frac{1}{N^{2}} |U(\boldsymbol{\theta}_{i}, \boldsymbol{\theta}_{i}) - U(\boldsymbol{\theta}'_{i}, \boldsymbol{\theta}'_{i})| + \frac{2}{N^{2}} \sum_{j \leq N, j \neq i} |U(\boldsymbol{\theta}_{i}, \boldsymbol{\theta}_{j}) - U(\boldsymbol{\theta}'_{i}, \boldsymbol{\theta}_{j})| \\
& \leq \frac{K}{N} (\|\boldsymbol{\theta}_{i} - \boldsymbol{\theta}'_{i}\|_{2} \wedge 1).
\end{aligned} (7.32)$$

Then, Eq. (7.31) follows immediately.

Lemma 7.4. Under the assumptions of Theorem 3, we have,

$$\max_{k \in [0, T/\varepsilon] \cap \mathbb{N}} \left| R_N(\overline{\boldsymbol{\theta}}^{k\varepsilon}) - R(\rho_{k\varepsilon}) \right| \le K\sqrt{1/N} \cdot \left( \sqrt{D + \log(N(T/\varepsilon \vee 1))} + z \right)$$
 (7.33)

with probability at least  $1 - e^{-z^2}$ .

*Proof.* By Eq. (7.32) and by Azuma-Höeffding inequality and union bound, we get

$$\max_{k \in [0, T/\varepsilon] \cap \mathbb{N}} \left| R_N(\overline{\boldsymbol{\theta}}^{k\varepsilon}) - \mathbb{E}R_N(\overline{\boldsymbol{\theta}}^{k\varepsilon}) \right| \le K\sqrt{1/N} \cdot \left( \sqrt{D + \log(N(T/\varepsilon \vee 1))} + z \right)$$
 (7.34)

with probability at least  $1 - e^{-z^2}$ . The claim follows since

$$\left| \mathbb{E} R_N(\overline{\boldsymbol{\theta}}^t) - R(\rho_t) \right| = \frac{1}{N} \left| \int U(\boldsymbol{\theta}, \boldsymbol{\theta}) \, \rho_t(\mathrm{d}\boldsymbol{\theta}) - \int U(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) \, \rho_t(\mathrm{d}\boldsymbol{\theta}_1) \, \rho_t(\mathrm{d}\boldsymbol{\theta}_2) \right| \le \frac{K}{N}. \tag{7.35}$$

The proof of the theorem follows from a straightforward application of Lemma 7.2, 7.3, 7.4. The proof for any bounded Lipschitz function f follows the same argument as Lemma 7.3, 7.4. As a result, for any sequence  $(N, \varepsilon = \varepsilon_N)$  such that  $N \to \infty$  and  $\varepsilon_N \to 0$  with  $N/\log(N/\varepsilon_N) \to \infty$  and  $\varepsilon_N \log(N/\varepsilon_N) \to 0$ , we have  $\hat{\rho}_{\lfloor k/\varepsilon \rfloor}^{(N)}$  converges weakly to  $\rho_t$  almost surely immediately.

#### 7.2 Proof of Theorem 3: Generalization to $\beta < \infty$

Here we generalize the proof given in the previous section to noisy SGD at finite temperature  $\beta < \infty$ . Since the proof follows the same scheme as in the noiseless case, we will limit ourselves to describing the differences.

Throughout this section we assume that conditions A1, A2, A3 hold. We also let

$$\Psi_{\lambda}(\boldsymbol{\theta}; \rho) = \frac{\lambda}{2} \|\boldsymbol{\theta}\|_{2}^{2} + V(\boldsymbol{\theta}) + \int U(\boldsymbol{\theta}, \boldsymbol{\theta}') \rho(\boldsymbol{\theta}') d\boldsymbol{\theta}'$$
(7.36)

for some  $\lambda \leq 1$ . Further we assume  $\rho_0$  is  $K_0^2$ -sub-Gaussian. Finally, we assume  $1 \leq \beta < \infty$ . For the reader's convenience, we reproduce here the form of the limiting PDE

$$\partial_t \rho_t(\boldsymbol{\theta}) = 2\xi(t) \nabla_{\boldsymbol{\theta}} \cdot \left[ \rho_t(\boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}} \Psi_{\lambda}(\boldsymbol{\theta}; \rho_t) \right] + 2\xi(t) / \beta \cdot \Delta_{\boldsymbol{\theta}} \rho_t(\boldsymbol{\theta}), \qquad (7.37)$$

which again should be interpreted in weak sense.

**Remark 7.2.** Recall conditionss A1, A2, A3 in the main text. By a modified argument of [Szn91, Theorem 1.1], conditions A1 and A3 are sufficient for the existence and uniqueness of solution of PDE (7.37) in weak sense. Section 10 provides further information of this PDE, including a proof of existence and uniqueness.

As in the noiseless case, there is an equivalent formulation of this PDE as a fixed point distribution for the following nonlinear dynamics, which is an integration form of a stochastic differential equation,

$$\overline{\boldsymbol{\theta}}_{i}^{t} = \boldsymbol{\theta}_{i}^{0} + 2 \int_{0}^{t} \xi(s) \, \boldsymbol{G}(\overline{\boldsymbol{\theta}}_{i}^{s}; \rho_{s}) \, \mathrm{d}s + \int_{0}^{t} \sqrt{2\xi(s)/\beta} \, \mathrm{d}\boldsymbol{W}_{i}(s) \,, \tag{7.38}$$

$$\rho_s = P_{\overline{\theta}_i^s}, \tag{7.39}$$

where  $\{W_i(s)\}_{s\geq 0}$  for  $i\leq N$  are independent *D*-dimensional Brownian motions, and  $G(\theta;\rho)\equiv -\nabla \Psi_{\lambda}(\theta;\rho)$ . The assumptions on  $U, V, \lambda$ , and  $\xi$  ensures that this nonlinear dynamics has a unique continuous solution.

This nonlinear dynamics should be compared with the noisy SGD dynamics [11] in the main text that can be written as follows for  $k \in \mathbb{N}$ :

$$\boldsymbol{\theta}_{i}^{k} = \boldsymbol{\theta}_{i}^{0} + 2\varepsilon \sum_{\ell=0}^{k-1} \xi(\ell\varepsilon) \, \boldsymbol{F}_{i}(\boldsymbol{\theta}^{\ell}; \boldsymbol{z}_{\ell}) + \int_{0}^{k\varepsilon} \sqrt{2\xi([s])/\beta} \, \mathrm{d}\boldsymbol{W}_{i}(s) \,, \tag{7.40}$$

where

$$F_i(\boldsymbol{\theta}; \boldsymbol{z}_k) = -\lambda \boldsymbol{\theta}_i + (y_k - \hat{y}(\boldsymbol{x}_k; \boldsymbol{\theta})) \nabla_{\boldsymbol{\theta}_i} \sigma_*(\boldsymbol{x}_k; \boldsymbol{\theta}_i), \quad \boldsymbol{\theta} = (\boldsymbol{\theta}_i)_{i \le N} \in \mathbb{R}^{D \times N}.$$
 (7.41)

It is convenient to collect some standard estimates about the solution of the stochastic differential equation (7.38).

**Lemma 7.5.** Assume  $\rho_0$  is  $K_0^2$ -sub-Gaussian,  $\xi(s)$  and  $G(\mathbf{0}; \rho_s)$  are  $K_0$ -bounded,  $G(\boldsymbol{\theta}; \rho_s)$  is  $K_0$ -Lipschitz in  $\boldsymbol{\theta}$ , and  $\beta \geq 1$ . Let  $(\overline{\boldsymbol{\theta}}_i^t)_{t\geq 0}$  for  $i \leq N$  be the solution of (7.38) with independent

initialization  $(\theta_i^0)_{i\leq N} \sim \rho_0$ . Let  $(\rho_t)_{t\geq 0}$  be the solution of PDE (7.37). Then there exists a constant K depending uniquely on  $K_0$ , such that

$$\mathbb{P}\left(\sup_{i \le N} \sup_{t \in [0,T]} \|\overline{\theta}_i^t\|_2 \le Ke^{KT} [\sqrt{D + \log N} + z]\right) \ge 1 - e^{-z^2},\tag{7.42}$$

and

$$\mathbb{P}\left(\sup_{i\leq N}\sup_{k\in[0,T/\varepsilon]\cap\mathbb{N}}\sup_{u\in[0,\varepsilon]}\|\overline{\boldsymbol{\theta}}_{i}^{k\varepsilon+u}-\overline{\boldsymbol{\theta}}_{i}^{k\varepsilon}\|_{2}\leq Ke^{KT}\left[\sqrt{D+\log(N(T/\varepsilon\vee1))}+z\right]\sqrt{\varepsilon}\right)\geq 1-e^{-z^{2}}, (7.43)$$

and for any  $t, h \ge 0, t + h \le T$ ,

$$d_{\text{BL}}(\rho_t, \rho_{t+h}) \le W_2(\rho_t, \rho_{t+h}) \le Ke^{KT} \sqrt{Dh}.$$
 (7.44)

*Proof.* We decompose the proof into three parts.

Part (a). First, note that for any D-dimensional  $K_0^2$ -sub-Gaussian random vector X, we have

$$\mathbb{E}_{\mathbf{X}}[\exp\{\tau \|\mathbf{X}\|_{2}^{2}/2\}] = \mathbb{E}_{\mathbf{X},\mathbf{G}}[\exp\{\tau \langle \mathbf{G}, \mathbf{X} \rangle\}] \le \mathbb{E}_{\mathbf{G}}[\exp\{\tau K_{0}^{2} \|\mathbf{G}\|_{2}^{2}\}/2] = (1 - \tau K_{0}^{2})^{-D/2}.$$
(7.45)

Note that  $(\theta_i^0)_{i\leq N}\sim \rho_0$  independently, and  $\rho_0$  is  $K_0^2$ -sub-Gaussian. Therefore

$$\mathbb{P}(\|\boldsymbol{\theta}_i^0\|_2 \ge u) \le \mathbb{E}[\exp(\tau\|\boldsymbol{\theta}_i\|_2^2/2)]/\exp\{\tau z^2/2\} \le (1 - \tau K_0^2)^{-D/2}\exp\{-\tau u^2/2\}.$$

Taking union bound over  $i \leq N$  gives

$$\mathbb{P}\left(\max_{i < N} \|\boldsymbol{\theta}_i^0\|_2 \ge u\right) \le (1 - \tau K_0^2)^{-D/2} \exp\{-\tau u^2/2 + \log N\}.$$

Taking  $\tau = 1/(2K_0^2)$  and  $u = 2K_0(\sqrt{D + \log N} + z)$ , we get

$$\mathbb{P}\Big(\max_{i \le N} \|\boldsymbol{\theta}_i^0\|_2 \ge 2K_0(\sqrt{D + \log N} + z)\Big) \le \exp\{-z^2\}.$$
 (7.46)

Then we define  $\mathbf{W}_{\xi,i}(t) \equiv \int_0^t \sqrt{2\xi(s)} \, \mathrm{d}\mathbf{W}_i(s)$ . We have  $\mathrm{Var}(W^j_{\xi,i}(t)) = \int_0^t 2\xi(s) \, \mathrm{d}s \leq 2K_0 t$  for  $j \leq D$ . Note  $\exp\{\tau \|\mathbf{W}_{\xi,i}(t)\|_2^2\}$  is a submartingale, due to Doob's martingale inequality, we have

$$\mathbb{P}\left(\sup_{t < T} \|\boldsymbol{W}_{\xi,i}(t)\|_{2} \ge u\right) \le \mathbb{E}[\exp\{\tau\|\boldsymbol{W}_{\xi,i}(T)\|_{2}^{2}/2\}] \cdot \exp\{-\tau u^{2}/2\} \le (1 - 2K_{0}T\tau)^{-D/2}\exp\{-\tau u^{2}/2\}.$$

Taking union bound over  $i \leq N$  gives

$$\mathbb{P}\Big(\max_{i < N} \sup_{t < T} \|\boldsymbol{W}_{\xi,i}(t)\|_2 \ge u\Big) \le (1 - 2K_0T\tau)^{-D/2} \exp\{-\tau u^2/2 + \log N\}.$$

Taking  $\tau = 1/(4K_0T)$  and  $u = 4\sqrt{K_0T}(\sqrt{D + \log N} + z)$ , we get

$$\mathbb{P}\Big(\max_{i \le N} \sup_{t \le T} \|\mathbf{W}_{\xi,i}(t)\|_2 \ge 4\sqrt{K_0 T} (\sqrt{D + \log N} + z)\Big) \le \exp\{-z^2\}.$$
 (7.47)

By noting that  $\xi(s)$ ,  $G(0; \rho_s)$  are  $K_0$ -bounded, and  $G(\theta; \rho_s)$  is  $K_0$ -Lipschitz in  $\theta$ , according to Eq. (7.38), there exists some constant K depending on  $K_0$ , such that

$$\Delta_i(t) \le K \int_0^t \Delta_i(s) ds + K[W/\sqrt{\beta} + \Theta],$$

where  $\Delta_i(t) \equiv \sup_{s \leq t} \|\overline{\boldsymbol{\theta}}_i^s\|_2$ ,  $W \equiv \max_{i \leq N} \sup_{t \leq T} \|\boldsymbol{W}_{\xi,i}(t)\|_2$ , and  $\Theta \equiv \max_{i \leq N} \|\boldsymbol{\theta}_i^0\|_2$ . Due to Gronwall's inequality, we have

$$\Delta_i(T) \le K \exp(KT)[W/\sqrt{\beta} + \Theta].$$

The high probability bound (7.42) holds by noting the high probability bound for  $\Theta$  and W in Eq. (7.46) and (7.47).

Part (b). Define  $\Delta_i(h; k, \varepsilon) = \sup_{0 \le u \le h} \|\overline{\boldsymbol{\theta}}_i^{k\varepsilon+u} - \overline{\boldsymbol{\theta}}_i^{k\varepsilon}\|_2$ . By noting that  $\xi(s)$ ,  $\boldsymbol{G}(\mathbf{0}; \rho_s)$  are  $K_0$ -bounded, and  $\boldsymbol{G}(\boldsymbol{\theta}; \rho_s)$  is  $K_0$ -Lipschitz in  $\boldsymbol{\theta}$ , according to Eq. (7.38), we have

$$\Delta_{i}(h; k, \varepsilon) \leq K \left[ \sup_{s < T} \|\overline{\boldsymbol{\theta}}_{i}^{s}\|_{2} + 1 \right] h + \frac{1}{\sqrt{\beta}} \sup_{0 < u < h} \|\boldsymbol{W}_{\xi, i, k}(u)\|_{2}, \tag{7.48}$$

where  $W_{\xi,i,k}(u) \equiv \int_{k\varepsilon}^{k\varepsilon+u} \sqrt{2\xi(s)} \,dW_i(s)$ . Similar to the bound Eq. (7.47), we have

$$\mathbb{P}\Big(\max_{i\leq N}\sup_{k\in[0,T/\varepsilon]\cap\mathbb{N}}\sup_{0\leq u\leq h}\|\boldsymbol{W}_{\xi,i,k}(u)\|_{2}\leq 4\sqrt{K_{0}h}\Big(\sqrt{D+\log(N(T/\varepsilon\vee1))}+z\Big)\Big)\geq 1-e^{-z^{2}}.$$
(7.49)

Plugging the bound Eq. (7.42) and Eq. (7.49) into Eq. (7.48), we have

$$\max_{i \le N} \sup_{k \in [0, T/\varepsilon] \cap \mathbb{N}} \Delta_i(h; k, \varepsilon) \le K e^{KT} \left[ \sqrt{D + \log N} + z \right] h + K \left( \sqrt{D + \log(N(T/\varepsilon \vee 1))} + z \right) \sqrt{h}$$

$$\le K e^{KT} \left[ \sqrt{D + \log(N(T/\varepsilon \vee 1))} + z \right] \sqrt{h}$$

with probability at least  $1 - e^{-z^2}$ .

Part (c). Equation (7.44) holds directly by noting that

$$W_2(\rho_t, \rho_{t+h})^2 \le \mathbb{E}\{\|\overline{\boldsymbol{\theta}}^t - \overline{\boldsymbol{\theta}}^{t+h}\|_2^2\}$$

and applying a integration over z in a modified version of Eq. (7.43) without union bound over  $i \leq N$  and  $k \in [0, T/\varepsilon] \cap \mathbb{N}$ .

As in the noiseless case, the key step consists in bounding the difference between the nonlinear dynamics and the SGD dynamics.

**Lemma 7.6.** Under the assumptions of Theorem 3, there exists a constant K depending uniquely on  $K_0, K_1, K_2, K_3$ , such that for any  $T \ge 0$ , we have

$$\max_{i \le N} \sup_{k \in [0, T/\varepsilon] \cap \mathbb{N}} \|\boldsymbol{\theta}_i^k - \overline{\boldsymbol{\theta}}_i^{k\varepsilon}\|_2 \le Ke^{KT} \cdot \sqrt{1/N \vee \varepsilon} \cdot \left[\sqrt{D + \log(N(T/\varepsilon \vee 1))} + z\right]$$
(7.50)

with probability at least  $1 - e^{-z^2}$ .

*Proof.* We take the difference of Eqs. (7.40) and (7.38), for  $t \in \mathbb{N}\varepsilon \cap [0,T]$ :

$$\|\boldsymbol{\theta}_{i}^{t/\varepsilon} - \overline{\boldsymbol{\theta}}_{i}^{t}\|_{2} \leq 2 \|\int_{0}^{t} \left[\xi(s) \boldsymbol{G}(\overline{\boldsymbol{\theta}}_{i}^{s}; \rho_{s}) - \xi([s]) \boldsymbol{G}(\overline{\boldsymbol{\theta}}_{i}^{[s]}; \rho_{[s]})\right] ds \|_{2}$$

$$+ 2 \int_{0}^{t} \|\xi([s]) \boldsymbol{G}(\overline{\boldsymbol{\theta}}_{i}^{[s]}; \rho_{[s]}) - \xi([s]) \boldsymbol{G}(\boldsymbol{\theta}_{i}^{\lfloor s/\varepsilon \rfloor}; \rho_{[s]}) \|_{2} ds$$

$$+ 2 \|\varepsilon \sum_{k=0}^{t/\varepsilon - 1} \xi(k\varepsilon) \left\{ \boldsymbol{F}_{i}(\boldsymbol{\theta}^{k}; \boldsymbol{z}_{k+1}) - \boldsymbol{G}(\boldsymbol{\theta}_{i}^{k}; \rho_{k\varepsilon}) \right\} \|_{2}$$

$$+ \|\int_{0}^{t} \left(\sqrt{2\xi(s)/\beta} - \sqrt{2\xi([s])/\beta}\right) d\boldsymbol{W}_{i}(s) \|_{2}$$

$$\equiv 2E_{1}^{i}(t) + 2E_{2}^{i}(t) + 2E_{3}^{i}(t) + E_{4}^{i}(t).$$

$$(7.51)$$

Terms  $E_2^i(t)$ ,  $E_3^i(t)$  can be bounded the same as in Lemma 7.2, i.e., Eq. (7.17) and (7.26), by noting

that the replacement of  $\Psi$  by  $\Psi_{\lambda}$  does not affect these estimates. To bound  $E_4^i(t)$ , notice that  $\mathbf{W}_{\xi,i} \equiv \int_0^T \left(\sqrt{2\xi(s)} - \sqrt{2\xi([s])}\right) \mathrm{d}\mathbf{W}_i(s)$  is a Gaussian random vector,  $\mathbf{W}_{\xi,i} \sim \mathsf{N}(\mathbf{0}, \tau^2 \mathbf{I}_D)$ , where, using the Lipschitz continuity of  $\xi$ ,

$$\tau^2 = \int_0^T \left( \sqrt{2\xi(s)} - \sqrt{2\xi([s])} \right)^2 ds \le K T \varepsilon.$$

By Gaussian concentration

$$\mathbb{P}(\|\boldsymbol{W}_{\xi,i}\|_2 \ge (\sqrt{D} + z)\tau) \le e^{-z^2/2},$$

and therefore by applying Doob's inequality to the submartingale  $t \mapsto E_4^i(t)$ , we get

$$\mathbb{P}\Big(\max_{s < T} E_4^i(s) \ge K(\sqrt{D} + z)\sqrt{T\varepsilon}\Big) \le e^{-z^2/2},$$

and hence

$$\mathbb{P}\left(\max_{i \le N} \max_{s \le T} E_4^i(s) \le K(\sqrt{D + \log N} + z)\sqrt{T\varepsilon}\right) \ge 1 - e^{-z^2/2}.\tag{7.52}$$

We need to modify the proof of Lemma 7.2 to bound terms  $E_1^i(t)$ .

$$E_{1}^{i}(t) \leq \left\| \int_{0}^{t} \left[ \xi(s) - \xi([s]) \right] \boldsymbol{G}(\overline{\boldsymbol{\theta}}_{i}^{s}; \rho_{s}) ds \right\|_{2} + \left\| \int_{0}^{t} \xi([s]) \left[ \boldsymbol{G}(\overline{\boldsymbol{\theta}}_{i}^{s}; \rho_{s}) - \boldsymbol{G}(\overline{\boldsymbol{\theta}}_{i}^{s}; \rho_{[s]}) \right] ds \right\|_{2}$$

$$+ \left\| \int_{0}^{t} \xi([s]) \left[ \boldsymbol{G}(\overline{\boldsymbol{\theta}}_{i}^{s}; \rho_{[s]}) - \boldsymbol{G}(\overline{\boldsymbol{\theta}}_{i}^{[s]}; \rho_{[s]}) \right] ds \right\|_{2}$$

$$\equiv E_{1,A}^{i}(t) + E_{1,B}^{i}(t) + E_{1,C}^{i}(t).$$

$$(7.53)$$

To bound the first term  $E_{1,A}^i(t)$ , due to the Lipschitz property of  $G(\theta;\rho)$  and the boundedness of  $G(0; \rho)$ , with probability at least  $1 - e^{-z^2}$ , we have for all  $i \leq N$  and  $t \leq T$ ,

$$E_{1,A}^{i}(t) \leq TK\varepsilon \cdot \sup_{s \in [0,T]} \|\boldsymbol{G}(\overline{\boldsymbol{\theta}}_{i}^{s}; \rho_{s})\|_{2} \leq TK\varepsilon \cdot \left[K \sup_{s \in [0,T]} \|\overline{\boldsymbol{\theta}}_{i}^{s}\|_{2} + K\right]$$

$$\leq Ke^{KT} \left[\sqrt{D + \log N} + z\right]\varepsilon. \tag{7.54}$$

Here the last inequality is due to Eq. (7.42) in Lemma 7.5.

To bound the second term  $E_{1,B}^i(t)$ , using the fact that  $\nabla_1 U$  is bounded Lipschitz, we have for all  $i \leq N$  and  $t \leq T$ ,

$$E_{1,B}^{i}(t) \leq TK \cdot \sup_{\boldsymbol{\theta} \in \mathbb{R}^{D}} \|\nabla_{1}U(\boldsymbol{\theta}; \rho_{s}) - \nabla_{1}U(\boldsymbol{\theta}; \rho_{[s]})\|_{2} \leq TK^{2} \cdot d_{\mathrm{BL}}(\rho_{s}, \rho_{[s]}) \leq Ke^{KT}\sqrt{D\varepsilon}.$$
 (7.55)

Here the last inequality is due to Eq. (7.44) in Lemma 7.5.

To bound the third term  $E_{1,C}^i(t)$ , with probability at least  $1 - e^{-z^2}$ , we have for all  $i \leq N$  and  $t \leq T$ ,

$$E_{1,C}^{i}(t) \leq TK \cdot \sup_{s \in [0,T]} \left\| \boldsymbol{G}(\overline{\boldsymbol{\theta}}_{i}^{s}; \rho_{[s]}) - \boldsymbol{G}(\overline{\boldsymbol{\theta}}_{i}^{[s]}; \rho_{[s]}) \right\|_{2}$$

$$\leq TK^{2} \cdot \sup_{s \in [0,T]} \left\| \overline{\boldsymbol{\theta}}_{i}^{s} - \overline{\boldsymbol{\theta}}_{i}^{[s]} \right\|_{2} \leq Ke^{KT} \left[ \sqrt{D + \log(N(T/\varepsilon \vee 1))} + z \right] \sqrt{\varepsilon}.$$

$$(7.56)$$

Here the last inequality is due to Eq. (7.43) in Lemma 7.5.

As a result, combining Eq. (7.17), (7.26), (7.27), (7.51), (7.52), (7.54), (7.55), and (7.56), defining

$$\Delta(t; N, \varepsilon) \equiv \max_{i \le N} \sup_{k \in [0, T/\varepsilon] \cap \mathbb{N}} \|\boldsymbol{\theta}_i^k - \overline{\boldsymbol{\theta}}_i^{k\varepsilon}\|_2, \qquad (7.57)$$

we get

$$\Delta(t; N, \varepsilon) \le K \int_0^t \Delta(s; N, \varepsilon) ds + \frac{Kt}{N} + E(T), \tag{7.58}$$

where

$$E(T) = Ke^{KT} \cdot \sqrt{1/N \vee \varepsilon} \cdot \left[ \sqrt{D + \log(N(T/\varepsilon \vee 1))} + z \right]. \tag{7.59}$$

Applying Gronwall's inequality gives the desired result.

The generalization of Theorem 3 to  $\beta < \infty$  follows from this lemma exactly as in the previous section.

#### 7.3 Proof of Proposition 2: Monotonicity of the risk

By simple algebra, we have

$$R(\rho_{t+h}) - R(\rho_t) = 2 \int \Psi(\boldsymbol{\theta}; \rho_t) (\rho_{t+h} - \rho_t) (d\boldsymbol{\theta}) + \langle U, (\rho_{t+h} - \rho_t)^{\otimes 2} \rangle.$$
 (7.60)

By Lemma 7.1,  $t \mapsto \rho_t$  is Lipschitz continuous in Wasserstein distance  $W_2(\rho_{t_1}, \rho_{t_2}) \leq K|t_1 - t_2|$ . Hence, we get

$$R(\rho_{t+h}) - R(\rho_t) = 2 \int \Psi(\boldsymbol{\theta}; \rho_t) (\rho_{t+h} - \rho_t) (d\boldsymbol{\theta}) + O(h^2)$$
(7.61)

$$= -4\xi(t) \int \|\nabla \Psi(\boldsymbol{\theta}; \rho_t)\|_2^2 \rho_t(d\boldsymbol{\theta}) h + o(h), \qquad (7.62)$$

where in the second step we used Eq. (7.3). This immediately implies that  $R(\rho_t)$  is non-increasing in t.

Let  $\rho$  be a fixed point of Eq. (7.1). Since  $\partial_t R(\rho_t)|_{\rho_0=\rho}=0$ , the above formula implies

$$\int \|\nabla \Psi(\boldsymbol{\theta}; \rho)\|_{2}^{2} \rho(\mathrm{d}\boldsymbol{\theta}) = 0, \qquad (7.63)$$

and therefore  $\rho$  is supported in the set of  $\boldsymbol{\theta}$ 's such that  $\nabla \Psi(\boldsymbol{\theta}; \rho) = \mathbf{0}$ .

Vice versa, if this is the case, setting  $\rho_0 = \rho$ , Eq. (7.3) implies  $\partial_t \langle \varphi, \rho_t \rangle = 0$ , then  $\rho_t \equiv \rho_0$  is a fixed point.

### 7.4 A general continuity result

It is useful to notice that the solution  $(\rho_t)_{t\geq 0}$  of the PDE (7.1) is continuous with respect to changes in  $V(\cdot)$ ,  $U(\cdot,\cdot)$ . Namely, we consider the following two PDEs:

$$\partial_t \rho_t(\boldsymbol{\theta}) = 2\xi(t)\nabla \cdot \left[\rho_t(\boldsymbol{\theta})\nabla \Psi(\boldsymbol{\theta}; \rho_t)\right],\tag{7.64}$$

$$\partial_t \tilde{\rho}_t(\boldsymbol{\theta}) = 2\xi(t)\nabla \cdot \left[ \tilde{\rho}_t(\boldsymbol{\theta})\nabla \tilde{\Psi}(\boldsymbol{\theta}; \tilde{\rho}_t) \right], \tag{7.65}$$

where

$$\Psi(\boldsymbol{\theta}; \rho) = V(\boldsymbol{\theta}) + \int U(\boldsymbol{\theta}, \boldsymbol{\theta}') \, \rho(\mathrm{d}\boldsymbol{\theta}') \,, \tag{7.66}$$

$$\tilde{\Psi}(\boldsymbol{\theta}; \tilde{\rho}) = \tilde{V}(\boldsymbol{\theta}) + \int \tilde{U}(\boldsymbol{\theta}, \boldsymbol{\theta}') \, \tilde{\rho}(\mathrm{d}\boldsymbol{\theta}') \,. \tag{7.67}$$

**Lemma 7.7.** Let assumptions A1, A3 hold both for V, U and  $\tilde{V}, \tilde{U}$ , and consider the solutions of Eqs. (7.64) and (7.65) with initial conditions  $\rho_0$ ,  $\tilde{\rho}_0$ . Then there exists  $K < \infty$  depending only on the constants  $K_1$ ,  $K_3$  in the assumptions (independent of D), such that

$$\sup_{t \in [0,T]} d_{\text{BL}}(\rho_t, \tilde{\rho}_t) \le K e^{KT} \cdot \left[ d_{\text{BL}}(\rho_0, \tilde{\rho}_0) + \varepsilon_0 \right], \tag{7.68}$$

where

$$\varepsilon_0 \equiv \sup_{\boldsymbol{\theta}, \boldsymbol{\theta}' \in \mathbb{R}^D} \left[ \|\nabla V(\boldsymbol{\theta}) - \nabla \tilde{V}(\boldsymbol{\theta})\|_2 + \|\nabla_1 U(\boldsymbol{\theta}, \boldsymbol{\theta}') - \nabla_1 \tilde{U}(\boldsymbol{\theta}, \boldsymbol{\theta}')\|_2 \right]. \tag{7.69}$$

*Proof.* The proof adapts the argument used to establish uniqueness in [Szn91]. Without loss of generality, we fix  $\xi(t) \equiv 1/2$ . We further denote by K generic constants depending on  $K_1$ ,  $K_3$ .

The assumption of bounded Lipschitz  $\nabla V$  and  $\nabla U$  implies that  $\nabla \Psi(\boldsymbol{\theta}; \rho)$  is K-bounded Lipschitz with respect to argument  $(\boldsymbol{\theta}, \rho)$ , that is,

$$\left\| \nabla \Psi(\boldsymbol{\theta}_1; \rho_1) - \nabla \Psi(\boldsymbol{\theta}_2; \rho_2) \right\|_2 \le K \left[ \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_2 \wedge 1 + d_{\text{BL}}(\rho_1, \rho_2) \right]. \tag{7.70}$$

The assumption of bounded Lipschitz  $\nabla \tilde{V}$  and  $\nabla \tilde{U}$  implies that  $\nabla \tilde{\Psi}(\boldsymbol{\theta}; \rho)$  is K-bounded Lipschitz. Under these conditions, according to [Szn91, Theorem 1.1], there is existence and uniqueness of PDE (7.64) and (7.65). We denote their solutions at time t to be  $\rho_t, \tilde{\rho}_t \in \mathcal{P}(\mathbb{R}^D)$  respectively.

Let  $\gamma_0 \in \mathscr{P}(\mathbb{R}^D \times \mathbb{R}^D)$  be a coupling of  $\rho_0$  and  $\tilde{\rho}_0$  that achieves  $2d_{\mathrm{BL}}(\rho_0, \tilde{\rho}_0)$ . Given these fixed  $(\rho_t)_{t>0}$  and  $(\tilde{\rho}_t)_{t>0}$ , consider the nonlinear dynamics

$$\boldsymbol{\theta}^t = \boldsymbol{\theta}^0 - \int_0^t \nabla \Psi(\boldsymbol{\theta}^s; \rho_s) \mathrm{d}s, \tag{7.71}$$

$$\tilde{\boldsymbol{\theta}}^t = \tilde{\boldsymbol{\theta}}^0 - \int_0^t \nabla \tilde{\Psi}(\tilde{\boldsymbol{\theta}}^s; \tilde{\rho}_s) ds, \tag{7.72}$$

with initialization  $(\boldsymbol{\theta}^0, \tilde{\boldsymbol{\theta}}^0) \sim \gamma_0$ . As implied by [Szn91, Theorem 1.1], since we have  $\boldsymbol{\theta}^0 \sim \rho_0$ ,  $\tilde{\boldsymbol{\theta}}^0 \sim \tilde{\rho}_0$ , it follows that  $\boldsymbol{\theta}_t \sim \rho_t$ ,  $\tilde{\boldsymbol{\theta}}_t \sim \tilde{\rho}_t$ , and therefore

$$d_{\mathrm{BL}}(\rho_t, \tilde{\rho}_t) \le 2 \int \left( \|\boldsymbol{\theta}^t - \tilde{\boldsymbol{\theta}}^t\|_2 \wedge 1 \right) \gamma_0(\mathrm{d}\boldsymbol{\theta}^0, \mathrm{d}\tilde{\boldsymbol{\theta}}^0). \tag{7.73}$$

Taking the difference of Eqs. (7.71) and (7.72), for any  $(\boldsymbol{\theta}^0, \tilde{\boldsymbol{\theta}}^0) \in \text{supp}(\gamma_0)$ ,

$$\|\boldsymbol{\theta}^{t} - \tilde{\boldsymbol{\theta}}^{t}\|_{2} \leq \int_{0}^{t} \|\nabla \Psi(\boldsymbol{\theta}^{s}; \rho_{s}) - \nabla \tilde{\Psi}(\tilde{\boldsymbol{\theta}}^{s}; \tilde{\rho}_{s})\|_{2} ds + \|\boldsymbol{\theta}^{0} - \tilde{\boldsymbol{\theta}}^{0}\|_{2}$$

$$\leq \int_{0}^{t} \|\nabla \Psi(\boldsymbol{\theta}^{s}; \rho_{s}) - \nabla \Psi(\tilde{\boldsymbol{\theta}}^{s}; \tilde{\rho}_{s})\|_{2} ds + \int_{0}^{t} \|\nabla \Psi(\tilde{\boldsymbol{\theta}}^{s}; \tilde{\rho}_{s}) - \nabla \tilde{\Psi}(\tilde{\boldsymbol{\theta}}^{s}; \tilde{\rho}_{s})\|_{2} ds + \|\boldsymbol{\theta}^{0} - \tilde{\boldsymbol{\theta}}^{0}\|_{2}$$

$$\equiv E_{1}(t) + E_{2}(t) + \|\boldsymbol{\theta}^{0} - \tilde{\boldsymbol{\theta}}^{0}\|_{2}.$$

$$(7.74)$$

Using bound (7.70), the first term  $E_1(t)$  is simply bounded by

$$E_1(t) \le K \int_0^t \left[ \|\boldsymbol{\theta}^s - \tilde{\boldsymbol{\theta}}^s\|_2 \wedge 1 + d_{\mathrm{BL}}(\rho_s, \tilde{\rho}_s) \right] \cdot \mathrm{d}s. \tag{7.75}$$

To bound the second term  $E_2(t)$ , we have

$$E_{2}(t) \leq t \times \sup_{\boldsymbol{\theta} \in \mathbb{R}^{D}, \rho \in \mathscr{P}(\mathbb{R}^{D})} \|\nabla \Psi(\boldsymbol{\theta}; \rho) - \nabla \tilde{\Psi}(\boldsymbol{\theta}; \rho)\|_{2}$$

$$\leq t \times \sup_{\boldsymbol{\theta}, \boldsymbol{\theta}' \in \mathbb{R}^{D}} \left[ \|\nabla V(\boldsymbol{\theta}) - \nabla \tilde{V}(\boldsymbol{\theta})\|_{2} + \|\nabla_{1} U(\boldsymbol{\theta}, \boldsymbol{\theta}') - \nabla_{1} \tilde{U}(\boldsymbol{\theta}, \boldsymbol{\theta}')\|_{2} \right] = t \cdot \varepsilon_{0}, \tag{7.76}$$

with the definition of  $\varepsilon_0$  given by Equation (7.69).

Combining Equation (7.74), (7.75), and (7.76), we have

$$\|\boldsymbol{\theta}^t - \tilde{\boldsymbol{\theta}}^t\|_2 \wedge 1 \le K \int_0^t \left[ \|\boldsymbol{\theta}^s - \tilde{\boldsymbol{\theta}}^s\|_2 \wedge 1 + d_{\text{BL}}(\rho_s, \tilde{\rho}_s) \right] \cdot ds + t \cdot \varepsilon_0 + \|\boldsymbol{\theta}^0 - \tilde{\boldsymbol{\theta}}^0\|_2 \wedge 1. \tag{7.77}$$

Averaging the above inequality over  $(\theta^0, \tilde{\theta}^0) \sim \gamma_0$ , and using inequality (7.73), we have

$$\int \|\boldsymbol{\theta}^t - \tilde{\boldsymbol{\theta}}^t\|_2 \wedge 1 \cdot d\gamma_0 \le 2d_{\text{BL}}(\rho_0, \tilde{\rho}_0) + 3K \int_0^t \left[ \int \|\boldsymbol{\theta}^s - \tilde{\boldsymbol{\theta}}^s\|_2 \wedge 1 \cdot d\gamma_0 \right] \cdot ds + t \cdot \varepsilon_0.$$
 (7.78)

Using Gronwall's inequality, for any  $t \in \mathbb{R}$ , we have

$$\int \|\boldsymbol{\theta}^t - \tilde{\boldsymbol{\theta}}^t\|_2 \wedge 1 \cdot \gamma_0(\mathrm{d}\boldsymbol{\theta}^0, \mathrm{d}\tilde{\boldsymbol{\theta}}^0) \le K \exp(Kt) \cdot [d_{\mathrm{BL}}(\rho_0, \tilde{\rho}_0) + \varepsilon_0].$$

Applying Equation (7.73), the result follows.

### 7.5 Some properties of the solution of the PDE (7.1)

In this section we prove four lemmas on the properties of the solution of the PDE (7.1), under conditions A1 and A3. All of these facts are quite standard, but we provide complete proofs for them for reader's convenience.

We will use several times the following notations. Let  $\rho_t$  be a solution of the PDE (7.1) with initialization  $\rho_0$ . Let  $(\boldsymbol{\theta}^t)_{t\geq 0}$  be the solution of the ordinary differential equation (ODE)

$$\dot{\boldsymbol{\theta}}^t = -2\xi(t)\nabla\Psi(\boldsymbol{\theta}^t; \rho_t), \qquad (7.79)$$

with initial condition  $\boldsymbol{\theta}^0$ . Without loss of generality, we will assume  $\xi(t) = 1/2$  throughout this section. If  $\boldsymbol{\theta}^0 \sim \rho_0$ , then for any  $t \geq 0$ , we have  $\boldsymbol{\theta}^t \sim \rho_t$ . We will denote by  $\boldsymbol{\varphi}^t : \mathbb{R}^D \to \mathbb{R}^D$  the map between initial conditions of this ODE, and its state at time t (i.e.  $\boldsymbol{\varphi}^t(\boldsymbol{\theta}^0) = \boldsymbol{\theta}^t$ ). Since  $\nabla \Psi(\cdot; \rho_t)$  is bounded and Lipschitz continuous, it follows that  $\boldsymbol{\varphi}^t$  is a homeomorphism on its image by Picard's theorem.

With these notations,  $\rho_t$  is the push forward of  $\rho_0$  under  $\varphi^t$ :  $\rho_t = \varphi_*^t \rho_0$ . In other words, for any Borel set B,  $\rho_t(\varphi^t(B)) = \rho_0(B)$ .

**Lemma 7.8.** Assume conditions A1, A3 hold. Let  $(\rho_t)_{t\geq 0}$  be the solution of the PDE (7.1) with initialization  $\rho_0$ . Let  $\Omega \subseteq \mathbb{R}^D$  be a Borel set. Suppose  $\varphi^t(\Omega) \subseteq \Omega$ , then we have  $\rho_t(\Omega) \geq \rho_0(\Omega)$ .

*Proof.* The lemma holds immediately by noting that 
$$\rho_t(\Omega) \geq \rho_t(\varphi^t(\Omega)) = \rho_0(\Omega)$$
.

**Lemma 7.9.** Assume conditions A1, A3 hold. Further assume there exists a constant  $K < \infty$  such that

$$|\partial_i \Psi(\boldsymbol{\theta}; \rho)| \le K \cdot \theta_i, \tag{7.80}$$

for any  $\boldsymbol{\theta} \in (0, \infty)^D$  and  $\rho \in \mathcal{P}([0, \infty]^D)$ . Let  $(\rho_t)_{t \geq 0}$  be the solution of the PDE (7.1) with initial condition  $\rho_0$  with  $\rho_0((0, \infty)^D) = 1$ . Then for any  $t < \infty$ ,  $\rho_t((0, \infty)^D) = 1$ .

*Proof.* According to Eqs. (7.80) and (7.79), we have for  $i \in [d]$ ,

$$\theta_i^0 \cdot \exp\{-Kt\} \le \theta_i^t \le \theta_i^0 \cdot \exp\{Kt\}. \tag{7.81}$$

Denote

$$\Omega_k(t) = [1/k \cdot \exp\{-Kt\}, k \cdot \exp\{Kt\}]^D.$$
 (7.82)

Then according to (7.81), we have  $\varphi^t(\Omega_k(0)) \subseteq \Omega_k(t)$ . Note  $\Omega_k(t)$  is increasing in k for fixed t, and  $\bigcup_k \Omega_k(t) = \bigcup_k \Omega_k(0) = (0, \infty)^D$ . Hence,

$$\rho_t((0,\infty)^D) = \lim_{k \to \infty} \rho_t(\Omega_k(t)) \ge \lim_{k \to \infty} \rho_t(\boldsymbol{\varphi}^t(\Omega_k(0))) = \lim_{k \to \infty} \rho_0(\Omega_k(0)) = \rho_0((0,\infty)^D) = 1. \quad (7.83)$$

**Lemma 7.10.** Let  $(\rho_t)_{t>0}$  be a continuous curve in a compact metric space  $(\Omega,d)$ . Denoting

$$\mathcal{S}_* \equiv \{ \rho_* \in \Omega : \exists (t_k)_{k \ge 1}, \lim_{k \to \infty} t_k = \infty, s.t., \lim_{k \to \infty} d(\rho_{t_k}, \rho_*) = 0 \}$$

to be the set of all limiting points of  $(\rho_t)_{t>0}$ . Then  $S_*$  is a connected compact set.

*Proof.* First, it is easy to see that  $S_*$  should be closed. Note that  $\Omega$  is a compact space, then  $S_*$  should be a compact set. If  $S_* = \{\rho_*\}$  is a singleton, this lemma holds automatically. Therefore, we would like to consider the case when  $S_*$  is not a singleton.

For any  $\rho_1, \rho_2 \in \mathcal{S}_*$ , and  $d(\rho_1, \rho_2) > 0$ . We would like to show  $\rho_1$  and  $\rho_2$  are connected in  $\mathcal{S}_*$ .

We use proof by contradiction. Now suppose  $\rho_1$  and  $\rho_2$  are not connected. Define  $\mathcal{A} \subseteq \mathcal{S}_*$  to be the maximal connected subset of  $\mathcal{S}_*$  containing  $\rho_1$ . It is easy to see that  $\mathcal{A}$  is compact. It is also easy to see that its complement  $\mathcal{B} \equiv \mathcal{S}_* \setminus \mathcal{A}$  is also a compact set, and  $\rho_2 \in \mathcal{B}$ . As a result, we have  $\mathcal{A} \cup \mathcal{B} = \mathcal{S}_*$ ,  $\mathcal{A} \cap \mathcal{B} = \emptyset$ , and  $\rho_1 \in \mathcal{A}$ ,  $\rho_2 \in \mathcal{B}$ .

Note that  $\Omega$  is a metric space, so it satisfies T4 separation axiom. Since  $\mathcal{A}$  and  $\mathcal{B}$  are closed sets and  $\mathcal{A} \cap \mathcal{B} = \emptyset$ , there exists an open set  $\mathcal{O}$ , such that  $\mathcal{A} \subseteq \mathcal{O}$ ,  $\mathcal{O} \cap \mathcal{B} = \emptyset$ . Hence,  $\partial \mathcal{O} \subseteq \mathcal{S}^c_*$ .

Note that  $\rho_1$  and  $\rho_2$  are limiting points of  $(\rho_t)_{t\geq 0}$  which is a continuous curve in  $\Omega$ . Therefore, it must cross the boundary  $\partial \mathcal{O}$  infinite times. That is, there is a sequence  $(t_k)_{k\geq 1}$  of time with  $\lim_{k\to\infty} t_k = \infty$ , such that  $\rho_{t_k} \in \partial \mathcal{O}$ . But since  $\partial \mathcal{O}$  is compact, there exists a limiting point  $\rho_* \in \partial \mathcal{O}$ , so that a subsequence of sequence  $\rho_{t_k}$  converges to  $\rho_*$ . Therefore,  $\rho_*$  should be a limiting point of  $(\rho_t)_{t\geq 0}$ . But this contradict with  $\partial \mathcal{O} \subseteq \mathcal{S}_*^c$ .

**Lemma 7.11.** Under the assumptions of A1 and A3, further assume that U, V are twice continuous differentiable, and that  $\rho_0$  has density with respect to the Lebesgue measure, bounded by  $M_0$ . Then  $\rho_t$  also has a density, bounded by  $M_t = K M_0 \exp\{KDt\}$  (where K depends on the constants in the assumptions).

*Proof.* Let  $J(\theta;t)$  for the Jacobian of  $\varphi^t(\cdot)$  at  $\theta^0 = \theta$ . Then Eq. (7.79) implies that  $J(\theta;t)$  satisfies

$$\frac{\mathrm{d}}{\mathrm{d}t} \boldsymbol{J}(\boldsymbol{\theta};t) = -\nabla^2 \Psi(\boldsymbol{\varphi}^t(\boldsymbol{\theta}); \rho_t) \, \boldsymbol{J}(\boldsymbol{\theta};t) \,, \tag{7.84}$$

with initial condition  $J(\theta;0) = I_D$ . This implies

$$\frac{\mathrm{d}}{\mathrm{d}t}\lambda_{\min}(\boldsymbol{J}(\boldsymbol{\theta};t)) \ge -\|\nabla^2\Psi(\boldsymbol{\varphi}^t(\boldsymbol{\theta});\rho_t)\|_{\mathrm{op}}\lambda_{\min}(\boldsymbol{J}(\boldsymbol{\theta};t)). \tag{7.85}$$

Therefore, using the fact that  $\|\nabla^2 \Psi(\boldsymbol{\theta}; \rho_t)\|_{\text{op}}$  is K-bounded, we obtain  $\lambda_{\min}(\boldsymbol{J}(\boldsymbol{\theta}; t)) \geq \exp(-Kt)$ . Finally, since  $\varphi^t$  is a diffeomorphism, we have

$$\rho_t(\boldsymbol{\theta}) = \rho_0((\boldsymbol{\varphi}^t)^{-1}(\boldsymbol{\theta})) \left| \det(\boldsymbol{J}((\boldsymbol{\varphi}^t)^{-1}(\boldsymbol{\theta});t)) \right|^{-1}$$
(7.86)

$$\leq \rho_0((\boldsymbol{\varphi}^t)^{-1}(\boldsymbol{\theta})) \exp(KDt).$$
 (7.87)

This completes the proof.

#### 7.6 Proof of Theorems 6: Stability conditions

In this section, we will prove the stability result in Theorem 6. Throughout the proof we can assume, without loss of generality,  $\xi(t) = 1/2$ . Indeed  $\xi(t)$  amounts just of a change of time. Further we introduce the matrix  $\mathbf{H}_1 = \mathbf{H}_1(\delta_{\theta_*}) \in \mathbb{R}^{D \times D}$  by

$$\boldsymbol{H}_{1}(\delta_{\boldsymbol{\theta}_{*}}) = \nabla^{2}V(\boldsymbol{\theta}_{*}) + \nabla_{1,1}^{2}U(\boldsymbol{\theta}_{*}, \boldsymbol{\theta}_{*}) + \nabla_{1,2}^{2}U(\boldsymbol{\theta}_{*}, \boldsymbol{\theta}_{*}), \qquad (7.88)$$

$$= \boldsymbol{H}_0(\delta_{\boldsymbol{\theta}_*}) + \nabla_{1,2}^2 U(\boldsymbol{\theta}_*, \boldsymbol{\theta}_*), \qquad (7.89)$$

where  $\mathbf{H}_0 \equiv \mathbf{H}_0(\delta_{\boldsymbol{\theta}_*}) = \nabla^2 V(\boldsymbol{\theta}_*) + \nabla^2_{1,1} U(\boldsymbol{\theta}_*, \boldsymbol{\theta}_*)$ . Notice that

$$\langle \boldsymbol{u}, \nabla_{1,2}^2 U(\boldsymbol{\theta}_*, \boldsymbol{\theta}_*) \boldsymbol{u} \rangle = \mathbb{E}\{\langle \boldsymbol{u}, \nabla_{\boldsymbol{\theta}} \sigma_*(\boldsymbol{x}; \boldsymbol{\theta}_*) \rangle^2\},$$
 (7.90)

and therefore  $\nabla_{1,2}^2 U(\boldsymbol{\theta}_*, \boldsymbol{\theta}_*) \succeq \mathbf{0}$ , whence  $\boldsymbol{H}_1 \succeq \boldsymbol{H}_0$ .

We first establish the condition for  $\rho_* = \delta_{\theta_*}$  to be a fixed point. Note that  $\Psi(\theta; \rho_*) = V(\theta) + U(\theta, \theta_*)$  and supp $(\rho_*) = \{\theta_*\}$ . Hence the condition [20] in the main text is satisfied if and only if  $\nabla_{\theta} \Psi(\theta; \rho_*)|_{\theta=\theta_*} = \mathbf{0}$ , i.e.  $\nabla V(\theta_*) + \nabla_1 U(\theta_*, \theta_*) = \mathbf{0}$ .

To establish the stability result of Theorem 6, the following lemma provides a key estimate.

**Lemma 7.12.** Under the assumptions of Theorem 6, let  $\lambda \equiv \lambda_{\min}(\mathbf{H}_0) > 0$ . Then there exists  $r_1, \varepsilon_1, \gamma > 0$  such that the following hold

(i) If  $supp(\rho) \subseteq B(\theta_*; r_1) \equiv \{\theta : \|\theta - \theta_*\|_2 \le r_1\}$ , then,

$$\int \langle (\boldsymbol{\theta} - \boldsymbol{\theta}_*), \nabla \Psi(\boldsymbol{\theta}; \rho) \rangle \, \rho(\mathrm{d}\boldsymbol{\theta}) \ge \frac{\lambda}{2} \int \|\boldsymbol{\theta} - \boldsymbol{\theta}_*\|_2^2 \, \rho(\mathrm{d}\boldsymbol{\theta}) \,. \tag{7.91}$$

(ii) If  $\int \|\boldsymbol{\theta} - \boldsymbol{\theta}_*\|_2^2 \rho(\mathrm{d}\boldsymbol{\theta}) \leq \varepsilon_1^2$  and  $\mathrm{supp}(\rho) \subseteq \mathsf{B}(\boldsymbol{\theta}_*; r_1)$ , then for any  $\boldsymbol{\theta} \in \mathsf{B}(\boldsymbol{\theta}_*; r_1) \setminus \mathsf{B}(\boldsymbol{\theta}_*; r_1/2)$ ,

$$\langle (\boldsymbol{\theta} - \boldsymbol{\theta}_*), \nabla \Psi(\boldsymbol{\theta}; \rho) \rangle \ge \gamma > 0.$$
 (7.92)

*Proof.* Note that

$$\nabla^2 \Psi(\boldsymbol{\theta}; \rho) = \nabla^2 V(\boldsymbol{\theta}) + \int \nabla_1^2 U(\boldsymbol{\theta}, \boldsymbol{\theta}') \, \rho(\mathrm{d}\boldsymbol{\theta}') \,. \tag{7.93}$$

Since  $\nabla^2 V(\boldsymbol{\theta})$  is continuous and  $\nabla_1^2 U(\boldsymbol{\theta}, \boldsymbol{\theta}')$  is bounded continuous, it follows that  $\boldsymbol{\theta} \mapsto \nabla^2 \Psi(\boldsymbol{\theta}; \rho)$  is continuous, and  $\rho \mapsto \nabla^2 \Psi(\boldsymbol{\theta}; \rho)$  is continuous in the weak topology, and in fact  $(\boldsymbol{\theta}, \rho) \mapsto \nabla^2 \Psi(\boldsymbol{\theta}; \rho)$  is continuous in the product topology.

Further, we have

$$\nabla^2 \Psi(\boldsymbol{\theta}_*; \rho_*) = \nabla^2 V(\boldsymbol{\theta}_*) + \nabla^2_{11} U(\boldsymbol{\theta}_*, \boldsymbol{\theta}_*) = \boldsymbol{H}_0.$$
 (7.94)

Since  $\mathbf{H}_0 \succ \mathbf{0}$  strictly, for any  $\delta > 0$  we can choose  $r_1 = r_1(\delta) > 0$  such that

$$\nabla^2 \Psi(\boldsymbol{\theta}; \rho) \succeq (1 - \delta) \, \boldsymbol{H}_0 \,, \tag{7.95}$$

$$\|\nabla_{12}^2 U(\boldsymbol{\theta}_*, \boldsymbol{\theta}) - \nabla_{12}^2 U(\boldsymbol{\theta}_*, \boldsymbol{\theta}_*)\|_{\text{op}} \le \delta,$$

$$(7.96)$$

for all  $\theta \in \mathsf{B}(\theta_*; r_1)$ , and  $\rho$  such that  $\mathsf{supp}(\rho) \subseteq \mathsf{B}(\theta_*; r_1)$ . If these conditions hold

$$\langle (\boldsymbol{\theta} - \boldsymbol{\theta}_*), \nabla \Psi(\boldsymbol{\theta}; \rho) \rangle = \langle (\boldsymbol{\theta} - \boldsymbol{\theta}_*), \nabla \Psi(\boldsymbol{\theta}; \rho) - \nabla \Psi(\boldsymbol{\theta}_*; \rho) \rangle + \langle (\boldsymbol{\theta} - \boldsymbol{\theta}_*), \nabla \Psi(\boldsymbol{\theta}_*; \rho) \rangle$$
(7.97)

$$= \langle (\boldsymbol{\theta} - \boldsymbol{\theta}_*), \nabla^2 \Psi(\tilde{\boldsymbol{\theta}}; \rho) (\boldsymbol{\theta} - \boldsymbol{\theta}_*) \rangle + \langle (\boldsymbol{\theta} - \boldsymbol{\theta}_*), \nabla \Psi(\boldsymbol{\theta}_*; \rho) \rangle$$
 (7.98)

$$\geq (1 - \delta) \langle (\boldsymbol{\theta} - \boldsymbol{\theta}_*), \boldsymbol{H}_0 (\boldsymbol{\theta} - \boldsymbol{\theta}_*) \rangle + \langle (\boldsymbol{\theta} - \boldsymbol{\theta}_*), \nabla \Psi(\boldsymbol{\theta}_*; \rho) \rangle. \tag{7.99}$$

In order to bound the second term, note that, since  $\nabla \Psi(\boldsymbol{\theta}_*; \rho_*) = \mathbf{0}$ ,

$$\nabla \Psi(\boldsymbol{\theta}_*; \rho) = \int \left[ \nabla_1 U(\boldsymbol{\theta}_*, \boldsymbol{\theta}') - \nabla_1 U(\boldsymbol{\theta}_*, \boldsymbol{\theta}_*) \right] \rho(\mathrm{d}\boldsymbol{\theta}') = \nabla_{12}^2 U(\boldsymbol{\theta}_*, \boldsymbol{\theta}_*) \boldsymbol{\mu} + \boldsymbol{\xi}, \qquad (7.100)$$

$$\boldsymbol{\mu} = \int (\boldsymbol{\theta} - \boldsymbol{\theta}_*) \, \rho(\mathrm{d}\boldsymbol{\theta}) \,, \tag{7.101}$$

$$\boldsymbol{\xi} = \int \left[ \nabla_1 U(\boldsymbol{\theta}_*, \boldsymbol{\theta}') - \nabla_1 U(\boldsymbol{\theta}_*, \boldsymbol{\theta}_*) - \nabla_{12}^2 U(\boldsymbol{\theta}_*, \boldsymbol{\theta}_*) (\boldsymbol{\theta}' - \boldsymbol{\theta}_*) \right] \rho(\mathrm{d}\boldsymbol{\theta}') \,. \tag{7.102}$$

Substituting in Eq. (7.99), we obtain

$$\langle (\boldsymbol{\theta} - \boldsymbol{\theta}_*), \nabla \Psi(\boldsymbol{\theta}; \rho) \rangle \ge (1 - \delta) \langle (\boldsymbol{\theta} - \boldsymbol{\theta}_*), \boldsymbol{H}_0(\boldsymbol{\theta} - \boldsymbol{\theta}_*) \rangle + \langle (\boldsymbol{\theta} - \boldsymbol{\theta}_*), (\boldsymbol{H}_1 - \boldsymbol{H}_0) \boldsymbol{\mu} \rangle + \langle (\boldsymbol{\theta} - \boldsymbol{\theta}_*), \boldsymbol{\xi} \rangle.$$
(7.103)

By the intermediate value theorem, for any  $v \in \mathbb{R}^D$ , there exists  $\tilde{\theta} = \tilde{\theta}(v, \theta) \in [\theta_*, \theta]$  such that

$$\langle \boldsymbol{v}, \boldsymbol{\xi} \rangle = \int \langle \boldsymbol{v}, [\nabla_{12}^2 U(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_*) - \nabla_{12}^2 U(\boldsymbol{\theta}_*, \boldsymbol{\theta}_*)] (\boldsymbol{\theta} - \boldsymbol{\theta}_*) \rangle \rho(\mathrm{d}\boldsymbol{\theta})$$
(7.104)

$$\geq -\int \|\boldsymbol{v}\|_{2} \|\nabla_{12}^{2} U(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_{*}) - \nabla_{12}^{2} U(\boldsymbol{\theta}_{*}, \boldsymbol{\theta}_{*})\|_{\mathrm{op}} \|\boldsymbol{\theta} - \boldsymbol{\theta}_{*}\|_{2} \rho(\mathrm{d}\boldsymbol{\theta})$$
(7.105)

$$\geq -\delta \|\boldsymbol{v}\|_{2} \int \|\boldsymbol{\theta} - \boldsymbol{\theta}_{*}\|_{2} \rho(\mathrm{d}\boldsymbol{\theta}) \tag{7.106}$$

$$\geq -\delta \|\boldsymbol{v}\|_2 \sqrt{\mathsf{Tr}(\boldsymbol{Q}) + \|\boldsymbol{\mu}\|_2^2} \tag{7.107}$$

$$\geq -\delta \|\boldsymbol{v}\|_2 \sqrt{\mathsf{Tr}(\boldsymbol{Q})} - \delta \|\boldsymbol{v}\|_2 \|\boldsymbol{\mu}\|_2, \tag{7.108}$$

where  $\mathbf{Q} = \int (\boldsymbol{\theta} - \mu)(\boldsymbol{\theta} - \mu)^{\mathsf{T}} \rho(\mathrm{d}\boldsymbol{\theta})$  is the covariance of  $(\boldsymbol{\theta} - \boldsymbol{\theta}_*)$ .

Let now consider the claim at point (i). Integrating Eq. (7.103) with respect to  $\rho(d\theta)$ , we get

$$\int \langle (\boldsymbol{\theta} - \boldsymbol{\theta}_*), \nabla \Psi(\boldsymbol{\theta}; \rho) \rangle \rho(\mathrm{d}\boldsymbol{\theta}) \ge (1 - \delta) \langle \boldsymbol{H}_0, \boldsymbol{Q} + \boldsymbol{\mu} \boldsymbol{\mu}^\mathsf{T} \rangle + \langle \boldsymbol{\mu}, (\boldsymbol{H}_1 - \boldsymbol{H}_0) \boldsymbol{\mu} \rangle + \langle \boldsymbol{\mu}, \boldsymbol{\xi} \rangle$$
(7.109)

$$\geq (1 - \delta)\langle \boldsymbol{H}_0, \boldsymbol{Q} \rangle + \langle \boldsymbol{\mu}, (\boldsymbol{H}_1 - \delta \boldsymbol{H}_0) \boldsymbol{\mu} \rangle - \delta \|\boldsymbol{\mu}\|_2 \sqrt{\mathsf{Tr}(\boldsymbol{Q})} - \delta \|\boldsymbol{\mu}\|_2^2$$
 (7.110)

$$\geq (1 - \delta)\langle \boldsymbol{H}_0, \boldsymbol{Q} \rangle + \langle \boldsymbol{\mu}, (\boldsymbol{H}_1 - \delta \boldsymbol{H}_0) \boldsymbol{\mu} \rangle - \frac{3\delta}{2} \|\boldsymbol{\mu}\|_2^2 - \frac{\delta}{2} \operatorname{Tr}(\boldsymbol{Q})$$
 (7.111)

$$= \langle (1 - \delta) \mathbf{H}_0 - \frac{\delta}{2} \mathbf{I}, \mathbf{Q} \rangle + \langle \boldsymbol{\mu}, (\mathbf{H}_1 - \delta \mathbf{H}_0 - \frac{3\delta}{2} \mathbf{I}) \boldsymbol{\mu} \rangle.$$
 (7.112)

By choosing  $\delta$  sufficiently small, we can ensure that  $(1 - \delta)\mathbf{H}_0 - (\delta/2)\mathbf{I} \succeq \lambda_{\min}(\mathbf{H}_0)\mathbf{I}/2$ ,  $\mathbf{H}_1 - \delta\mathbf{H}_0 - (3\delta/2)\mathbf{I} \succeq \lambda_{\min}(\mathbf{H}_1)\mathbf{I}/2$ , and therefore

$$\int \langle (\boldsymbol{\theta} - \boldsymbol{\theta}_*), \nabla \Psi(\boldsymbol{\theta}; \rho) \rangle \, \rho(\mathrm{d}\boldsymbol{\theta}) \ge \frac{1}{2} \lambda_{\min}(\boldsymbol{H}_0) \, \mathsf{Tr}(\boldsymbol{Q}) + \frac{1}{2} \lambda_{\min}(\boldsymbol{H}_1) \, \|\boldsymbol{\mu}\|_2^2, \tag{7.113}$$

which yields the claim (7.91).

Next consider point (ii). In this case, Eq. (7.107) implies

$$\langle (\boldsymbol{\theta} - \boldsymbol{\theta}_*), \boldsymbol{\xi} \rangle \ge -\delta \varepsilon_1 \|\boldsymbol{\theta} - \boldsymbol{\theta}_*\|_2.$$
 (7.114)

Substituting in Eq. (7.103), and using  $\|\boldsymbol{\mu}\|_2 \leq \varepsilon_1$ , we get

$$\langle (\boldsymbol{\theta} - \boldsymbol{\theta}_*), \nabla \Psi(\boldsymbol{\theta}; \rho) \rangle \ge (1 - \delta) \langle \boldsymbol{H}_0, (\boldsymbol{\theta} - \boldsymbol{\theta}_*)^{\otimes 2} \rangle - \varepsilon_1 (\lambda_{\max}(\boldsymbol{H}_1) + \lambda_{\max}(\boldsymbol{H}_0) + \delta) \|\boldsymbol{\theta} - \boldsymbol{\theta}_*\|_2$$

$$\ge (1 - \delta) \lambda \left(\frac{r_1}{2}\right)^2 - \varepsilon_1 (\lambda_{\max}(\boldsymbol{H}_1) + \lambda_{\max}(\boldsymbol{H}_0) + \delta) r_1. \tag{7.115}$$

This is strictly positive for all  $\varepsilon_1$  small enough, hence implying the claim (7.92).

We are now in position of proving Theorem 6.

Proof of Theorem 6. Let  $r_0 = \min(r_1/2, \varepsilon_1/2)$  and assume, without loss of generality  $t_0 = 0$ , so that  $\operatorname{supp}(\rho_0) \subseteq \mathsf{B}(\boldsymbol{\theta}_*; r_0)$ . We also define

$$T_1 \equiv \inf \left\{ t : \int \|\boldsymbol{\theta} - \boldsymbol{\theta}_*\|_2^2 \rho_t(\mathrm{d}\boldsymbol{\theta}) > \varepsilon_1^2 \right\},$$
 (7.116)

$$T_2 \equiv \inf \left\{ t : \operatorname{supp}(\rho_t) \not\subseteq \mathsf{B}(\boldsymbol{\theta}_*; r_1) \right\},$$
 (7.117)

$$T_* \equiv \min(T_1, T_2). \tag{7.118}$$

As usual, we adopt the convention that the infimum of an empty set is equal to  $+\infty$ .

Define  $\varphi_1(\boldsymbol{\theta}) = h(\|\boldsymbol{\theta} - \boldsymbol{\theta}_*\|_2)$ , with h to be an non-decreasing function with

$$h(r) = \begin{cases} 0 & \text{if } r < r_1/2, \\ \text{smooth intropolation} & \text{if } r_1/2 \le r < 5r_1/8, \\ 2r/r_1 - 1 & \text{if } 5r_1/8 \le r < 7r_1/8, \\ \text{smooth intropolation} & \text{if } 7r_1/8 \le r < r_1, \\ 1 & \text{if } r \ge r_1. \end{cases}$$
(7.119)

For any  $t < T_*$ , the PDE (7.1) implies

$$\partial_t \langle \varphi_1, \rho_t \rangle = -\int \langle \nabla \varphi_1(\boldsymbol{\theta}), \nabla \Psi(\boldsymbol{\theta}; \rho_t) \rangle \rho_t(d\boldsymbol{\theta})$$
 (7.120)

$$= -\frac{2}{r_1} \int h'(\|\boldsymbol{\theta} - \boldsymbol{\theta}_*\|_2) \langle \frac{(\boldsymbol{\theta} - \boldsymbol{\theta}_*)}{\|\boldsymbol{\theta} - \boldsymbol{\theta}_*\|_2}, \nabla \Psi(\boldsymbol{\theta}; \rho_t) \rangle \rho_t(d\boldsymbol{\theta})$$
(7.121)

$$\leq -\frac{4\gamma}{r_1^2} \rho_t \Big( \mathsf{B}(\boldsymbol{\theta}_*; 7r_1/8) \setminus \mathsf{B}(\boldsymbol{\theta}_*; 5r_1/8) \Big) , \tag{7.122}$$

where, in the last inequality, we used Lemma 7.12.(ii). Next, define

$$\varphi_2(\boldsymbol{\theta}) = \frac{1}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}_*\|_2^2. \tag{7.123}$$

Applying again Eq. (7.1), we get, for  $t \leq T_*$ ,

$$\partial_t \langle \varphi_2, \rho_t \rangle = -\int \langle \nabla \varphi_2(\boldsymbol{\theta}), \nabla \Psi(\boldsymbol{\theta}; \rho_t) \rangle \rho_t(d\boldsymbol{\theta})$$
 (7.124)

$$= -\int \langle (\boldsymbol{\theta} - \boldsymbol{\theta}_*), \nabla \Psi(\boldsymbol{\theta}; \rho_t) \rangle \rho_t(d\boldsymbol{\theta})$$
 (7.125)

$$\leq -\lambda \left\langle \varphi_2, \rho_t \right\rangle. \tag{7.126}$$

Together the last two bounds imply  $T_* = \infty$ . Indeed assume by contradiction  $T_* < \infty$ . Then either  $T_1 \le T_2$ ,  $T_1 < \infty$ , or  $T_2 < T_1$ ,  $T_2 < \infty$ .

Consider the first case:  $T_1 \leq T_2$ ,  $T_1 < \infty$ . Since  $\langle \rho_{T_1}, \varphi_2 \rangle \geq \varepsilon_1^2$  but  $\langle \rho_0, \varphi_2 \rangle \leq r_0^2 \leq \varepsilon_1^2/4$ , there exists  $t < T_*$  such that  $\partial_t \langle \rho_0, \varphi_2 \rangle > 0$ . However this contradicts Eq. (7.126). Consider then the second case:  $T_2 < T_1$ ,  $T_2 < \infty$ . This implies  $\langle \rho_{T_2}, \varphi_1 \rangle > 0$ , but on the other hand  $\langle \rho_0, \varphi_1 \rangle = 0$ . Hence, there exists  $t < T_*$  such that  $\partial_t \langle \rho_0, \varphi_1 \rangle > 0$ . However this contradicts Eq. (7.122).

We conclude that  $T_* = \infty$  and hence we can apply Eq. (7.126) for any t, thus obtaining  $\partial_t \langle \varphi_2, \rho_t \rangle \leq -\lambda \langle \varphi_2, \rho_t \rangle$  and hence  $\langle \varphi_2, \rho_t \rangle \leq (r_0^2/2)e^{-\lambda t}$ , which concludes the proof.

# 7.7 Proof of Theorem 7: Instability conditions

In this section we will prove the instability result of Theorem 7. Throughout the section, we assume  $\xi(t) \equiv 1/2$ . We will use several times the nonlinear dynamics, defined for  $\rho_t$  a solution of Eq. (7.1) with initial condition  $\rho_0$ :

$$\dot{\boldsymbol{\theta}}^t = -\nabla \Psi(\boldsymbol{\theta}^t; \rho_t) \,. \tag{7.127}$$

**Lemma 7.13.** Let  $\nu$  be a probability measure on  $\mathbb{R}^d$ , absolutely continuous with respect to the Lebesgue measure, with density bounded by M, and let  $\mathbf{u} \in \mathbb{R}^d$  be a unit vector. Further assume that, for some  $\mathbf{x}_0 \in \mathbb{R}^d$ , r > 0, we have  $\nu(\mathsf{B}(\mathbf{x}_0; r)) \geq 1 - \varepsilon$ , with  $0 < \varepsilon < 1/20$ . Then there exists a coupling  $\gamma$  of  $\nu$  with itself (i.e. a probability distribution on  $\mathbb{R}^d \times \mathbb{R}^d$  with marginals  $\int \gamma(\cdot, d\mathbf{x}) = \int \gamma(d\mathbf{x}, \cdot) = \nu(\cdot)$ ) and a constant L = L(d, r, M) such that the following holds. If  $(\mathbf{x}_1, \mathbf{x}_2) \sim \gamma$ , then

$$\gamma\left(\langle \boldsymbol{u}, \boldsymbol{x}_1 - \boldsymbol{x}_2 \rangle \ge \frac{1}{L}; \ \boldsymbol{P}_{\boldsymbol{u}}^{\perp}(\boldsymbol{x}_1 - \boldsymbol{x}_2) = \boldsymbol{0}\right) \ge \frac{9}{10},$$
 (7.128)

where  $m{P}_{m{u}}^{\perp} = m{I} - m{u}m{u}^{\mathsf{T}}$  is the projector orthogonal to vector  $m{u}$ .

Proof. First consider the case d=1: in this case, the assumption  $\nu(\mathsf{B}(\boldsymbol{x}_0;r)) \geq 1-\varepsilon$  is not required. Denote by F the distribution function associated to  $\nu$  (i.e.  $F(x) \equiv \nu((-\infty,x])$ ). By assumption F is differentiable with  $F'(x) \leq M$ . In order to construct the desired coupling, let Z be a random variable uniformly distributed in [0,1]. For a small constant  $\xi_0 > 0$ , define the random variables  $(X_1, X_2)$  by letting

$$X_1 = F^{-1}(Z), (7.129)$$

$$X_2 = \begin{cases} F^{-1}(Z - \xi_0) & \text{if } Z > \xi_0, \\ F^{-1}(Z + 1 - \xi_0) & \text{if } Z < \xi_0. \end{cases}$$
 (7.130)

(Note that  $X_2$  is not defined for  $Z = \xi_0$  but this is a zero-probability event.) On the event  $\{Z > \xi_0\}$  (which has probability  $1 - \xi_0$ ), we have, for some  $W \in [X_1, X_2]$ ,

$$\xi_0 = F'(W)(X_1 - X_2) \le M(X_1 - X_2).$$
 (7.131)

By choosing  $\xi_0$  small enough, this proves the claim for d=1.

Consider next d > 1 and assume without loss of generality  $u = e_1$ .

Let  $\overline{\nu}(\cdot) = \nu(\cdot | \mathbf{X} \in \mathsf{B}(\mathbf{x}_0; r))$ ,  $\mathbf{X}_a^b \equiv (X_a, \dots, X_b)$ , and denote by  $f_{1|[2,d]}$  the density of  $\overline{\nu}(X_1 \in \mathcal{X}_2^n)$ , and by  $f_{[a,b]}$  the density of  $\overline{\nu}(\mathbf{X}_a^b \in \cdot)$ . We then have

$$f_{1|[2,d]}(x_1|\boldsymbol{x}_2^d) = \frac{f_{[1,d]}(\boldsymbol{x}_1^d)}{f_{[2,d]}(\boldsymbol{x}_2^d)} \le \frac{M}{f_{[2,d]}(\boldsymbol{x}_2^d)}.$$
 (7.132)

Further, we have

$$\overline{\nu}(\{\boldsymbol{x}: f_{[2,d]}(\boldsymbol{x}_2^d) \le \Delta\}) = \int \mathbf{1}_{f_{[2,d]}(\boldsymbol{x}_2^d) \le \Delta} f_{[2,d]}(\boldsymbol{x}_2^d) \, \mathrm{d}\boldsymbol{x}_2^d$$
 (7.133)

$$\leq \Delta \int_{\mathsf{B}((\boldsymbol{x}_0)_2^d;r)} \mathrm{d}\boldsymbol{x}_2^d \leq C_d \Delta r^{d-1} \,. \tag{7.134}$$

In order to construct the coupling, we sample  $Z \sim \nu$ . If  $Z \notin \mathsf{B}(\boldsymbol{x}_0;r)$ , then we take  $X_1 = X_2 = Z$ . If  $Z \in \mathsf{B}(\boldsymbol{x}_0;r)$  and  $\max_{x_1} f_{1|[2,d]}(x_1|\boldsymbol{Z}_2^d) > M/\Delta$ , we also take  $X_1 = X_2 = Z$ . Otherwise we have  $Z \in \mathsf{B}(\boldsymbol{x}_0;r)$  and  $\max_{x_1} f_{1|[2,d]}(x_1|\boldsymbol{Z}_2^d) \leq M/\Delta$ , then we sample  $(X_{1,1},X_{2,1})$  from the coupling developed in the case d=1 applied to  $f_{1|[2,d]}(\cdot|\boldsymbol{Z}_2^d)$ , and set  $X_1 = (X_{1,1},\boldsymbol{Z}_2^d)$ ,  $X_2 = (X_{2,1},\boldsymbol{Z}_2^d)$ . Now define  $\gamma$  to be the joint distribution of  $X_1, X_2$ . Then  $\gamma$  is a coupling of  $\nu$  with itself.

The above analysis yields

$$\gamma\left(\langle \boldsymbol{u}, \boldsymbol{X}_1 - \boldsymbol{X}_2 \rangle \ge \frac{\xi_0 \Delta}{M}; \quad \boldsymbol{P}_{\boldsymbol{u}}^{\perp}(\boldsymbol{X}_1 - \boldsymbol{X}_2) = \boldsymbol{0}\right) \ge 1 - \xi_0 - C_d \Delta r^{d-1} - \varepsilon. \tag{7.135}$$

Hence, we can choose  $\Delta, \xi_0$  small enough so that the claim (7.128) holds.

For any  $u \in \mathbb{R}$ , define the level set  $\tilde{\mathcal{L}}(u)$ ,

$$\tilde{\mathcal{L}}(u) \equiv \{ \boldsymbol{\theta} \in \mathbb{R}^D : \ \Psi(\boldsymbol{\theta}; \rho_*) \le u \}. \tag{7.136}$$

According to the notation of Theorem 7, we have  $\mathcal{L}(\eta) = \tilde{\mathcal{L}}(\Psi(\theta_*; \rho_*) - \eta)$  for any  $\eta \in \mathbb{R}$ .

**Lemma 7.14.** For any  $u \in \mathbb{R}$ ,  $\Delta > 0$  such that  $\partial \hat{\mathcal{L}}(u_0)$  is compact for all  $u_0 \in (u - \Delta, u)$ , there exists  $\varepsilon_{0,\#} > 0$  such that the following holds. Let  $(\rho_t)_{t \geq t_0}$  be a solution of the PDE (7.1) such that  $d_{\mathrm{BL}}(\rho_t, \rho_*) \leq \varepsilon_{0,\#}$  for all  $t \geq t_0$ . Let  $(\boldsymbol{\theta}^t)_{t \geq t_0}$  be a solution of the ODE (7.127) with  $\Psi(\boldsymbol{\theta}^{t_0}; \rho_*) \leq u - \Delta$ . Then  $\Psi(\boldsymbol{\theta}^t; \rho_*) \leq u$  for all  $t \geq t_0$ .

Proof. By Sard's theorem [GP10], there exists  $u_0 \in (u - \Delta, u)$  such that the boundary  $\partial \tilde{\mathcal{L}}(u_0)$  contains no critical points of  $\Psi(\cdot; \rho_*)$ . If we define  $g_0 = \min_{\boldsymbol{\theta} \in \partial \tilde{\mathcal{L}}(u_0)} \|\nabla \Psi(\boldsymbol{\theta}; \rho_*)\|_2$ , the minimum is achieved by compactness, and therefore we have  $g_0 > 0$  strictly. Notice that by the differentiability assumptions on V and U,  $\partial \tilde{\mathcal{L}}(u_0)$  is a  $C^1$  submanifold of  $\mathbb{R}^D$ , with  $\nabla \Psi(\boldsymbol{\theta}; \rho_*)$  orthogonal to  $\partial \tilde{\mathcal{L}}_0(u_0)$  and directed toward the exterior. Further, as observed already above,

$$\|\nabla \Psi(\boldsymbol{\theta}; \rho_t) - \nabla \Psi(\boldsymbol{\theta}; \rho_*)\|_2 = \left\| \int \nabla_{\boldsymbol{\theta}} U(\boldsymbol{\theta}; \boldsymbol{\theta}') (\rho_t - \rho_*) (\mathrm{d}\boldsymbol{\theta}') \right\|_2$$
 (7.137)

$$\leq K d_{\mathrm{BL}}(\rho_t, \rho_*) \leq K \varepsilon_{0,\#}. \tag{7.138}$$

By choosing  $\varepsilon_{0,\#}$  small enough, we can ensure  $\|\nabla \Psi(\boldsymbol{\theta}; \rho_t) - \nabla \Psi(\boldsymbol{\theta}; \rho_*)\|_2 \le g_0/3$  for all  $\boldsymbol{\theta}$  and all  $t \ge t_0$ .

Assume by contradiction that  $\Psi(\boldsymbol{\theta}^{t_1}; \rho_*) > u$  for some  $t_1 \geq t_0$ , and let  $t_* = \sup\{t \leq t_1 : \Psi(\boldsymbol{\theta}^t; \rho_*) \leq u_0\}$ . Note that, by continuity of the trajectory,  $\boldsymbol{\theta}^{t_*} \in \partial \tilde{\mathcal{L}}(u_0)$ . We then must have

$$0 \le \frac{\mathrm{d}}{\mathrm{d}t} \Psi(\boldsymbol{\theta}^{t_*}; \rho_*) = -\langle \nabla \Psi(\boldsymbol{\theta}^{t_*}; \rho_{t_*}), \nabla \Psi(\boldsymbol{\theta}^{t_*}; \rho_*) \rangle$$
 (7.139)

$$\leq -\|\nabla \Psi(\boldsymbol{\theta}^{t_*}; \rho_*)\|_2^2 + \|\nabla \Psi(\boldsymbol{\theta}^{t_*}; \rho_*)\|_2 \|\nabla \Psi(\boldsymbol{\theta}^{t_*}; \rho_{t_*}) - \nabla \Psi(\boldsymbol{\theta}^{t_*}; \rho_*)\|_2$$
 (7.140)

$$\leq -\frac{2}{3}g_0 \|\nabla \Psi(\boldsymbol{\theta}^{t_*}; \rho_*)\|_2, \tag{7.141}$$

which leads to a contradiction since  $\boldsymbol{\theta}^{t_*} \in \partial \tilde{\mathcal{L}}(u_0)$  and hence  $\|\nabla \Psi(\boldsymbol{\theta}^{t_*}; \rho_*)\|_2 > 0$ .

To prove Theorem 7, let now assume by contradiction that  $\rho_t \Rightarrow \rho_* = p_* \delta_{\theta_*} + (1 - p_*) \tilde{\rho}_*$  weakly. Then for any  $\varepsilon_0, r_0 > 0$  (to be chosen below), we can find  $t_0 = t_0(\varepsilon_0, r_0)$  such that

$$d_{\mathrm{BL}}(\rho_t, \rho_*) \le \varepsilon_0, \quad |\rho_t(\mathsf{B}(\boldsymbol{\theta}_*; r_0)) - p_*| \le \varepsilon_0$$
 (7.142)

for all  $t \geq t_0$ . Let  $\overline{\rho}_{t_0}$  be the conditional probability measure of  $\rho_{t_0}$  given  $\boldsymbol{\theta} \in \mathsf{B}(\boldsymbol{\theta}_*; r_0)$ . By Lemma 7.11,  $\overline{\rho}_{t_0}$  has a density upper bounded by a constant  $M = M(\varepsilon_0, t_0)$  (note that  $\overline{\rho}_{t_0}(S) \leq \rho_{t_0}(S)/(p_* - \varepsilon_0)$ ).

Set  $\mathbf{H}_0 = \mathbf{H}_0(\rho_*) = \nabla^2 \Psi(\boldsymbol{\theta}_*; \rho_*)$ . Since  $\boldsymbol{\theta}_*$  is a critical point of  $\boldsymbol{\theta} \mapsto \Psi(\boldsymbol{\theta}; \rho_*)$ , for any  $\delta > 0$ , we can find  $r_1(\delta) > 0$  such that

$$\boldsymbol{\theta} \in \mathsf{B}(\boldsymbol{\theta}_*; r_1) \Rightarrow \|\nabla^2 \Psi(\boldsymbol{\theta}; \rho_*) - \boldsymbol{H}_0\|_{\mathrm{op}} \leq \frac{\delta}{2}, \|\nabla \Psi(\boldsymbol{\theta}_*; \rho_*)\|_2 = 0.$$
 (7.143)

As shown in the proof of Theorem 6, the function  $(\theta, \rho) \mapsto \nabla^2 \Psi(\theta; \rho)$  is continuous when the space of probability distributions  $\rho$  is endowed with the weak topology. Analogously  $\rho \mapsto \nabla \Psi(\theta_*; \rho)$  is continuous in the weak topology. Hence for this  $\delta > 0$  and  $r_1(\delta) > 0$ , there exists  $\varepsilon_{0,*}(\delta, r_1) > 0$  small enough such that, the following inequalities hold

$$\boldsymbol{\theta} \in \mathsf{B}(\boldsymbol{\theta}_*; r_1), \ d_{\mathrm{BL}}(\rho, \rho_*) \le \varepsilon_{0,*} \Rightarrow \ \left\| \nabla^2 \Psi(\boldsymbol{\theta}; \rho) - \boldsymbol{H}_0 \right\|_{\mathrm{op}} \le \delta, \quad \left\| \nabla \Psi(\boldsymbol{\theta}_*; \rho) \right\|_2 \le \delta^2 r_1 / 2.$$

$$(7.144)$$

Let us emphasize that  $r_1$  depends on  $\delta$  but can be taken to be independent of  $\varepsilon_0$ . Further, by an application of the intermediate value theorem, for all  $\theta \in \mathsf{B}(\theta_*; r_1)$ ,

$$\left|\Psi(\boldsymbol{\theta}; \rho_*) - \Psi(\boldsymbol{\theta}_*; \rho_*) - \frac{1}{2} \langle (\boldsymbol{\theta} - \boldsymbol{\theta}_*), \boldsymbol{H}_0(\boldsymbol{\theta} - \boldsymbol{\theta}_*) \rangle \right| \leq \frac{1}{2} \delta \|\boldsymbol{\theta} - \boldsymbol{\theta}_*\|_2^2.$$
 (7.145)

For  $r_0 < r_1$ ,  $\boldsymbol{\theta}^{t_0} \in \mathsf{B}(\boldsymbol{\theta}_*; r_0)$ , we let  $(\boldsymbol{\theta}^t)_{t \geq t_0}$  be the solution of Eq. (7.127) with this initial condition. We then define

$$t_{\text{exit}}(\boldsymbol{\theta}^{t_0}, r_1) = \inf \left\{ t \ge t_0 : \boldsymbol{\theta}^t \notin \mathsf{B}(\boldsymbol{\theta}_*; r_1) \right\}, \tag{7.146}$$

$$t_{\text{return}}(\boldsymbol{\theta}^{t_0}, r_0, r_1) = \inf \{ t > t_{\text{exit}}(\boldsymbol{\theta}^{t_0}, r_1) : \boldsymbol{\theta}^t \in \mathsf{B}(\boldsymbol{\theta}_*; r_0) \}.$$
 (7.147)

**Lemma 7.15.** Under the conditions of Theorem 7, there exists  $r_1 > 0$  and  $\varepsilon_{0,*} > 0$  such that, for all  $r_0 \leq r_1$ ,  $\varepsilon_0 \leq \varepsilon_{0,*}$ , there exists  $T_{\rm UB}(\varepsilon_0, r_0, r_1, t_0)$  such that the following happens. If  $d_{\rm BL}(\rho_t, \rho_*) \leq \varepsilon_0$  and  $|\rho_t(\mathsf{B}(\boldsymbol{\theta}_*; r_0)) - p_*| \leq \varepsilon_0$  for all  $t \geq t_0$  for some  $t_0$ , then

$$\rho_{t_0}\Big(\big\{\boldsymbol{\theta}^{t_0} \in \mathsf{B}(\boldsymbol{\theta}_*; r_0) : t_{\mathrm{exit}}(\boldsymbol{\theta}^{t_0}, r_1) \le T_{\mathrm{UB}}(\varepsilon_0, r_0, r_1, t_0)\big\}\Big) \ge \frac{1}{3} p_*. \tag{7.148}$$

*Proof.* Let  $\boldsymbol{u}$  be an eigenvector of  $\boldsymbol{H}_0$  corresponding to the eigenvalue  $\lambda_{\min}(\boldsymbol{H}_0) = -\lambda_1$ . By condition B1 of Theorem 7, we have  $\lambda_1 > 0$ . Let  $-\lambda_2$  denote the second smallest eigenvalue (which can be positive). We further denote by  $\boldsymbol{P} \in \mathbb{R}^{D \times D}$  the orthogonal projector onto the eigenspace corresponding to  $\lambda_{\min}(\boldsymbol{H}_0)$  and by  $\boldsymbol{P}_{\perp} = \boldsymbol{I} - \boldsymbol{P}$  the projector onto the orthogonal subspace.

We fix a  $\delta \leq (\lambda_1 - \lambda_2)/10$ . Then we choose  $r_1 > 0$  and  $\varepsilon_{0,*} > 0$  such that Eq. (7.144) holds, with an additional requirement that  $\varepsilon_{0,*} < p_*/10$ . We will prove this lemma with this choice of  $r_1$  and  $\varepsilon_{0,*}$ .

We always denote  $(\boldsymbol{\theta}_i^t)_{t \geq t_0}$  to be the solution of Eq. (7.127) with initial condition  $\boldsymbol{\theta}_i^{t_0}$ , for i = 1, 2. First we claim that, for  $0 < \delta \leq (\lambda_1 - \lambda_2)/10$ , assuming

$$\|\nabla^2 \Psi(\boldsymbol{\theta}; \rho_t) - \boldsymbol{H}_0\|_{\text{op}} \le \delta, \quad \forall t \ge t_0, \quad \forall \boldsymbol{\theta} \in \mathsf{B}(\boldsymbol{\theta}_*; r_1),$$
 (7.149)

then for any  $\boldsymbol{\theta}_1^{t_0}, \boldsymbol{\theta}_2^{t_0} \in \mathsf{B}(\boldsymbol{\theta}_*; r_1)$  with  $\boldsymbol{P}_{\perp}(\boldsymbol{\theta}_1^{t_0} - \boldsymbol{\theta}_2^{t_0}) = \boldsymbol{0}$ , we have

$$\|\boldsymbol{\theta}_{1}^{t} - \boldsymbol{\theta}_{2}^{t}\|_{2} \ge \|\boldsymbol{\theta}_{1}^{t_{0}} - \boldsymbol{\theta}_{2}^{t_{0}}\|_{2} e^{\lambda_{1}(t - t_{0})/2}$$

$$(7.150)$$

 $\text{for all } t \in [t_0, t_{\text{exit}}(\boldsymbol{\theta}_1^{t_0}, r_1) \wedge t_{\text{exit}}(\boldsymbol{\theta}_2^{t_0}, r_1)].$ 

For now we assume this claim holds. Fix  $r_0 \leq r_1$  and  $\varepsilon_0 \leq \varepsilon_{0,*}$ . Define  $\gamma$  to be the coupling of Lemma 7.13 corresponding to  $\boldsymbol{u}$  which is the eigenvector corresponding to the least eigenvalue of  $\boldsymbol{H}_0$ , and  $\nu = \overline{\rho}_{t_0}$  which is the conditional measure of  $\rho_{t_0}$  given  $\boldsymbol{\theta}^{t_0} \in \mathsf{B}(\boldsymbol{\theta}_*; r_0)$ . Note  $\overline{\rho}_{t_0}$  has a density upper bounded by a constant  $M = M(\varepsilon_0, t_0)$ . By Lemma 7.13, we have  $\gamma(\mathcal{E}) \geq 9/10$ , where

$$\mathcal{E} \equiv \left\{ (\boldsymbol{\theta}_1^{t_0}, \boldsymbol{\theta}_2^{t_0}) \in \mathsf{B}(\boldsymbol{\theta}_*; r_0) \times \mathsf{B}(\boldsymbol{\theta}_*; r_0) : \langle \boldsymbol{u}, \boldsymbol{\theta}_1^{t_0} - \boldsymbol{\theta}_2^{t_0} \rangle \ge \frac{1}{Z}; \ \boldsymbol{P}_{\boldsymbol{u}}^{\perp}(\boldsymbol{\theta}_1^{t_0} - \boldsymbol{\theta}_2^{t_0}) = \boldsymbol{0} \right\}$$
(7.151)

for some  $Z = Z(\varepsilon_0, r_0, t_0) > 0$ . Now we take  $(\boldsymbol{\theta}_1^{t_0}, \boldsymbol{\theta}_2^{t_0}) \in \mathcal{E}$ . Note the assumption of this lemma gives  $d_{\text{BL}}(\rho_t, \rho_*) \leq \varepsilon_0 \leq \varepsilon_{0,*}$  for all  $t \geq t_0$ . According to Eq. (7.144), we have Eq. (7.149) holds, and due to this claim, we have  $\|\boldsymbol{\theta}_1^t - \boldsymbol{\theta}_2^t\|_2 \geq (1/Z)e^{\lambda_1(t-t_0)/2}$  for all  $t \in [t_0, t_{\text{exit}}(\boldsymbol{\theta}_1^{t_0}, r_1) \wedge t_{\text{exit}}(\boldsymbol{\theta}_2^{t_0}, r_1)]$ .

due to this claim, we have  $\|\boldsymbol{\theta}_1^t - \boldsymbol{\theta}_2^t\|_2 \ge (1/Z)e^{\lambda_1(t-t_0)/2}$  for all  $t \in [t_0, t_{\text{exit}}(\boldsymbol{\theta}_1^{t_0}, r_1) \land t_{\text{exit}}(\boldsymbol{\theta}_2^{t_0}, r_1)]$ . Define  $T_{\text{UB}}(\varepsilon_0, r_0, r_1, t_0) = (2/\lambda_1)\log(2Z\,r_1)$ . Then for  $t > T_{\text{UB}}$ , we have  $\|\boldsymbol{\theta}_1^t - \boldsymbol{\theta}_2^t\|_2 \ge 2r_1$ . This is impossible if  $\boldsymbol{\theta}_1^t, \boldsymbol{\theta}_2^t \in \mathsf{B}(\boldsymbol{\theta}_*; r_1)$  and hence we deduce  $(t_{\text{exit}}(\boldsymbol{\theta}_1^{t_0}, r_1) \land t_{\text{exit}}(\boldsymbol{\theta}_2^{t_0}, r_1)) \le T_{\text{UB}}$  for all  $(\boldsymbol{\theta}_1^{t_0}, \boldsymbol{\theta}_2^{t_0}) \in \mathcal{E}$ .

Therefore, we get

$$\begin{split} &\frac{9}{10} \leq \gamma(\mathcal{E}) \leq \gamma \Big( \Big\{ (\boldsymbol{\theta}_1^{t_0}, \boldsymbol{\theta}_2^{t_0}) \in \mathsf{B}(\boldsymbol{\theta}_*; r_0) \times \mathsf{B}(\boldsymbol{\theta}_*; r_0) \ : \ t_{\mathrm{exit}}(\boldsymbol{\theta}_1^{t_0}, r_1) \wedge t_{\mathrm{exit}}(\boldsymbol{\theta}_2^{t_0}, r_1) \leq T_{\mathrm{UB}} \Big\} \Big) \\ &\leq \gamma \Big( \Big\{ \boldsymbol{\theta}_1^{t_0} \in \mathsf{B}(\boldsymbol{\theta}_*; r_0) \ : \ t_{\mathrm{exit}}(\boldsymbol{\theta}_1^{t_0}, r_1) \leq T_{\mathrm{UB}} \Big\} \Big) + \gamma \Big( \Big\{ \boldsymbol{\theta}_2^{t_0} \in \mathsf{B}(\boldsymbol{\theta}_*; r_0) \ : \ t_{\mathrm{exit}}(\boldsymbol{\theta}_2^{t_0}, r_1) \leq T_{\mathrm{UB}} \Big\} \Big) \\ &= 2 \, \overline{\rho}_{t_0} \Big( \Big\{ \boldsymbol{\theta}^{t_0} \in \mathsf{B}(\boldsymbol{\theta}_*; r_0) \ : \ t_{\mathrm{exit}}(\boldsymbol{\theta}^{t_0}, r_1) \leq T_{\mathrm{UB}} \Big\} \Big) \,. \end{split}$$

Denoting by S the event in the last expression, we obtain  $\rho_{t_0}(S) \ge (p_* - \varepsilon_0)\overline{\rho}_{t_0}(S) \ge (9/20)(p_* - \varepsilon_0) \ge p_*/3$  by noting that  $\varepsilon_0 < p_*/10$ .

**Proof of the claim.** Define the quantities

$$x_{\parallel}(t) = \|\boldsymbol{P}(\boldsymbol{\theta}_{1}^{t} - \boldsymbol{\theta}_{2}^{t})\|_{2}^{2},$$
 (7.152)

$$x_{\perp}(t) = \|\mathbf{P}_{\perp}(\mathbf{\theta}_1^t - \mathbf{\theta}_2^t)\|_2^2.$$
 (7.153)

We then have, for  $t \in [t_0, t_{\text{exit}}(\boldsymbol{\theta}_1^{t_0}, r_1) \wedge t_{\text{exit}}(\boldsymbol{\theta}_2^{t_0}, r_1)],$ 

$$\begin{split} \dot{x}_{\parallel}(t) &= 2 \langle \boldsymbol{P}(\boldsymbol{\theta}_{1}^{t} - \boldsymbol{\theta}_{2}^{t}), -\nabla \Psi(\boldsymbol{\theta}_{1}^{t}; \rho_{t}) + \nabla \Psi(\boldsymbol{\theta}_{2}^{t}; \rho_{t}) \rangle \\ &\stackrel{(a)}{=} 2 \langle \boldsymbol{P}(\boldsymbol{\theta}_{1}^{t} - \boldsymbol{\theta}_{2}^{t}), -\nabla^{2} \Psi(\tilde{\boldsymbol{\theta}}^{t}; \rho_{t}) (\boldsymbol{\theta}_{1}^{t} - \boldsymbol{\theta}_{2}^{t}) \rangle \\ &= -2 \langle (\boldsymbol{\theta}_{1}^{t} - \boldsymbol{\theta}_{2}^{t}), \boldsymbol{P} \nabla^{2} \Psi(\tilde{\boldsymbol{\theta}}^{t}; \rho_{t}) \boldsymbol{P} (\boldsymbol{\theta}_{1}^{t} - \boldsymbol{\theta}_{2}^{t}) \rangle - 2 \langle (\boldsymbol{\theta}_{1}^{t} - \boldsymbol{\theta}_{2}^{t}), \boldsymbol{P} \nabla^{2} \Psi(\tilde{\boldsymbol{\theta}}^{t}; \rho_{t}) \boldsymbol{P}_{\perp} (\boldsymbol{\theta}_{1}^{t} - \boldsymbol{\theta}_{2}^{t}) \rangle \\ &\geq 2 (\lambda_{1} - \delta) \|\boldsymbol{P}(\boldsymbol{\theta}_{1}^{t} - \boldsymbol{\theta}_{2}^{t})\|_{2}^{2} - 2\delta \|\boldsymbol{P}(\boldsymbol{\theta}_{1}^{t} - \boldsymbol{\theta}_{2}^{t})\|_{2} \|\boldsymbol{P}_{\perp} (\boldsymbol{\theta}_{1}^{t} - \boldsymbol{\theta}_{2}^{t})\|_{2} \\ &\geq 2 (\lambda_{1} - \delta) x_{\parallel}(t) - \delta(x_{\parallel}(t) + x_{\perp}(t)), \end{split}$$

where in (a) we used the intermediate value theorem (with  $\tilde{\boldsymbol{\theta}}^t$  a point between  $\boldsymbol{\theta}_1^t$  and  $\boldsymbol{\theta}_2^t$ ), and in (b) we used Eq. (7.149).

Proceeding analogously for  $x_{\perp}(t)$ , we get (for a new choice of  $\tilde{\theta}^t$ )

$$\dot{x}_{\perp}(t) = 2\langle \boldsymbol{P}_{\perp}(\boldsymbol{\theta}_{1}^{t} - \boldsymbol{\theta}_{2}^{t}), -\nabla\Psi(\boldsymbol{\theta}_{1}^{t}; \rho_{t}) + \nabla\Psi(\boldsymbol{\theta}_{2}^{t}; \rho_{t})\rangle 
= 2\langle \boldsymbol{P}_{\perp}(\boldsymbol{\theta}_{1}^{t} - \boldsymbol{\theta}_{2}^{t}), -\nabla^{2}\Psi(\tilde{\boldsymbol{\theta}}^{t}; \rho_{t})(\boldsymbol{\theta}_{1}^{t} - \boldsymbol{\theta}_{2}^{t})\rangle 
= -2\langle (\boldsymbol{\theta}_{1}^{t} - \boldsymbol{\theta}_{2}^{t}), \boldsymbol{P}_{\perp}\nabla^{2}\Psi(\tilde{\boldsymbol{\theta}}^{t}; \rho_{t})\boldsymbol{P}_{\perp}(\boldsymbol{\theta}_{1}^{t} - \boldsymbol{\theta}_{2}^{t})\rangle - 2\langle (\boldsymbol{\theta}_{1}^{t} - \boldsymbol{\theta}_{2}^{t}), \boldsymbol{P}_{\perp}\nabla^{2}\Psi(\tilde{\boldsymbol{\theta}}^{t}; \rho_{t})\boldsymbol{P}(\boldsymbol{\theta}_{1}^{t} - \boldsymbol{\theta}_{2}^{t})\rangle 
\leq 2(\lambda_{2} + \delta)\|\boldsymbol{P}_{\perp}(\boldsymbol{\theta}_{1}^{t} - \boldsymbol{\theta}_{2}^{t})\|_{2}^{2} + 2\delta\|\boldsymbol{P}(\boldsymbol{\theta}_{1}^{t} - \boldsymbol{\theta}_{2}^{t})\|_{2}\|\boldsymbol{P}_{\perp}(\boldsymbol{\theta}_{1}^{t} - \boldsymbol{\theta}_{2}^{t})\|_{2} 
\leq 2(\lambda_{2} + \delta)x_{\parallel}(t) + \delta(x_{\parallel}(t) + x_{\perp}(t)).$$

Summarizing, we obtained the inequalities

$$\dot{x}_{\parallel}(t) \ge (2\lambda_1 - 3\delta)x_{\parallel}(t) - \delta x_{\perp}(t),$$
 (7.154)

$$\dot{x}_{\perp}(t) \le \delta x_{\parallel}(t) + (2\lambda_2 + 3\delta)x_{\perp}(t). \tag{7.155}$$

The matrix of coefficients on the right-hand side is

$$\mathbf{A} = \begin{pmatrix} 2\lambda_1 - 3\delta & -\delta \\ \delta & 2\lambda_2 + 3\delta \end{pmatrix} . \tag{7.156}$$

This has a (un-normalized) left eigenvectors (1, -v), (-v, 1) with eigenvalues  $\xi_{\pm}$  given by:

$$v = \frac{1}{\delta} \left[ \lambda_1 - \lambda_2 - 3\delta - \sqrt{(\lambda_1 - \lambda_2 - 3\delta)^2 - \delta^2} \right],\tag{7.157}$$

$$\xi_{\pm} = \lambda_1 + \lambda_2 \pm \sqrt{(\lambda_1 - \lambda_2 - 3\delta)^2 - \delta^2}. \tag{7.158}$$

Note we took  $\delta < (\lambda_1 - \lambda_2)/10$ , we have v > 0, and  $\xi_+ \geq \lambda_1$ .

Multiplying the inequalities (7.154), (7.155) by (1, -v), we thus obtain

$$\frac{\mathrm{d}}{\mathrm{d}t} \left( x_{\parallel}(t) - v \, x_{\perp}(t) \right) \ge \xi_{+} \left( x_{\parallel}(t) - v \, x_{\perp}(t) \right). \tag{7.159}$$

Since we assumed  $x_{\perp}(t_0) = 0$ , whence, for all  $t \in [t_0, t_{\text{exit}}(\boldsymbol{\theta}_1^{t_0}, r_1) \land t_{\text{exit}}(\boldsymbol{\theta}_2^{t_0}, r_1)]$ , we have

$$x_{\parallel}(t) \ge x_{\parallel}(t) - v \, x_{\perp}(t) \ge x_{\parallel}(t_0) \, e^{\xi_{+}(t - t_0)} \ge x_{\parallel}(t_0) \, e^{\lambda_1(t - t_0)}.$$
 (7.160)

We next strengthen the last lemma and prove that trajectories that exit  $\mathsf{B}(\boldsymbol{\theta}_*; r_1)$  do not re-enter  $\mathsf{B}(\boldsymbol{\theta}_*; r_0)$ .

**Lemma 7.16.** Under the conditions of Theorem 7, there exists  $r_{0,*}, r_1 > 0$  (with  $r_{0,*} < r_1$ ) and  $\varepsilon_{0,*} > 0$  such that, for all  $r_0 \le r_{0,*}$ ,  $\varepsilon_0 \le \varepsilon_{0,*}$ , there exists  $T_{\text{UB}}(\varepsilon_0, r_0, r_1, t_0)$  such that the following happens. If  $d_{\text{BL}}(\rho_t, \rho_*) \le \varepsilon_0$  and  $|\rho_t(\mathsf{B}(\boldsymbol{\theta}_*; r_0)) - p_*| \le \varepsilon_0$  for all  $t \ge t_0$  for some  $t_0$ , then

$$\rho_{t_0}\Big(\big\{\boldsymbol{\theta}^{t_0} \in \mathsf{B}(\boldsymbol{\theta}_*; r_0) : \ t_{\text{exit}}(\boldsymbol{\theta}^{t_0}, r_1) \le T_{\text{UB}}(\varepsilon_0, r_0, r_1, t_0), \ t_{\text{return}}(\boldsymbol{\theta}^{t_0}, r_0, r_1) = \infty\big\}\Big) \ge \frac{1}{3} p_*. \quad (7.161)$$

*Proof.* Let  $P_+$  be the projector onto the eigenspace of  $-H_0$  corresponding to positive eigenvalues, and  $P_-$  the projector onto the subspace corresponding to negative eigenvalues, and let

 $\lambda_0 \equiv \min_{i \leq D} |\lambda_i(\boldsymbol{H}_0)|$  to be the least absolute value of eigenvalue of  $\boldsymbol{H}_0$ . By condition B1 of Theorem 7, we have  $\lambda_0 > 0$ . Let  $\lambda_{\max}$  denote the largest absolute value of eigenvalue of  $\boldsymbol{H}_0$ .

Fix a  $\delta$  such that  $0 < \delta \le \min\{\lambda_0/(1 + \lambda_0 + \lambda_{\max}), \sqrt{\lambda_0/\lambda_{\max}}, \lambda_1 - \lambda_2, 1\}/10$ , where  $\lambda_1, \lambda_2$  are as defined in Lemma 7.15. Next we choose  $r_1$  as per Lemma 7.15, and we further require  $\lambda_0 r_1^2 \le \eta_0$ , where  $\eta_0$  is as per condition B3 in the statement of Theorem 7. We take  $\varepsilon_{0,*}$  to be the minimum of the parameter  $\varepsilon_{0,*}$  as per Lemma 7.15 and the parameter  $\varepsilon_{0,\#}$  as per Lemma 7.14, where in Lemma 7.14, we choose  $u = \Psi(\boldsymbol{\theta}_*; \rho_*) - \lambda_0 r_1^2/8$ , and  $\Delta = \lambda_0 r_1^2/8$ . Then we will choose smaller  $r_1$  and  $\varepsilon_{0,*}$  so that Eq. (7.144) holds. Finally, we take  $r_{0,*} = \delta r_1 < r_1$ . We will prove this lemma with this choice of  $r_1$ ,  $\varepsilon_{0,*}$ , and  $r_{0,*}$ , and with the same function  $T_{\text{UB}}$  as per Lemma 7.15.

Define

$$t_*(\boldsymbol{\theta}^{t_0}; r_1, \delta) \equiv \sup \left\{ t \in (t_0, t_{\text{exit}}(\boldsymbol{\theta}^{t_0}, r_1)) : \|\boldsymbol{\theta}_1^t - \boldsymbol{\theta}_*\|_2 < \delta r_1 \right\}, \tag{7.162}$$

and define

$$z_{+}(t) = \|P_{+}(\theta^{t} - \theta_{*})\|_{2}^{2}, \tag{7.163}$$

$$z_{-}(t) = \|\mathbf{P}_{-}(\boldsymbol{\theta}^{t} - \boldsymbol{\theta}_{*})\|_{2}^{2}. \tag{7.164}$$

We bound the evolution of these quantities following the same argument used above for  $x_{\parallel}(t)$ ,  $x_{\perp}(t)$ . Namely

$$\begin{split} \dot{z}_{+}(t) = & 2\langle \boldsymbol{P}_{+}(\boldsymbol{\theta}^{t} - \boldsymbol{\theta}_{*}), -\nabla\Psi(\boldsymbol{\theta}^{t}; \rho_{t}) + \nabla\Psi(\boldsymbol{\theta}_{*}; \rho_{t})\rangle - 2\langle \boldsymbol{P}_{+}(\boldsymbol{\theta}^{t} - \boldsymbol{\theta}_{*}), \nabla\Psi(\boldsymbol{\theta}_{*}; \rho_{t})\rangle \\ = & - 2\langle \boldsymbol{P}_{+}(\boldsymbol{\theta}^{t} - \boldsymbol{\theta}_{*}), \nabla^{2}\Psi(\tilde{\boldsymbol{\theta}}^{t}; \rho_{t})(\boldsymbol{\theta}^{t} - \boldsymbol{\theta}_{*})\rangle - 2\langle \boldsymbol{P}_{+}(\boldsymbol{\theta}^{t} - \boldsymbol{\theta}_{*}), \nabla\Psi(\boldsymbol{\theta}_{*}; \rho_{t})\rangle \\ = & - 2\langle (\boldsymbol{\theta}^{t} - \boldsymbol{\theta}_{*}), \boldsymbol{P}_{+}\nabla^{2}\Psi(\tilde{\boldsymbol{\theta}}^{t}; \rho_{t})\boldsymbol{P}_{+}(\boldsymbol{\theta}^{t} - \boldsymbol{\theta}_{*})\rangle \\ & - 2\langle (\boldsymbol{\theta}^{t} - \boldsymbol{\theta}_{*}), \boldsymbol{P}_{+}\nabla^{2}\Psi(\tilde{\boldsymbol{\theta}}^{t}; \rho_{t})\boldsymbol{P}_{-}(\boldsymbol{\theta}^{t} - \boldsymbol{\theta}_{*})\rangle - 2\langle \boldsymbol{P}_{+}(\boldsymbol{\theta}^{t} - \boldsymbol{\theta}_{*}), \nabla\Psi(\boldsymbol{\theta}_{*}; \rho_{t})\rangle \\ \geq & 2(\lambda_{0} - \delta)\|\boldsymbol{P}_{+}(\boldsymbol{\theta}^{t} - \boldsymbol{\theta}_{*})\|_{2}^{2} - 2\delta\|\boldsymbol{P}_{+}(\boldsymbol{\theta}^{t} - \boldsymbol{\theta}_{*})\|_{2}\|\boldsymbol{P}_{-}(\boldsymbol{\theta}_{1}^{t} - \boldsymbol{\theta}_{2}^{t})\|_{2} - \delta^{2}r_{1}\|\boldsymbol{P}_{+}(\boldsymbol{\theta}^{t} - \boldsymbol{\theta}_{*})\|_{2} \\ \geq & (2\lambda_{0} - 3\delta)z_{+}(t) - \delta z_{-}(t) - \delta^{2}r_{1}\sqrt{z_{+}(t)} \,. \end{split}$$

For  $t \in [t_*(\boldsymbol{\theta}^{t_0}; r_1, \delta), t_{\text{exit}}(\boldsymbol{\theta}^{t_0}; r_1)]$ , we have  $\sqrt{z_+(t) + z_-(t)} \ge \delta r_1$ . Using the inequality  $\sqrt{a(a+b)} \le a+b$  holding for non-negative a and b, we have

$$\dot{z}_{+}(t) \ge (2\lambda_0 - 3\delta)z_{+}(t) - \delta z_{-}(t) - \delta^2 r_1 \sqrt{z_{+}(t)}$$
(7.165)

$$\geq (2\lambda_0 - 3\delta)z_+(t) - \delta z_-(t) - \delta \sqrt{z_+(t)(z_+(t) + z_-(t))}$$
(7.166)

$$\geq (2\lambda_0 - 3\delta)z_+(t) - \delta z_-(t) - \delta z_+(t) - \delta z_-(t) \tag{7.167}$$

$$\geq (2\lambda_0 - 4\delta)z_+(t) - 2\delta z_-(t). \tag{7.168}$$

Proceeding analogously for  $z_{-}$ , we arrive at the inequalities

$$\dot{z}_{+}(t) \ge (2\lambda_0 - 4\delta)z_{+}(t) - 2\delta z_{-}(t), \qquad (7.169)$$

$$\dot{z}_{-}(t) \le 2\delta z_{+}(t) - (2\lambda_0 - 4\delta)z_{-}(t), \tag{7.170}$$

for  $t \in [t_*(\boldsymbol{\theta}^{t_0}; r_1, \delta), t_{\text{exit}}(\boldsymbol{\theta}^{t_0}; r_1)]$ . The matrix of coefficients on the right-hand side has a left eigenvector of the form (-w, 1) with corresponding eigenvalue  $-\tilde{\xi}$ , whereby  $\tilde{\xi} = \sqrt{\lambda_0^2 - 4\delta^2}$  and

 $w = (\lambda_0 - \sqrt{\lambda_0^2 - 4\delta^2})/(2\delta)$ . In particular, since  $\delta < \lambda_0/10$ , we have  $\tilde{\xi} \ge \lambda_0/2 > 0$  and w > 0. Multiplying the above inequalities by (-w, 1), we get

$$\frac{\mathrm{d}}{\mathrm{d}t} \left( -wz_{+}(t) + z_{-}(t) \right) \le -\tilde{\xi} \left( -wz_{+}(t) + z_{-}(t) \right), \tag{7.171}$$

and therefore, for all  $t \in [t_*(\boldsymbol{\theta}^{t_0}; r_1, \delta), t_{\text{exit}}(\boldsymbol{\theta}^{t_0}; r_1)], z_-(t) \leq w z_+(t) + e^{-\tilde{\xi}t} (-w z_+(0) + z_-(0)) \leq w z_+(t) + \delta^2 r_1^2$ . In particular, for  $t = t_{\text{exit}}(\boldsymbol{\theta}^{t_0}; r_1)$ , using  $z_+(t_{\text{exit}}) + z_-(t_{\text{exit}}) = r_1^2$ , we finally obtain

$$\|P_{+}(\theta^{t_{\text{exit}}} - \theta_{*})\|_{2}^{2} \ge r_{1}^{2} \left(\frac{1 - \delta^{2}}{1 + w}\right) \ge r_{1}^{2} (1 - \delta),$$
 (7.172)

$$\left\| \mathbf{P}_{-}(\boldsymbol{\theta}^{t_{\text{exit}}} - \boldsymbol{\theta}_{*}) \right\|_{2}^{2} \le r_{1}^{2} \delta. \tag{7.173}$$

Using Eq. (7.145), we obtain

$$\Psi(\boldsymbol{\theta}^{t_{\text{exit}}}; \rho_*) \le \Psi(\boldsymbol{\theta}_*; \rho_*) + \frac{1}{2} \langle (\boldsymbol{\theta} - \boldsymbol{\theta}_*), \boldsymbol{H}_0(\boldsymbol{\theta} - \boldsymbol{\theta}_*) \rangle + \frac{1}{2} \delta r_1^2$$
(7.174)

$$\leq \Psi(\boldsymbol{\theta}_*; \rho_*) - \frac{1}{2}\lambda_0 \|\boldsymbol{P}_+(\boldsymbol{\theta}^{t_{\text{exit}}} - \boldsymbol{\theta}_*)\|_2^2 + \frac{1}{2}\lambda_{\text{max}} \|\boldsymbol{P}_-(\boldsymbol{\theta}^{t_{\text{exit}}} - \boldsymbol{\theta}_*)\|_2^2 + \frac{1}{2}\delta r_1^2 \qquad (7.175)$$

$$\leq \Psi(\boldsymbol{\theta}_*; \rho_*) - \frac{1}{2}\lambda_0 r_1^2 + \frac{1}{2}(1 + \lambda_0 + \lambda_{\max})\delta r_1^2.$$
 (7.176)

Since  $\delta \leq \lambda_0/(10(1+\lambda_0+\lambda_{\max}))$ , we can ensure that  $\Psi(\boldsymbol{\theta}^{t_{\text{exit}}}; \rho_*) \leq \Psi(\boldsymbol{\theta}_*; \rho_*) - \lambda_0 r_1^2/4$ . By Lemma 7.14, since  $d_{\text{BL}}(\rho_t, \rho_*) \leq \varepsilon_{0,*} \leq \varepsilon_{0,\#}$  for all  $t \geq t_0$ , we have  $\Psi(\boldsymbol{\theta}^t; \rho_*) \leq \Psi(\boldsymbol{\theta}_*; \rho_*) - \lambda_0 r_1^2/8$  for all  $t \geq t_{\text{exit}}(\boldsymbol{\theta}^{t_0}; r_1)$ . Note for all  $\boldsymbol{\theta} \in \mathsf{B}(\boldsymbol{\theta}_*; \delta r_1)$ , we have  $\Psi(\boldsymbol{\theta}; \rho_*) \geq \Psi(\boldsymbol{\theta}_*; \rho_*) - \lambda_{\max} \delta^2 r_1^2/2$ . Since  $\delta \leq \sqrt{\lambda_0/\lambda_{\max}}/10$ , we have  $\boldsymbol{\theta}^t \notin \mathsf{B}(\boldsymbol{\theta}_*; \delta r_1)$  for all  $t \geq t_{\text{exit}}(\boldsymbol{\theta}^{t_0}; r_1)$ .

This implies that, for any  $\boldsymbol{\theta}^{t_0} \in \mathsf{B}(\boldsymbol{\theta}_*; r_0)$  for  $r_0 \leq r_{0,*}$  with  $t_{\mathrm{exit}}(\boldsymbol{\theta}^{t_0}, r_1) \leq T_{\mathrm{UB}}(\varepsilon_0, r_0, r_1, t_0) < \infty$ , it will never return to  $\mathsf{B}(\boldsymbol{\theta}_*; r_0)$ . This gives the desired result.

Finally we upper bound the probability that  $\boldsymbol{\theta}^t \in \mathsf{B}(\boldsymbol{\theta}_*; r_0)$  for some  $t > t_0$ , given that  $\boldsymbol{\theta}^{t_0} \not\in \mathsf{B}(\boldsymbol{\theta}_*; r_0)$ . We define

$$t_{\text{enter}}(\boldsymbol{\theta}^{t_0}, r_0) = \inf \{ t \ge t_0 : \boldsymbol{\theta}^t \in \mathsf{B}(\boldsymbol{\theta}_*; r_0) \}.$$
 (7.177)

**Lemma 7.17.** Under the conditions of Theorem 7, for any  $\eta > 0$ , there exists  $r_{0,*} > 0$  and  $\varepsilon_{0,*} > 0$  such that, for all  $r_0 \leq r_{0,*}$ ,  $\varepsilon_0 \leq \varepsilon_{0,*}$ , the following happens. If  $d_{BL}(\rho_t, \rho_*) \leq \varepsilon_0$  and  $|\rho_t(\mathsf{B}(\boldsymbol{\theta}_*; r_0)) - p_*| \leq \varepsilon_0$  for all  $t \geq t_0$  for some  $t_0$ , then

$$\rho_{t_0}\Big(\big\{\boldsymbol{\theta}^{t_0} \not\in \mathsf{B}(\boldsymbol{\theta}_*; r_0) : t_{\text{enter}}(\boldsymbol{\theta}^{t_0}, r_0) = \infty\big\}\Big) \ge 1 - p_* - \eta. \tag{7.178}$$

Proof. Due to condition B2 of Theorem 7, we can choose  $u_1$  with  $\Psi(\boldsymbol{\theta}_*; \rho_*) - \eta_0 < u_1 < \Psi(\boldsymbol{\theta}_*; \rho_*)$  (where  $\eta_0$  is as per condition B3 of Theorem 7) such that  $\rho_*(\tilde{\mathcal{L}}(u_1)) \geq 1 - p_* - \eta/2$  (recall the notation  $\tilde{\mathcal{L}}$  defined as Eq. (7.136)). By taking  $\varepsilon_{0,*}$  small enough, and since  $\boldsymbol{\theta} \mapsto \Psi(\boldsymbol{\theta}; \rho_*)$  is Lipschitz continuous, we can also choose  $u_2 \in (u_1, \Psi(\boldsymbol{\theta}_*; \rho_*))$  such that  $\rho_{t_0}(\tilde{\mathcal{L}}(u_2)) \geq 1 - p_* - \eta$ . Fix  $u_3 \in (u_2, \Psi(\boldsymbol{\theta}_*, \rho_*))$ . Applying Lemma 7.14, we can further reduce  $\varepsilon_{0,*}$ , so that for any initialization  $\boldsymbol{\theta}^{t_0} \in \tilde{\mathcal{L}}(u_2)$ , we have  $\boldsymbol{\theta}^t \in \tilde{\mathcal{L}}(u_3)$  for all t. Further, by continuity of  $\Psi(\cdot; \rho_*)$ , we can choose  $r_{0,*}$ 

small enough so that  $\mathsf{B}(\boldsymbol{\theta}_*; r_{0,*}) \cap \tilde{\mathcal{L}}(u_3) = \emptyset$ , whence

$$\rho_{t_0} \Big( \big\{ \boldsymbol{\theta}^{t_0} \not\in \mathsf{B}(\boldsymbol{\theta}_*; r_0) : t_{\text{enter}}(\boldsymbol{\theta}^{t_0}, r_0) = \infty \big\} \Big)$$
 (7.179)

$$\geq \rho_{t_0} \left( \left\{ \boldsymbol{\theta}^{t_0} : \ \Psi(\boldsymbol{\theta}^{t_0}; \rho_*) < u_2, \ t_{\text{enter}}(\boldsymbol{\theta}^{t_0}, r_0) = \infty \right\} \right)$$
 (7.180)

$$= \rho_{t_0} \Big( \{ \boldsymbol{\theta}^{t_0} : \ \Psi(\boldsymbol{\theta}^{t_0}; \rho_*) < u_2 \} \Big) \ge 1 - p_* - \eta.$$
 (7.181)

The proof of Theorem 7 follows immediately from Lemma 7.16 and Lemma 7.17. Indeed, let  $\eta = p_*/10$ . Take  $\varepsilon_0 \leq \min\{\varepsilon_{0,*}, p_*/10\}$  where  $\varepsilon_{0,*}$  is the minimum of  $\varepsilon_{0,*}$  as per Lemma 7.16 and 7.17. Take  $r_1$  as per Lemma 7.16. Take  $r_0 \leq \min\{r_{0,*}, r_1\}$  where  $r_{0,*}$  is the minimum of  $r_{0,*}$  as per Lemma 7.16 and 7.17. With this choice of  $\varepsilon_0$  and  $r_0$ , there exists  $t_0 > 0$  such that Eq. (7.142) holds for all  $t \geq t_0$ . Setting  $t_* = T_{\text{UB}}(\varepsilon_0, r_0, r_1, t_0) \geq t_0$  with  $T_{\text{UB}}$  given in Lemma 7.16. Denoting by  $\mathbb{P}_{t_0, \rho_{t_0}}$  be the probability distribution over trajectories of (7.127) with  $\boldsymbol{\theta}^{t_0} \sim \rho_{t_0}$ , we have

$$\begin{split} \rho_{t_*}(\mathsf{B}(\boldsymbol{\theta}_*;r_0)) = & \mathbb{P}_{t_0,\rho_{t_0}} \big( \boldsymbol{\theta}^{t_*} \in \mathsf{B}(\boldsymbol{\theta}_*;r_0) \big) \\ = & \mathbb{P}_{t_0,\rho_{t_0}} \big( \boldsymbol{\theta}^{t_0} \in \mathsf{B}(\boldsymbol{\theta}_*;r_0); \ \boldsymbol{\theta}^{t_*} \in \mathsf{B}(\boldsymbol{\theta}_*;r_0) \big) + \mathbb{P}_{t_0,\rho_{t_0}} \big( \boldsymbol{\theta}^{t_0} \not\in \mathsf{B}(\boldsymbol{\theta}_*;r_0); \ \boldsymbol{\theta}^{t_*} \in \mathsf{B}(\boldsymbol{\theta}_*;r_0) \big) \\ \leq & \mathbb{P}_{t_0,\rho_{t_0}} \big( \boldsymbol{\theta}^{t_0} \in \mathsf{B}(\boldsymbol{\theta}_*;r_0) \big) - \mathbb{P}_{t_0,\rho_{t_0}} \big( \boldsymbol{\theta}^{t_0} \in \mathsf{B}(\boldsymbol{\theta}_*;r_0); \ t_{\mathrm{exit}}(\boldsymbol{\theta}^{t_0};r_1) < t_*, t_{\mathrm{return}}(\boldsymbol{\theta}^{t_0};r_0) = \infty \big) \\ + & \mathbb{P}_{t_0,\rho_{t_0}} \big( \boldsymbol{\theta}^{t_0} \not\in \mathsf{B}(\boldsymbol{\theta}_*;r_0) \big) - \mathbb{P}_{t_0,\rho_{t_0}} \big( \boldsymbol{\theta}^{t_0} \not\in \mathsf{B}(\boldsymbol{\theta}_*;r_0); \ t_{\mathrm{enter}}(\boldsymbol{\theta}^{t_0},r_0) = \infty \big) \\ \leq & 1 - \frac{1}{3} p_* - (1 - p_* - \eta) = 2p_*/3 + \eta \,. \end{split}$$

Since we also had  $\rho_t(\mathsf{B}(\theta_*; r_0)) \ge p_* - \varepsilon_0$  for all  $t \ge t_0$ , note  $\eta, \varepsilon_0 \le p_*/10$ , we reached a contradiction.

# 8 Centered isotropic Gaussians

In this section we consider the centered isotropic Gaussians example discussed in the main text. That is, we assume the joint law of (y, x) to be as follows:

With probability 1/2: y = +1,  $\boldsymbol{x} \sim N(\boldsymbol{0}, (1 + \Delta)^2 \boldsymbol{I}_d)$ .

With probability 1/2: y = -1,  $\boldsymbol{x} \sim N(\boldsymbol{0}, (1 - \Delta)^2 \boldsymbol{I}_d)$ .

We assume  $0 < \Delta < 1$ , and choose  $\sigma_*(\boldsymbol{x}; \boldsymbol{\theta}_i) = \sigma(\langle \boldsymbol{x}, \boldsymbol{w}_i \rangle)$  for some activation function  $\sigma$ . Define  $q(r) \equiv \mathbb{E}\{\sigma(rG)\}$  for  $G \sim \mathsf{N}(0,1)$ . We assume  $\sigma(\cdot)$  satisfies the following conditions S0 - S4:

So  $x \mapsto \sigma(x)$  is bounded, non-decreasing, Lipschitz continuous. Its weak derivative  $x \mapsto \sigma'(x)$  is Lipschitz in a neighborhood of 0.

S1 q is analytic on  $(0,\infty)$  with  $\sup_{r\in[0,\infty]}q''(r)<\infty$ .

S2 q'(r) > 0 for all  $r \in (0, \infty)$ , with  $\sup_{r \in [0, \infty]} q'(r) < \infty$ , and  $\lim_{r \to 0} q'(r) = \lim_{r \to \infty} q'(r) = 0$ .

 $\mathsf{S3} \ -\infty < q(0+) < -1, \ 1 < q(+\infty) < \infty, \ \mathrm{and} \ -1 < (q(0+) + q(+\infty))/2 < 1.$ 

S4 Letting  $Z(r) \equiv q'(\tau_- r)/q'(\tau_+ r)$  for some  $\tau_+ > \tau_- > 0$  we have Z'(r) > 0 for all  $r \in (0, \infty)$ .

Note that condition S1 and part of S2 are implied by S0, but we list them here for conveniency. Some of these assumptions can be relaxed at the cost of extra technical work. In the interest of simplicity, we prefer to avoid being overly general.

As our running example we will use

$$\sigma(t) = \begin{cases}
s_1 & \text{if } t \le t_1, \\
(s_2(t - t_1) + s_1(t_2 - t))/(t_2 - t_1) & \text{if } t \in (t_1, t_2), \\
s_2 & \text{if } t \ge t_2.
\end{cases}$$
(8.1)

In particular, we choose  $s_1 = -2.5$ ,  $s_2 = 7.5$ ,  $t_1 = 0.5$ ,  $t_2 = 1.5$  in our simulations. In section 8.5, we check that this choice satisfies the above assumptions.

Throughout this section, we set  $\tau_{\pm} = (1 \pm \Delta)$  and  $q_{+}(r) = q(\tau_{+}r)$ ,  $q_{-}(r) = q(\tau_{-}r)$ . Also, we will assume  $\xi(t) = 1/2$ , since other choices of  $\xi(\cdot)$  merely amounts to a time reparametrization.

Before analyzing our model, we introduce the function space and space of probability measures we will work on. We equip the set  $[0,\infty]$  with a metric  $\bar{d}$ , where  $\bar{d}(x,y) = |1/(1+x)-1/(1+y)|$  for any  $x,y \in [0,\infty]$ . Then  $([0,\infty],\bar{d})$  is a compact metric space, and we will still denote it by  $[0,\infty]$  for simplicity in notations. We denote  $C_b([0,\infty])$  to be the set of bounded continuous functions on  $[0,\infty]$ , where continuity is defined using the topology generated by  $\bar{d}$ . More explicitly, we have isomorphism

$$C_b([0,\infty]) \simeq \{ f \in C([0,\infty)) : \exists f(+\infty) \equiv \lim_{r \to +\infty} f(r), \sup_{r \in [0,\infty]} f(r) < \infty \}.$$
(8.2)

Because of condition S2 and S3, we have  $q, q' \in C_b([0, \infty])$ .

Let  $\mathscr{P}([0,\infty])$  be the set of probability measures on  $[0,\infty]$ . Due to Prokhorov's theorem, there exists a complete metric  $\bar{d}_{\mathscr{P}}$  on  $\mathscr{P}([0,\infty])$  equivalent to the topology of weak convergence, so that  $(\mathscr{P}([0,\infty]),\bar{d}_{\mathscr{P}})$  is a compact metric space. In this section, we will denote by  $\overline{\mathscr{P}}=\mathscr{P}([0,\infty])$ .

#### 8.1 Statics

Since the distribution of x is invariant under rotations for each of the two classes, so are the functions

$$V(\mathbf{w}) = v(\|\mathbf{w}\|_2), \qquad U(\mathbf{w}_1, \mathbf{w}_2) = u_0(\|\mathbf{w}_1\|_2, \|\mathbf{w}_2\|_2, \langle \mathbf{w}_1, \mathbf{w}_2 \rangle).$$
 (8.3)

These take the form

$$v(r) = -\frac{1}{2}q(\tau_{+}r) + \frac{1}{2}q(\tau_{-}r), \quad q(t) = \mathbb{E}\{\sigma(tG)\}$$
(8.4)

$$u_0(r_1, r_2, r_1 r_2 \cos \alpha) = \frac{1}{2} \mathbb{E} \{ \sigma(\tau_+ r_1 G_1) \sigma(\tau_+ r_2 G_2) \} + \frac{1}{2} \mathbb{E} \{ \sigma(\tau_- r_1 G_1) \sigma(\tau_+ r_2 G_2) \},$$
 (8.5)

where expectations are with respect to standard normals  $G, G_1, G_2 \sim N(0, 1)$ , with  $(G_1, G_2)$  jointly Gaussian and  $\mathbb{E}\{G_1G_2\} = \cos \alpha$ .

In order to minimize  $R(\rho)$ , it is sufficient to restrict ourselves to distributions that are invariant under rotations. Indeed, for any probability distribution  $\rho$  on  $\mathbb{R}^d$ , we can define its symmetrization by letting, for any Borel set  $Q \subseteq \mathbb{R}^d$ ,

$$\rho_s(Q) \equiv \int \rho(\mathbf{R} Q) \ \mu_{\text{Haar}}(\mathrm{d}\mathbf{R}), \tag{8.6}$$

where  $\mu_{\text{Haar}}$  is the Haar measure over the group of orthogonal rotations. Since  $\rho \mapsto R(\rho)$  is convex,  $R(\rho_s) \leq R(\rho)$ .

We therefore restrict ourselves to  $\rho$ 's that are invariant under rotations. In other words, under  $\rho$ , the vector  $\boldsymbol{w}$  is uniformly random conditional on  $\|\boldsymbol{w}\|_2$ . We denote by  $\overline{\rho}$  the probability distribution of  $\|\boldsymbol{w}\|_2$  when  $\boldsymbol{w} \sim \rho$  and we let  $\overline{R}_d(\overline{\rho})$  denote the resulting risk. We then have

$$\overline{R}_d(\overline{\rho}) = 1 + 2 \int v(r) \,\overline{\rho}(\mathrm{d}r) + \int u_d(r_1, r_2) \,\overline{\rho}(\mathrm{d}r_1) \,\overline{\rho}(\mathrm{d}r_2) \,, \tag{8.7}$$

$$u_d(r_1, r_2) = \mathbb{E}[u_0(r_1, r_2, r_1 r_2 \cos \Theta)]. \tag{8.8}$$

where  $\Theta \sim (1/Z_d) \sin^{d-2} \theta \cdot \mathbf{1} \{ \theta \in [0, \pi] \} d\theta$ .

As  $d \to \infty$ , we have  $\lim_{d \to \infty} u_d(r_1, r_2) = u_\infty(r_1, r_2)$  (uniformly over compact sets), with

$$u_{\infty}(r_1, r_2) = \frac{1}{2} \Big[ q(\tau_+ r_1) q(\tau_+ r_2) + q(\tau_- r_1) q(\tau_+ r_2) \Big], \tag{8.9}$$

and the risk function converges to

$$\overline{R}_{\infty}(\overline{\rho}) = \frac{1}{2} \left( 1 - \int q(\tau_{+}r) \,\overline{\rho}(\mathrm{d}r) \right)^{2} + \frac{1}{2} \left( 1 + \int q(\tau_{-}r) \,\overline{\rho}(\mathrm{d}r) \right)^{2} \,. \tag{8.10}$$

We also define

$$\psi_d(r;\overline{\rho}) = v(r) + \int u_d(r,r')\,\overline{\rho}(\mathrm{d}r')\,. \tag{8.11}$$

For  $d = \infty$ , we have the simpler expression

$$\psi_{\infty}(r;\overline{\rho}) = \lambda_{+}(\overline{\rho}) \cdot q_{+}(r) + \lambda_{-}(\overline{\rho}) \cdot q_{-}(r), \tag{8.12}$$

$$\lambda_{+}(\overline{\rho}) = \frac{1}{2} [\langle q_{+}, \overline{\rho} \rangle - 1], \tag{8.13}$$

$$\lambda_{-}(\overline{\rho}) = \frac{1}{2} [\langle q_{-}, \overline{\rho} \rangle + 1]. \tag{8.14}$$

The following theorem provides a characterization of global minimizers of  $\overline{R}_d(\overline{\rho})$ .

**Proposition 4** (Lemma 1 in the main text). For any  $d \leq \infty$ , define

$$\psi_d(r;\overline{\rho}) \equiv v(r) + \int u_d(r,r') \,\overline{\rho}(\mathrm{d}r') \,. \tag{8.15}$$

Then

- 1.  $\overline{\rho}_*$  is a global minimizer of  $\overline{R}_d(\overline{\rho})$  if and only if  $\operatorname{supp}(\overline{\rho}_*) \subseteq \operatorname{arg\,min}_r \psi_d(r; \overline{\rho}_*)$ .
- 2. In particular,  $\overline{\rho}_* = \delta_{r_*}$  is a global minimizer or  $\overline{R}_d(\overline{\rho})$  if and only if  $v(r) + u_d(r, r_*) \ge v(r_*) + u(r_*, r_*)$  for all r.

*Proof.* Point 1 is essentially a special case of the second part of Proposition 1 in the main text (cf. Eq. (6.7)) and follows by the same argument. Point 2 is follows by taking  $\overline{\rho}_* = \delta_{r_*}$ .

Given the last result, it is interesting to understand whether the optimal radial distribution  $\overline{\rho}_*$  is a single point mass or not. Under the ansatz  $\overline{\rho} = \delta_r$  (a single point mass at radius r) we obtain an effective risk  $\overline{R}_d^{(1)}(r) \equiv \overline{R}_d(\delta_r)$  defined by  $\overline{R}_d^{(1)}(r) = 1 + 2v(r) + u_d(r,r)$ , which is plotted in Figure 11.6 for the case of our running example (8.1), and  $\Delta = 0.4$ .

Let  $r_* = r_*(\Delta, d)$  be the minimizer of  $\overline{R}_d^{(1)}(r)$ , and define, for  $d \leq \infty$ ,

$$\Delta_d = \sup \left\{ \Delta : \ v(r) + u_d(r, r_*) \ge v(r_*) + u_d(r_*, r_*), \ \forall r \ge 0 \right\}.$$
 (8.16)

In the case  $d = \infty$ , the minimization problem simplifies further. Either the minimum risk is 0, or it is achieved at a point mass  $\overline{\rho}_* = \delta_{r_*}$ .

**Theorem 8.** Consider  $d = \infty$ . Recall  $\overline{\mathscr{P}} = \mathscr{P}([0,\infty])$ . In this case  $\Delta_{\infty}$  defined as per Eq. (8.16) is such that  $\Delta_{\infty} \in (0,1)$ . Further

- 1. For  $\Delta < \Delta_{\infty}$ ,  $\inf_{\overline{\rho} \in \overline{\mathscr{P}}} \overline{R}_{\infty}(\overline{\rho}) > 0$  and the unique global minimizer of risk function  $\overline{R}_{\infty}(\overline{\rho})$  is a point mass located at some  $r_*(\Delta) \in (0,\infty)$ .
- 2. For  $\Delta \geq \Delta_{\infty}$ , all global minimizers of risk function  $\overline{R}_{\infty}(\overline{\rho})$  have risk zero, and there exists a global minimizer that has compact support bounded away from 0.

Proof of Theorem 8. Recall the definitions  $q_+(r) = q(\tau_+ r)$  and  $q_-(r) = q(\tau_- r)$ . Further, we define the set  $\Gamma \subseteq [0,1]$  by

$$\Gamma = \{ \Delta : \exists r \in (0, +\infty), \text{ s.t., } q_+(r) \ge 1 \text{ and } q_-(r) \le -1 \}.$$
 (8.17)

According to condition S3, for  $\Delta = 1$ , we have  $q_{-}(r) = q(0) < -1$  and  $q_{+}(+\infty) = q(+\infty) > +1$ . Since q is continuous, it is easy to see that there exists an  $\varepsilon > 0$ , such that  $[1 - \varepsilon, 1] \subseteq \Gamma$ . Further, for  $\Delta = 0$  we have  $q_{+}(r) = q_{-}(r)$ . By continuity, there exists an  $\varepsilon > 0$ , such that  $[0, \varepsilon] \in [0, 1] \setminus \Gamma$ . Since q is an increasing function, we have

$$\Gamma = [\Delta_{\infty}, 1], \qquad \Delta_{\infty} = \inf_{\Delta \in \Gamma} \Delta.$$
 (8.18)

By the remarks above, we have  $0 < \Delta_{\infty} < 1$ . Notice that this definition does not coincide with the one in Eq. (8.16). However, the proof below (together with Proposition 4) implies that the two definitions actually coincide.

Part (1):  $\Delta < \Delta_{\infty}$ .

Step 1. Prove that  $\inf_{\overline{\rho} \in \overline{\mathscr{P}}} \overline{R}_{\infty}(\overline{\rho}) > 0$  as  $\Delta < \Delta_{\infty}$ .

First, we consider the optimization problem

$$f_* \equiv \sup_{\overline{\rho} \in \overline{\mathscr{P}}} \left\{ \langle q_+, \overline{\rho} \rangle - 1 \quad \text{s.t.} \quad \langle q_-, \overline{\rho} \rangle \le -1 \right\}.$$
 (8.19)

We claim that, for  $\Delta < \Delta_{\infty}$  we have  $f_* < 0$ . Indeed, for any  $\lambda \in [0, +\infty)$ , we have the following upper bound

$$f_* \le \sup_{\overline{\rho} \in \overline{\mathscr{P}}} \{ L(\overline{\rho}, \lambda) \equiv \langle q_+, \overline{\rho} \rangle - 1 - \lambda (\langle q_-, \overline{\rho} \rangle + 1) \}.$$
 (8.20)

Since  $q_+ - \lambda q_- \in C_b([0, +\infty])$ , then  $L(\cdot, \lambda)$  is continuous in  $\overline{\rho}$  in weak topology. By the compactness of  $\overline{\mathscr{P}}$ , the supremum of  $L(\cdot, \lambda)$  is attained by some  $\overline{\rho}_{\lambda} \in \overline{\mathscr{P}}$ . This  $\overline{\rho}_{\lambda}$  should satisfy

$$\operatorname{supp}(\overline{\rho}_{\lambda}) \subseteq \operatorname{argmax}_{r \in [0,+\infty]} \{q_{+}(r) - \lambda q_{-}(r)\}.$$

Let  $h(r) \equiv q_{+}(r) - \lambda q_{-}(r)$ . Note the supremum of h should either satisfy

$$h'(r) = q'_{+}(r) - \lambda q'_{-}(r) = 0, \tag{8.21}$$

for  $r \in (0, \infty)$ , or the supremum should be attained at the boundary 0 or  $+\infty$ . According to condition S4,  $[q'_{-}(r)/q'_{+}(r)]' > 0$  for  $r \in (0, \infty)$ , the equation (8.21) has at most one solution  $r_* \in (0, \infty)$ .

Assume that there exists  $r_* \in (0, \infty)$  such that  $h'(r_*) = 0$ . Then we have h'(r) > 0 for  $0 < r < r_*$ , and h'(r) < 0 for  $r_* < r < +\infty$ , whence  $\operatorname{supp}(\overline{\rho}_{\lambda}) = \{r_*\}$ . If h'(r) = 0 does not have a solution in  $(0, \infty)$ , the only supremum of h(r) could be achieved at 0 or  $+\infty$ . Therefore,  $\operatorname{supp}(\overline{\rho}_{\lambda}) = \{0\}$  or  $\operatorname{supp}(\overline{\rho}_{\lambda}) = \{+\infty\}$ . This concludes that, for any  $\lambda \in [0, +\infty)$ ,  $\operatorname{sup}_{\overline{\rho} \in \overline{\mathscr{P}}} L(\overline{\rho}, \lambda)$  is achieved by a point mass. Therefore, we have

$$f_* \le \inf_{\lambda \in [0, +\infty)} \sup_{r \in [0, +\infty]} \{q_+(r) - 1 - \lambda (q_-(r) + 1)\} = q_+(q_-^{-1}(-1)) - 1.$$

For  $\Delta < \Delta_{\infty}$ , the right hand side of the above inequality is less than 0. Therefore, we cannot have a probability distribution  $\overline{\rho}$  such that  $\langle q_+, \overline{\rho} \rangle = 1$  and  $\langle q_-, \overline{\rho} \rangle = -1$ . The infimum of the risk cannot be 0.

### Step 2. Show that the global minimizer should be a delta function for $\Delta < \Delta_{\infty}$ .

According to Proposition 1, the global minimizer  $\overline{\rho}_* \in \overline{\mathscr{P}}$  should satisfy

$$\operatorname{supp}(\overline{\rho}_*) \subseteq \arg\min_{r \in [0,+\infty]} \psi_{\infty}(r; \overline{\rho}_*),$$

with  $\psi_{\infty}$  given in Eq. (8.12).

As proved in the last step, as  $\Delta < \Delta_{\infty}$ , we cannot have both  $\lambda_{+}(\overline{\rho}_{*}) = 0$  and  $\lambda_{-}(\overline{\rho}_{*}) = 0$ . The argument given above also implies that  $\psi_{\infty}(r; \overline{\rho}_{*})$  is minimized at a unique point, and hence the support of  $\overline{\rho}_{*}$  should be a single point. This proves the first part of the theorem.

### Part (2): $\Delta \geq \Delta_{\infty}$ .

For  $\Delta \geq \Delta_{\infty}$ , there exists r > 0, such that  $q(\tau_+ r) \geq 1$ , and  $q(\tau_- r) \leq -1$ . Therefore, there exists  $r_* > 0$  such that  $q(\tau_+ r_*) - 1 = -1 - q(\tau_- r_*) = \varepsilon_* \geq 0$ . Consider the following probability measure on  $[0, +\infty]$ ,

$$\overline{\rho}_* = \frac{1}{1 + \varepsilon_*} \delta_{r_*} + \frac{\varepsilon_*}{(1 + \varepsilon_*)(q(+\infty) - q(0))} [q(+\infty)\delta_0 - q(0)\delta_{+\infty}].$$

It can be checked that  $\overline{R}_{\infty}(\overline{\rho}_*) = 0$ .

We would like to show further that there exists a global minimizer that is compactly supported. We construct this global minimizer as following. First, define

$$r_0 = \inf\{r : q_-(r) \ge -1\}.$$

Then we know that  $q_-(r_0)=-1$  and  $q_+(r_0)\geq 1$ . Now for any  $0\leq r\leq r_0$ , define  $u(r)=q_-^{-1}(-2-q_-(r))$ . According to condition S3, we have  $-1<[q(0)+q(+\infty)]/2<1$ , then u(r) is

well defined on  $[0, r_0]$ . It is easy to see that  $u(r_0) = r_0$ , and  $[q_-(r) + q_-(u(r))]/2 = -1$  for any  $0 \le r \le r_0$ . Now we consider the function  $z(r) = [q_+(r) + q_+(u(r))]/2 - 1$ . Note that  $z(r_0) > 0$ , and  $z(0) \le [q(0) + q(\infty)]/2 - 1 < 0$ . Therefore, there exists  $r_*$  satisfying  $0 < r_* \le r_0$  such that  $z(r_*) = 0$ . Consider the following probability measure on  $(0, +\infty)$ ,

$$\overline{\rho}_* = \frac{1}{2} [\delta_{r_*} + \delta_{u(r_*)}].$$

It is easy to see that  $\overline{R}_{\infty}(\overline{\rho}_*) = 0$ .

# 8.2 Dynamics: Fixed points

We specialize the general evolution (7.1) to the present case. Assuming  $\rho_0$  to be spherically symmetric, then  $\rho_t$  is spherically symmetric for any  $t \geq 0$ . We let  $\overline{\rho}_t$  denote the distribution of  $\|\boldsymbol{w}\|_2$  when  $\boldsymbol{w} \sim \rho_t$ . This satisfies the following PDE:

$$\partial_t \overline{\rho}_t(r) = 2\xi(t)\partial_r \left[ \overline{\rho}_t(r)\partial_r \psi_d(r; \overline{\rho}_t) \right]. \tag{8.22}$$

We will view this as an evolution in the space of probability distribution on the completed half-line  $\mathcal{P}([0,\infty])$ .

In analogy with Proposition 2, we can prove the following characterization of fixed points.

**Proposition 5.** A distribution  $\overline{\rho} \in \mathscr{P}([0,\infty])$  is a fixed point of the PDE (8.22) if and only if

$$\operatorname{supp}(\overline{\rho}) \subseteq \{ r \in [0, \infty] : \partial_r \psi_d(r; \overline{\rho}) = 0 \}. \tag{8.23}$$

Notice, in particular, global minimizers of  $\overline{R}_d(\overline{\rho})$  are fixed points of this evolution, but not vice-versa. The next result classifies fixed points.

**Theorem 9.** Consider  $d = \infty$  and recall the definition of  $\lambda_{+}(\overline{\rho})$  and  $\lambda_{-}(\overline{\rho})$  given by Eqs. (8.13) and (8.14). Then the fixed points of the PDE (8.22) (i.e. the probability measures  $\overline{\rho} \in \mathscr{P}([0,\infty])$  satisfying (8.23)) are of one of the following types

- (a) A fixed point with zero risk.
- (b) A point mass  $\overline{\rho}_{r_*} = \delta_{r_*}$  at some location  $r_* \notin \{0, +\infty\}$ , but not of type (a).
- (c) A mixture of the type  $\overline{\rho} = a_0 \delta_0 + a_\infty \delta_{+\infty} + a \delta_{r_*}$ , but not of type (a) or (b).

For  $\Delta < \Delta_{\infty}$ , the PDE has a unique fixed point of type (b), with  $\lambda_{+}(\overline{\rho}_{*}) < 0$  and  $\lambda_{-}(\overline{\rho}_{*}) > 0$ ; it has no type-(a) fixed points; it has possibly fixed points of type (c).

For  $\Delta > \Delta_{\infty}$ , the PDE has some fixed points of type (b), with  $\lambda_{+}(\overline{\rho}_{*}) > 0$  and  $\lambda_{-}(\overline{\rho}_{*}) < 0$ ; it also has some type-(a) fixed points; it has possibly fixed points of type (c).

For  $\Delta = \Delta_{\infty}$ , the PDE has a unique fixed point of type (a) which is also a delta function at some location  $r_*$ , and no type (b) fixed points; it has possibly fixed points of type (c).

*Proof.* We use the characterization of fixed points in Proposition 5. Recall that  $\psi_{\infty}(r; \overline{\rho}_*)$  is defined as in Equation (8.12). The derivative  $\partial_r \psi_{\infty}(r; \overline{\rho})$  gives

$$\partial_r \psi_{\infty}(r; \overline{\rho}) = \lambda_+(\overline{\rho}) q'_+(r) + \lambda_-(\overline{\rho}) q'_-(r). \tag{8.24}$$

If a fixed point has  $\lambda_+(\overline{\rho}_*) = \lambda_-(\overline{\rho}_*) = 0$ , then  $\overline{R}_{\infty}(\overline{\rho}_*) = 0$ . This is type-(a) fixed point. Consider then the case  $(\lambda_+(\overline{\rho}_*), \lambda_-(\overline{\rho}_*)) \neq (0,0)$ . For the same reason as in the proof of Theorem 8, we conclude that  $\partial_r \psi_{\infty}(r; \overline{\rho}_*)$  has at most three zeros, two of which are located at 0 and  $+\infty$ . This proves that all fixed points are of type (a), (b) or (c).

We already proved in Theorem 8 that, for  $\Delta < \Delta_{\infty}$ ,  $\inf_{\overline{\rho}} \overline{R}_{\infty}(\overline{\rho}) > 0$ . Therefore, for  $\Delta < \Delta_{\infty}$ , there is no type (a) fixed points.

We next prove that, as  $\Delta < \Delta_{\infty}$ , fixed point of type (b) is always unique. The location of the delta fixed point should satisfy

$$\partial_r \psi_{\infty}(r_*; \delta_{r_*}) = [q'_+(r_*)(q_+(r_*) - 1) + q'_-(r_*)(q_-(r_*) + 1)]/2 = 0. \tag{8.25}$$

Note that  $\partial_r \psi_\infty(r_*; \delta_{r_*}) < 0$  for r > 0 small enough, and  $\partial_r \psi_\infty(r_*; \delta_{r_*}) > 0$  for r large enough, whence this equation has at least one solution  $r_* \in (0, \infty)$ . In order to prove that it has a unique solution in  $(0, +\infty)$ , define  $r_+ \equiv \inf\{r : q_+(r) \ge 1\}$  and  $r_- \equiv \inf\{r : q_-(r) \ge -1\}$ . Note that  $q'_+(r_*), q'_-(r_*) > 0$  and that, in order to satisfy Eq. (8.25), the terms  $\lambda_+(\delta_{r_*}) = 1/2 \cdot (q_+(r_*) - 1)$  and  $\lambda_-(\delta_{r_*}) = 1/2 \cdot (q_-(r_*) + 1)$  must have opposite signs. For  $\Delta < \Delta_\infty$ , we must have  $\lambda_+(\delta_{r_*}) < 0$  and  $\lambda_-(\delta_{r_*}) > 0$ , and all stationary points should be within  $[r_-, r_+]$ . Note that  $q'_-(r)/q'_+(r)$  is strictly increasing, and  $[1 - q_+(r)]/[1 + q_-(r)]$  is decreasing on  $[r_-, r_+]$ . Therefore, the fixed point of type  $\delta_{r_*}$  with  $r_* \in (0, \infty)$  is unique.

For  $\Delta > \Delta_{\infty}$ , we must have  $\lambda_{+}(\overline{\rho}_{*}) > 0$  and  $\lambda_{-}(\overline{\rho}_{*}) < 0$ , and all solutions should be within  $[r_{+}, r_{-}]$ . There could possibly be multiple fixed points of type  $\delta_{r_{*}}$  with  $r_{*} \in [r_{+}, r_{-}]$ .

If  $\Delta = \Delta_{\infty}$ , it is easy to see that,  $\overline{\rho}_* = \delta_{r_*}$  at some  $r_* \in (0, \infty)$  is the unique fixed point with zero risk, and the unique fixed point as a point mass.

# 8.3 Dynamics: Convergence to global minimum for $d = \infty$

In this section, denote  $\mathscr{P}_{\text{good}}$  to be

$$\mathscr{P}_{\text{good}} = \{ \overline{\rho}_0 \in \mathscr{P}((0, \infty)) : \overline{R}_{\infty}(\overline{\rho}_0) < 1, \overline{\rho}_0 \text{ has bounded density on } (0, \infty) \}. \tag{8.26}$$

We then prove that the  $d=\infty$  dynamics converges to a global minimizer from any initialization in  $\mathscr{P}_{\text{good}}$ .

**Theorem 10.** Consider the PDE (8.22) for  $d = \infty$ , with initialization  $\overline{\rho}_0 \in \mathscr{P}_{good}$ . It has a unique solution  $(\overline{\rho}_t)_{t \geq 0}$ , such that

$$\lim_{t\to +\infty} \overline{R}_{\infty}(\overline{\rho}_t) = \inf_{\overline{\rho}\in \overline{\mathscr{P}}} \overline{R}_{\infty}(\overline{\rho})\,.$$

*Proof.* Without loss of generality, we assume  $\xi(t) = 1/2$ . First we show the existence and uniqueness of solution of the PDE.

#### Step 1. Existence and uniqueness of solution. Mass $\overline{\rho}_t((0,\infty)) = 1$ for all t.

According to conditions S1 - S3, q(r), q'(r), and q''(r) are uniformly bounded on  $[0, \infty]$ . Recall that

$$v(r) = 1/2 \cdot [q_{-}(r) - q_{+}(r)],$$
  

$$u_{\infty}(r_{1}, r_{2}) = 1/2 \cdot [q_{+}(r_{1})q_{+}(r_{2}) + q_{-}(r_{1})q_{-}(r_{2})].$$

Hence v'(r),  $\partial_1 u_{\infty}(r_1, r_2)$ , v''(r),  $\partial_{11}^2 u_{\infty}(r_1, r_2)$ ,  $\partial_{12}^2 u_{\infty}(r_1, r_2)$  are uniformly bounded. Recall we further assumed  $\xi(t) \equiv 1/2$ . Therefore, conditions A1 and A3 are satisfied with D = 1, V = v, and

U = u. By Remark 7.1, there is the existence and uniqueness of solution of PDE (8.22) for  $d = \infty$ . Denote this solution to be  $(\overline{\rho}_t)_{t\geq 0}$ .

Recall the formula of  $\partial_r \psi_{\infty}(r; \overline{\rho})$  given in Equation (8.24), it is easy to see that the assumption of Lemma 7.9 is satisfied with d=1 and  $\Psi=\psi_{\infty}$ . Hence, we have  $\overline{\rho}_t((0,\infty))=1$  for any  $t<\infty$ .

### Step 2. Classify the limiting set $S_*$ .

Recall the definition of  $(\mathscr{P}([0,+\infty]), \bar{d}_{\mathscr{P}})$  at the beginning of Section 8. Since  $(\mathscr{P}([0,+\infty]), \bar{d}_{\mathscr{P}})$  is a compact metric space, and  $(\overline{\rho}_t)_{t\geq 0}$  is a continuous curve in this space, then there exists a subsequence  $(t_k)_{k\geq 1}$  of times, such that  $(\overline{\rho}_{t_k})_{k\geq 1}$  converges in metric  $\bar{d}_{\mathscr{P}}$  to a probability distribution  $\overline{\rho}_* \in \mathscr{P}([0,+\infty])$ .

Analogously to Proposition 2 (using Eq. (8.22)), we have

$$\partial_t \overline{R}_{\infty}(\overline{\rho}_t) = -\int [\partial_r \psi_{\infty}(r; \overline{\rho}_t)]^2 \overline{\rho}_t(\mathrm{d}r).$$

Since  $\overline{R}_{\infty}(\overline{\rho}_t) \geq 0$ , we have

$$\lim_{t \to +\infty} \int [\partial_r \psi_\infty(r; \overline{\rho}_t)]^2 \, \overline{\rho}_t(\mathrm{d}r) = 0.$$

Recall the definition of  $\lambda_{+}(\overline{\rho})$  and  $\lambda_{-}(\overline{\rho})$  given by Eq. (8.13) and (8.14). Since  $q \in C_b([0,\infty])$ , we have

$$\lim_{k \to \infty} \lambda_{+}(\overline{\rho}_{t_{k}}) = \lambda_{+}(\overline{\rho}_{*}), \quad \lim_{k \to \infty} \lambda_{-}(\overline{\rho}_{t_{k}}) = \lambda_{-}(\overline{\rho}_{*}). \tag{8.27}$$

Note  $\partial_r \psi_\infty(r; \overline{\rho})$  is given by Eq. (8.24), and  $q' \in C_b([0, +\infty])$ , hence

$$\lim_{k \to +\infty} \langle [\partial_r \psi_{\infty}(\,\cdot\,; \overline{\rho}_{t_k})]^2, \overline{\rho}_{t_k} \rangle = \langle [\partial_r \psi_{\infty}(\,\cdot\,; \overline{\rho}_*)]^2, \overline{\rho}_* \rangle,$$

which implies

$$\langle [\partial_r \psi_\infty(\,\cdot\,;\overline{\rho}_*)]^2,\overline{\rho}_* \rangle = 0.$$

In other words, any limiting point  $\overline{\rho}_*$  of the PDE is a fixed point of the PDE (8.22).

Note  $\overline{R}_{\infty}(\overline{\rho}) = 1/2 \cdot [\lambda_{+}(\overline{\rho})^{2} + \lambda_{-}(\overline{\rho})^{2}]$ , we have

$$\lim_{k \to +\infty} \overline{R}_{\infty}(\overline{\rho}_{t_k}) = \overline{R}_{\infty}(\overline{\rho}_*).$$

Note  $\overline{R}_{\infty}(\overline{\rho}_t)$  is decreasing with t, hence

$$\lim_{t \to +\infty} \overline{R}_{\infty}(\overline{\rho}_t) = \overline{R}_{\infty}(\overline{\rho}_*).$$

Let  $S_* = S_*(\overline{\rho}_0)$  be the set of all limiting points of the  $(\overline{\rho}_t)_{t\geq 0}$ ,

$$\mathcal{S}_* = \{ \overline{\rho}_* \in \mathscr{P}([0,\infty]) : \exists (t_k)_{k \ge 1}, \lim_{k \to \infty} t_k = +\infty, s.t., \lim_{k \to \infty} \overline{d}_{\mathscr{P}}(\overline{\rho}_*, \overline{\rho}_{t_k}) = 0 \}.$$

Due to Lemma 7.10,  $S_*$  is a connected compact set. Since  $\overline{R}_{\infty}(\overline{\rho}_t)$  is decreasing as t increases, we have  $\overline{R}_{\infty}(\overline{\rho}_*) \equiv \overline{R}_*$  is a constant for all  $\overline{\rho}_* \in S_*$ . Since we assumed  $\overline{R}_{\infty}(\overline{\rho}_0) < 1$ , and  $\overline{R}_{\infty}(\overline{\rho}_t)$  is decreasing in t, we have  $\overline{R}_* < 1$ .

Let  $\overline{\rho}_*$  be a fixed point of PDE such that  $\lambda_+(\overline{\rho}_*) \geq 0, \lambda_-(\overline{\rho}_*) \geq 0$  or  $\lambda_+(\overline{\rho}_*) \leq 0, \lambda_-(\overline{\rho}_*) \leq 0$  but not both  $\lambda_+(\overline{\rho}_*)$  and  $\lambda_-(\overline{\rho}_*)$  equal 0. In this case, according to Eq. (8.24),  $\partial_r \psi_\infty(r; \overline{\rho}_*)$  must

be strictly increasing or strictly decreasing in r. Since  $\operatorname{supp}(\overline{\rho}_*) \subseteq \{r \in [0, \infty] : \partial_r \psi_\infty(r; \overline{\rho}_*) = 0\}$ ,  $\overline{\rho}_*$  must be a combination of two delta functions located at 0 and  $+\infty$ , i.e.,  $\overline{\rho}_* = a_0 \delta_0 + (1 - a_0) \delta_\infty$ . But for a fixed point of this type, it is easy to see that  $\overline{R}_\infty(\overline{\rho}_*) \geq 1$ . Such fixed points  $\overline{\rho}_*$  cannot be one of the limiting points of the PDE since  $\overline{R}_\infty(\overline{\rho}_0) < 1$ .

Let L be a mapping  $L: \mathscr{P}([0,+\infty]) \to \mathbb{R}^2$ ,  $\overline{\rho} \mapsto (\lambda_+(\overline{\rho}), \lambda_-(\overline{\rho}))$ . The above argument implies that for any  $\overline{\rho}_0 \in \mathscr{P}_{\text{good}}$ , we have

$$L(\mathcal{S}_*(\overline{\rho}_0)) \cap (\{(\lambda_+, \lambda_-) : \lambda_+ \ge 0, \lambda_- \ge 0, \text{ or } \lambda_+ \le 0, \lambda_- \le 0\} \setminus \{(0, 0)\}) = \emptyset.$$

Since  $S_*$  is a connected set,  $L(S_*)$  should also be a connected set. Further notice that  $\overline{R}_{\infty}(\overline{\rho}_*) = 1/2 \cdot [\lambda_+(\overline{\rho}_*)^2 + \lambda_-(\overline{\rho}_*)^2]$ , and  $\overline{R}_{\infty}(\overline{\rho}_1) = \overline{R}_{\infty}(\overline{\rho}_2)$  for any  $\overline{\rho}_1, \overline{\rho}_2 \in S_*$ . Therefore, we can only have  $L(S_*) \subseteq \mathcal{P}_2 \equiv \{(\lambda_+, \lambda_-) : \lambda_+ > 0, \lambda_- < 0\}$ , or  $L(S_*) \subseteq \mathcal{P}_1 \equiv \{(\lambda_+, \lambda_-) : \lambda_+ < 0, \lambda_- > 0\}$ , or  $L(S_*) = \{(0, 0)\}$ .

# Step 3. Finish the proof using two claims.

We make the following two claims.

Claim (1). If  $L(S_*) \subseteq \mathcal{P}_1$ , then for any  $\overline{\rho}_* \in S_*$ , we have  $\overline{\rho}_*((0,\infty)) = 1$ .

Claim (2). We cannot have  $L(\mathcal{S}_*) \subseteq \mathcal{P}_2$ .

Here we assume these two claims hold, and use them to prove our results. For  $\Delta < \Delta_{\infty}$ , we proved in Theorem 9 that, there is not a fixed point such that  $L(\overline{\rho}_*) = (0,0)$ . Therefore, we cannot have  $L(\mathcal{S}_*) = \{(0,0)\}$ . Due to Claim (2), we cannot have  $L(\mathcal{S}_*) \subseteq \mathcal{P}_2$ . Hence, we must have  $L(\mathcal{S}_*) \subseteq \mathcal{P}_1$ . According to Theorem 9, for  $\Delta < \Delta_{\infty}$ , the only fixed point of PDE with  $\overline{\rho}_*((0,\infty)) = 1$  is a point mass at some location  $r_*$ . Furthermore, this delta function fixed point is unique and is also the global minimizer of the risk. Therefore, we conclude that, as  $\Delta < \Delta_{\infty}$ , the PDE will converge to this global minimizer.

For  $\Delta \geq \Delta_{\infty}$ , according to Claim (1), if  $\overline{\rho}_*$  is a limiting point such that  $L(\overline{\rho}_*) \in \mathcal{P}_1$ , then  $\overline{\rho}_*((0,\infty)) = 1$ . According to Theorem 9, a fixed point  $\overline{\rho}_*$  with  $\overline{\rho}_*((0,\infty)) = 1$  and  $L(\overline{\rho}_*) \neq (0,0)$  must be a point mass at some location  $r_*$ , with  $L(\overline{\rho}_*) \in \mathcal{P}_2$ . Therefore, we cannot have  $L(\mathcal{S}_*) \subseteq \mathcal{P}_1$ . Claim (2) also tells us that we cannot have  $L(\mathcal{S}_*) \subseteq \mathcal{P}_2$ . Hence, we must have  $L(\mathcal{S}_*) = \{(0,0)\}$ . In this case, all the points in the set  $\mathcal{S}_*$  have risk 0. Therefore, we conclude that, as  $\Delta \geq \Delta_{\infty}$ , the PDE will converge to some limiting set with risk 0.

#### Step 4. Proof of the two claims.

We are left with the task of proving the two claims above. Before that, we introduce some useful notations. Recall  $Z(r)=q'_-(r)/q'_+(r)$  for  $r\in(0,+\infty)$ . According to condition S4, Z'(r)>0 for  $r\in(0,+\infty)$ . This implies that  $Z(0+)\equiv Z_0\geq 0$  and  $Z(+\infty)\equiv Z_\infty\leq\infty$  exist. We rewrite  $\partial_r\psi_\infty(r;\overline{\rho})$  as

$$\partial_r \psi_{\infty}(r; \overline{\rho}) = \lambda_+(\overline{\rho}) q'_+(r) + \lambda_-(\overline{\rho}) q'_-(r) = \lambda_-(\overline{\rho}) q'_+(r) [\lambda_+(\overline{\rho})/\lambda_-(\overline{\rho}) + Z(r)]. \tag{8.28}$$

Proof of Claim (1). If  $L(\mathcal{S}_*) \subseteq \mathcal{P}_1$ , then for any  $\overline{\rho}_* \in \mathcal{S}_*$ , we have  $\overline{\rho}_*(\{0,\infty\}) = 0$ .

Assume  $L(S_*) \subseteq \mathcal{P}_1$ . Then, we must have  $L(S_*) \subseteq \mathcal{P}_1 \cap \{(\lambda_+, \lambda_-) : Z_0 < -\lambda_+/\lambda_- < Z_\infty\}$ . Otherwise suppose there exists  $\overline{\rho}_* \in S_*$ , such that  $-\lambda_+(\overline{\rho}_*)/\lambda_-(\overline{\rho}_*) \ge Z_\infty$  or  $-\lambda_+(\overline{\rho}_*)/\lambda_-(\overline{\rho}_*) \le Z_0$ , according to Eq. (8.28),  $\psi_\infty(r; \overline{\rho}_*)$  must be strictly increasing or strictly decreasing in r. Since  $\sup(\overline{\rho}_*) \subseteq \{r \in [0,\infty] : \partial_r \psi_\infty(r; \overline{\rho}_*) = 0\}$ , then  $\overline{\rho}_*$  must be a combination of two delta functions

located at 0 and  $+\infty$ . But such  $\overline{\rho}_*$  must have  $\overline{R}_{\infty}(\overline{\rho}_*) \geq 1$ , and thus  $\overline{\rho}_*$  cannot be a limiting point of the PDE. Hence the claim that  $L(\mathcal{S}_*) \subseteq \mathcal{P}_1 \cap \{(\lambda_+, \lambda_-) : Z_0 < -\lambda_+/\lambda_- < Z_{\infty}\}$  holds.

Since  $S_*$  is a compact set, and L is a continuous map, then  $L(S_*)$  is a compact set. Therefore, there must exist  $\varepsilon_0 > 0$ , so that for any  $\overline{\rho}_* \in S_*$ , we have  $Z_0 + 3\varepsilon_0 < -\lambda_+(\overline{\rho}_*)/\lambda_-(\overline{\rho}_*) < Z_\infty - 3\varepsilon_0$ . For this  $\varepsilon_0 > 0$ , since  $S_*$  contains all the limiting points of PDE starting from  $\overline{\rho}_0$ , there exists  $t_0$  large enough, so that as  $t \geq t_0$ , we have  $Z_0 + 2\varepsilon_0 < -\lambda_+(\overline{\rho}_t)/\lambda_-(\overline{\rho}_t) < Z_\infty - 2\varepsilon_0$ , and  $\lambda_+(\overline{\rho}_t) < 0$ ,  $\lambda_-(\overline{\rho}_t) > 0$ . For the same  $\varepsilon_0$ , since Z(r) is continuous at 0 and  $+\infty$ , there exists  $0 < r_0 < r_\infty < \infty$ , so that  $Z(r) < Z_0 + \varepsilon_0$  for  $r \in (0, r_0)$ , and  $Z(r) > Z_\infty - \varepsilon_0$  for  $r \in (r_\infty, \infty)$ . Therefore, for any  $t \geq t_0$ ,  $\partial_r \psi_\infty(r; \overline{\rho}_t) < 0$  for any  $r \in (0, r_0)$ , and  $\partial_r \psi_\infty(r; \overline{\rho}_t) > 0$  for any  $r \in (r_\infty, +\infty)$ .

As a result, according to the equation (8.28), we must have  $\partial_r \psi_\infty(r; \overline{\rho}_t) < 0$  for any  $r \in (0, r_0)$  and  $t \ge t_0$ , and  $\partial_r \psi_\infty(r; \overline{\rho}_t) > 0$  for any  $r \in (r_\infty, \infty)$  and  $t \ge t_0$ .

Due to Lemma 7.9,  $\overline{\rho}_{t_0}((0,\infty)) = 1$ . Denoting  $\Omega_k = [1/k, k]$ , then  $\lim_{k\to\infty} \overline{\rho}_{t_0}(\Omega_k) = 1$ . With this choice of  $\Omega_k$ , for any  $k \geq \{r_\infty, 1/r_0\}$ , and for any  $t \geq t_0$ , we have  $\langle \partial_r \psi_\infty(r; \overline{\rho}_t), \boldsymbol{n}(r) \rangle > 0$  for  $r \in \partial \Omega_k$  where  $\boldsymbol{n}(r)$  is the normal vector point outside  $\Omega_k$ . Therefore, if we consider the ODE

$$\dot{r}(t) = -\partial \psi_{\infty}(r(t); \overline{\rho}_t). \tag{8.29}$$

starting with  $r(t_0) \in \Omega_k$ , r(t) cannot leak outside  $\Omega_k$  from either boundaries of  $\Omega_k$ , and we must have  $r(t) \in \Omega_k$  for any  $t \ge t_0$ . Due to Lemma 7.8,  $\overline{\rho}_t(\Omega_k) \ge \overline{\rho}_{t_0}(\Omega_k)$  for any  $t \ge t_0$ . As a result, we conclude that for any  $\overline{\rho}_* \in \mathcal{S}_*$ ,

$$\overline{\rho}_*(\cup_k \Omega_k) \ge \lim_{k \to \infty} \overline{\rho}_*(\Omega_k) \ge \lim_{k \to \infty} \overline{\rho}_{t_0}(\Omega_k) = 1. \tag{8.30}$$

Note  $\bigcup_k \Omega_k = (0, \infty)$ . This gives  $\overline{\rho}_*(\{0, \infty\}) = 0$ , which proves Claim (1).

# Proof of Claim (2), step (1). If $L(S_*) \subseteq \mathcal{P}_2$ , then $S_*$ must be a singleton.

In the case  $L(S_*) \subseteq \mathcal{P}_2$ , the argument is similar to the proof of Claim (1), and hence will be presented in a synthetic form. First, we must have  $L(S_*) \subseteq \mathcal{P}_2 \cap \{(\lambda_+, \lambda_-) : Z_0 < -\lambda_+/\lambda_- < Z_\infty\}$ . Therefore, there must exist  $\varepsilon_0 > 0$ , so that for any  $\overline{\rho}_* \in S_*$ , we have  $Z_0 + 3\varepsilon_0 < -\lambda_+(\overline{\rho}_*)/\lambda_-(\overline{\rho}_*) < Z_\infty - 3\varepsilon_0$ . For this  $\varepsilon_0 > 0$ , there exists  $t_0$  large enough, so that as  $t \geq t_0$ , we have  $Z_0 + 2\varepsilon_0 < -\lambda_+(\overline{\rho}_t)/\lambda_-(\overline{\rho}_t) < Z_\infty - 2\varepsilon_0$ , and  $\lambda_+(\overline{\rho}_t) > 0$ ,  $\lambda_-(\overline{\rho}_t) < 0$ . Further, there exists  $0 < r_0 < r_\infty < \infty$ , so that  $\partial_r \psi_\infty(r; \overline{\rho}_t) > 0$  for any  $r \in (0, r_0)$  and  $t \geq t_0$ , and  $\partial_r \psi_\infty(r; \overline{\rho}_t) < 0$  for any  $r \in (r_\infty, \infty)$  and  $t \geq t_0$ .

Therefore, if we consider the ODE (8.29) starting with  $r(t_0) \in [0, r_0)$ , we must have  $r(t) \in [0, r_0)$  for any  $t \geq t_0$ ; if we start with  $r(t_0) \in (r_\infty, \infty]$ , we must have  $r(t) \in (r_\infty, \infty]$  for any  $t \geq t_0$ . Due to Lemma 7.8,  $\{\overline{\rho}_t([0,r))\}_{t\geq t_0}$  for  $0 < r \leq r_0$  and  $\{\overline{\rho}_t((r,+\infty])\}_{t\geq t_0}$  for  $r \geq r_\infty$  must be non-decreasing in t. According to Theorem 9, we can express  $\overline{\rho}_* \in \mathcal{S}_*$  in the form  $\overline{\rho}_* = a_0(\overline{\rho}_*)\delta_0 + a_\infty(\overline{\rho}_*)\delta_\infty + a(\overline{\rho}_*)\delta_{r_*}$ . By the stated monotonicity property, for any  $\overline{\rho}_1, \overline{\rho}_2 \in \mathcal{S}_*$ , it holds that  $a_0(\overline{\rho}_1) = a_0(\overline{\rho}_2)$ ,  $a_\infty(\overline{\rho}_1) = a_\infty(\overline{\rho}_2)$ , and hence  $a(\overline{\rho}_1) = a(\overline{\rho}_2)$ . We denote them in short as  $a_0, a_\infty$ , and a.

For any such fixed point  $\overline{\rho}_* \in \mathcal{S}_*$ , since we must have  $\operatorname{supp}(\overline{\rho}_*) \subseteq \{r : \partial_r \psi_\infty(r; \overline{\rho}_*) = 0\}$ ,  $r_* \in (0, +\infty)$  should be a solution of  $\phi(r) = 0$  where

$$\phi(r) = (a_0 q(0) + a_\infty q_\infty + a q_+(r) - 1) q'_+(r) + (a_0 q(0) + a_\infty q_\infty + a q_-(r) + 1) q'_-(r).$$

By condition S1, the function  $\phi(r)$  is analytic, and it is not constant. Therefore, the set of all its zeros  $\{r_*^i\}_{i\in\mathbb{N}}\subseteq(0,+\infty)$  is a countable set, and it does not have accumulation points in  $(0,+\infty)$ . Furthermore, according to Lemma 7.10, the limiting set  $\mathcal{S}_*$  should be a connected compact set

with respect to the metric  $\bar{d}_{\mathscr{P}}$ . Therefore, the limiting set could only be a singleton. That is,  $\mathcal{S}_* = \{a_0\delta_0 + a_\infty\delta_\infty + a\delta_{r_*}\}$  for some  $r_*$ .

Proof of Claim (2), step (2). If  $\overline{\rho}_*$  is a fixed point with  $L(\overline{\rho}_*) \in \mathcal{P}_2$ , then  $\overline{\rho}_*$  is unstable.

We apply Theorem 7 to  $\overline{\rho}_* = a_0 \delta_0 + a_\infty \delta_\infty + a \delta_{r_*}$ . We will check the conditions of Theorem 7 to show that this type of fixed point is unstable.

First we check condition B1. Since  $[q'_{-}(r)/q'_{+}(r)]' > 0$  and  $q'_{+}(r) > 0$  for  $r \in (0, +\infty)$ , we have

$$q''_{-}(r_*)q'_{+}(r_*) - q''_{+}(r_*)q'_{-}(r_*) > 0. (8.31)$$

Note the stationary condition of the PDE implies

$$\partial_r \psi(r_*; \overline{\rho}_*) = \lambda_+(\overline{\rho}_*) q'_+(r_*) + \lambda_-(\overline{\rho}_*) q'_-(r_*) = 0, \tag{8.32}$$

and  $\lambda_{+}(\overline{\rho}_{*}) > 0$ ,  $\lambda_{-}(\overline{\rho}_{*}) < 0$ . Combined with the equation above, we have

$$\partial_r^2 \psi_{\infty}(r_*; \overline{\rho}_*) = \lambda_+(\overline{\rho}_*) q''_+(r_*) + \lambda_-(\overline{\rho}_*) q''_-(r_*) = [q'_+(r_*) q''_-(r_*) - q'_-(r_*) q''_+(r_*)] \cdot \lambda_-(\overline{\rho}_*) / q'_+(r_*) < 0.$$
(8.33)

This verifies condition B1 of Theorem 7.

Second, since  $\lambda_{+}(\overline{\rho}_{*}) > 0$  and  $\lambda_{-}(\overline{\rho}_{*}) < 0$ , according to Equation (8.28), we must have  $\partial_{r}\psi_{\infty}(r;\overline{\rho}_{*}) > 0$  for  $r \in (0,r_{*})$ , and  $\partial_{r}\psi_{\infty}(r;\overline{\rho}_{*}) < 0$  for  $r \in (r_{*},\infty)$ . Therefore, we have  $\psi_{\infty}(0;\overline{\rho}_{*}) < \psi_{\infty}(r_{*};\overline{\rho}_{*})$  and  $\psi_{\infty}(+\infty;\overline{\rho}_{*}) < \psi_{\infty}(r_{*};\overline{\rho}_{*})$ . Note  $\mathcal{L}(\eta) \equiv \{r : \psi_{\infty}(r;\overline{\rho}_{*}) \le \psi_{\infty}(r_{*};\overline{\rho}_{*}) - \eta\}$ . For any  $\eta > 0$  small enough,  $\overline{\rho}_{*}(\mathcal{L}(\eta)) = 1 - a$ , which verifies condition B2. It is also easy to see that, for any  $\eta > 0$ ,  $\partial \mathcal{L}(\eta)$  is a compact set, hence condition B3 holds. Note that we assumed further that  $\overline{\rho}_{0}$  has a bounded density with respect to Lebesgue measure, all the assumptions of Theorem 7 are satisfied. Theorem 7 implies that the PDE cannot converge to  $\overline{\rho}_{*}$ . As a result, we conclude that we cannot have  $\mathcal{L}(\mathcal{S}_{*}(\overline{\rho}_{0})) \subseteq \mathcal{P}_{2}$  for  $\overline{\rho}_{0} \in \mathscr{P}_{\text{good}}$ . This proves Claim (2).

### 8.4 Proof of Theorem 1

The key step consists in proving that the dynamics for large but finite d is well approximated by the dynamics at  $d = \infty$ . The key estimate is provided by the next lemma.

**Lemma 8.1.** Assume  $\sigma$  satisfies condition S0, recall the definition of  $u_d$  and  $u_{\infty}$  given by Equation (8.8) and (8.9). Then we have

$$\lim_{d \to \infty} \sup_{r_1, r_2 \in [0, \infty)} |u_d(r_1, r_2) - u_\infty(r_1, r_2)| = 0,$$

and

$$\lim_{d\to\infty} \sup_{r_1,r_2\in[0,\infty)} |\partial_1 u_d(r_1,r_2) - \partial_1 u_\infty(r_1,r_2)| = 0.$$

*Proof.* Recall that  $u_d$  is given by

$$u_d(r_1, r_2) = 1/2 \cdot [u_{d,1}(r_1, r_2) + u_{d,2}(r_1, r_2)],$$
  

$$u_{d,1}(r_1, r_2) = \mathbb{E}[\sigma(r_1(1 + \Delta)G_1)\sigma(r_2(1 + \Delta)(G_1\cos\Theta + G_2\sin\Theta))],$$
  

$$u_{d,2}(r_1, r_2) = \mathbb{E}[\sigma(r_1(1 - \Delta)G_1)\sigma(r_2(1 - \Delta)(G_1\cos\Theta + G_2\sin\Theta))],$$

where  $(G_1, G_2) \sim \mathsf{N}(0, \mathbf{I}_2)$ , and  $\Theta \sim (1/Z_d)\sin(\theta)^{d-2} \cdot \mathbf{1}\{\theta \in [0, \pi]\}d\theta$  are mutually independent. Define  $G_3 = G_1\cos\Theta + G_2\sin\Theta$ , then

$$|u_{d,1}(r_1, r_2) - u_{\infty,1}(r_1, r_2)|$$

$$= |\mathbb{E}[\sigma(r_1(1+\Delta)G_1)[\sigma(r_2(1+\Delta)G_3) - \sigma(r_2(1+\Delta)G_2)]]|$$

$$\leq ||\sigma||_{\infty} \mathbb{E}[|\sigma(r_2(1+\Delta)G_3) - \sigma(r_2(1+\Delta)G_2)|],$$
(8.34)

and

$$\begin{aligned} &|\partial_{1}u_{d,1}(r_{1}, r_{2}) - \partial_{1}u_{\infty,1}(r_{1}, r_{2})| \\ &= |\mathbb{E}[(1+\Delta)G_{1} \cdot \sigma'(r_{1}(1+\Delta)G_{1})[\sigma(r_{2}(1+\Delta)G_{3}) - \sigma(r_{2}(1+\Delta)G_{2})]]| \\ &\leq (1+\Delta)\|\sigma'\|_{\infty}\mathbb{E}[G_{1}^{2}]^{1/2}\mathbb{E}[[\sigma(r_{2}(1+\Delta)G_{3}) - \sigma(r_{2}(1+\Delta)G_{2})]^{2}]^{1/2} \\ &\leq (1+\Delta)\|\sigma'\|_{\infty}(2\|\sigma\|_{\infty}^{1/2}) \cdot \mathbb{E}[|\sigma(r_{2}(1+\Delta)G_{3}) - \sigma(r_{2}(1+\Delta)G_{2})|]^{1/2}. \end{aligned}$$
(8.35)

According to condition S0,  $\|\sigma'\|_{\infty}$  and  $\|\sigma\|_{\infty}$  are bounded, it is sufficient to bound the following quantity uniformly for  $r \in [0, \infty)$ 

$$T(r) \equiv 1/2 \cdot \mathbb{E}\{|\sigma(rG_2) - \sigma(rG_3)|\} = \mathbb{E}\{[\sigma(rG_2) - \sigma(rG_3)] \mathbf{1}_{G_2 > G_3}\}.$$
(8.36)

We claim that, for any  $a \in \mathbb{R}$ ,

$$\mathbb{P}(G_3 \le a, G_2 \ge a) \le \mathbb{P}(G_3 \le 0, G_2 \ge 0) = \mathbb{E}[|\pi/2 - \Theta|/(2\pi)]. \tag{8.37}$$

Assuming this claim holds, let us show that it implies the desired bound on T(r). We have

$$T(r) = \mathbb{E}\left\{ \int_{\mathbb{R}} \sigma'(t) \, \mathbf{1}_{rG_2 \ge t \ge rG_3} \, \mathrm{d}t \right\} = \int_{\mathbb{R}} \sigma'(t) \, \mathbb{P}\left\{ G_2 \ge t/r \ge G_3 \right\} \, \mathrm{d}t$$
$$\leq \sup_{a \in \mathbb{R}} \mathbb{P}(G_3 \le a, G_2 \ge a) \, \int_{\mathbb{R}} \sigma'(t) \, \mathrm{d}t \le 2 \|\sigma\|_{\infty} \cdot \mathbb{E}[|\pi/2 - \Theta|/(2\pi)] \, .$$

Note that  $\cos(\Theta) \stackrel{\mathrm{d}}{=} Z_1/\|\boldsymbol{Z}\|_2$  for  $\boldsymbol{Z} \sim \mathsf{N}(0,\boldsymbol{I}_d)$  and hence  $\mathbb{E}\{|\Theta - \pi/2|\} \leq K/\sqrt{d}$  for a universal constant K. We therefore obtain

$$\sup_{r} |T(r)| \le (K/\pi) \|\sigma\|_{\infty} / \sqrt{d}. \tag{8.38}$$

We are left with the task of proving Eq. (8.37).

Denote  $X = G_2$  and  $Y = G_3$  for simplicity in notations. Note that  $(X,Y) \stackrel{\mathrm{d}}{=} (Y,X) \stackrel{\mathrm{d}}{=} (-X,-Y)$ . It follows that we can assume, without loss of generality, a > 0. We have

$$\begin{split} \mathbb{P}(Y \leq a, X \geq a) = & \mathbb{P}(Y \leq 0, X \geq a) + \mathbb{P}(0 \leq Y \leq a, X \geq a), \\ \mathbb{P}(Y \leq 0, X \geq 0) = & \mathbb{P}(Y \leq 0, X \geq a) + \mathbb{P}(Y \leq 0, 0 \leq X \leq a), \end{split}$$

suffice to prove that

$$\mathbb{P}(0 \le Y \le a, X \ge a) \le \mathbb{P}(Y \le 0, 0 \le X \le a).$$

Define U = (X - Y)/2, V = (X + Y)/2, and  $\mathcal{A}_1 = \{0 \le Y \le a, X \ge a\}$ ,  $\mathcal{A}_2 = \{Y \le 0, 0 \le X \le a\}$ . It is easy to see that  $[U|\Theta = \theta]$  and  $[V|\Theta = \theta]$  are independent normal random variables. Therefore, it is sufficient to show  $\mathbb{P}(\mathcal{A}_1|U = u, \Theta = \theta) \le \mathbb{P}(\mathcal{A}_2|U = u, \Theta = \theta)$  for  $u \ge 0$  and  $\theta \in [0, \pi]$  (as u < 0, both conditional probability equal 0).

Fix an  $u \geq 0$  and  $\theta \in [0, \pi]$ . Consider the closed interval  $\mathcal{I}_i = \mathcal{I}_i(u) \subseteq \mathbb{R}$  for i = 1, 2, with definition  $\mathcal{I}_i(u) \equiv \{v : \{U = u, V = v, \Theta = \theta\} \subseteq \mathcal{A}_i\}$ . Then  $\mathbb{P}(\mathcal{A}_i|U = u, \Theta = \theta) = \int_{\mathcal{I}_i(u)} p_{V|\Theta}(v|\theta) \mathrm{d}v$ , where  $p_{V|\Theta}(v|\theta)$  is the density of  $[V|\Theta = \theta]$  at v. It is not hard to see that every element in  $\mathcal{I}_1$  is greater or equal to a/2, and every element in  $\mathcal{I}_2$  is less or equal to a/2; in the meanwhile,  $\mathcal{I}_1$  and  $\mathcal{I}_2$  are symmetric with respect to a/2. Note that  $[V|\Theta = \theta]$  is a Gaussian random variable with zero mean, therefore  $p_{V|\Theta}(a/2 + s|\theta) \leq p_{V|\Theta}(a/2 - s|\theta)$  for any  $s \geq 0$  and  $\theta \in [0, \pi]$ . This implies that  $\mathbb{P}(\mathcal{A}_1|U = u, \Theta = \theta) \leq \mathbb{P}(\mathcal{A}_2|U = u, \Theta = \theta)$ , for any  $u \geq 0$  and  $\theta \in [0, \pi]$ .

**Lemma 8.2.** Let  $y \sim \text{Unif}(\{-1,+1\})$ ,  $[\boldsymbol{x}|y=+1] \sim \mathsf{N}(\mathbf{0},\boldsymbol{\Sigma}_+)$ ,  $[\boldsymbol{x}|y=-1] \sim \mathsf{N}(0,\boldsymbol{\Sigma}_-)$  with  $\tau_-^2 \boldsymbol{I}_D \leq \boldsymbol{\Sigma}_+, \boldsymbol{\Sigma}_- \leq \tau_+^2 \boldsymbol{I}_D$  for some  $0 < \tau_- < \tau_+ < \infty$ . Assume that the activation function  $\sigma$  satisfies condition S0. Define

$$V(\boldsymbol{\theta}) = -\mathbb{E}[y \,\sigma(\langle \boldsymbol{x}, \boldsymbol{\theta} \rangle)],$$

$$U(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) = \mathbb{E}[\sigma(\langle \boldsymbol{x}, \boldsymbol{\theta}_1 \rangle)\sigma(\langle \boldsymbol{x}, \boldsymbol{\theta}_2 \rangle)].$$
(8.39)

Then assumptions A2 and A3 are satisfied.

*Proof.* Note that x is sub-Gaussian, and by condition S0 we have  $\sigma'$  is bounded, then  $\nabla_{\theta}\sigma(\langle x,\theta\rangle) = \sigma'(\langle x,\theta\rangle)x$  is also sub-Gaussian (with sub-Gaussian parameter independent of D). Condition S0 also gives that  $\sigma$  is bounded, therefore assumption A2 is satisfied.

To verify assumption A3, it is sufficient to check that  $\nabla V$ ,  $\nabla_1 U$ ,  $\nabla_{12}^2 U$ ,  $\nabla^2 V$ , and  $\nabla_{11}^2 U$  are uniformly bounded in  $\ell_2$  norm (for the gradients) or operator norm (for the Hessians). For any unit vector  $\boldsymbol{n}$ , we have

$$\langle \nabla V(\boldsymbol{\theta}), \boldsymbol{n} \rangle = -\mathbb{E}[y\sigma'(\langle \boldsymbol{x}, \boldsymbol{\theta} \rangle)\langle \boldsymbol{x}, \boldsymbol{n} \rangle],$$
 (8.40)

$$\langle \nabla_1 U(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2), \boldsymbol{n} \rangle = \mathbb{E}[\sigma'(\langle \boldsymbol{x}, \boldsymbol{\theta}_1 \rangle) \langle \boldsymbol{x}, \boldsymbol{n} \rangle \sigma(\langle \boldsymbol{x}, \boldsymbol{\theta}_2 \rangle)], \tag{8.41}$$

$$\langle \nabla_{12}^2 U(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2), \boldsymbol{n}^{\otimes 2} \rangle = \mathbb{E}[\sigma'(\langle \boldsymbol{x}, \boldsymbol{\theta}_1 \rangle) \langle \boldsymbol{x}, \boldsymbol{n} \rangle^2 \sigma'(\langle \boldsymbol{x}, \boldsymbol{\theta}_2 \rangle)]. \tag{8.42}$$

Since  $\|\sigma\|_{\infty}$ ,  $\|\sigma'\|_{\infty} < \infty$ , applying Cauchy-Schwarz inequality, we have  $\nabla V$ ,  $\nabla_1 U$ ,  $\nabla_{12}^2 U$  are uniformly bounded.

It is difficult to bound  $\nabla^2 V$  and  $\nabla_1^2 U$  directly because  $\sigma'$  may not be differentiable. We will use a longer argument to bound them.

First, for a bounded-Lipschitz function f, and for  $g \in \{1, \sigma\}$ , define

$$W_{f,g}(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) = \mathbb{E}_{\boldsymbol{G}}[f(\langle \boldsymbol{\theta}_1, \boldsymbol{G} \rangle)g(\langle \boldsymbol{\theta}_2, \boldsymbol{G} \rangle)], \tag{8.43}$$

where  $G \sim N(0, I_d)$ . Since we have  $\tau_-^2 I_D \leq \Sigma_+, \Sigma_- \leq \tau_+^2 I_D$  for some  $0 < \tau_- < \tau_+ < \infty$ , in order to bound  $\nabla^2 V$  and  $\nabla_1^2 U$ , it is sufficient to bound  $\nabla_1^2 W_{\sigma,1}$  and  $\nabla_1^2 W_{\sigma,\sigma}$ .

Since  $\sigma'$  is  $K_0$ -Lipschitz on  $[-2\delta_0, 2\delta_0]$  for some  $\delta_0 > 0$  and  $K_0 < \infty$ , then, there exists a function  $\sigma_0 : \mathbb{R} \to \mathbb{R}$ , so that  $\sigma_0$  is non-decreasing and K-bounded-Lipschitz,  $\sigma'_0$  is K-bounded-Lipschitz, and  $\sigma_0(r) = \sigma(r)$  for  $r \in [-\delta_0, \delta_0]$ . For this  $\sigma_0$ , a second weak derivative exists and  $|\sigma''_0| \leq K$ . Hence

$$\langle \nabla_1^2 W_{\sigma_0, g}(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2), \boldsymbol{n}^{\otimes 2} \rangle = \mathbb{E}[\sigma_0''(\langle \boldsymbol{\theta}_1, \boldsymbol{G} \rangle) \langle \boldsymbol{G}, \boldsymbol{n} \rangle^2 g(\langle \boldsymbol{\theta}_2, \boldsymbol{G} \rangle)]$$
(8.44)

is uniformly bounded for g=1 or  $g=\sigma$ . Let  $h=\sigma-\sigma_0$ , then h=0 for  $r\in [-\delta_0,\delta_0]$ , and h is K-bounded-Lipschitz for some constant K. It is sufficient to bound  $\nabla_1^2 W_{h,g}$  for  $g\in \{1,\sigma\}$ .

Since G is Gaussian, using Stein's formula, for any unit vector n, we have

$$\langle \nabla_{1}W_{h,g}(\boldsymbol{\theta}_{1},\boldsymbol{\theta}_{2}),\boldsymbol{n}\rangle = \mathbb{E}[h'(\langle\boldsymbol{\theta}_{1},\boldsymbol{G}\rangle)\langle\boldsymbol{n},\boldsymbol{G}\rangle g(\langle\boldsymbol{\theta}_{2},\boldsymbol{G}\rangle)]$$

$$= \underbrace{\frac{1}{\|\boldsymbol{\theta}_{1}\|_{2}^{2}}\mathbb{E}[h(\langle\boldsymbol{\theta}_{1},\boldsymbol{G}\rangle)\langle\boldsymbol{\theta}_{1},\boldsymbol{G}\rangle\langle\boldsymbol{n},\boldsymbol{G}\rangle g(\langle\boldsymbol{\theta}_{2},\boldsymbol{G}\rangle)]}_{E_{1}(\boldsymbol{\theta}_{1},\boldsymbol{\theta}_{2},\boldsymbol{n})} - \underbrace{\frac{1}{\|\boldsymbol{\theta}_{1}\|_{2}^{2}}\mathbb{E}[h(\langle\boldsymbol{\theta}_{1},\boldsymbol{G}\rangle)\langle\boldsymbol{\theta}_{1},\boldsymbol{n}\rangle g(\langle\boldsymbol{\theta}_{2},\boldsymbol{G}\rangle)]}_{E_{2}(\boldsymbol{\theta}_{1},\boldsymbol{\theta}_{2},\boldsymbol{n})} (8.45)$$

$$- \underbrace{\frac{1}{\|\boldsymbol{\theta}_{1}\|_{2}^{2}}\mathbb{E}[h(\langle\boldsymbol{\theta}_{1},\boldsymbol{G}\rangle)\langle\boldsymbol{n},\boldsymbol{G}\rangle g'(\langle\boldsymbol{\theta}_{2},\boldsymbol{G}\rangle)\langle\boldsymbol{\theta}_{2},\boldsymbol{\theta}_{1}\rangle]}_{E_{3}(\boldsymbol{\theta}_{1},\boldsymbol{\theta}_{2},\boldsymbol{n})}.$$

Taking directional derivatives of  $E_1$  and  $E_2$ , we have

$$\langle \nabla_{1} E_{1}(\boldsymbol{\theta}_{1}, \boldsymbol{\theta}_{2}, \boldsymbol{n}), \boldsymbol{n} \rangle = \underbrace{\frac{1}{\|\boldsymbol{\theta}_{1}\|_{2}^{2}} \mathbb{E}[h'(\langle \boldsymbol{\theta}_{1}, \boldsymbol{G} \rangle) \langle \boldsymbol{\theta}_{1}, \boldsymbol{G} \rangle \langle \boldsymbol{n}, \boldsymbol{G} \rangle^{2} g(\langle \boldsymbol{\theta}_{2}, \boldsymbol{G} \rangle)]}_{E_{11}} + \underbrace{\frac{1}{\|\boldsymbol{\theta}_{1}\|_{2}^{2}} \mathbb{E}[h(\langle \boldsymbol{\theta}_{1}, \boldsymbol{G} \rangle) \langle \boldsymbol{n}, \boldsymbol{G} \rangle^{2} g(\langle \boldsymbol{\theta}_{2}, \boldsymbol{G} \rangle)]}_{E_{12}} - \underbrace{\frac{2\langle \boldsymbol{\theta}_{1}, \boldsymbol{n} \rangle}{\|\boldsymbol{\theta}_{1}\|_{2}^{4}} \mathbb{E}[h(\langle \boldsymbol{\theta}_{1}, \boldsymbol{G} \rangle) \langle \boldsymbol{\theta}_{1}, \boldsymbol{G} \rangle \langle \boldsymbol{n}, \boldsymbol{G} \rangle g(\langle \boldsymbol{\theta}_{2}, \boldsymbol{G} \rangle)]}_{E_{13}}, \tag{8.46}$$

and

$$\langle \nabla_{1} E_{2}(\boldsymbol{\theta}_{1}, \boldsymbol{\theta}_{2}, \boldsymbol{n}), \boldsymbol{n} \rangle = \underbrace{\frac{1}{\|\boldsymbol{\theta}_{1}\|_{2}^{2}} \mathbb{E}[h'(\langle \boldsymbol{\theta}_{1}, \boldsymbol{G} \rangle) \langle \boldsymbol{\theta}_{1}, \boldsymbol{n} \rangle \langle \boldsymbol{G}, \boldsymbol{n} \rangle g(\langle \boldsymbol{\theta}_{2}, \boldsymbol{G} \rangle)]}_{E_{21}} + \underbrace{\frac{1}{\|\boldsymbol{\theta}_{1}\|_{2}^{2}} \mathbb{E}[h(\langle \boldsymbol{\theta}_{1}, \boldsymbol{G} \rangle) g(\langle \boldsymbol{\theta}_{2}, \boldsymbol{G} \rangle)]}_{E_{22}} - \underbrace{\frac{2\langle \boldsymbol{\theta}_{1}, \boldsymbol{n} \rangle}{\|\boldsymbol{\theta}_{1}\|_{2}^{4}} \mathbb{E}[h(\langle \boldsymbol{\theta}_{1}, \boldsymbol{G} \rangle) \langle \boldsymbol{\theta}_{1}, \boldsymbol{n} \rangle g(\langle \boldsymbol{\theta}_{2}, \boldsymbol{G} \rangle)]}_{E_{22}}.$$

$$(8.47)$$

To bound  $E_{11}$ , note h'(r) = 0 for  $r \in (-\delta_0, \delta_0)$ , and  $|h'(r)| \leq K$  for  $r \in \mathbb{R}$ , we have

$$E_{11} \leq \frac{K}{\|\boldsymbol{\theta}_{1}\|_{2}} \mathbb{E}\left[\mathbf{1}\{|\langle\boldsymbol{\theta}_{1},\boldsymbol{G}\rangle| \geq \delta_{0}\} \cdot |\langle\boldsymbol{\theta}_{1}/\|\boldsymbol{\theta}_{1}\|_{2},\boldsymbol{G}\rangle| \cdot \langle\boldsymbol{n},\boldsymbol{G}\rangle^{2} |g(\langle\boldsymbol{\theta}_{2},\boldsymbol{G}\rangle)|\right]$$

$$\leq \frac{K}{\|\boldsymbol{\theta}_{1}\|_{2}} \cdot \mathbb{P}(|\langle\boldsymbol{\theta}_{1},\boldsymbol{G}\rangle| \geq \delta_{0})^{1/2} \cdot \{\mathbb{E}[(\langle\boldsymbol{\theta}_{1}/\|\boldsymbol{\theta}_{1}\|_{2},\boldsymbol{G}\rangle^{2}\langle\boldsymbol{n},\boldsymbol{G}\rangle^{4}g(\langle\boldsymbol{\theta}_{2},\boldsymbol{G}\rangle))^{2}]\}^{1/2}.$$
(8.48)

Take  $r = \|\boldsymbol{\theta}_1\|_2$ , then

$$1/\|\boldsymbol{\theta}_1\|_2 \cdot \mathbb{P}(|\langle \boldsymbol{\theta}_1, \boldsymbol{G} \rangle| \ge \delta_0)^{1/2} \le 1/r \cdot \exp\{-\delta_0^2/(4r^2)\}$$
(8.49)

is uniformly bounded for  $r \in [0, \infty]$ . Hence  $E_{11}$  is uniformly bounded. Using a similar argument, we can show that each terms  $E_{12}$ ,  $E_{13}$ ,  $E_{21}$ ,  $E_{22}$ , and  $E_{23}$  are uniformly bounded.

Now we look at  $\nabla_1 E_3(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{n})$ . We have

$$\langle \nabla_{1} E_{3}(\boldsymbol{\theta}_{1}, \boldsymbol{\theta}_{2}, \boldsymbol{n}), \boldsymbol{n} \rangle = \underbrace{\frac{1}{\|\boldsymbol{\theta}_{1}\|_{2}^{2}} \mathbb{E}[h'(\langle \boldsymbol{\theta}_{1}, \boldsymbol{G} \rangle) \langle \boldsymbol{n}, \boldsymbol{G} \rangle^{2} g'(\langle \boldsymbol{\theta}_{2}, \boldsymbol{G} \rangle) \langle \boldsymbol{\theta}_{2}, \boldsymbol{\theta}_{1} \rangle]}_{E_{31}} + \underbrace{\frac{1}{\|\boldsymbol{\theta}_{1}\|_{2}^{2}} \mathbb{E}[h(\langle \boldsymbol{\theta}_{1}, \boldsymbol{G} \rangle) \langle \boldsymbol{n}, \boldsymbol{G} \rangle g'(\langle \boldsymbol{\theta}_{2}, \boldsymbol{G} \rangle) \langle \boldsymbol{\theta}_{2}, \boldsymbol{n} \rangle]}_{E_{32}} - \underbrace{\frac{2\langle \boldsymbol{\theta}_{1}, \boldsymbol{n} \rangle}{\|\boldsymbol{\theta}_{1}\|_{2}^{4}} \mathbb{E}[h(\langle \boldsymbol{\theta}_{1}, \boldsymbol{G} \rangle) \langle \boldsymbol{n}, \boldsymbol{G} \rangle g'(\langle \boldsymbol{\theta}_{2}, \boldsymbol{G} \rangle) \langle \boldsymbol{\theta}_{2}, \boldsymbol{\theta}_{1} \rangle]}_{E_{33}}.$$

$$(8.50)$$

In order to bound  $E_{32}$ , we apply Stein's formula to get

$$E_{32} = \frac{\langle \boldsymbol{\theta}_{2}, \boldsymbol{n} \rangle}{\|\boldsymbol{\theta}_{1}\|_{2}^{2} \|\boldsymbol{\theta}_{2}\|_{2}^{2}} \Big\{ \mathbb{E}[h(\langle \boldsymbol{\theta}_{1}, \boldsymbol{G} \rangle) \langle \boldsymbol{n}, \boldsymbol{G} \rangle g(\langle \boldsymbol{\theta}_{2}, \boldsymbol{G} \rangle) \langle \boldsymbol{\theta}_{2}, \boldsymbol{G} \rangle] \\ - \mathbb{E}[h(\langle \boldsymbol{\theta}_{1}, \boldsymbol{G} \rangle) \langle \boldsymbol{n}, \boldsymbol{\theta}_{2} \rangle g(\langle \boldsymbol{\theta}_{2}, \boldsymbol{G} \rangle)] - \mathbb{E}[h'(\langle \boldsymbol{\theta}_{1}, \boldsymbol{G} \rangle) \langle \boldsymbol{\theta}_{1}, \boldsymbol{\theta}_{2} \rangle \langle \boldsymbol{n}, \boldsymbol{G} \rangle g(\langle \boldsymbol{\theta}_{2}, \boldsymbol{G} \rangle)] \Big\}.$$

$$(8.51)$$

For each terms above, we can bound them using the same argument as for bounding  $E_{11}$ . Similarly, we can bound  $E_{33}$ . We cannot apply directly Stein's formula to  $E_{31}$  similar to what we did for  $E_{32}$ , because  $h' = \sigma' - \sigma'_0$  may not have weak derivative. However, recall that h'(r) = 0 for  $r \in [-\delta_0, \delta_0]$  and h' is K-bounded. Therefore, we can find a function  $h_0 : \mathbb{R} \to \mathbb{R}$ , such that  $|h'(r)| \le h_0(r)$  for  $r \in \mathbb{R}$ ,  $h_0(r) = 0$  for  $r \in [-\delta_0/2, \delta_0/2]$ , and  $h_0$  is K-bounded-Lipschitz (for some larger constant K). Hence, recalling that  $g'(r) \ge 0$ , we get

$$E_{31} \leq \frac{1}{\|\boldsymbol{\theta}_1\|_2} \mathbb{E}[h_0(\langle \boldsymbol{\theta}_1, \boldsymbol{G} \rangle) \langle \boldsymbol{n}, \boldsymbol{G} \rangle^2 g'(\langle \boldsymbol{\theta}_2, \boldsymbol{G} \rangle) \|\boldsymbol{\theta}_2\|_2]. \tag{8.52}$$

We can apply Stein's formula to the right hand side of the last equation. Using the same argument as above, we obtain that  $E_{31}$  is uniformly bounded.

As a result,  $\nabla^2 V$  and  $\nabla_1^2 U$  are uniformly bounded. Therefore, assumption A3 is satisfied.

We are now in position to prove Theorem 1.

Proof of Theorem 1. First we consider PDE (8.22) for  $d = \infty$ . We fix an initial radial density  $\overline{\rho}_0 \in \mathscr{P}_{\text{good}}$ . Due to Theorem 10, for any  $\eta > 0$ , there exists  $T = T(\eta, \overline{\rho}_0, \Delta) > 0$ , so that the solution  $(\overline{\rho}_t^{\infty})_{t\geq 0}$  of PDE (8.22) for  $d = \infty$  with initialization  $\overline{\rho}_0$  satisfies

$$\overline{R}_{\infty}(\overline{\rho}_t^{\infty}) \le \inf_{\overline{\rho} \in \overline{\mathscr{P}}} \overline{R}_{\infty}(\overline{\rho}) + \eta/5$$

for any  $t \geq T$ .

Then we consider the general PDE

$$\partial_t \rho_t(\boldsymbol{\theta}) = 2\xi(t) \nabla \cdot \left[ \rho_t(\boldsymbol{\theta}) \nabla \Psi(\boldsymbol{\theta}; \rho_t) \right], \tag{8.53}$$

with initialization  $\rho_0$  the distribution of  $r\mathbf{n}$ , where  $(r, \mathbf{n}) \sim \overline{\rho}_0 \times \text{Unif}(\mathbb{S}^{d-1})$ . Due to Lemma 8.2, we have the existence and uniqueness of the solution of PDE (8.53), and let  $(\rho_t)_{t\geq 0}$  be the solution. Let  $\overline{\rho}_t^d$  be the radial marginal distribution of  $\rho_t$ . It is easy to see that  $(\overline{\rho}_t^d)_{t\geq 0}$  is the unique solution of (8.22) for d finite.

Now, we would like to bound the distance of  $\overline{\rho}_t^d$  and  $\overline{\rho}_t^\infty$  using Lemma 7.7. We take D=1,  $V=v,\ U=u_d,\ \tilde{V}=v,\ \tilde{U}=u_\infty$  in Lemma 7.7. Let  $\varepsilon_0(d)$  be defined as in Eq. (7.69). Due to Lemma 8.1, we have  $\varepsilon_0(d)\to 0$  as  $d\to\infty$ . Therefore, according to Lemma 7.7, we have  $\lim_{d\to\infty}\sup_{t\le 10T}d_{\rm BL}(\overline{\rho}_t^d,\overline{\rho}_t^\infty)=0$ . Further note that  $\overline{R}_\infty$  is uniformly continuous with respect to  $\overline{\rho}$  in bounded-Lipschitz distance. Therefore, there exists  $d_0=d_0(\eta,\overline{\rho}_0,\Delta)$  large enough, so that for  $d\ge d_0$  we have

$$|\overline{R}_{\infty}(\overline{\rho}_t^d) - \overline{R}_{\infty}(\overline{\rho}_t^{\infty})| \le \eta/5.$$

for any t < 10T.

Next we would like to bound the difference of  $\overline{R}_{\infty}(\overline{\rho})$  and  $\overline{R}_{d}(\overline{\rho})$  for any  $\overline{\rho}$ . Note

$$|\overline{R}_{\infty}(\overline{\rho}) - \overline{R}_{d}(\overline{\rho})| \le \int |u_{d}(r_{1}, r_{2}) - u_{\infty}(r_{1}, r_{2})| \,\overline{\rho}(\mathrm{d}r_{1})\overline{\rho}(\mathrm{d}r_{2}). \tag{8.54}$$

By Lemma 8.1, there exists  $d_0 = d_0(\eta, \Delta)$  large enough, so that for  $d \geq d_0$ , we have

$$\sup_{\overline{\rho} \in \overline{\mathscr{P}}} |\overline{R}_{\infty}(\overline{\rho}) - \overline{R}_{d}(\overline{\rho})| \le \eta/5. \tag{8.55}$$

Finally, let  $(\boldsymbol{\theta}^k)_{k\geq 1}$  be the trajectory of SGD, with step size  $s_k = \varepsilon \xi(k\varepsilon)$ , and initialization  $\boldsymbol{w}_i^0 \sim_{iid} \rho_0$  for  $i \leq N$ . We apply Theorem 3 to bound the difference of the law of trajectory of SGD and the solution of PDE (8.53). The assumptions of Theorem 3 are verified by Lemma 8.2. As a consequence, there exists constant K (which depend uniquely on the constants in assumptions A1 A2 A3), such that for any  $t \leq 10T$ , we have

$$R_N(\boldsymbol{\theta}^{\lfloor t/\varepsilon \rfloor}) - \overline{R}_d(\overline{\rho}_t^d) \le Ke^{10KT} \cdot \operatorname{err}_{N,d}(z).$$

with probability  $1 - e^{-z^2}$ , where

$$\operatorname{err}_{N,d}(z) = \sqrt{1/N \vee \varepsilon} \cdot \left\lceil \sqrt{d + \log(N(1/\varepsilon \vee 1))} + z \right\rceil.$$

As a consequence, for any  $\delta > 0$ , there exists  $C_0 = C_0(\delta, \eta, \overline{\rho}_0, \Delta)$ , so that as  $N, 1/\varepsilon \geq C_0 d$  and  $\varepsilon \geq 1/N^{10}$ , for any  $t \leq 10T$ , we have

$$R_N(\boldsymbol{\theta}^{\lfloor t/\varepsilon \rfloor}) - \overline{R}_d(\overline{\rho}_t^d) \le \eta/5$$

with probability at least  $1 - \delta$ .

Therefore, the trajectory  $\boldsymbol{\theta}^{\lfloor t/\varepsilon \rfloor}$  of SGD as  $t \in [T, 10T]$  satisfies

$$R_{N}(\boldsymbol{\theta}^{\lfloor t/\varepsilon \rfloor}) \leq \overline{R}_{d}(\overline{\rho}_{t}^{d}) + \eta/5 \leq \overline{R}_{\infty}(\overline{\rho}_{t}^{d}) + 2\eta/5 \leq \overline{R}_{\infty}(\overline{\rho}_{t}^{\infty}) + 3\eta/5$$

$$\leq \inf_{\overline{\rho} \in \overline{\mathscr{P}}} \overline{R}_{\infty}(\overline{\rho}) + 4\eta/5 \leq \inf_{\overline{\rho} \in \overline{\mathscr{P}}} \overline{R}_{d}(\overline{\rho}) + \eta = \inf_{\rho \in \mathscr{P}(\mathbb{R}^{d})} R(\rho) + \eta$$

$$\leq \inf_{\boldsymbol{\theta} \in \mathbb{R}^{d \times N}} R_{N}(\boldsymbol{\theta}) + \eta$$

with probability at least  $1 - \delta$ . This gives the desired result.

#### 8.5 Checking conditions S0–S4 for the running example

**Lemma 8.3.** Consider the activation function  $\sigma$  with definition in Equation (8.1), with  $s_1 < s_2$ ,  $s_1 < -1$ ,  $(s_1 + s_2)/2 > 1$ ,  $(3s_1 + s_2)/4 \in (-1,1)$ ,  $0 < t_1 < t_2$ . For  $r \in (0,+\infty)$ , define  $q(r) = \mathbb{E}_G[\sigma(rG)]$  where  $G \sim \mathsf{N}(0,1)$ . Then conditions SO-S4 hold.

**Remark 8.1.** The requirements of Lemma 8.3 are not restrictive. An example of parameters that satisfies all conditions gives  $s_1 = -2.5$ ,  $s_2 = 7.5$ ,  $t_1 = 0.5$ ,  $t_2 = 1.5$ .

*Proof.* It is straightforward to see that condition S0 holds. To show condition S1, denote by  $\sigma'(r)$  the weak derivative of  $\sigma(r)$ , we calculate the function q'(r) for r > 0 explicitly,

$$q'(r) = \mathbb{E}[\sigma'(rG)G] = \frac{s_2 - s_1}{t_2 - t_1} \int_{\mathbb{R}} \mathbf{1}\{rx \in [t_1, t_2]\} \cdot \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{x^2}{2}\right\} \cdot x \cdot dx$$

$$= \frac{s_2 - s_1}{\sqrt{2\pi}(t_2 - t_1)} \left\{ \exp\left[-\frac{t_1^2}{2r^2}\right] - \exp\left[-\frac{t_2^2}{2r^2}\right] \right\}.$$
(8.56)

Since  $s_1 < s_2$  and  $0 < t_1 < t_2$ , it is easy to see that q'(r) is analytic on  $(0, \infty)$ , and hence q(r) is analytic on  $(0, \infty)$ . Differentiating q'(r) in Eq. (8.56), it is easy to see that  $\lim_{r\to\infty} q''(r) = 0$ , and q''(0+) = 0. Hence, we have  $\sup_{r \in [0,+\infty]} q''(r) < \infty$ . Then condition S1 holds.

q''(0+) = 0. Hence, we have  $\sup_{r \in [0,+\infty]} q''(r) < \infty$ . Then condition S1 holds. Since  $s_2 > s_1$ ,  $0 < t_1 < t_2$ , we have q'(r) > 0 for  $r \in (0,+\infty)$ ,  $\lim_{r \to \infty} q'(r) = 0$ , and q'(0+) = 0. Hence, we have  $\sup_{r \in [0,+\infty]} q'(r) < \infty$ . Then condition S2 holds. Note that  $q(0) = \sigma(0) = s_1 < -1$ , and  $q(+\infty) = (s_1 + s_2)/2 > 1$ . In addition,  $[q(0) + q(+\infty)]/2 = (3s_1 + s_2)/4 \in (-1,1)$ . Therefore, condition S3 holds.

Finally, we show that condition S4 holds. Define  $p(r) = \exp[-t_1^2/(2r^2)] - \exp[-t_2^2/(2r^2)]$ , which is a positively scaled version of q'(r). To show that for  $r \in (0, \infty)$ ,

$$[q'(\tau_{-}r)/q'(\tau_{+}r)]' = [\tau_{-} \cdot q''(\tau_{-}r)q'(\tau_{+}r) - \tau_{+} \cdot q'(\tau_{-}r)q''(\tau_{+}r)]/[q'(\tau_{+}r)]^{2} > 0,$$

we only need to show that for  $r \in (0, \infty)$ 

$$F_1(r) \equiv \tau_- \cdot p'(\tau_- r) p(\tau_+ r) - \tau_+ \cdot p'(\tau_+ r) p(\tau_- r) > 0.$$

We have

$$\begin{split} F_1(r) &= + \, 1/(\tau_-^2 r^3) \cdot \{t_1^2 \exp[-t_1^2/(2\tau_-^2 r^2)] - t_2^2 \exp[-t_2^2/(2\tau_-^2 r^2)]\} \\ & \times \{ \exp[-t_1^2/(2\tau_+^2 r^2)] - \exp[-t_2^2/(2\tau_+^2 r^2)]\} \\ & - \, 1/(\tau_+^2 r^3) \cdot \{t_1^2 \exp[-t_1^2/(2\tau_+^2 r^2)] - t_2^2 \exp[-t_2^2/(2\tau_+^2 r^2)]\} \\ & \times \{ \exp[-t_1^2/(2\tau_-^2 r^2)] - \exp[-t_2^2/(2\tau_-^2 r^2)]\}. \end{split}$$

Define  $x \equiv t_2^2/(2\tau_+^2r^2) > 0$ ,  $s \equiv \tau_+^2/\tau_-^2 > 1$ ,  $0 < c \equiv t_1^2/t_2^2 < 1$ , we have

$$F_{1}(r) = + t_{2}^{2}/(\tau_{+}^{2}r^{3}) \cdot \{cs \cdot \exp[-xsc] - s\exp[-xs]\} \cdot \{\exp[-xc] - \exp[-x]\}$$

$$- t_{2}^{2}/(\tau_{+}^{2}r^{3}) \cdot \{c \cdot \exp[-xc] - \exp[-x]\} \cdot \{\exp[-xsc] - \exp[-xs]\}$$

$$= t_{2}^{2}/(\tau_{+}^{2}r^{3})\{(cs - c)\exp[-xc - xsc] + (c - s)\exp[-xs - xc]$$

$$+ (1 - cs)\exp[-x - xsc] + (s - 1)\exp[-x - xs]\}$$

$$= t_{2}^{2}/(\tau_{+}^{2}r^{3})\exp\{-x - xsc\}\{(cs - c)\exp[x - xc]$$

$$+ (c - s)\exp[x - xs - xc + xsc] + (1 - cs) + (s - 1)\exp[xsc - xs]\}.$$

Define

$$F_2(x; s, c) = (cs - c) \exp[x - xc] + (c - s) \exp[x - xs - xc + xsc] + (1 - cs) + (s - 1) \exp[xsc - xs].$$

It is sufficient to show that  $F_2(x; s, c) > 0$  for x > 0, s > 1, and 0 < c < 1. Note that  $F_2(0+; s, c) = 0$ . Hence it is sufficient to show that  $\partial_x F_2(x; s, c) > 0$  for x > 0.

We have

$$\partial_x F_2(x; s, c) = c(s-1)(1-c) \exp[x-xc] + (s-c)(s-1)(1-c) \exp[x-xs-xc+xsc] + (s-1)s(c-1) \exp[xsc-xs]$$

$$= (s-1)(1-c) \exp[xsc-xs] \{c \cdot \exp[x-xc-xsc+xs] + (s-c) \exp[x-xc] - s\}.$$

Define

$$F_3(x; s, c) = c \cdot \exp[x - xc - xsc + xs] + (s - c) \exp[x - xc] - s.$$

Note that s > 1 and  $0 \le c < 1$ ,  $F_3(0+; s, c) = 0$ . It is therefore sufficient to show that  $\partial_x F_3(x; s, c) > 0$  for x > 0.

We have

$$\partial_x F_3(x; s, c) = c(1-c)(1+s) \exp[x - xc - xsc + xs] + (s-c)(1-c) \exp[x - xc].$$

Since 0 < c < 1, s > 1, and x > 0, we have  $\partial_x F_3(x; s, c) > 0$ , and hence condition S4 holds.

# 9 Centered anisotropic Gaussians

In this section we consider the centered anisotropic Gaussian example discussed in the main text. That is, we assume the joint law of (y, x) to be as follows:

With probability 1/2: y = +1,  $\boldsymbol{x} \sim N(0, \boldsymbol{\Sigma}_{+})$ .

With probability 1/2: y = -1,  $\boldsymbol{x} \sim \mathsf{N}(\boldsymbol{0}, \boldsymbol{\Sigma}_{-})$ .

We will assume  $\Sigma_+$ ,  $\Sigma_+$  to be diagonalizable in the same orthonormal basis, and to differ only on a subspace of dimension  $s_0$ . We want to study whether and how the neural network will identify this subspace of relevant features. Without loss of generality, we can assume that the eigenvalues correspond to the standard basis. In order to focus on the simplest possible model of this type, we will choose:

$$\Sigma_{+} = \operatorname{Diag}\left(\underbrace{(1+\Delta)^{2}, \dots, (1+\Delta)^{2}}_{s_{0}}, \underbrace{1, \dots, 1}_{d-s_{0}}\right), \tag{9.1}$$

$$\Sigma_{-} = \operatorname{Diag}\left(\underbrace{(1-\Delta)^{2}, \dots, (1-\Delta)^{2}}_{s_{0}}, \underbrace{1, \dots, 1}_{d-s_{0}}\right). \tag{9.2}$$

We assume  $0 < \Delta < 1$ . As in the previous section, we choose  $\sigma_*(\boldsymbol{x}; \boldsymbol{\theta}_i) = \sigma(\langle \boldsymbol{x}, \boldsymbol{w}_i \rangle)$  for some activation function  $\sigma$ . Define  $q(r) \equiv \mathbb{E}\{\sigma(rG)\}$  for  $G \sim \mathsf{N}(0,1)$ . We assume  $\sigma(\cdot)$  satisfies conditions S0 - S4 stated at the beginning of Section 8. We will still use the specific  $\sigma$  in Eq. (8.1) as our running example.

Throughout this section, we assume  $s_0 = \gamma \cdot d$  for some fixed  $0 < \gamma < 1$ . Therefore, as  $d \to \infty$ , we have  $s_0 = \gamma \cdot d \to \infty$  and  $d - s_0 = (1 - \gamma) \cdot d \to \infty$ . For any  $\mathbf{w} \in \mathbb{R}^d$ , we denote  $\mathbf{w}_1 \in \mathbb{R}^{s_0}$  and  $\mathbf{w}_2 \in \mathbb{R}^{d-s_0}$  by writing  $\mathbf{w} = (\mathbf{w}_1, \mathbf{w}_2)$ . We denote  $\tau_+ = 1 + \Delta$  and  $\tau_- = 1 - \Delta$ . Then we have  $0 < \tau_- < 1 < \tau_+ < 2$ . Denote  $q_+(r) = q(\tau_+ r)$  and  $q_-(r) = q(\tau_- r)$ . For any  $\mathbf{a} = (a_1, a_2) \in \mathbb{R}^2$ , denote

$$r_{+}(\boldsymbol{a}) = (\tau_{+}^{2}a_{1}^{2} + a_{2}^{2})^{1/2}, \quad r_{-}(\boldsymbol{a}) = (\tau_{-}^{2}a_{1}^{2} + a_{2}^{2})^{1/2}.$$
 (9.3)

Before analyzing our model, we introduce the function space and space of probability measures we will work on. Let  $E_2 \equiv [0, +\infty)^2 \cup \{\infty\}$ . Note there is a bijection  $\iota$  between  $E_2$  and  $\mathbb{S}^2 \cap \{(x, y, x) \in \mathbb{R}^3 : x, y \geq 0\}$ . Indeed, for any  $\mathbf{r} = (r_1, r_2) \in [0, +\infty)^2$ , consider the line crossing  $(r_1, r_2, 0)$  and (0, 0, 1). This line intersects with  $\mathbb{S}^2$  at two points. One intersection point is (0, 0, 1), and we denote the other intersection point as  $\iota(\mathbf{r})$ . Moreover, let  $\iota(\infty) = (0, 0, 1)$ . With this bijection  $\iota$ , we equip  $E_2$  with a metric  $\bar{d}$  induced by the usual round metric on  $\mathbb{S}^2$ . Then  $(E_2, \bar{d})$  is a compact metric space, and we will still denote it as  $E_2$  for simplicity in notations. We denote  $C_b(E_2)$  to be the set of bounded continuous functions on  $E_2$ , where continuity is defined using the topology generated by  $\bar{d}$ . More explicitly, we have isomorphism

$$C_b(E_2) \simeq \{ f \in C([0,\infty)^2) : \exists f(\infty) \equiv \lim_{\|\boldsymbol{r}\|_2 \to \infty} f(\boldsymbol{r}), \sup_{\boldsymbol{r} \in E_2} f(\boldsymbol{r}) < \infty \}.$$
 (9.4)

Because of condition S2 and S3, we have  $q \circ r_+, q \circ r_-, q' \circ r_+, q' \circ r_- \in C_b(E_2)$ .

Let  $\mathscr{P}(E_2)$  be the set of probability measures on  $E_2$ . Due to Prokhorov's theorem, there exists a complete metric  $\bar{d}_{\mathscr{P}}$  on  $\mathscr{P}(E_2)$  equivalent to the topology of weak convergence, so that  $(\mathscr{P}(E_2), \bar{d}_{\mathscr{P}})$  is a compact metric space. In this section, we will denote by  $\overline{\mathscr{P}} = \mathscr{P}(E_2)$ .

### 9.1 Statics

Since the distribution of x is invariant under rotations in first  $s_0$  coordinates, and invariant under rotations in last  $d - s_0$  coordinates, so are the functions

$$V(\mathbf{a}) = v(\|\mathbf{a}_1\|_2, \|\mathbf{a}_2\|_2), \tag{9.5}$$

$$U(\boldsymbol{a}, \boldsymbol{b}) = u_0(\|\boldsymbol{a}_1\|_2, \|\boldsymbol{b}_1\|_2, \langle \boldsymbol{a}_1, \boldsymbol{b}_1 \rangle, \|\boldsymbol{a}_2\|_2, \|\boldsymbol{b}_2\|_2, \langle \boldsymbol{a}_2, \boldsymbol{b}_2 \rangle). \tag{9.6}$$

These take the form

$$v(a_1, a_2) = -\frac{1}{2}q(r_+(a_1, a_2)) + \frac{1}{2}q(r_-(a_1, a_2)), \quad q(t) = \mathbb{E}\{\sigma(tG)\}$$

and

 $u_0(a_1, b_1, a_1b_1\cos\alpha, a_2, b_2, a_2b_2\cos\beta)$ 

$$= \frac{1}{2} \mathbb{E} \left\{ \sigma(\tau_{+} a_{1} F_{1} + a_{2} G_{1}) \sigma(\tau_{+} b_{1} F_{2} + b_{2} G_{2}) \right\} \\ + \frac{1}{2} \mathbb{E} \left\{ \sigma(\tau_{-} a_{1} F_{1} + a_{2} G_{1}) \sigma(\tau_{-} b_{1} F_{2} + b_{2} G_{2}) \right\},$$

where expectations are with respect to standard normals  $G, F_1, F_2, G_1, G_2 \sim N(0, 1)$ , with  $(F_1, F_2)$  independent of  $(G_1, G_2)$ . Moreover,  $(F_1, F_2)$  are jointly Gaussian,  $(G_1, G_2)$  are jointly Gaussian, and covariance  $\mathbb{E}\{F_1F_2\} = \cos \alpha$ ,  $\mathbb{E}\{G_1G_2\} = \cos \beta$ .

In order to minimize  $R(\rho)$ , it is sufficient to restrict ourselves to distributions that are invariant under product of rotations. Indeed, for any probability distribution  $\rho$  on  $\mathbb{R}^d$ , we can define its symmetrization by letting, for any Borel set  $Q_1 \subseteq \mathbb{R}^{s_0}$ ,  $Q_2 \subseteq \mathbb{R}^{d-s_0}$ ,

$$\rho_s(Q_1 \times Q_2) \equiv \int \rho((\mathbf{R}_1 Q_1) \times (\mathbf{R}_2 Q_2)) \ \mu_{\text{Haar}}(\mathrm{d}\mathbf{R}_1) \mu_{\text{Haar}}(\mathrm{d}\mathbf{R}_2), \qquad (9.7)$$

where  $\mu_{\text{Haar}}$  is the Haar measure over the group of orthogonal rotations. Since  $\rho \mapsto R(\rho)$  is convex,  $R(\rho_s) \leq R(\rho)$ .

We therefore restrict ourselves to  $\rho$ 's that are invariant under product of rotations. In other words, under  $\rho$ , the vector  $\boldsymbol{w} = (\boldsymbol{w}_1, \boldsymbol{w}_2) \in \mathbb{R}^d$  is sampled as following:  $\boldsymbol{w}_1 \in \mathbb{R}^{s_0}$  is uniformly random conditional on  $\|\boldsymbol{w}_1\|_2$ , and  $\boldsymbol{w}_2 \in \mathbb{R}^{d-s_0}$  is uniformly random conditional on  $\|\boldsymbol{w}_2\|_2$ . We denote by  $\overline{\rho} \in \mathscr{P}(E_2)$  the probability distribution of  $(\|\boldsymbol{w}_1\|_2, \|\boldsymbol{w}_2\|_2)$  when  $\boldsymbol{w} \sim \rho$  and we let  $\overline{R}_d(\overline{\rho})$  denote the corresponding risk. We then have

$$\overline{R}_d(\overline{\rho}) = 1 + 2 \int v(r_1, r_2) \,\overline{\rho}(d\mathbf{r}) + \int u_d(a_1, a_2, b_1, b_2) \,\overline{\rho}(d\mathbf{a}) \,\overline{\rho}(d\mathbf{b}), \qquad (9.8)$$

and

$$u_d(a_1, a_2, b_1, b_2) = \mathbb{E}_{\Theta_1, \Theta_2}[u_0(a_1, b_1, a_1b_1 \cos \Theta_1, a_2, b_2, a_2b_2 \cos \Theta_2)], \tag{9.9}$$

where  $\Theta_1 \sim (1/Z_{s_0}) \sin^{s_0-2} \theta \cdot \mathbf{1} \{\theta \in [0, \pi]\} d\theta$  and  $\Theta_2 \sim (1/Z_{d-s_0}) \sin^{d-s_0-2} \theta \cdot \mathbf{1} \{\theta \in [0, \pi]\} d\theta$  are independent.

As  $d \to \infty$ , we have  $\lim_{d \to \infty} u_d(a_1, a_2, b_1, b_2) = u_\infty(a_1, a_2, b_1, b_2)$ , with

$$u_{\infty}(a_1, a_2, b_1, b_2) = \frac{1}{2} \left[ q(r_+(a_1, a_2)) q(r_+(b_1, b_2)) + q(r_-(a_1, a_2)) q(r_-(b_1, b_2)) \right], \tag{9.10}$$

and the risk function converges to (for  $\mathbf{a} = (a_1, a_2)$ )

$$\overline{R}_{\infty}(\overline{\rho}) = \frac{1}{2} \left( 1 - \int q(r_{+}(\boldsymbol{a})) \, \overline{\rho}(\mathrm{d}\boldsymbol{a}) \right)^{2} + \frac{1}{2} \left( 1 + \int q(r_{-}(\boldsymbol{a})) \, \overline{\rho}(\mathrm{d}\boldsymbol{a}) \right)^{2}. \tag{9.11}$$

We also define

$$\psi_d(\boldsymbol{a}; \overline{\rho}) = v(\boldsymbol{a}) + \int u_d(\boldsymbol{a}, \boldsymbol{b}) \, \overline{\rho}(\mathrm{d}\boldsymbol{b}) \,. \tag{9.12}$$

For  $s_0 = \gamma \cdot d$  with  $0 < \gamma < 1$  and  $d \to \infty$ , we have the simpler expression

$$\psi_{\infty}(\boldsymbol{a}; \overline{\rho}) = \lambda_{+}(\overline{\rho}) \cdot q(r_{+}(\boldsymbol{a})) + \lambda_{-}(\overline{\rho}) \cdot q(r_{-}(\boldsymbol{a})), \tag{9.13}$$

$$\lambda_{+}(\overline{\rho}) = \frac{1}{2} [\langle q \circ r_{+}, \overline{\rho} \rangle - 1], \tag{9.14}$$

$$\lambda_{-}(\overline{\rho}) = \frac{1}{2} [\langle q \circ r_{-}, \overline{\rho} \rangle + 1]. \tag{9.15}$$

The following theorem provides a characterization of the global minimizers of  $\overline{R}_{\infty}(\overline{\rho})$ .

**Theorem 11.** Consider  $d = \infty$ . Recall  $\overline{\mathscr{P}} = \mathscr{P}(E_2)$  where  $E_2 \equiv [0, +\infty)^2 \cup \{\infty\}$ . Then there exists  $\Delta_{\infty} \in (0, 1)$ , such that

- 1. For  $\Delta < \Delta_{\infty}$ ,  $\inf_{\overline{\rho} \in \overline{\mathscr{P}}} \overline{R}_{\infty}(\overline{\rho}) > 0$  and the unique global minimizer of risk function  $\overline{R}_{\infty}(\overline{\rho})$  is a point mass located at  $(r_*,0)$  for some  $r_* = r_*(\Delta) \in (0,\infty)$ .
- 2. For  $\Delta \geq \Delta_{\infty}$ , all global minimizers of risk function  $\overline{R}_{\infty}(\overline{\rho})$  have risk zero, and there exists a global minimizer that has finite support.

Proof. Throughout the proof, we will denote  $\overline{R}_{\infty}^{(1)}: \mathscr{P}([0,\infty]) \to \mathbb{R}$  as the risk function defined as in Eq. (8.10), and  $\overline{R}_{\infty}^{(2)}: \mathscr{P}(E_2) \to \mathbb{R}$  as the risk function defined as in Eq. (9.11). Recall the definition  $\tau_+ = 1 + \Delta$ ,  $\tau_- = 1 - \Delta$ ,  $q_+(r) = q(\tau_+ r)$ ,  $q_-(r) = q(\tau_- r)$ ,  $r_+(\boldsymbol{a}) = (\tau_+^2 a_1^2 + a_2^2)^{1/2}$ , and  $r_-(\boldsymbol{a}) = (\tau_-^2 a_1^2 + a_2^2)^{1/2}$  for  $\boldsymbol{a} = (a_1, a_2) \in E_2$ .

Suppose  $\overline{\rho}_2^* \in \arg\min_{\overline{\rho}_2 \in \mathscr{P}(E_2)} \overline{R}_{\infty}^{(2)}(\overline{\rho}_2)$ . Then we must have  $\langle q \circ r_+, \overline{\rho}_2^* \rangle \leq 1$  and  $\langle q \circ r_-, \overline{\rho}_2^* \rangle \geq -1$ . Indeed, if either  $\langle q \circ r_+, \overline{\rho}_2^* \rangle > 1$  or  $\langle q \circ r_-, \overline{\rho}_2^* \rangle < -1$ , since  $q(+\infty) > 1$  and q(0) < -1, the distribution  $\overline{\rho}_2' = a_0 \delta_0 + a_\infty \delta_\infty + (1 - a_0 - a_\infty) \overline{\rho}_2^*$  with appropriate choice of  $a_0$  and  $a_\infty$  will give a lower risk.

This  $\overline{\rho}_2^* \in \mathscr{P}(E_2)$  induces a  $\overline{\rho}_1 \in \mathscr{P}([0,\infty])$  as follows: for any Borel set  $B \subseteq [0,\infty]$ ,  $\overline{\rho}_1(B) = \overline{\rho}_2^*(\{r \in E_2 : \|r\|_2 \in B\})$ . For this  $\overline{\rho}_1$ , it is easy to see that  $\langle q_-, \overline{\rho}_1 \rangle \leq \langle q \circ r_-, \overline{\rho}_2^* \rangle$  and  $\langle q_+, \overline{\rho}_1 \rangle \geq \langle q \circ r_+, \overline{\rho}_2^* \rangle$ , and the equalities hold if and only if  $\overline{\rho}_2^*(E_1) = 1$ , where  $E_1 \equiv ([0, +\infty) \times \{0\}) \cup \{\infty\}$ . Since  $q(+\infty) > 1$  and q(0) < -1, we can take  $\overline{\rho}_1^* = a_0 \delta_0 + a_\infty \delta_\infty + (1 - a_0 - a_\infty) \overline{\rho}_1$  with appropriate choice of  $a_0$  and  $a_\infty$ , so that  $\langle q \circ r_+, \overline{\rho}_2^* \rangle \leq \langle q_+, \overline{\rho}_1^* \rangle \leq 1$  and  $\langle q \circ r_-, \overline{\rho}_2^* \rangle \geq \langle q_-, \overline{\rho}_1^* \rangle \geq -1$ . Therefore, we always have  $\inf_{\overline{\rho}_1 \in \mathscr{P}([0,\infty])} \overline{R}_\infty^{(1)}(\overline{\rho}_1) \leq \inf_{\overline{\rho}_2 \in \mathscr{P}(E_2)} \overline{R}_\infty^{(2)}(\overline{\rho}_2)$ , and  $\overline{\rho}_2^*(E_1) = 1$  for any  $\overline{\rho}_2^* \in \arg \min_{\overline{\rho}_2 \in \mathscr{P}(E_2)} \overline{R}_\infty^{(2)}(\overline{\rho}_2)$ . Note that  $\overline{R}_\infty^{(2)}(\overline{\rho}_1 \times \delta_0) = \overline{R}_\infty^{(1)}(\overline{\rho}_1)$  for any  $\overline{\rho}_1 \in \mathscr{P}([0,\infty])$ . Hence, we must have  $\inf_{\overline{\rho}_1 \in \mathscr{P}([0,\infty])} \overline{R}_\infty^{(1)}(\overline{\rho}_1) = \inf_{\overline{\rho}_2 \in \mathscr{P}(E_2)} \overline{R}_\infty^{(2)}(\overline{\rho}_2)$ .

Due to the above argument, we reduced our analysis to the centered isotropic Gaussians case. All the conclusions can be proved using the same argument as in the proof of Theorem 8.

### 9.2 Dynamics: Fixed points

We specialize the general evolution (7.1) to the present case. Assuming  $\rho_0$  to be invariant with respect to products of orthogonal transformations, the same happens for  $\rho_t$ . We let  $\overline{\rho}_t \in \mathscr{P}(E_2)$  denote the distribution of  $(\|\boldsymbol{w}_1\|_2, \|\boldsymbol{w}_2\|_2)$  when  $\boldsymbol{w} \sim \rho_t$ . Then  $\overline{\rho}_t$  satisfies the following PDE:

$$\partial_t \overline{\rho}_t(\mathbf{r}) = 2\xi(t)\nabla \cdot \left[\overline{\rho}_t(\mathbf{r})\nabla \psi_d(\mathbf{r}; \overline{\rho}_t)\right]. \tag{9.16}$$

We will view this as an evolution in the space of probability distribution on  $\overline{\mathscr{P}} = \mathscr{P}(E_2)$ . In analogy with Proposition 2, we can prove the following characterization of fixed points.

**Proposition 6.** A distribution  $\overline{\rho} \in \overline{\mathscr{P}}$  is a fixed point of the PDE (9.16) if and only if

$$\operatorname{supp}(\overline{\rho}) \subseteq \{ r \in E_2 : \nabla_r \psi_d(r; \overline{\rho}) = \mathbf{0} \}. \tag{9.17}$$

Notice, in particular, global minimizers of  $\overline{R}_d(\overline{\rho})$  are fixed points of this evolution, but not vice-versa. The next result classifies fixed points.

**Theorem 12.** Consider  $d = \infty$ , and recall the definition of  $\lambda_{+}(\overline{\rho})$  and  $\lambda_{-}(\overline{\rho})$  given by Eq. (9.15) and (9.14). Then the fixed points of the PDE (9.16) (i.e. the probability measures  $\overline{\rho} \in \overline{\mathscr{P}}$  satisfying (9.17)) must be of one of the following types

- (a) A fixed point with zero risk.
- (b) A point mass  $\overline{\rho}_{r_*} = \delta_{(r_*,0)}$  at some location  $(r_*,0)$  with  $r_* \notin \{0,+\infty\}$ , but not of type (a).
- (c) A mixture of the type  $\overline{\rho} = a_0 \delta_0 + a_\infty \delta_\infty + a_1 \delta_{(r_{*1},0)} + a_2 \overline{\rho}_2$  with  $\operatorname{supp}(\overline{\rho}_2) \subseteq \{0\} \times (0,\infty)$ , but not of type (b) and (a).

For  $\Delta < \Delta_{\infty}$ , the PDE has a unique fixed point of type (b), with  $\lambda_{+}(\overline{\rho}_{*}) < 0$  and  $\lambda_{-}(\overline{\rho}_{*}) > 0$ ; it has no type-(a) fixed points; it has possibly fixed points of type (c).

For  $\Delta > \Delta_{\infty}$ , the PDE has some fixed points of type (b), with  $\lambda_{+}(\overline{\rho}_{*}) > 0$  and  $\lambda_{-}(\overline{\rho}_{*}) < 0$ ; it also has some type-(a) fixed points; it has possibly fixed points of type (c).

For  $\Delta = \Delta_{\infty}$ , the PDE has a unique fixed point of type (a) which is also a delta function at some location  $(r_{*1}, 0)$ , and no type (b) fixed points; it has possibly fixed points of type (c).

*Proof.* We use the characterization of fixed points in Proposition 6. Recall that  $\psi_{\infty}(\mathbf{r}; \overline{\rho}_{*})$  is defined as in Eq. (9.13). The gradient  $\nabla \psi_{\infty}(\mathbf{r}; \overline{\rho})$  is given by

$$\partial_{r_1} \psi_{\infty}(\boldsymbol{r}; \overline{\rho}) = \lambda_+(\overline{\rho}) q'(r_+(\boldsymbol{r})) \tau_+^2 r_1 / r_+(\boldsymbol{r}) + \lambda_-(\overline{\rho}) q'(r_-(\boldsymbol{r})) \tau_-^2 r_1 / r_-(\boldsymbol{r}), 
\partial_{r_2} \psi_{\infty}(\boldsymbol{r}; \overline{\rho}) = \lambda_+(\overline{\rho}) q'(r_+(\boldsymbol{r})) r_2 / r_+(\boldsymbol{r}) + \lambda_-(\overline{\rho}) q'(r_-(\boldsymbol{r})) r_2 / r_-(\boldsymbol{r}).$$
(9.18)

If a fixed point  $\overline{\rho}_*$  gives  $\lambda_+(\overline{\rho}_*) = \lambda_-(\overline{\rho}_*) = 0$ , then  $\overline{R}_{\infty}(\overline{\rho}_*) = 0$ . This is type-(a) fixed point. Consider then the case  $(\lambda_+(\overline{\rho}_*), \lambda_-(\overline{\rho}_*)) \neq (0, 0)$ .

Suppose  $\overline{\rho}_*((0,+\infty)^2) > 0$ . Since q'(r) > 0 and  $\tau_+ > 1 > \tau_-$ , in order for  $\nabla \psi_\infty(r; \overline{\rho}_*) = \mathbf{0}$  for some  $r \in (0,+\infty)^2$ , we must have  $(\lambda_+(\overline{\rho}_*),\lambda_-(\overline{\rho}_*)) = (0,0)$ . Therefore, as  $\overline{\rho}_*$  is a fixed point with  $(\lambda_+(\overline{\rho}_*),\lambda_-(\overline{\rho}_*)) \neq (0,0)$ , we must have  $\overline{\rho}_*((0,+\infty)^2) = 0$ . That is, we can write  $\overline{\rho}_* = a_0\delta_0 + a_\infty\delta_\infty + a_1\overline{\rho}_1 + a_2\overline{\rho}_2$ , with  $\sup(\overline{\rho}_1) \in (0,\infty) \times \{0\}$ , and  $\sup(\overline{\rho}_2) \in \{0\} \times (0,\infty)$ .

The solutions of  $\nabla \psi_{\infty}((r_1, r_2); \overline{\rho}_*) = 0$  with  $r_2 = 0$  are of the form  $\mathbf{0}$ ,  $(r_{*1}, 0)$ , and  $\infty$ . Therefore,  $\overline{\rho}_1 = \delta_{(r_{*1}, 0)}$  for some  $r_{*1} \in (0, \infty)$ . Hence, as  $\overline{\rho}_*$  is not a type-(a) stationary point, it must be a type-(b) or type-(c) stationary point.

This proves that all fixed points are of type (a), (b), or (c). The remaining claims follows the same argument as the proof of Theorem 9.

# 9.3 Dynamics: Convergence to global minimum for $d = \infty$

In this section, denote  $\mathscr{P}_{good}$  to be

$$\mathscr{P}_{\text{good}} = \{ \overline{\rho}_0 \in \mathscr{P}((0, \infty)^2) : \overline{R}_{\infty}(\overline{\rho}_0) < 1 \}. \tag{9.19}$$

We then prove that the  $d=\infty$  dynamics converges to a global minimizer from any initialization  $\overline{\rho}_0\in\mathscr{P}_{\mathrm{good}}$ .

**Theorem 13.** Consider the PDE (9.16) for  $d = \infty$ , with initialization  $\overline{\rho}_0 \in \mathscr{P}_{good}$ . It has a unique solution  $(\overline{\rho}_t)_{t\geq 0}$ , such that

$$\lim_{t\to +\infty} \overline{R}_{\infty}(\overline{\rho}_t) = \inf_{\overline{\rho}\in \overline{\mathscr{P}}} \overline{R}_{\infty}(\overline{\rho}).$$

*Proof.* Without loss of generality, we assume  $\xi(t) = 1/2$ . First we show the existence and uniqueness of solution of the PDE.

Step 1. Existence and uniqueness of solution. Mass  $\overline{\rho}_t((0,\infty)^2) = 1$  for all t.

According to conditions S1 - S3, q(r), q'(r), and q''(r) are uniformly bounded on  $[0,\infty]$ . Note

$$\begin{split} v(\boldsymbol{r}) = & 1/2 \cdot [q(r_{-}(\boldsymbol{r})) - q(r_{+}(\boldsymbol{r}))], \\ u_{\infty}(\boldsymbol{r}_{1}, \boldsymbol{r}_{2}) = & 1/2 \cdot [q(r_{+}(\boldsymbol{r}_{1}))q(r_{+}(\boldsymbol{r}_{2})) + q(r_{-}(\boldsymbol{r}_{1}))q(r_{-}(\boldsymbol{r}_{2}))]. \end{split}$$

Then  $\nabla v(\boldsymbol{r})$ ,  $\nabla_1 u_{\infty}(\boldsymbol{r}_1, \boldsymbol{r}_2)$ ,  $\nabla^2 v(\boldsymbol{r})$ ,  $\nabla^2_{11} u_{\infty}(\boldsymbol{r}_1, \boldsymbol{r}_2)$ ,  $\nabla^2_{12} u_{\infty}(\boldsymbol{r}_1, \boldsymbol{r}_2)$  are uniformly bounded. Therefore, conditions A1 and A3 are satisfied with D=2, V=v, and U=u. Then, there is the existence and uniqueness of solution of PDE (9.16) for  $d=\infty$ . Denote this solution to be  $(\overline{\rho}_t)_{t>0}$ .

Recall the expression for  $\nabla \psi_{\infty}(\mathbf{r}; \overline{\rho})$  in Eq. (9.18). It is easy to see that the assumption of Lemma 7.9 is satisfied with d=2 and  $\Psi=\psi_{\infty}$ . Hence, we have  $\overline{\rho}_t((0,\infty)^2)=1$  for any fixed  $t<\infty$ .

# Step 2. Classify the limiting set $S_*$ .

Recall the definition of  $(\mathscr{P}(E_2), \bar{d}_{\mathscr{P}})$  at the beginning of Section 9. Since  $(\mathscr{P}(E_2), \bar{d}_{\mathscr{P}})$  is a compact metric space, and  $(\overline{\rho}_t)_{t\geq 0}$  is a continuous curve in this space, then there exists a subsequence  $(t_k)_{k\geq 1}$  of times, such that  $(\overline{\rho}_{t_k})_{k\geq 1}$  converges in metric  $\bar{d}_{\mathscr{P}}$  to a probability distribution  $\overline{\rho}_* \in \mathscr{P}(E_2)$ . For any  $\overline{\rho}_0 \in \mathscr{P}_{\text{good}}$ , let  $\mathcal{S}_* = \mathcal{S}_*(\overline{\rho}_0)$  be the set of limiting points of the PDE,

$$\mathcal{S}_* = \{ \overline{\rho}_* \in \mathscr{P}(E_2) : \exists (t_k)_{k \ge 1}, \lim_{k \to \infty} t_k = +\infty, s.t., \lim_{k \to \infty} \overline{d}_{\mathscr{P}}(\overline{\rho}_*, \overline{\rho}_{t_k}) = 0 \}.$$

Analogous to the proof of Theorem 10, we have the following properties for  $S_*$ :

- 1.  $S_*$  is connected and compact.
- 2. For any  $\overline{\rho}_* \in \mathcal{S}_*$ ,  $\overline{\rho}_*$  is a fixed point of PDE.
- 3. For any  $\overline{\rho}_* \in \mathcal{S}_*$ ,  $\overline{R}_{\infty}(\overline{\rho}_*) = \overline{R}_* < 1$ .

Recall the definition of  $\lambda_{+}(\overline{\rho}_{*})$  and  $\lambda_{-}(\overline{\rho}_{*})$  given by Equation (9.14) and (9.15). Let  $\overline{\rho}_{*}$  be a fixed point of PDE such that  $\lambda_{+}(\overline{\rho}_{*}) \geq 0$ ,  $\lambda_{-}(\overline{\rho}_{*}) \geq 0$  or  $\lambda_{+}(\overline{\rho}_{*}) \leq 0$ ,  $\lambda_{-}(\overline{\rho}_{*}) \leq 0$  but not both  $\lambda_{+}(\overline{\rho}_{*})$  and  $\lambda_{-}(\overline{\rho}_{*})$  equal 0. In this case, according to Eq. (9.18), both  $\partial_{r_{1}}\psi_{\infty}(r;\overline{\rho}_{*})$  and  $\partial_{r_{2}}\psi_{\infty}(r;\overline{\rho}_{*})$  must be strictly positive or strictly negative. Since  $\operatorname{supp}(\overline{\rho}_{*}) \subseteq \{r \in E_{2} : \nabla_{r}\psi_{\infty}(r;\overline{\rho}_{*}) = \mathbf{0}\}$ ,  $\overline{\rho}_{*}$  must be a combination of two delta functions located at  $\mathbf{0}$  and  $\infty$ , i.e.,  $\overline{\rho}_{*} = a_{0}\delta_{\mathbf{0}} + (1 - a_{0})\delta_{\infty}$ . But for a fixed point like this, it is easy to see that  $\overline{R}_{\infty}(\overline{\rho}_{*}) \geq 1$ . Such fixed points  $\overline{\rho}_{*}$  cannot be one of the limiting points of the PDE since  $\overline{R}_{\infty}(\overline{\rho}_{0}) < 1$ .

Let L be a mapping  $L: \mathscr{P}(E_2) \to \mathbb{R}^2$ ,  $\overline{\rho} \mapsto (\lambda_+(\overline{\rho}), \lambda_-(\overline{\rho}))$ . The above argument concludes that for any  $\overline{\rho}_0 \in \mathscr{P}_{\text{good}}$ , we have

$$L(\mathcal{S}_*(\overline{\rho}_0)) \cap (\{(\lambda_+, \lambda_-) : \lambda_+ \ge 0, \lambda_- \ge 0, \text{ or } \lambda_+ \le 0, \lambda_- \le 0\} \setminus \{(0, 0)\}) = \emptyset.$$

Since  $S_*$  is a connected set,  $L(S_*)$  should also be a connected set. Further notice that  $\overline{R}_{\infty}(\overline{\rho}_*) = 1/2 \cdot [\lambda_+(\overline{\rho}_*)^2 + \lambda_-(\overline{\rho}_*)^2]$ , and  $\overline{R}_{\infty}(\overline{\rho}_1) = \overline{R}_{\infty}(\overline{\rho}_2)$  for any  $\overline{\rho}_1, \overline{\rho}_2 \in S_*$ . Therefore, we can only have  $L(S_*) \subseteq \mathcal{P}_2 \equiv \{(\lambda_+, \lambda_-) : \lambda_+ > 0, \lambda_- < 0\}$ , or  $L(S_*) \subseteq \mathcal{P}_1 \equiv \{(\lambda_+, \lambda_-) : \lambda_+ < 0, \lambda_- > 0\}$ , or  $L(S_*) = \{(0, 0)\}$ .

#### Step 3. Finish the proof using two claims.

We make the following two claims.

Claim (1). If 
$$L(S_*) \subseteq \mathcal{P}_1$$
, then for any  $\overline{\rho}_* \in S_*$ , we have  $\overline{\rho}_*((0,\infty) \times \{0\}) = 1$ .

Claim (2). We cannot have  $L(\mathcal{S}_*) \subseteq \mathcal{P}_2$ .

Here we assume these two claims holds, and use it to prove our results. For  $\Delta < \Delta_{\infty}$ , we proved in Theorem 12 that, there is no fixed point such that  $L(\overline{\rho}_*) = (0,0)$ . Therefore, we cannot have  $L(\mathcal{S}_*) = \{(0,0)\}$ . Due to Claim (2), we cannot have  $L(\mathcal{S}_*) \subseteq \mathcal{P}_2$ . Hence, we must have  $L(\mathcal{S}_*) \subseteq \mathcal{P}_1$ .

According to Theorem 12, for  $\Delta < \Delta_{\infty}$ , the only fixed point of PDE with  $\overline{\rho}_*((0,\infty) \times \{0\}) = 1$  is a point mass at some location  $\mathbf{r}_* = (r_{*1}, 0)$ . Furthermore, this delta function fixed point is unique and is also the global minimizer of the risk. Therefore, we conclude that, for  $\Delta < \Delta_{\infty}$ , the PDE will converge to this global minimizer.

For  $\Delta \geq \Delta_{\infty}$ , according to Claim (1), if  $\overline{\rho}_*$  is a limiting point such that  $L(\overline{\rho}_*) \in \mathcal{P}_1$ , then  $\overline{\rho}_*((0,\infty) \times \{0\}) = 1$ . According to Theorem 12, a fixed point  $\overline{\rho}_*$  with  $\overline{\rho}_*((0,\infty) \times \{0\}) = 1$  and  $L(\overline{\rho}_*) \neq (0,0)$  must be a point mass at some location  $r_* = (r_{*1},0)$ , with  $L(\overline{\rho}_*) \in \mathcal{P}_2$ . Therefore, we cannot have  $L(\mathcal{S}_*) \subseteq \mathcal{P}_1$ . Claim (2) also tells us that we cannot have  $L(\mathcal{S}_*) \subseteq \mathcal{P}_2$ . Hence, we must have  $L(\mathcal{S}_*) = \{(0,0)\}$ . In this case, all the points in the set  $\mathcal{S}_*$  have risk 0. Therefore, we conclude that, as  $\Delta \geq \Delta_{\infty}$ , the PDE will converge to some limiting set with risk 0.

#### Step 4. Proof of the two claims.

We are left with the task of proving the two claims above. Before that, we introduce some useful notions used in the proof. Define Z(r) for  $r \in E_2$ ,

$$Z(\mathbf{r}) \equiv [q'(r_{-}(\mathbf{r}))r_{-}(\mathbf{r})]/[q'(r_{+}(\mathbf{r}))r_{+}(\mathbf{r})]. \tag{9.20}$$

Define  $Z_l(r) \equiv Z((r, lr))$  for  $r, l \in [0, \infty]$ . Then we have

$$Z_l(r) = \left[ q'((\tau_-^2 + l^2)^{1/2}r) / q'((\tau_+^2 + l^2)^{1/2}r) \right] \cdot \left[ (\tau_-^2 + l^2)^{1/2} / (\tau_+^2 + l^2)^{1/2} \right]. \tag{9.21}$$

According to condition S4, for any fixed  $l \in [0, \infty]$ ,  $Z_l(r)$  is increasing in r.

Recall the formula of  $\nabla_{\mathbf{r}}\psi_{\infty}(\mathbf{r};\overline{\rho})$  given by Equation (9.18). Define

$$\chi_{\text{nm}}(\mathbf{r}; \overline{\rho}) \equiv \langle \nabla_{\mathbf{r}} \psi_{\infty}(\mathbf{r}; \overline{\rho}), \mathbf{r} / || \mathbf{r} ||_{2} \rangle, \tag{9.22}$$

$$\chi_{\rm tg}(\mathbf{r};\overline{\rho}) \equiv \langle \nabla_{\mathbf{r}} \psi_{\infty}(\mathbf{r};\overline{\rho}), (-r_2, r_1) / \|\mathbf{r}\|_2 \rangle. \tag{9.23}$$

Then we have

$$\chi_{\text{nm}}(\boldsymbol{r};\overline{\rho}) = \lambda_{+}(\overline{\rho})q'(r_{+}(\boldsymbol{r}))r_{+}(\boldsymbol{r})/\|\boldsymbol{r}\|_{2} + \lambda_{-}(\overline{\rho})q'(r_{-}(\boldsymbol{r}))r_{-}(\boldsymbol{r})/\|\boldsymbol{r}\|_{2},$$

$$= \lambda_{-}(\overline{\rho})q'(r_{+}(\boldsymbol{r}))r_{+}(\boldsymbol{r})/\|\boldsymbol{r}\|_{2} \cdot [\lambda_{+}(\overline{\rho})/\lambda_{-}(\overline{\rho}) + Z(\boldsymbol{r})],$$

$$(9.24)$$

and

$$\chi_{\text{tg}}(\boldsymbol{r}; \overline{\rho}) = [+\lambda_{+}(\overline{\rho})(1 - \tau_{+}^{2})q'(r_{+}(\boldsymbol{r}))/r_{+}(\boldsymbol{r}) + \lambda_{-}(\overline{\rho})(1 - \tau_{-}^{2})q'(r_{-}(\boldsymbol{r}))/r_{-}(\boldsymbol{r})] \times r_{1}r_{2}/\|\boldsymbol{r}\|_{2}.$$

$$(9.25)$$

**Proof of Claim** (1). If  $L(S_*) \subseteq \mathcal{P}_1$ , then for any  $\overline{\rho}_* \in S_*$ , we have  $\overline{\rho}_*((0,\infty) \times \{0\}) = 1$ .

Assume  $L(S_*) \subseteq \mathcal{P}_1$ . There must exist  $t_0$  large enough, so that as  $t \geq t_0$ , we have  $\lambda_+(\overline{\rho}_t) < 0$ , and  $\lambda_-(\overline{\rho}_t) > 0$ . Therefore, we must have  $\chi_{\rm tg}(\boldsymbol{r}; \overline{\rho}_t) > 0$  for any  $\boldsymbol{r} \in (0, \infty)^2$ . We denote

$$\Gamma_k \equiv \{ r \in [0, \infty)^2 : r_2 \le k \cdot r_1 \}.$$
 (9.26)

Consider the ODE

$$\dot{\boldsymbol{r}}(t) = -\nabla_{\boldsymbol{r}}\psi_{\infty}(\boldsymbol{r}(t); \overline{\rho}_t), \tag{9.27}$$

starting with  $r(t_0) \in \Gamma_k$  for some  $k \in (0, \infty)$ , we claim  $r(t) \in \Gamma_k$  for any  $t \geq t_0$ . Indeed, for any  $r \in \partial \Gamma_k \cap \{r : r_2 = kr_1 > 0\}$ , its normal vector pointing outside  $\Gamma_k$  gives  $n(r) = (-r_2, r_1)/||r||_2$ , and hence  $\langle \nabla_r \psi_\infty(r; \overline{\rho}), n(r) \rangle = \chi_{tg}(r; \overline{\rho}_t) > 0$ . Therefore, r(t) cannot leak outside  $\Gamma_k$  from this

boundary. Further note that r(t) cannot reach the boundary  $([0,\infty)\times\{0\})\cup\{\infty\}$  for any finite time t. This proves the claim that  $r(t)\in\Gamma_k$  for any  $t\geq t_0$ .

According to Lemma 7.8, we have  $\rho_t(\Gamma_k) \geq \rho_{t_0}(\Gamma_k)$  for any  $k \in (0, \infty)$ . Furthermore, according to Lemma 7.9,  $\overline{\rho}_{t_0}((0, \infty)^2) = 1$ , hence  $\lim_{k \to \infty} \overline{\rho}_{t_0}(\Gamma_k) = 1$ . Therefore, for any  $\overline{\rho}_* \in \mathcal{S}_*$ , we must have

$$\overline{\rho}_*(\{0\} \times (0, \infty)) \le \lim_{k \to \infty} \overline{\rho}_*([0, \infty)^2 \setminus \Gamma_k) \le \lim_{k \to \infty} \overline{\rho}_{t_0}([0, \infty)^2 \setminus \Gamma_k) = 0. \tag{9.28}$$

Theorem 12 implies that for any such fixed point  $\overline{\rho}_*$ , we have  $\operatorname{supp}(\overline{\rho}_*) \subseteq ([0,\infty) \times \{0\}) \cup \{\infty\}$ .

In this case, we claim  $L(S_*) \subseteq \mathcal{P}_1 \cap \{(\lambda_+, \lambda_-) : Z_0(0) < -\lambda_+/\lambda_- < Z_0(\infty)\}$ . Indeed, suppose there exists  $\overline{\rho}_* \in S_*$ , such that  $-\lambda_+(\overline{\rho}_*)/\lambda_-(\overline{\rho}_*) \ge Z_0(\infty)$  or  $-\lambda_-(\overline{\rho}_*)/\lambda_-(\overline{\rho}_*) \le Z_0(0)$ , according to Equation (9.24),  $\chi_{\text{nm}}((r,0);\overline{\rho}_*)$  must be strictly positive or strictly negative. However, we know  $\sup(\overline{\rho}_*) \in \{r : \nabla \psi_\infty(r;\overline{\rho}_*) = \mathbf{0}\}$ . Hence,  $\overline{\rho}_*$  should be a combination of two delta functions located at  $\mathbf{0}$  and  $\infty$ . Such fixed point  $\overline{\rho}_*$  has risk  $\overline{R}_\infty(\overline{\rho}_*) \ge 1$ , hence  $\overline{\rho}_*$  cannot be a limiting point of the PDE. Hence the claim holds.

Since  $S_*$  is a compact set, and L is a continuous map, then  $L(S_*)$  is a compact set. Therefore, there must exist  $\varepsilon_0 > 0$ , so that for any  $\overline{\rho}_* \in S_*$ , we have  $Z_0(0) + 3\varepsilon_0 < -\lambda_+(\overline{\rho}_*)/\lambda_-(\overline{\rho}_*) < Z_0(\infty) - 3\varepsilon_0$ . For this  $\varepsilon_0 > 0$ , we take  $t_0$  large enough, so that for  $t \geq t_0$ , we have  $Z_0(0) + 2\varepsilon_0 < -\lambda_+(\overline{\rho}_t)/\lambda_-(\overline{\rho}_t) < Z_0(\infty) - 2\varepsilon_0$ , and  $\lambda_+(\overline{\rho}_t) < 0$ ,  $\lambda_-(\overline{\rho}_t) > 0$ .

According to the conditions S0 - S4 on q(r), for any fixed l,  $Z_l(r)$  is an increasing function of r, and for any fixed r,  $Z_l(r)$  is continuous in l. Therefore, for the fixed  $\varepsilon_0 > 0$ , there exists  $0 < r_0 < r_\infty < \infty$  and b > 0, such that

$$\sup_{r \in [0, r_0]} \sup_{l \in [0, b]} Z_l(r) < Z_0(0) + \varepsilon_0, \tag{9.29}$$

$$\inf_{r \in [r_{\infty}, \infty]} \inf_{l \in [0, b]} Z_l(r) > Z_0(\infty) - \varepsilon_0. \tag{9.30}$$

As a result, for any  $t \ge t_0$ , we have

$$\chi_{\text{nm}}(\boldsymbol{r}; \overline{\rho}_t) < 0, \quad \forall \boldsymbol{r} \in \mathsf{B}(\boldsymbol{0}; r_0) \cap \Gamma_b,$$
  
$$\chi_{\text{nm}}(\boldsymbol{r}; \overline{\rho}_t) > 0, \quad \forall \boldsymbol{r} \in \mathsf{B}(\boldsymbol{0}; r_\infty)^c \cap \Gamma_b,$$

$$(9.31)$$

where  $\Gamma_{(\cdot)}$  is defined as in Equation (9.26).

According to Lemma 7.9,  $\overline{\rho}_{t_0}((0,\infty)^2)=1$ . Define

$$O_k = \Gamma_k \cap \mathsf{B}(\mathbf{0}; k) \cap \mathsf{B}(\mathbf{0}; 1/k)^c. \tag{9.32}$$

We have  $O_k$  is increasing in k, and  $\bigcup_k O_k \supset (0, \infty)^2$ . Hence  $\lim_{k\to\infty} \overline{\rho}_{t_0}(O_k) = 1$ . Now we fix a parameter k.

Recall the formula for  $\chi_{\text{nm}}$  and  $\chi_{\text{tg}}$  given by Equation (9.24) and (9.25). It is easy to see that, there exists  $0 < u_{k1}, u_{k2} < \infty$  depending on  $(b, k, \tau_+, \tau_-, Z_0(0), Z_0(\infty), \varepsilon_0)$ , such that for any  $\mathbf{r} \in (0, \infty)^2$  with  $b \cdot r_1 \leq r_2 \leq k \cdot r_1$ , and  $t \geq t_0$ , we have

$$\chi_{\text{tg}}(\mathbf{r}; \overline{\rho}_t) \ge u_{k1} |\lambda_+(\overline{\rho}_t)| q'(r_+(\mathbf{r})) > 0, \tag{9.33}$$

$$|\chi_{\text{nm}}(\mathbf{r}; \overline{\rho}_t)| \le u_{k2} |\lambda_+(\overline{\rho}_t)| q'(r_+(\mathbf{r})) < \infty, \tag{9.34}$$

and hence

$$|\chi_{\rm nm}(\mathbf{r}; \overline{\rho}_t)|/\chi_{\rm tg}(\mathbf{r}; \overline{\rho}_t) \le u_{k2}/u_{k1} \equiv u_k < \infty.$$
 (9.35)

Consider the following spiral curve  $\mathbf{r}_k^{\infty}(s) = (r_{k1}^{\infty}(s), r_{k2}^{\infty}(s))$ , with

$$r_{k1}^{\infty}(s) = k \cdot \cos(\arctan(k) - s) \exp\{2u_k s\},$$
  

$$r_{k2}^{\infty}(s) = k \cdot \sin(\arctan(k) - s) \exp\{2u_k s\},$$
(9.36)

and another spiral curve  $\mathbf{r}_k^0(s) = (r_{k1}^0(s), r_{k2}^0(s))$ , with

$$r_{k1}^{0}(s) = 1/k \cdot \cos(\arctan(k) - s) \exp\{-2u_{k}s\},\$$

$$r_{k2}^{0}(s) = 1/k \cdot \sin(\arctan(k) - s) \exp\{-2u_{k}s\},\$$
(9.37)

for  $s \in [0, s_{k*}]$  with  $s_{k*} = \arctan(k) - \arctan(b)$ .

Because of inequality (9.35), along the curve  $\mathbf{r}_k^{\infty}(s)$ , denoting  $\mathbf{n}(\mathbf{r}_k^{\infty}(s))$  to be its normal vector with  $[\mathbf{n}(\mathbf{r}_k^{\infty}(s))]_2 > 0$ , we have for any  $t \geq t_0$  and  $s \in [0, s_{k*}]$ ,

$$\langle \nabla \psi_{\infty}(\boldsymbol{r}_{k}^{\infty}(s); \overline{\rho}_{t}), \boldsymbol{n}(\boldsymbol{r}_{k}^{\infty}(s)) \rangle > 0.$$
 (9.38)

Along the curve  $r_k^0(s)$ , denoting  $n(r_k^0(s))$  to be its normal vector with  $[n(r_k^0(s))]_2 > 0$ , we have for any  $t \ge t_0$  and  $s \in [0, s_{k*}]$ ,

$$\langle \nabla \psi_{\infty}(\boldsymbol{r}_{k}^{0}(s); \overline{\rho}_{t}), \boldsymbol{n}(\boldsymbol{r}_{k}^{0}(s)) \rangle > 0,$$
 (9.39)

Define the set  $\Omega_k$  to be

$$\Omega_{k} = \Gamma_{k} \cap \mathsf{B}(\mathbf{0}; k \cdot \exp\{2u_{k}s_{k*}\}) \cap \mathsf{B}(\mathbf{0}; 1/k \cdot \exp\{-2u_{k}s_{k*}\})^{c} 
\cap \{\boldsymbol{r} : \exists s \in [0, s_{k*}], s.t., r_{1} = r_{k1}^{\infty}(s), r_{2} \geq r_{k2}^{\infty}(s)\}^{c} 
\cap \{\boldsymbol{r} : \exists s \in [0, s_{k*}], s.t., r_{1} = r_{k1}^{0}(s), r_{2} \geq r_{k2}^{0}(s)\}^{c}.$$
(9.40)

Consider the ODE (9.27) starting with  $\mathbf{r}(t_0) \in \Omega_k$  for any  $k \geq \{r_\infty, 1/r_0\}$ , we claim  $\mathbf{r}(t) \in \Omega_k$  for any  $t \geq t_0$ . Indeed, combining Eq. (9.31), (9.33), (9.39), and (9.38), for any  $\mathbf{r} \in \partial \Omega_k \setminus (([0,\infty) \times \{0\}) \cup \{\infty\})$  and  $t \geq t_0$ , the gradient  $\nabla \psi_\infty(\mathbf{r}; \overline{\rho}_t)$  pointing outside  $\Omega_k$ . Therefore,  $\mathbf{r}(t)$  cannot leak outside  $\Gamma_k$  from this boundary. Further note that  $\mathbf{r}(t)$  cannot reach the boundary  $([0,\infty) \times \{0\}) \cup \{\infty\}$  for any finite time t. This proves the claim that  $\mathbf{r}(t) \in \Omega_k$  for any  $t \geq t_0$ . According to Lemma 7.8,  $\overline{\rho}_t(\overline{\Omega}_k) \geq \overline{\rho}_{t_0}(\overline{\Omega}_k)$  for any  $k \geq \{r_\infty, 1/r_0\}$  and  $t \geq t_0$ .

Recall the definition of  $O_k$  given by Equation (9.32). Note that  $O_k \subseteq \Omega_k$ , and  $\lim_{k\to\infty} \overline{\rho}_{t_0}(\overline{O}_k) = 1$ , which implies  $\lim_{k\to\infty} \overline{\rho}_{t_0}(\overline{\Omega}_k) = 1$ . Hence, for any  $\overline{\rho}_* \in \mathcal{S}_*$ ,

$$\overline{\rho}_*(\cup_k \overline{\Omega}_k) \ge \lim_{k \to \infty} \overline{\rho}_*(\overline{\Omega}_k) \ge \lim_{k \to \infty} \overline{\rho}_{t_0}(\overline{\Omega}_k) = 1. \tag{9.41}$$

It is easy to see that  $\bigcup_k \overline{\Omega}_k = (0, \infty) \times [0, \infty)$ . Combining with the fact that  $\overline{\rho}_*((0, \infty)^2) = 0$  for any  $\overline{\rho}_* \in \mathcal{S}_*$ , claim (1) holds.

# Proof of Claim (2). We cannot have $L(S_*) \subseteq \mathcal{P}_2$ .

In the case  $L(S_*) \subseteq \mathcal{P}_2$ , the argument is similar to the proof of Claim (1), and hence will be presented in a synthetic form. First, there exists  $t_0$  large enough, so that as  $t \geq t_0$ , we have  $\lambda_+(\overline{\rho}_t) > 0$ , and  $\lambda_-(\overline{\rho}_t) < 0$ . Then  $\chi_{\rm tg}(\boldsymbol{r}; \overline{\rho}_t) < 0$  for any  $\boldsymbol{r} \in (0, \infty)^2$ . Letting

$$\Gamma_k \equiv \{ r \in [0, \infty)^2 : r_1 \le k \cdot r_2 \},$$
(9.42)

According to the same argument as in the proof of Claim (1), we have  $\rho_t(\Gamma_k) \geq \rho_{t_0}(\Gamma_k)$  for any  $k \in (0, \infty)$  and  $t \geq t_0$ . As a result, we have  $\sup(\overline{\rho}_*) \subseteq (\{0\} \times [0, \infty)) \cup \{\infty\}$ .

However, the fixed point  $\overline{\rho}_*$  with support on  $(\{0\} \times [0,\infty)) \cup \{\infty\}$  has risk  $\overline{R}_{\infty}(\overline{\rho}_*) \geq 1$ . Therefore, we cannot have  $L(\mathcal{S}_*) \subseteq \mathcal{P}_2$ . This proves claim (2).

### 9.4 Dynamics: Proof of Theorem 2

We will prove that the dynamics for large but finite d is well approximated by the dynamics at  $d = \infty$ . The key estimate is provided by the next lemma.

**Lemma 9.1.** Assume  $\sigma$  satisfies condition S0, recall the definition of  $u_d$  and  $u_{\infty}$  given by Equation (9.9) and (9.10). Assuming  $k = \gamma \cdot d$  for some  $\gamma \in (0,1)$ , then we have

$$\lim_{d\to\infty} \sup_{\boldsymbol{a},\boldsymbol{b}\in[0,\infty)^2} |u_d(\boldsymbol{a},\boldsymbol{b}) - u_\infty(\boldsymbol{a},\boldsymbol{b})| = 0.$$

and

$$\lim_{d\to\infty} \sup_{\boldsymbol{a},\boldsymbol{b}\in[0,\infty)^2} \|\nabla_{\boldsymbol{a}} u_d(\boldsymbol{a},\boldsymbol{b}) - \nabla_{\boldsymbol{a}} u_\infty(\boldsymbol{a},\boldsymbol{b})\|_2 = 0.$$

*Proof.* We rewrite  $u_d$  here as

$$u_{d}(\boldsymbol{a}, \boldsymbol{b}) = 1/2 \cdot [u_{d,1}(\boldsymbol{a}, \boldsymbol{b}) + u_{d,2}(\boldsymbol{a}, \boldsymbol{b})],$$

$$u_{d,1}(\boldsymbol{a}, \boldsymbol{b}) = \mathbb{E}[\sigma(\tau_{+}a_{1}F_{1} + a_{2}G_{1})\sigma(\tau_{+}b_{1}(F_{1}\cos\Theta_{1} + F_{2}\sin\Theta_{1}) + b_{2}(G_{1}\cos\Theta_{2} + G_{2}\sin\Theta_{2}))],$$

$$u_{d,2}(\boldsymbol{a}, \boldsymbol{b}) = \mathbb{E}[\sigma(\tau_{-}a_{1}F_{1} + a_{2}G_{1})\sigma(\tau_{-}b_{1}(F_{1}\cos\Theta_{1} + F_{2}\sin\Theta_{1}) + b_{2}(G_{1}\cos\Theta_{2} + G_{2}\sin\Theta_{2}))],$$

where

$$(F_1, F_2, G_1, G_2) \sim N(0, I_4),$$
 (9.43)

$$\Theta_1 \sim (1/Z_{s_0})\sin(\theta)^{s_0-2}\mathbf{1}\{\theta \in [0,\pi]\}d\theta,$$
 (9.44)

$$\Theta_2 \sim (1/Z_{d-s_0}) \sin(\theta)^{d-s_0-2} \mathbf{1} \{ \theta \in [0, \pi] \} d\theta,$$
 (9.45)

are mutually independent.

Define  $F_3 = F_1 \cos \Theta_1 + F_2 \sin \Theta_1$ ,  $G_3 = G_1 \cos \Theta_2 + G_2 \sin \Theta_2$ , then

$$|u_{d,1}(\boldsymbol{a},\boldsymbol{b}) - u_{\infty,1}(\boldsymbol{a},\boldsymbol{b})|$$

$$= |\mathbb{E}\{\sigma(\tau_{+}a_{1}F_{1} + a_{2}G_{1})[\sigma(\tau_{+}b_{1}F_{3} + b_{2}G_{3}) - \sigma(\tau_{+}b_{1}F_{2} + b_{2}G_{2})]\}|$$

$$\leq ||\sigma||_{\infty} \cdot \mathbb{E}\{|\sigma(\tau_{+}b_{1}F_{3} + b_{2}G_{3}) - \sigma(\tau_{+}b_{1}F_{2} + b_{2}G_{2})|\},$$

$$(9.46)$$

and

$$\begin{aligned} &|\partial_{a_{1}}u_{d,1}(\boldsymbol{a},\boldsymbol{b}) - \partial_{a_{1}}u_{\infty,1}(\boldsymbol{a},\boldsymbol{b})| \\ &= &|\mathbb{E}\{\tau_{+}F_{1} \cdot \sigma'(\tau_{+}a_{1}F_{1} + a_{2}G_{1})[\sigma(\tau_{+}b_{1}F_{3} + b_{2}G_{3}) - \sigma(\tau_{+}b_{1}F_{2} + b_{2}G_{2})]\}| \\ &\leq \tau_{+} \|\sigma'\|_{\infty} \mathbb{E}[F_{1}^{2}]^{1/2} \mathbb{E}\{[\sigma(\tau_{+}b_{1}F_{3} + b_{2}G_{3}) - \sigma(\tau_{+}b_{1}F_{2} + b_{2}G_{2})]^{2}\}^{1/2} \\ &\leq \tau_{+} \|\sigma'\|_{\infty} (2\|\sigma\|_{\infty}^{1/2}) \cdot \mathbb{E}\{|\sigma(\tau_{+}b_{1}F_{3} + b_{2}G_{3}) - \sigma(\tau_{+}b_{1}F_{2} + b_{2}G_{2})|\}^{1/2}. \end{aligned}$$

$$(9.47)$$

We have similar bounds for  $|\partial_{a_2}u_{d,1}(\boldsymbol{a},\boldsymbol{b}) - \partial_{a_2}u_{\infty,1}(\boldsymbol{a},\boldsymbol{b})|$ .

According to condition S0,  $\|\sigma'\|_{\infty}$  and  $\|\sigma\|_{\infty}$  are bounded, it is sufficient to bound the following quantity uniformly for  $r \in [0, \infty)$  and  $\boldsymbol{a} \in \mathbb{S}^1$ ,

$$T(r, \boldsymbol{a}) \equiv 1/2 \cdot \mathbb{E}\{|\sigma(rH_2) - \sigma(rH_3)|\} = \mathbb{E}\{[\sigma(rH_2) - \sigma(rH_3)] \,\mathbf{1}_{H_2 > H_3}\}, \tag{9.48}$$

where

$$H_2 = H_2(\mathbf{a}) = [\tau_+ a_1 F_2 + a_2 G_2] / [\tau_+^2 a_1^2 + a_2^2]^{1/2}, \tag{9.49}$$

$$H_3 = H_3(\mathbf{a}) = [\tau_+ a_1 F_3 + a_2 G_3] / [\tau_+^2 a_1^2 + a_2^2]^{1/2}.$$
(9.50)

We denote  $\Theta_3 = \Theta_3(\mathbf{a}) = \arcsin\{\mathbb{E}[H_2H_3|\Theta_1,\Theta_2]\}$ . It is easy to see that  $H_2, H_3 \sim \mathsf{N}(0,1)$  with

$$\sin(\Theta_3) = \mathbb{E}[H_2 H_3 | \Theta_1, \Theta_2] = [\tau_+^2 a_1^2 \sin \Theta_1 + a_2^2 \sin \Theta_2] / [\tau_+^2 a_1^2 + a_2^2]. \tag{9.51}$$

Using the same argument as in the proof of Theorem 8.1, we have for any  $z \in \mathbb{R}$ 

$$\mathbb{P}(H_3 \le z, H_2 \ge z) \le \mathbb{P}(H_3 \le 0, H_2 \ge 0) = \mathbb{E}[|\pi/2 - \Theta_3|/(2\pi)]. \tag{9.52}$$

Hence, we have

$$T(r, \boldsymbol{a}) = \mathbb{E}\left\{ \int_{\mathbb{R}} \sigma'(t) \, \mathbf{1}_{rH_2 \ge t \ge rH_3} \, \mathrm{d}t \right\} = \int_{\mathbb{R}} \sigma'(t) \, \mathbb{P}\left\{ H_2 \ge t/r \ge H_3 \right\} \, \mathrm{d}t$$
$$\leq \sup_{z \in \mathbb{R}} \mathbb{P}(H_3 \le z, H_2 \ge z) \, \int_{\mathbb{R}} \sigma'(t) \, \mathrm{d}t \le 2 \|\sigma\|_{\infty} \cdot \mathbb{E}[|\pi/2 - \Theta_3|/(2\pi)] \, .$$

Note that  $\cos(\Theta_1) \stackrel{\mathrm{d}}{=} Y_1/\|\boldsymbol{Y}\|_2$ , for  $\boldsymbol{Y} \sim \mathsf{N}(0, \boldsymbol{I}_{s_0})$ , and  $\cos(\Theta_2) \stackrel{\mathrm{d}}{=} Z_1/\|\boldsymbol{Z}\|_2$ , for  $\boldsymbol{Z} \sim \mathsf{N}(0, \boldsymbol{I}_{d-s_0})$ . Hence, there exists a universal constant K, such that  $\mathbb{E}\{|\Theta_1 - \pi/2|\} \leq K/\sqrt{s_0}$ ,  $\mathbb{E}\{|\Theta_2 - \pi/2|\} \leq K/\sqrt{d-s_0}$ .

Note the relationship of  $\Theta_3 = \Theta_3(a)$  with  $(\Theta_1, \Theta_2)$  is given by Eq. (9.51), which yields

$$\sin(\Theta_3(\boldsymbol{a})) \ge \min\{\sin\Theta_1, \sin\Theta_2\},\tag{9.53}$$

hence

$$|\pi/2 - \Theta_3(\mathbf{a})| \le \max\{|\pi/2 - \Theta_1|, |\pi/2 - \Theta_2|\}.$$
 (9.54)

As a result,

$$\sup_{\boldsymbol{a} \in \mathbb{S}^1} \mathbb{E}\{|\Theta_3(\boldsymbol{a}) - \pi/2|\} \le K \cdot \max\{1/\sqrt{s_0}, 1/\sqrt{d - s_0}\}. \tag{9.55}$$

We therefore obtain

$$\sup_{r \in \mathbb{R}, \boldsymbol{a} \in \mathbb{S}^1} |T(r, \boldsymbol{a})| \le K/\pi \cdot \|\sigma\|_{\infty} \cdot \max\{1/\sqrt{s_0}, 1/\sqrt{d - s_0}\}. \tag{9.56}$$

The lemma holds by noting that as  $d \to \infty$ , we have  $s_0 \to \infty$  and  $d - s_0 \to \infty$ .

Proof of Theorem 2. Recall the definition of  $\overline{R}_{\infty}$  given by Eq. (9.11), and R given by Eq. (6.2). Recall the set of good initialization given by

$$\mathscr{P}_{\text{good}} = \{ \overline{\rho}_0 \in \mathscr{P}((0, \infty)) : \lim_{d \to \infty} R(\overline{\rho} \times \text{Unif}(\mathbb{S}^{d-1})) < 1 \}.$$

Define  $\mathscr{P}^1_{\mathrm{good}}$  and  $\mathscr{P}^2_{\mathrm{good}}$  to be

$$\mathscr{P}_{\text{good}}^1 = \{ \overline{\rho}_0^1 \in \mathscr{P}((0,\infty)) : \overline{R}_{\infty}(\overline{\rho}_0^2) < 1, \text{ where } \overline{\rho}_0^2 \sim (\gamma^{1/2}u, (1-\gamma)^{1/2}u) \text{ with } u \sim \overline{\rho}_0^1 \}, \quad (9.57)$$

$$\mathscr{P}_{\text{good}}^2 = \{ \overline{\rho}_0^2 \in \mathscr{P}((0,\infty)^2) : \overline{R}_\infty(\overline{\rho}_0^2) < 1 \}. \tag{9.58}$$

With this definition, it is easy to see that  $\mathscr{P}^1_{\text{good}} = \mathscr{P}_{\text{good}}$ .

For any  $\overline{\rho}_0^1 \in \mathscr{P}_{\mathrm{good}}^1$ , let  $u \sim \overline{\rho}_0^1$ ,  $Y_1 \sim \chi^2(\gamma \cdot d)$ , and  $Y_2 \sim \chi^2((1-\gamma) \cdot d)$  be independent. We take  $u_{d1} = u \cdot [Y_1/(Y_1+Y_2)]^{1/2}$ ,  $u_{d2} = u \cdot [Y_2/(Y_1+Y_2)]^{1/2}$ ,  $u_d = (u_{d1},u_{d2})$ ,  $u_{\infty 1} = u \cdot [s_0/d]^{1/2} = u \cdot \gamma^{1/2}$ ,  $u_{\infty 2} = u \cdot [(d-s_0)/d]^{1/2} = u \cdot (1-\gamma)^{1/2}$ , and  $u_{\infty} = (u_{\infty 1},u_{\infty 2})$ . Denote  $\overline{\rho}_0^{2,d}$  to be the distribution of  $u_d$ , and  $\overline{\rho}_0^{2,\infty}$  to be the distribution of  $u_\infty$ . Then we have  $\overline{\rho}_0^{2,\infty} \in \mathscr{P}_{\mathrm{good}}^2$ . Further, if we sample  $(r,n) \sim \overline{\rho}_0^1 \times \mathrm{Unif}(\mathbb{S}^{d-1})$  and  $(r,n_1,n_2) \sim \overline{\rho}_0^{2,d} \times \mathrm{Unif}(\mathbb{S}^{k-1}) \times \mathrm{Unif}(\mathbb{S}^{d-k-1})$ , then  $rn \stackrel{\mathrm{d}}{=} (r_1n_1,r_2n_2)$ . Here we bound  $d_{\mathrm{BL}}(\overline{\rho}_0^{2,d},\overline{\rho}_0^{2,\infty})$ . Note the joint distribution of  $u_d$  and  $u_\infty$  is a coupling of  $\overline{\rho}_0^{2,d}$  and  $\overline{\rho}_0^{2,\infty}$ , hence

$$d_{\text{BL}}(\overline{\rho}_0^{2,d}, \overline{\rho}_0^{2,\infty}) \leq \mathbb{E}[\|\boldsymbol{u}_d - \boldsymbol{u}_\infty\|_2 \wedge 1]$$

$$= \mathbb{E}[\{u[((Y_1/(Y_1 + Y_2))^{1/2} - \gamma^{1/2})^2 + ((Y_2/(Y_1 + Y_2))^{1/2} - (1 - \gamma)^{1/2})^2]^{1/2}\} \wedge 1]. \tag{9.59}$$

It is easy to see that  $\lim_{d\to\infty} Y_1/(Y_1+Y_2) = \gamma$  almost surely. Bounded convergence theorem implies that  $\lim_{d\to\infty} d_{\rm BL}(\overline{\rho}_0^{2,d},\overline{\rho}_0^{2,\infty}) = 0$ .

Now we consider the PDE (9.16) for  $d = \infty$ . We fix its initialization  $\overline{\rho}_0^{2,\infty} \in \mathscr{P}_{\text{good}}^2$  induced by  $\overline{\rho}_0^1 \in \mathscr{P}_{\text{good}}^1$ . Denote the solution of PDE (9.16) to be  $(\overline{\rho}_t^{\infty})_{t\geq 0}$ . Due to Theorem 13, for any  $\eta > 0$ , there exists  $T = T(\eta, \overline{\rho}_0^1, \gamma, \Delta) > 0$ , so that its solution  $(\overline{\rho}_t^{\infty})_{t\geq 0}$  satisfies

$$\overline{R}_{\infty}(\overline{\rho}_t^{\infty}) \le \inf_{\overline{\rho} \in \mathscr{P}(E_2)} \overline{R}_{\infty}(\overline{\rho}) + \eta/5$$

for any  $t \geq T$ .

Then we consider the general PDE

$$\partial_t \rho_t(\boldsymbol{\theta}) = 2\xi(t) \nabla \cdot \left[ \rho_t(\boldsymbol{\theta}) \nabla \Psi(\boldsymbol{\theta}; \rho_t) \right], \tag{9.60}$$

with initialization  $\rho_0$  the distribution of  $r\boldsymbol{n}$ , where  $(r,\boldsymbol{n}) \sim \overline{\rho}_0^1 \times \mathrm{Unif}(\mathbb{S}^{d-1})$ . Due to Lemma 8.2 and Remark 7.1, we have the existence and uniqueness of the solution of PDE (9.60). We denote its solution to be  $(\rho_t)_{t\geq 0}$ . Let  $\overline{\rho}_t^d$  be the distribution of  $(\|\boldsymbol{w}_1\|_2, \|\boldsymbol{w}_2\|_2)$  with  $\boldsymbol{w} = (\boldsymbol{w}_1, \boldsymbol{w}_2) \sim \rho_t$ ,  $\boldsymbol{w}_1 \in \mathbb{R}^{s_0}$  and  $\boldsymbol{w}_2 \in \mathbb{R}^{d-s_0}$ . It is easy to see that  $(\overline{\rho}_t^d)_{t\geq 0}$  is the unique solution of (9.16) with initialization  $\overline{\rho}_0^{2,d}$ .

Now, we would like to bound the distance of  $\overline{\rho}_t^d$  and  $\overline{\rho}_t^\infty$  using Lemma 7.7. We take D=2, V=v,  $U=u_d, \tilde{V}=v, \tilde{U}=u_\infty$  in Lemma 7.7. Let  $\varepsilon_0(d)$  be as defined in Eq. (7.69). Due to Lemma 9.1, we have  $\lim_{d\to\infty}\varepsilon_0(d)=0$ . We also showed that  $\lim_{d\to\infty}d_{\rm BL}(\overline{\rho}_0^{2,d},\overline{\rho}_0^{2,\infty})=0$ . Therefore, according to Lemma 7.7, we have  $\lim_{d\to\infty}\sup_{t\le 10T}d_{\rm BL}(\overline{\rho}_t^{2,d},\overline{\rho}_t^{2,\infty})=0$ . Further note  $\overline{R}_\infty$  is uniformly continuous with respect to  $\overline{\rho}$  in bounded-Lipschitz distance. Therefore, there exists  $d_0=d_0(\eta,\overline{\rho}_0^1,\gamma,\Delta)$  large enough, so that for  $d\ge d_0$  we have

$$|\overline{R}_{\infty}(\overline{\rho}_t^d) - \overline{R}_{\infty}(\overline{\rho}_t^{\infty})| \le \eta/5.$$

for any  $t \leq 10T$ .

Then we would like to bound the difference of  $\overline{R}_{\infty}(\overline{\rho})$  and  $\overline{R}_{d}(\overline{\rho})$  for any  $\overline{\rho}$ . Note

$$|\overline{R}_{\infty}(\overline{\rho}) - \overline{R}_{d}(\overline{\rho})| \leq \int |u_{d}(\boldsymbol{a}, \boldsymbol{b}) - u_{\infty}(\boldsymbol{a}, \boldsymbol{b})| \overline{\rho}(\mathrm{d}\boldsymbol{a}) \overline{\rho}(\mathrm{d}\boldsymbol{b}). \tag{9.61}$$

By Lemma 9.1, there exists  $d_0 = d_0(\eta, \Delta)$  large enough, so that for  $d \geq d_0$ , we have

$$\sup_{\overline{\rho} \in \mathscr{P}(E_2)} |\overline{R}_{\infty}(\overline{\rho}) - \overline{R}_d(\overline{\rho})| \le \eta/5. \tag{9.62}$$

Finally, let  $(\theta^k)_{k\geq 1}$  be the trajectory of SGD, with step size  $s_k = \varepsilon \xi(k\varepsilon)$ , and initialization  $\mathbf{w}_i^0 \sim_{iid} \rho_0$  for  $i \leq N$ . We apply Theorem 3 to bound the difference of the law of trajectory of SGD and the solution of PDE (9.60). The assumptions of Theorem 3 are verified by Lemma 8.2. As a consequence, there exists constant K (which depend uniquely on the constants in assumptions A1 A2 A3), such that

$$R_N(\boldsymbol{\theta}^{\lfloor t/\varepsilon \rfloor}) - \overline{R}_d(\overline{\rho}_t^d) \le Ke^{10KT} \cdot \operatorname{err}_{N,d}(z).$$

with probability  $1 - e^{-z^2}$  for any  $t \le 10T$ , where

$$\mathrm{err}_{N,d}(z) = \sqrt{1/N \vee \varepsilon} \cdot \left\lceil \sqrt{D + \log(N(1/\varepsilon \vee 1))} + z \right\rceil.$$

As a consequence, for any  $\delta > 0$ , there exists  $C_0 = C_0(\delta, \eta, \overline{\rho}_0^1, \gamma, \Delta)$ , so that as  $N, 1/\varepsilon \geq C_0 d$  and  $\varepsilon \geq 1/N^{10}$ , for  $t \leq 10T$ , we have

$$R_N(\boldsymbol{\theta}^{\lfloor t/\varepsilon \rfloor}) - \overline{R}_d(\overline{\rho}_t^d) \le \eta/5$$

with probability at least  $1 - \delta$ .

Therefore, the trajectory  $\boldsymbol{\theta}^{\lfloor t/\varepsilon \rfloor}$  of SGD as  $t \in [T, 10T]$  satisfies

$$R_{N}(\boldsymbol{\theta}^{\lfloor t/\varepsilon \rfloor}) \leq \overline{R}_{d}(\overline{\rho}_{t}^{d}) + \eta/5 \leq \overline{R}_{\infty}(\overline{\rho}_{t}^{d}) + 2\eta/5 \leq \overline{R}_{\infty}(\overline{\rho}_{t}^{\infty}) + 3\eta/5$$

$$\leq \inf_{\overline{\rho} \in \overline{\mathscr{P}}} \overline{R}_{\infty}(\overline{\rho}) + 4\eta/5 \leq \inf_{\overline{\rho} \in \overline{\mathscr{P}}} \overline{R}_{d}(\overline{\rho}) + \eta = \inf_{\rho \in \mathscr{P}(\mathbb{R}^{d})} R(\rho) + \eta$$

$$\leq \inf_{\boldsymbol{\theta} \in \mathbb{R}^{d \times N}} R_{N}(\boldsymbol{\theta}) + \eta$$

with probability at least  $1 - \delta$ . This gives the desired result.

# 10 Finite temperature

We will states the lemma regarding statics properties of the finite temperature free energy in Section 10.1, and regarding dynamics properties in Section 10.2. We will prove Proposition 3, Theorem 4, and Theorem 5 in Section 10.3. Throughout Section 10.1 and 10.2, to distinguish the dimension of parameters with the generalized differential operator, we will denote the dimension of parameters by d instead of D. This should not be confused with the dimension of feature vectors, which never appears throughout this section.

We introduce the set K of admissible probability densities,

$$\mathcal{K} = \left\{ \rho : \mathbb{R}^d \to [0, +\infty) \text{ measurable } : \int_{\mathbb{R}^d} \rho(\boldsymbol{\theta}) d\boldsymbol{\theta} = 1, M(\rho) < \infty \right\}, \tag{10.1}$$

where

$$M(\rho) \equiv \int_{\mathbb{R}^d} \|\boldsymbol{\theta}\|_2^2 \rho(\boldsymbol{\theta}) d\boldsymbol{\theta}.$$
 (10.2)

Recall

$$R(\rho) = R_{\#} + 2 \int_{\mathbb{R}^d} V(\boldsymbol{\theta}) \rho(\boldsymbol{\theta}) d\boldsymbol{\theta} + \int_{\mathbb{R}^d \times \mathbb{R}^d} U(\boldsymbol{\theta}, \boldsymbol{\theta}') \rho(\boldsymbol{\theta}) \rho(\boldsymbol{\theta}') d\boldsymbol{\theta} d\boldsymbol{\theta}',$$
(10.3)

$$R_{\#} = \mathbb{E}\{y^2\}, \qquad V(\boldsymbol{\theta}) = -\mathbb{E}\{y\,\sigma_*(\boldsymbol{x};\boldsymbol{\theta})\},$$
 (10.4)

$$U(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) = \mathbb{E}\left\{\sigma_*(\boldsymbol{x}; \boldsymbol{\theta}_1)\sigma_*(\boldsymbol{x}; \boldsymbol{\theta}_2)\right\},\tag{10.5}$$

$$\Psi(\boldsymbol{\theta}; \rho) = V(\boldsymbol{\theta}) + \int_{\mathbb{R}^d} U(\boldsymbol{\theta}, \boldsymbol{\theta}') \, \rho(\boldsymbol{\theta}') d\boldsymbol{\theta}'.$$
 (10.6)

Let

$$R_{\lambda}(\rho) = \lambda M(\rho) + R(\rho), \tag{10.7}$$

$$\Psi_{\lambda}(\boldsymbol{\theta}; \rho) = \lambda/2 \cdot \|\boldsymbol{\theta}\|_{2}^{2} + V(\boldsymbol{\theta}) + \int_{\mathbb{R}^{d}} U(\boldsymbol{\theta}, \boldsymbol{\theta}') \rho(\boldsymbol{\theta}') d\boldsymbol{\theta}', \tag{10.8}$$

$$\operatorname{Ent}(\rho) = -\int_{\mathbb{R}^d} \rho(\boldsymbol{\theta}) \log \rho(\boldsymbol{\theta}) d\boldsymbol{\theta}, \tag{10.9}$$

$$F_{\beta,\lambda}(\rho) = 1/2 \cdot [\lambda M(\rho) + R(\rho)] - 1/\beta \cdot \text{Ent}(\rho). \tag{10.10}$$

#### 10.1 Statics

**Lemma 10.1.** For any  $\rho \in \mathcal{K}$ , we have

$$\operatorname{Ent}(\rho) \le \int_{\mathbb{R}^d} \rho(\boldsymbol{\theta}) \cdot |\min(\log \rho(\boldsymbol{\theta}), 0)| \cdot d\boldsymbol{\theta} \le 1 + M(\rho)/\sigma^2 + d \cdot \log(2\pi\sigma^2)$$
 (10.11)

for any  $\sigma^2 > 0$ .

*Proof.* Define  $\Omega_0 = \{ \theta : 1/(\sqrt{2\pi}\sigma)^d \cdot \exp\{-\|\theta\|_2^2/(2\sigma^2)\} \le \rho(\theta)^{1/2} \le 1 \}$ . Then we have

$$\operatorname{Ent}(\rho) = -\int_{\mathbb{R}^d} \rho(\boldsymbol{\theta}) \log \rho(\boldsymbol{\theta}) d\boldsymbol{\theta} \le \int_{\mathbb{R}^d} \rho(\boldsymbol{\theta}) \cdot |\min(\log \rho(\boldsymbol{\theta}), 0)| \cdot d\boldsymbol{\theta}$$
$$\le \int_{\Omega_0} \rho(\boldsymbol{\theta}) \cdot |\min(\log \rho(\boldsymbol{\theta}), 0)| \cdot d\boldsymbol{\theta} + \int_{\Omega_0^c} \rho(\boldsymbol{\theta}) \cdot |\min(\log \rho(\boldsymbol{\theta}), 0)| \cdot d\boldsymbol{\theta}.$$

The first term is bounded by

$$\int_{\Omega_0} \rho(\boldsymbol{\theta}) \cdot |\min(\log \rho(\boldsymbol{\theta}), 0)| \cdot d\boldsymbol{\theta} \leq \int_{\mathbb{R}^d} \rho(\boldsymbol{\theta}) [\|\boldsymbol{\theta}\|_2^2 / \sigma^2 + d \cdot \log(2\pi\sigma^2)] d\boldsymbol{\theta} = M(\rho) / \sigma^2 + d \cdot \log(2\pi\sigma^2).$$

Noting that  $|\rho \log \rho| \leq \sqrt{\rho}$  for any  $\rho \in [0,1]$ , the second term is bounded by

$$\begin{split} &\int_{\Omega_0^c} \rho(\boldsymbol{\theta}) \cdot |\min(\log \rho(\boldsymbol{\theta}), 0)| \cdot \mathrm{d}\boldsymbol{\theta} \leq \int_{\Omega_0^c} \rho(\boldsymbol{\theta})^{1/2} \mathbf{1} \{ \rho(\boldsymbol{\theta}) \leq 1 \} \mathrm{d}\boldsymbol{\theta} \\ \leq &\int_{\mathbb{R}^d} 1/(\sqrt{2\pi}\sigma)^d \cdot \exp\{-\|\boldsymbol{\theta}\|_2^2/(2\sigma^2)\} \mathrm{d}\boldsymbol{\theta} = 1. \end{split}$$

**Lemma 10.2.** Assume U and V are bounded-Lipschitz. Then for any  $\lambda > 0$  and  $0 < \beta < \infty$ ,  $F_{\beta,\lambda}(\rho)$  has a unique minimizer  $\rho_* \in \mathcal{K}$ . Moreover, we have

$$F_{\beta,\lambda}(\rho) \ge 1/2 \cdot R(\rho) + \lambda/4 \cdot M(\rho) - 1/\beta \cdot [1 + d \cdot \log(8\pi/(\beta\lambda))]. \tag{10.12}$$

*Proof.* First, by Lemma 10.1, we have

$$F_{\beta,\lambda}(\rho) = 1/2 \cdot R(\rho) + \lambda/2 \cdot M(\rho) - 1/\beta \cdot \text{Ent}(\rho)$$
  
 
$$\geq 1/2 \cdot R(\rho) + \lambda/2 \cdot M(\rho) - 1/\beta \cdot [1 + M(\rho)/\sigma^2 + d \cdot \log(2\pi\sigma^2)].$$

Taking  $\sigma^2 = 4/(\beta \lambda)$  gives Eq. (10.12).

The argument to show the existence and uniqueness of minimizer of  $F_{\beta,\lambda}$  is similar to the proof of [JKO98, Proposition 4.1], and we will just give a sketch here. Since U, V are bounded-Lipschitz, it follows that  $\rho \mapsto R(\rho)$  is continuous with respect to the topology of weak convergence in  $L^1(\mathbb{R}^d)$ . Fatou's lemma implies that M is lower semi-continuous. [JKO98, Proposition 4.1] shows the upper semi-continuity of Ent. Hence  $F_{\beta,\lambda}$  is lower semi-continuous. Note (as just shown)  $F_{\beta,\lambda}$  is lower bounded, there exists a sequence  $(\rho_k)_{k\geq 1} \subset \mathcal{K}$  such that  $\lim_{k\to\infty} F_{\beta,\lambda}(\rho_k) = \inf_{\rho\in\mathcal{K}} F_{\beta,\lambda}(\rho) > -\infty$ . By the same argument as [JKO98, Proposition 4.1], we can see that  $\{\int \max\{\rho_k \log \rho_k, 0\}\} d\theta\}_{k\geq 1}$  and  $\{M(\rho_k)\}_{k\geq 1}$  are uniformly upper bounded, and by de la Vallée-Poussin criterion, there exists  $\rho_* \in \mathcal{K}$  such that there is a subsequence of  $(\rho_k)_{k\geq 1}$  converges weakly to  $\rho_*$  in  $L^1(\mathbb{R}^d)$ . The lower semi-continuity of  $F_{\beta,\lambda}$  implies that  $\rho_*$  is the minimizer of  $F_{\beta,\lambda}$ . Uniqueness follows by noting that U is positive semi-definite, Ent is strongly concave, and  $\langle V, \rho \rangle$  and M are linear in  $\rho$ , so that  $F_{\beta,\lambda}$  is a strongly convex functional.

For any  $\rho \in \mathcal{K}$ , we call the following equation the Boltzmann fixed point condition

$$\rho(\boldsymbol{\theta}) = 1/Z(\beta, \lambda; \rho) \exp\{-\beta \Psi_{\lambda}(\boldsymbol{\theta}; \rho)\},$$

$$Z(\beta, \lambda; \rho) = \int \exp\{-\beta \Psi_{\lambda}(\boldsymbol{\theta}; \rho)\} d\boldsymbol{\theta}.$$
(10.13)

**Lemma 10.3.** Under the assumption of Lemma 10.2, the minimizer  $\rho_* \in \mathcal{K}$  of  $F_{\beta,\lambda}(\rho)$  satisfies the Boltzmann fixed point condition.

*Proof.* We denote  $\mu_0$  to be the Lebesgue measure on  $\mathbb{R}^d$ .

First, we show that  $\rho$  is positive almost everywhere. Let  $\rho_* \in \mathcal{K}$  be a minimizer of  $F(\rho)$ , and assume by contradiction that there exists a measurable set  $\Omega_0 \subset \mathbb{R}^d$ , such that  $\mu_0(\Omega_0) > 0$ , and  $\rho_*(\Omega_0) = 0$ . Without loss of generality, we assume that the support of  $\Omega_0$  is compact so that  $\mu_0(\Omega_0) < \infty$ , otherwise we can always consider the intersection of  $\Omega_0$  with a large ball. Define  $\rho_{\varepsilon} = (1 - \varepsilon)\rho_* + \varepsilon/\mu_0(\Omega_0) \cdot \mathbf{1}_{\Omega_0} \in \mathcal{K}$ . It is easy to see that there exists  $\varepsilon_0 > 0$  and  $C < \infty$ , such that  $|R_{\lambda}(\rho_*) - R_{\lambda}(\rho_{\varepsilon})| \leq C \cdot \varepsilon$ , and

$$\operatorname{Ent}(\rho_{\varepsilon}) = (1 - \varepsilon)\operatorname{Ent}(\rho_{*}) - (1 - \varepsilon)\log(1 - \varepsilon) + \varepsilon\log(\mu_{0}(\Omega_{0})/\varepsilon)$$
  
 
$$\geq \operatorname{Ent}(\rho_{*}) - C \cdot \varepsilon + \varepsilon\log(\mu_{0}(\Omega_{0})/\varepsilon)$$

for any  $\varepsilon < \varepsilon_0$ . As  $\varepsilon$  is sufficiently small, we have  $F_{\beta,\lambda}(\rho_{\varepsilon}) < F_{\beta,\lambda}(\rho_*)$ . This contradict with the fact that  $\rho_* \in \mathcal{K}$  is the minimizer of  $F_{\beta,\lambda}(\rho)$ .

Next we show that, for all  $\theta \in \mathbb{R}^d$ ,

$$\Psi_{\lambda}(\boldsymbol{\theta}; \rho_*) + 1/\beta \cdot \log \rho_*(\boldsymbol{\theta}) \equiv \gamma(\beta, \lambda; \rho_*)$$
(10.14)

for some constant  $\gamma(\beta, \lambda; \rho_*)$ .

Let  $\rho_* \in \mathcal{K}$  be the minimizer of  $F_{\beta,\lambda}(\rho)$ . Fix  $\varepsilon_0 > 0$  and define  $\Gamma_{\varepsilon_0} \equiv \{\boldsymbol{\theta} \in \mathbb{R}^d : \rho_*(\boldsymbol{\theta}) \geq \varepsilon_0\} \cap \mathsf{B}(\mathbf{0}; 1/\varepsilon_0)$ , and  $\mathcal{A}_{\varepsilon_0} \equiv \{v \in C^{\infty}(\mathbb{R}^d) : ||v||_{\infty} \leq 1, \operatorname{supp}(v) \subseteq \Gamma_{\varepsilon_0}, \int_{\mathbb{R}^d} v(\boldsymbol{\theta}) d\boldsymbol{\theta} = 0\}$ . For any  $v \in \mathcal{A}_{\varepsilon_0}$ , define  $\rho_{\varepsilon,v} = \rho + \varepsilon v$ . Note that, for  $-\varepsilon_0 < \varepsilon < \varepsilon_0$ , we have  $\rho_{\varepsilon,v} \in \mathcal{K}$ . Since  $\rho_*$  is the minimizer of  $F_{\beta,\lambda}(\rho)$ , we must have  $\lim_{\varepsilon \to 0+} [F_{\beta,\lambda}(\rho_{\varepsilon,v}) - F_{\beta,\lambda}(\rho_*)]/\varepsilon \geq 0$ . It can be easily verified that

$$\lim_{\varepsilon \to 0} [F_{\beta,\lambda}(\rho_{\varepsilon,v}) - F_{\beta,\lambda}(\rho_*)]/\varepsilon = \int_{\mathbb{R}^d} [\Psi_{\lambda}(\boldsymbol{\theta}; \rho_*) + 1/\beta \cdot \log \rho_*(\boldsymbol{\theta})] v(\boldsymbol{\theta}) d\boldsymbol{\theta},$$

which implies

$$\int_{\mathbb{R}^d} [\Psi_{\lambda}(\boldsymbol{\theta}; \rho_*) + 1/\beta \cdot \log \rho_*(\boldsymbol{\theta})] v(\boldsymbol{\theta}) d\boldsymbol{\theta} = 0$$
(10.15)

for any  $v \in \mathcal{A}_{\varepsilon_0}$ . This implies that Eq. (10.14) holds for any  $\boldsymbol{\theta} \in \Gamma_{\varepsilon_0}$ . But note that  $\mu_0(\mathbb{R}^d \setminus (\cup_{\varepsilon_0>0}\Gamma_{\varepsilon_0})) = 0$ . This implies that Eq. (10.14) holds almost surely.

Note we have  $\int \rho_*(\boldsymbol{\theta}) d\boldsymbol{\theta} = 1$ . Therefore, we must have  $\gamma(\beta, \lambda; \rho_*) = -1/\beta \cdot \log Z(\beta, \lambda; \rho_*)$ . This proves that  $\rho_*$  satisfies the Boltzmann fixed point condition.

**Lemma 10.4.** Under the assumption of Lemma 10.2, the Boltzmann fixed point condition has a unique solution in K.

*Proof.* The last two lemmas already imply that the Boltzmann fixed point condition has at least one solution. Assume  $\rho_1, \rho_2 \in K$  to be two such solutions. Then  $\rho_i$  is positive, and

$$\log Z(\beta, \lambda; \rho_i) = -\beta \Psi_{\lambda}(\boldsymbol{\theta}; \rho_i) - \log \rho_i(\boldsymbol{\theta}).$$

Therefore

$$0 = \int_{\mathbb{R}^d} [\log Z(\beta, \lambda; \rho_1) - \log Z(\beta, \lambda; \rho_2)] \cdot [\rho_1(\boldsymbol{\theta}) - \rho_2(\boldsymbol{\theta})] d\boldsymbol{\theta}$$
$$= -\beta \langle U, (\rho_1 - \rho_2)^{\otimes 2} \rangle - \int_{\mathbb{R}^d} \log(\rho_1(\boldsymbol{\theta}) / \rho_2(\boldsymbol{\theta})) [\rho_1(\boldsymbol{\theta}) - \rho_2(\boldsymbol{\theta})] d\boldsymbol{\theta}.$$

Note the right hand side does not equal 0 unless  $\rho_1 = \rho_2$ .

**Lemma 10.5.** Under the assumption of Lemma 10.2, and further assume condition A3 holds. Let  $\rho_*^{\beta,\lambda}$  be the minimizer of  $F_{\beta,\lambda}(\rho)$ . Then there is a constant K depending on the parameter  $K_3$  in condition A3, such that for any  $\beta \geq 1$ , we have

$$R(\rho_*^{\beta,\lambda}) \le \inf_{\rho \in \mathscr{P}(\mathbb{R}^d)} R_{\lambda}(\rho) + K(1+\lambda) [d\log(2+1/\lambda)]/\beta. \tag{10.16}$$

*Proof.* Fix a  $\rho \in \mathscr{P}(\mathbb{R}^d)$ . Let  $g_{\tau}(\boldsymbol{\theta})$  be the density for  $\mathsf{N}(\mathbf{0}, \tau^2 \mathbf{I}_d)$ . Denote  $\rho * g_{\tau}$  to be the convolution of  $\rho$  and  $g_{\tau}$ . Now we derive the formula for  $F_{\beta,\lambda}(\rho * g_{\tau})$ .

Let  $G, G_1, G_2 \sim N(0, I_d)$  be independent, we have

$$R(\rho * g_{\tau}) = R(\rho) + 2 \int \{\mathbb{E}[V(\boldsymbol{\theta} + \tau \boldsymbol{G})] - V(\boldsymbol{\theta})\} \rho(\mathrm{d}\boldsymbol{\theta})$$
$$+ \int \{\mathbb{E}[U(\boldsymbol{\theta}_1 + \tau \boldsymbol{G}_1, \boldsymbol{\theta}_2 + \tau \boldsymbol{G}_2)] - U(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)\} \rho(\mathrm{d}\boldsymbol{\theta}_1) \rho(\mathrm{d}\boldsymbol{\theta}_2).$$

Using the intermediate value theorem and Cauchy-Schwarz inequality, and noting that  $\nabla^2 V$  is  $K_3$ -bounded by condition A3, we have

$$\int \{V(\boldsymbol{\theta}) - \mathbb{E}[V(\boldsymbol{\theta} + \tau \boldsymbol{G})]\} \rho(\mathrm{d}\boldsymbol{\theta}) 
= \tau \int \mathbb{E}[\langle \nabla V(\boldsymbol{\theta}), \boldsymbol{G} \rangle] \rho(\mathrm{d}\boldsymbol{\theta}) + \frac{\tau^2}{2} \int \mathbb{E}[\langle \nabla^2 V(\tilde{\boldsymbol{\theta}}), \boldsymbol{G}^{\otimes 2} \rangle] \rho(\mathrm{d}\boldsymbol{\theta}) \leq \frac{\tau^2}{2} K_3 d,$$

We have similar bound for the U term. Therefore,

$$R(\rho * g_{\tau}) \le R(\rho) + 2\tau^2 K_3 d.$$
 (10.17)

For the term  $M(\rho * g_{\tau})$ , we have

$$M(\rho * g_{\tau}) = \int \mathbb{E}[\|\boldsymbol{\theta} + \tau \boldsymbol{G}\|_{2}^{2}]\rho(\mathrm{d}\boldsymbol{\theta}) = M(\rho) + \tau^{2}d.$$
 (10.18)

Next we give a lower bound for  $\operatorname{Ent}(\rho * g_{\tau})$ :

$$\text{Ent}(\rho * g_{\tau}) \ge \text{Ent}(g_{\tau}) = (d/2)\log(2\pi e \tau^2).$$
 (10.19)

As a result, taking  $\tau = 1/\beta$ , we have

$$F_{\beta,\lambda}(\rho_*^{\beta,\lambda}) \le (1/2)R_{\lambda}(\rho) + (2K_3 + \lambda)d/(2\beta^2) + d \cdot \log(2\pi e\beta^2)/(2\beta).$$
 (10.20)

Combining with Eq. (10.12), we have

$$R(\rho_*^{\beta,\lambda}) \le R_\lambda(\rho) + \frac{(2K_3 + \lambda)d}{\beta^2} + \frac{2}{\beta} + \frac{d \cdot \log(2\pi e \beta^2)}{\beta} - \frac{2d \cdot \log(\lambda\beta/(8\pi))}{\beta}$$
(10.21)

for any  $\rho \in \mathscr{P}(\mathbb{R}^d)$ . Hence, the theorem holds by taking infimum over  $\rho \in \mathscr{P}(\mathbb{R}^d)$ .

### 10.2 Dynamics

Recall that the finite-temperature distributional dynamics reads:

$$\partial_t \rho_t(\boldsymbol{\theta}) = 2\xi(t) \nabla_{\boldsymbol{\theta}} \cdot (\nabla_{\boldsymbol{\theta}} \Psi_{\lambda}(\boldsymbol{\theta}; \rho_t) \rho_t(\boldsymbol{\theta})) + 2\xi(t) / \beta \cdot \Delta_{\boldsymbol{\theta}} \rho_t(\boldsymbol{\theta}). \tag{10.22}$$

We say  $(\rho_t)_{t\geq 0} \subseteq \mathscr{P}(\mathbb{R}^d)$  is a weak solution of (10.22), if for any  $\zeta \in C_0^{\infty}(\mathbb{R} \times \mathbb{R}^d)$  (the space of smooth functions, decaying to 0 at infinity), we have

$$\int_{\mathbb{R}^{d}} \rho_{0}(\boldsymbol{\theta}) \zeta_{0}(\boldsymbol{\theta}) d\boldsymbol{\theta}$$

$$= -\int_{(0,\infty)\times\mathbb{R}^{d}} [\partial_{t} \zeta_{t}(\boldsymbol{\theta}) - 2\xi(t) \langle \nabla_{\boldsymbol{\theta}} \Psi_{\lambda}(\boldsymbol{\theta}; \rho_{t}), \nabla_{\boldsymbol{\theta}} \zeta_{t}(\boldsymbol{\theta}) \rangle + 2\xi(t) \Delta_{\boldsymbol{\theta}} \zeta_{t}(\boldsymbol{\theta})] \rho_{t}(d\boldsymbol{\theta}) dt$$
(10.23)

Notice that this notion of weak solution is equivalent to the one introduced earlier in Eq. (7.3), see for instance [San15, Proposition 4.2].

**Lemma 10.6.** Assume conditions A1, A2 and A3 hold. Let initialization  $\rho_0 \in \mathcal{K}$  so that  $F_{\beta,\lambda}(\rho_0) < \infty$ . Then, the weak solution  $(\rho_t)_{t\geq 0} \subseteq \mathscr{P}(\mathbb{R}^d)$  of PDE (10.23) exists and is unique. Moreover, for any fixed t,  $\rho_t \in \mathcal{K}$  is absolutely continuous with respect to the Lebesgue measure, and  $\operatorname{Ent}(\rho_t)$  and  $M(\rho_t)$  are uniformly bounded in t.

*Proof.* Without loss of generality, we assume  $\xi(t) \equiv 1/2$ .

We use the JKO scheme of [JKO98, Theorem 5.1] to show the existence, uniqueness, and absolute continuousness of solution of PDE (10.22). Since the proof is basically the same as the proof of [JKO98, Theorem 5.1], we will skip several details.

First, we consider the following discrete scheme. Let  $\overline{\rho}_0^h = \rho_0$ , and define  $\{\overline{\rho}_k^h\}_{k\in\mathbb{N}}$  recursively by

$$\overline{\rho}_{k+1}^h \in \arg\min_{\rho \in \mathcal{K}} \{ hF(\rho) + (1/2)W_2^2(\rho, \overline{\rho}_k^h) \}, \tag{10.24}$$

where  $W_2(\mu, \nu)$  is the Wasserstein distance between  $\mu, \nu \in \mathscr{P}(\mathbb{R}^d)$ , with definition

$$W_2^2(\mu,\nu) = \inf \Big\{ \int_{\mathbb{R}^d \times \mathbb{R}^d} \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_2^2 \gamma(\mathrm{d}\boldsymbol{\theta}_1,\mathrm{d}\boldsymbol{\theta}_2) : \gamma \text{ is a coupling of } \mu,\nu \Big\}.$$

For any  $\overline{\rho}_{k-1}^h$ , the optimization problem (10.24) has a unique minimizer  $\overline{\rho}_k^h \in \mathcal{K}$ , where the proof is basically the same as Lemma 10.2, by additionally noting that  $W_2^2(\rho, \overline{\rho}_{k-1}^h)$  as a function of  $\rho$  is lower bounded, lower semi-continuous, and convex over  $\rho \in \mathcal{K}$ .

Hence, we have a sequence of probability densities  $(\overline{\rho}_k^h)_{k\geq 0}$  with each  $\overline{\rho}_k^h \in \mathcal{K}$ . Now we define its interpolation  $\rho^h:(0,\infty)\times\mathbb{R}^d\to[0,\infty)$  by

$$\rho^h(t, \cdot) = \overline{\rho}_k^h$$
 for  $t \in [kh, (k+1)h)$  and  $k \in \mathbb{N}$ .

In the following, we will show that this  $\rho^h$  approximately satisfies PDE (10.23) in the weak form. Let  $\boldsymbol{\xi} \in C_0^{\infty}(\mathbb{R}^d, \mathbb{R}^d)$  be a smooth vector field with bounded support, and define the corresponding flux  $\{\Phi_{\tau}\}_{{\tau}\in\mathbb{R}}$  by

$$\partial_{\tau} \Phi_{\tau} = \boldsymbol{\xi} \circ \Phi_{\tau} \text{ for all } \tau \in \mathbb{R} \quad \text{ and } \quad \Phi_{0} = \text{id.}$$
 (10.25)

For any  $\tau \in \mathbb{R}$ , let the measure  $\nu_{\tau}$  to be the push forward of  $\overline{\rho}_k^h$  under  $\Phi_{\tau}$ . This means that

$$\int_{\mathbb{R}^d} \nu_{\tau}(\boldsymbol{\theta}) \zeta(\boldsymbol{\theta}) d\boldsymbol{\theta} = \int_{\mathbb{R}^d} \overline{\rho}_k^h(\boldsymbol{\theta}) \zeta(\Phi_{\tau}(\boldsymbol{\theta})) d\boldsymbol{\theta}, \quad \text{for all} \quad \zeta \in C(\mathbb{R}^d).$$
 (10.26)

Since  $\overline{\rho}_k^h$  is the minimizer of optimization problem (10.24), we have for each  $\tau > 0$ ,

$$\left(\frac{1}{2}W_2^2(\overline{\rho}_{k-1}^h, \nu_\tau) + hF(\nu_\tau)\right) - \left(\frac{1}{2}W_2^2(\overline{\rho}_{k-1}^h, \overline{\rho}_k^h) + hF(\overline{\rho}_k^h)\right) \ge 0. \tag{10.27}$$

Using the result in the proof of [JKO98, Theorem 5.1], and noting  $\nabla V$  is bounded Lipschitz, we have

$$\frac{\mathrm{d}}{\mathrm{d}\tau} [\langle V, \nu_{\tau} \rangle]_{\tau=0} = \int_{\mathbb{R}^d} \langle \nabla V(\boldsymbol{\theta}), \boldsymbol{\xi}(\boldsymbol{\theta}) \rangle \, \overline{\rho}_k^h(\boldsymbol{\theta}) \mathrm{d}\boldsymbol{\theta}, \tag{10.28}$$

$$\frac{\mathrm{d}}{\mathrm{d}\tau}[\mathrm{Ent}(\nu_{\tau})]_{\tau=0} = \int_{\mathbb{R}^d} \overline{\rho}_k^h(\boldsymbol{\theta}) \cdot \mathrm{div}(\boldsymbol{\xi}(\boldsymbol{\theta})) \mathrm{d}\boldsymbol{\theta}, \tag{10.29}$$

$$\lim \sup_{\tau \to 0+} \frac{1}{\tau} [M(\nu_{\tau}) - M(\overline{\rho}_{k}^{h})] \le \int_{\mathbb{R}^{d}} 2\langle \boldsymbol{\theta}, \boldsymbol{\xi}(\boldsymbol{\theta}) \rangle \, \overline{\rho}_{k}^{h}(\boldsymbol{\theta}) \mathrm{d}\boldsymbol{\theta}, \tag{10.30}$$

$$\lim \sup_{\tau \to 0+} \frac{1}{\tau} [W_2^2(\overline{\rho}_{k-1}^h, \nu_\tau) - W_2^2(\overline{\rho}_{k-1}^h, \overline{\rho}_k^h)] \le \int_{\mathbb{R}^d} 2\langle (\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2), \boldsymbol{\xi}(\boldsymbol{\theta}_1) \rangle \, p(\mathrm{d}\boldsymbol{\theta}_1, \mathrm{d}\boldsymbol{\theta}_2), \tag{10.31}$$

where p is an optimal coupling of  $\rho_k^h$  and  $\rho_{k-1}^h$  in Wasserstein metric. Further we have for any  $\zeta \in C_0^{\infty}(\mathbb{R}^d)$ ,

$$\left| \int_{\mathbb{R}^d} (\overline{\rho}_k^h - \overline{\rho}_{k-1}^h) \zeta d\boldsymbol{\theta} - \int_{\mathbb{R} \times \mathbb{R}} \langle \boldsymbol{\theta}_1 - \boldsymbol{\theta}_2, \nabla \zeta(\boldsymbol{\theta}_1) \rangle d\boldsymbol{p} \right| \le \frac{1}{2} \sup_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\nabla^2 \zeta(\boldsymbol{\theta})\|_{\text{op}} W_2^2(\overline{\rho}_k^h, \overline{\rho}_{k-1}^h). \tag{10.32}$$

We need to further calculate the derivative of  $\langle U, \nu_{\tau}^{\otimes 2} \rangle$  with respect to  $\tau$ . Note U is symmetric, we have

$$\frac{1}{\tau} [\langle U, \nu_{\tau}^{\otimes 2} \rangle - \langle U, (\overline{\rho}_{k}^{h})^{\otimes 2} \rangle] - 2 \int_{\mathbb{R}^{d} \times \mathbb{R}^{d}} \langle \nabla_{\boldsymbol{\theta}_{1}} U(\boldsymbol{\theta}_{1}, \boldsymbol{\theta}_{2}), \boldsymbol{\xi}(\boldsymbol{\theta}_{1}) \rangle \overline{\rho}_{k}^{h}(\boldsymbol{\theta}_{1}) \overline{\rho}_{k}^{h}(\boldsymbol{\theta}_{2}) d\boldsymbol{\theta}_{1} d\boldsymbol{\theta}_{2}$$

$$= \int_{\mathbb{R}^{d} \times \mathbb{R}^{d}} \{ \frac{1}{\tau} [U(\Phi_{\tau}(\boldsymbol{\theta}_{1}), \Phi_{\tau}(\boldsymbol{\theta}_{2})) - U(\Phi_{\tau}(\boldsymbol{\theta}_{1}), \boldsymbol{\theta}_{2})] - \langle \nabla_{\boldsymbol{\theta}_{2}} U(\Phi_{\tau}(\boldsymbol{\theta}_{1}), \boldsymbol{\theta}_{2}), \boldsymbol{\xi}(\boldsymbol{\theta}_{2}) \rangle \} \overline{\rho}_{k}^{h}(\boldsymbol{\theta}_{1}) \overline{\rho}_{k}^{h}(\boldsymbol{\theta}_{2}) d\boldsymbol{\theta}_{1} d\boldsymbol{\theta}_{2}$$

$$+ \int_{\mathbb{R}^{d} \times \mathbb{R}^{d}} \{ \frac{1}{\tau} [U(\Phi_{\tau}(\boldsymbol{\theta}_{1}), \boldsymbol{\theta}_{2}) - U(\boldsymbol{\theta}_{1}, \boldsymbol{\theta}_{2})] - \langle \nabla_{\boldsymbol{\theta}_{1}} U(\boldsymbol{\theta}_{1}, \boldsymbol{\theta}_{2}), \boldsymbol{\xi}(\boldsymbol{\theta}_{1}) \rangle \} \overline{\rho}_{k}^{h}(\boldsymbol{\theta}_{1}) \overline{\rho}_{k}^{h}(\boldsymbol{\theta}_{2}) d\boldsymbol{\theta}_{1} d\boldsymbol{\theta}_{2}$$

$$+ \int_{\mathbb{R}^{d} \times \mathbb{R}^{d}} [\langle \nabla_{\boldsymbol{\theta}_{2}} U(\Phi_{\tau}(\boldsymbol{\theta}_{1}), \boldsymbol{\theta}_{2}), \boldsymbol{\xi}(\boldsymbol{\theta}_{2}) \rangle - \langle \nabla_{\boldsymbol{\theta}_{2}} U(\boldsymbol{\theta}_{1}, \boldsymbol{\theta}_{2}), \boldsymbol{\xi}(\boldsymbol{\theta}_{2}) \rangle ] \overline{\rho}_{k}^{h}(\boldsymbol{\theta}_{1}) \overline{\rho}_{k}^{h}(\boldsymbol{\theta}_{2}) d\boldsymbol{\theta}_{1} d\boldsymbol{\theta}_{2}.$$

According to condition A3,  $\nabla_{\boldsymbol{\theta}_1} U(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)$  is Lipschitz in  $(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)$ , and note  $\boldsymbol{\xi}(\boldsymbol{\theta}) \in C_0^{\infty}(\mathbb{R}^d)$  is uniformly bounded, hence  $1/\tau \cdot [U(\Phi_{\tau}(\boldsymbol{\theta}_1), \boldsymbol{\theta}_2) - U(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)] - \langle \nabla_{\boldsymbol{\theta}_1} U(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2), \boldsymbol{\xi}(\boldsymbol{\theta}_1) \rangle$ ,  $1/\tau [U(\Phi_{\tau}(\boldsymbol{\theta}_1), \boldsymbol{\theta}_2) - U(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2), \boldsymbol{\xi}(\boldsymbol{\theta}_1) \rangle$ , and  $[\langle \nabla_{\boldsymbol{\theta}_2} U(\Phi_{\tau}(\boldsymbol{\theta}_1), \boldsymbol{\theta}_2), \boldsymbol{\xi}(\boldsymbol{\theta}_2) \rangle - \langle \nabla_{\boldsymbol{\theta}_2} U(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2), \boldsymbol{\xi}(\boldsymbol{\theta}_2) \rangle]$  converges to 0 for  $\tau \to 0+$ , uniformly over  $(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) \in \mathbb{R}^d \times \mathbb{R}^d$ . Therefore, we have

$$\frac{\mathrm{d}}{\mathrm{d}\tau} [\langle U, \nu_{\tau}^{\otimes 2} \rangle]_{\tau=0} = 2 \int_{\mathbb{R}^d \times \mathbb{R}^d} \langle \nabla_{\boldsymbol{\theta}_1} U(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2), \boldsymbol{\xi}(\boldsymbol{\theta}_1) \rangle \cdot \overline{\rho}_k^h(\boldsymbol{\theta}_1) \overline{\rho}_k^h(\boldsymbol{\theta}_2) \mathrm{d}\boldsymbol{\theta}_1 \mathrm{d}\boldsymbol{\theta}_2. \tag{10.33}$$

Combining Eq. (10.28) to (10.33), choosing  $\boldsymbol{\xi} = \nabla \zeta$  and  $\boldsymbol{\xi} = -\nabla \zeta$ , we have for any  $\zeta \in C_0^{\infty}(\mathbb{R})$ ,

$$\left| \int_{\mathbb{R}^d} \left\{ \frac{1}{h} (\overline{\rho}_k^h - \overline{\rho}_{k-1}^h) \zeta + (\langle \nabla_{\boldsymbol{\theta}} \Psi_{\lambda}(\boldsymbol{\theta}; \overline{\rho}_k^h), \nabla \zeta \rangle - \Delta \zeta) \overline{\rho}_k^h \right\} d\boldsymbol{\theta} \right| \leq \frac{1}{2} \sup_{\mathbb{R}^d} \|\nabla^2 \zeta\|_{\text{op}} \cdot \frac{1}{h} W_2^2 (\overline{\rho}_{k-1}^h, \overline{\rho}_k^h). \tag{10.34}$$

According to the estimates in [JKO98, Theorem 5.1], for any  $T < \infty$ , there exists a constant  $C < \infty$  such that for all  $N \in \mathbb{N}$  and all  $h \in (0,1]$  with  $Nh \leq T$ , there holds

$$\max \left\{ M(\overline{\rho}_N^h), \quad \int_{\mathbb{R}^d} \max \{ \overline{\rho}_N^h \log(\overline{\rho}_N^h), 0 \} d\boldsymbol{\theta}, \quad R(\overline{\rho}_N^h), \quad \frac{1}{h} \sum_{k=1}^N W_2^2(\overline{\rho}_k^h, \overline{\rho}_{k-1}^h) \right\} \le C. \tag{10.35}$$

As in [JKO98, Theorem 5.1], by de la Vallée-Poussin criterion, the second condition in Eq. (10.35) implies that there exists a measurable function  $(t, \theta) \mapsto \rho(t, \theta)$  and a sequence  $(h_s)_{s\geq 1}$  with  $\lim_{s\to\infty}h_s=0$ , such that  $(t,\theta)\mapsto \rho^{h_s}(t,\theta)$  converges to  $\rho$  weakly in  $L^1((0,T)\times\mathbb{R}^d)$  for all  $T<\infty$ . Eq. (10.35) also guarantees that  $\rho(t,\cdot)\in\mathcal{K}$  for almost every  $t\in(0,\infty)$ , and  $M(\rho),R(\rho)\in L^\infty((0,T))$  for all  $T<\infty$ . By Eq. (10.34) and (10.35), we have that  $\rho$  satisfies Eq. (10.23). Since this equation is not affected by changing  $\rho(t,\cdot)$  for a set of values of t with measure 0, we can ensure that the  $\rho(t,\cdot)\in\mathcal{K}$  for all t. Therefore,  $\rho$  is a solution of the weak form of PDE (10.23).

The uniqueness of solution of Eq. (10.23) can be proved using standard method from theory of elliptic-parabolic equations (see, for instance, [JKO98, Theorem 5.1]). In the proof of uniqueness we need the smoothness property of the solution, which is proved by Lemma 10.7.

**Lemma 10.7.** Assume conditions A1 - A4 hold. Let initialization  $\rho_0 \in \mathcal{K}$  with  $F_{\beta,\lambda}(\rho_0) < \infty$ . Denote the solution of PDE (10.22) to be  $(\rho_t)_{t\geq 0}$ . Then  $\rho_t(\boldsymbol{\theta})$  as a function of  $(t,\boldsymbol{\theta})$  is in  $C^{1,2}((0,\infty)\times\mathbb{R}^d)$ , where  $C^{1,2}((0,\infty)\times\mathbb{R}^d)$  is the function space of continuous function with continuous derivative in time, and second order continuous derivative in space.

Before proving this lemma, we give some notations in the following.

For any open set  $\Omega \subseteq \mathbb{R}^d$ , and  $1 \leq p \leq \infty$ , define  $L^p(\Omega)$  to be the Banach space consisting of all measurable functions on  $\Omega$  with a finite norm

$$||u||_{L^p(\Omega)} \equiv \left(\int_{\Omega} |u(\boldsymbol{\theta})|^p d\boldsymbol{\theta}\right)^{1/p}.$$
 (10.36)

We say  $u \in L^p_{loc}(\Omega)$  if for any compact subset  $\Omega' \subset \Omega$ , we have  $u \in L^p(\Omega')$ . We denote  $\|\cdot\|_{L^p(\mathbb{R}^d)}$  simply by  $\|\cdot\|_{L^p}$ .

For any nonnegative integer l and  $1 \leq p \leq \infty$ , we denote  $W_p^l(\Omega)$  to be the Banach space (Sobolev space) consisting of the elements of  $L^p(S)$  having generalized derivatives of all forms up to order l included, that are p'th power integrable on  $\Omega$ . The norm in  $W_p^l(\Omega)$  is defined by the equality

$$||u||_{L^{p}(\Omega)}^{(l)} = \sum_{j=0}^{l} \langle \langle u \rangle \rangle_{L^{p}(\Omega)}^{(j)}, \quad \langle \langle u \rangle \rangle_{L^{p}(\Omega)}^{(j)} = \sum_{|\boldsymbol{\alpha}|=j} ||D_{\boldsymbol{\theta}}^{\boldsymbol{\alpha}} u||_{L^{p}(\Omega)}, \tag{10.37}$$

where  $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_d)$  is a multi-index with  $|\boldsymbol{\alpha}| = \sum_{i=1}^d \alpha_i$ , and  $D_{\boldsymbol{\theta}}^{\boldsymbol{\alpha}} u = \partial^{|\boldsymbol{\alpha}|} u / \partial \theta_1^{\alpha_1} \cdots \partial \theta_d^{\alpha_d}$ . Let  $(t_1, t_2) \subseteq (0, T)$  be an open interval and  $\Omega \subseteq \mathbb{R}^d$  be an open set, in these three paragraphs

Let  $(t_1, t_2) \subseteq (0, T)$  be an open interval and  $\Omega \subseteq \mathbb{R}^d$  be an open set, in these three paragraphs we temporarily denote  $S = (t_1, t_2) \times \Omega$ . For any  $1 \leq r, p \leq \infty$ , define  $L^{r,p}(S)$  to be the Banach space consisting of all measurable functions on S with a finite norm

$$||u||_{L^{r,p}(S)} \equiv \left(\int_{t_1}^{t_2} \left(\int_{\Omega} |u(t,\boldsymbol{\theta})|^p d\boldsymbol{\theta}\right)^{r/p} dt\right)^{1/r}.$$
 (10.38)

We say  $u \in L^{r,p}_{loc}(S)$  if for any compact subset  $[t'_1, t'_2] \subset (t_1, t_2)$  and compact subset  $\Omega' \subset \Omega$ , we have  $u \in L^{r,p}([t'_1, t'_2] \times \Omega')$ . We will denote  $L^{p,p}(S)$  by  $L^p(S)$ , and  $L^{p,p}_{loc}(S)$  by  $L^p_{loc}(S)$ .

For nonnegative integer l and  $1 \leq p \leq \infty$ , we denote  $W_p^{2l,l}(S)$  to be the Banach space consisting of the elements of  $L^p(S)$  having generalized derivatives of the form  $D_t^r D_{\theta}^{\alpha}$  with r and  $\alpha$  satisfying the inequality  $2r + |\alpha| \leq 2l$ . The corresponding norm is defined by

$$||u||_{L^{p}(S)}^{(2l)} = \sum_{j=0}^{2l} \langle \langle u \rangle \rangle_{L^{p}(S)}^{(j)}, \quad \langle \langle u \rangle \rangle_{L^{p}(S)}^{(j)} = \sum_{|\alpha|+2r=j} ||D_{t}^{r} D_{\theta}^{\alpha} u||_{L^{p}(S)}.$$
 (10.39)

We denote  $C^{m,n}(S)$  to be the function space of continuous function with m continuous derivative in time, and n continuous derivatives in space. For example,  $u \in C^{1,2}(S)$  if and only if  $u, \partial_t u, \nabla_{\theta} u, \nabla_{\theta}^2 u \in C^{0,0}(S) \equiv C(S)$ . We say  $u \in C_c^{m,n}(S)$  if  $u \in C^{m,n}(S)$  and the support of u is compact. We will denote  $C^{n,n}(S)$  by  $C^n(S)$ , and  $C_c^{n,n}(S)$  by  $C_c^n(S)$ .

For any measurable functions f, g defined on  $\mathbb{R}^d$ , we denote f \* g to be their space convolution, which is a measurable function on  $\mathbb{R}^d$ , with

$$(f * g)(\boldsymbol{\theta}) = \int_{\mathbb{R}^d} f(\boldsymbol{\theta}') g(\boldsymbol{\theta} - \boldsymbol{\theta}') d\boldsymbol{\theta}'.$$
 (10.40)

For any measurable function u, v defined on  $\mathbb{R} \times \mathbb{R}^d$ , we denote  $u *_2 v$  to be their space and time convolution, which is a measurable function on  $\mathbb{R} \times \mathbb{R}^d$ , with

$$(u *2 v)(t, \boldsymbol{\theta}) = \int_{\mathbb{R}} dt' \int_{\mathbb{R}^d} u(t', \boldsymbol{\theta}') v(t - t', \boldsymbol{\theta} - \boldsymbol{\theta}') d\boldsymbol{\theta}'.$$
 (10.41)

If u, v are defined on a subset of  $\mathbb{R} \times \mathbb{R}^d$ , we define  $u *_2 v$  using their zero extensions.

We denote G to be the heat kernel, where for t > 0, we have

$$G(t, \boldsymbol{\theta}) = t^{-d/2} g(t^{-1/2} \boldsymbol{\theta}), \quad g(\boldsymbol{\theta}) = (2\pi)^{-d/2} \exp\{-1/2 \cdot \|\boldsymbol{\theta}\|_2^2\}.$$
 (10.42)

*Proof.* The proof is similar to the one of [JKO98, Theorem 5.1], so we will skip some details. Without loss of generality we can set  $\beta = 1$ , and  $\xi(t) = 1/2$  (different choices can be obtained by rescaling  $\Psi(\theta; \rho)$  and reparametrizing time).

Let  $E = (0, \infty) \times \mathbb{R}^d$ . With a slight abuse of notations, we denote  $\Psi(t, \boldsymbol{\theta}) = \Psi_{\lambda}(\boldsymbol{\theta}; \rho_t)$ . Since  $V \in C^4(\mathbb{R}^d)$ , and  $\nabla_1^k U$  are uniformly bounded for  $0 \le k \le 4$ , we have  $\nabla_{\boldsymbol{\theta}}^k \Psi \in L^{\infty}_{loc}(E)$  for  $0 \le k \le 4$ .

In the following, we will write  $\rho(t, \boldsymbol{\theta}) = \rho_t(\boldsymbol{\theta})$  for clarity. When we write  $\rho(t)$ , we regard it as a function in  $L^1(\mathbb{R}^d)$  at any fixed t. For other functions, we also use this convention.

# Step 1. Show that $\rho \in L^{\infty,p}_{loc}(E)$ .

Taking G to be the heat kernel, it is easy to see that

$$||G(t)||_{L^p} = t^{(\frac{1}{p}-1)\frac{d}{2}} ||g||_{L^p}, \quad ||\nabla G(t)||_{L^p} = t^{\frac{1}{p}\frac{d}{2} - \frac{d+1}{2}} ||\nabla g||_{L^p}.$$

Then for any  $\eta \in C_c^{\infty}(\mathbb{R}^d)$ , Duhamel's principle gives

$$\rho(t)\eta = \int_{\varepsilon}^{t} [\rho(s)(\Delta \eta - \langle \nabla \Psi(s), \nabla \eta \rangle)] * G(t - s) ds$$

$$+ \int_{\varepsilon}^{t} [\rho(s)(2\nabla \eta - \eta \nabla \Psi(s))] * \nabla G(t - s) ds + (\rho(\varepsilon)\eta) * G_{\varepsilon}(t)$$
(10.43)

for almost every  $0 \le \varepsilon < t < \infty$ , where \* denotes convolution in the  $\boldsymbol{\theta}$ -variables, and  $G_{\varepsilon}(t,\boldsymbol{\theta}) \equiv G(t-\varepsilon,\boldsymbol{\theta})$ . By Young's convolution inequality, we have  $||f*g||_{L^r} \le C||f||_{L^p}||g||_{L_q}$  for 1/p+1/q=1/r+1 and  $p,q,r \ge 1$ . For fixed t, we estimate the  $L^p(\mathbb{R}^d)$  norm of  $\rho(t)\eta$ , which gives

$$\begin{split} \|\rho(t)\eta\|_{L^{p}} &\leq \int_{\varepsilon}^{t} \|\rho(s)(\Delta\eta - \langle \nabla\Psi(s), \nabla\eta\rangle)\|_{L^{1}} \|G(t-s)\|_{L^{p}} \mathrm{d}s \\ &+ \int_{\varepsilon}^{t} \|\rho(t)(2\nabla\eta - \eta\nabla\Psi(t))\|_{L^{1}} \|\nabla G(t-s)\|_{L^{p}} \mathrm{d}s + \|\rho(\varepsilon)\eta\|_{L^{1}} \|G(t-\varepsilon)\|_{L^{p}} \\ &\leq & \sup_{s \in [\varepsilon,t]} \|\rho(s)(\Delta\eta - \langle \nabla\Psi(s), \nabla\eta\rangle)\|_{L^{1}} \|g\|_{L^{p}} \int_{0}^{t-\varepsilon} s^{(\frac{1}{p}-1)\frac{d}{2}} \mathrm{d}s \\ &+ & \sup_{s \in [\varepsilon,t]} \|\rho(s)(2\nabla\eta - \eta\nabla\Psi(s))\|_{L^{1}} \|\nabla g\|_{L^{p}} \int_{0}^{t-\varepsilon} s^{\frac{1}{p}\frac{d}{2} - \frac{d+1}{2}} \mathrm{d}s \\ &+ \|\rho(\varepsilon)\eta\|_{L^{1}} \|g\|_{L^{p}} (t-\varepsilon)^{(\frac{1}{p}-1)\frac{d}{2}} \end{split}$$

for almost every  $0 \le \varepsilon < t < \infty$ . For p < d/(d-1), the s-integrals are finite. Therefore, we have  $\rho \eta \in L^{\infty,p}((\delta,T) \times \mathbb{R}^d)$  for any  $\delta,T$  such that  $\varepsilon < \delta < T < \infty$ . Hence we have  $\rho \in L^{\infty,p}_{\mathrm{loc}}((0,\infty) \times \mathbb{R}^d)$ . Step 2. Show that  $\rho \in L^{\infty}_{\mathrm{loc}}((0,\infty) \times \mathbb{R}^d)$  using bootstrap.

In what follows, we let  $E \equiv (0, \infty) \times \mathbb{R}^d$ .

We can iteratively use the strategy in step 1 to show that  $\rho \in L^{\infty}_{loc}(E)$ . We will summarize our key estimates in Step 1 as follows. For any measurable function u defined on  $S = (\delta, T) \times \mathbb{R}^d$  for some  $0 \le \delta < T < \infty$ , we have

$$||u *_{2} G||_{L^{\infty,p_{o}}(S)} \le C||u||_{L^{\infty,p_{i}}(S)}, \tag{10.44}$$

$$||u *_{2} \nabla G||_{L^{\infty,p_{o}}(S)} \le C||u||_{L^{\infty,p_{i}}(S)}, \tag{10.45}$$

provided that the  $p_o, p_i$  satisfy the relations

$$1 \le p_i \le p_o, \quad d \cdot (1/p_i - 1/p_o) < 1.$$
 (10.46)

Here, C is a constant depends only on  $T, \delta$  and on  $p_i, p_o$ .

Define  $\varphi_1 \equiv \rho(\Delta \eta - \langle \nabla \Psi, \nabla \eta \rangle) \mathbf{1}\{t > \varepsilon\}, \ \varphi_2 \equiv \rho(2\nabla \eta - \eta \nabla \Psi) \mathbf{1}\{t > \varepsilon\}, \ \text{and} \ \psi \equiv \rho(\varepsilon)\eta$ . Then Eq. (10.43) reads

$$\rho \eta = \varphi_1 *_2 G + \varphi_2 *_2 \nabla G + \psi * G_{\varepsilon}. \tag{10.47}$$

Since  $\psi = \rho(\varepsilon)\eta \in L^1(\mathbb{R}^d)$ , the behavior of  $\psi * G_{\varepsilon}$  on  $S = (\delta, T) \times \mathbb{R}^d$  for  $\varepsilon < \delta < T < \infty$  will be extremely nice: for any generalized gradient  $D_t^r D^{\alpha} [\psi * G_{\varepsilon}]$ ,

$$||D_t^r D^{\alpha}[\psi * G_{\varepsilon}]||_{L^{\infty}(S)} \le ||\psi||_{L^1(\mathbb{R}^d)} ||D_t^r D^{\alpha} G_{\varepsilon}||_{L^{\infty}(S)} < \infty.$$

$$(10.48)$$

Hence  $D_t^r D^{\alpha}[\psi * G_{\varepsilon}] \in L^{\infty}(S)$ . From now on, we fix  $0 < \varepsilon < \delta < T < \infty$  and take  $S \equiv (\delta, T) \times \mathbb{R}^d$ . According to Eq. (10.47) we have

$$\|\rho\eta\|_{L^{\infty,p_o}(S)} \le \|\varphi_1 *_2 G\|_{L^{\infty,p_o}(S)} + \|\varphi_2 *_2 \nabla G\|_{L^{\infty,p_o}(S)} + \|\psi * G_{\varepsilon}\|_{L^{\infty,p_o}(S)}$$

$$\le C\{\|\varphi_1\|_{L^{\infty,p_i}(S)} + \|\varphi_2\|_{L^{\infty,p_i}(S)} + \|\psi\|_{L^1(\mathbb{R}^d)}\}$$
(10.49)

Now we assume  $\rho \in L^{\infty,p_i}_{loc}(E)$  for some  $p_i$ . Note  $\nabla \Psi \in L^{\infty}_{loc}(E)$  so that  $\max\{\|\varphi_1\|_{L^{\infty,p_i}(S)}, \|\varphi_2\|_{L^{\infty,p_i}(S)}\} \le C_{\eta}\|\rho\|_{L^{\infty,p_i}((\delta,T)\times\Omega_2)}$ , where  $\Omega_2 \supseteq \operatorname{supp}(\eta)$  is a compact set. As a result, for any  $\eta \in C^{\infty}_c(\mathbb{R}^d)$ , we have

$$\|\rho\|_{L^{\infty,p_o}((\delta,T)\times\Omega_1)} \le C_{\eta}(\|\rho\|_{L^{\infty,p_i}((\delta,T)\times\Omega_2)} + 1),$$
 (10.50)

where  $\Omega_1 \subseteq \text{supp}(\eta) \subseteq \Omega_2$ . Therefore,  $\rho \in L^{\infty,p_o}_{\text{loc}}(E)$ , where  $p_i, p_o$  satisfy Eq. (10.46).

Note there exists a sequence  $p_{i,l}, p_{o,l}$  for  $1 \le l \le k$  and  $k < \infty$ , so that  $p_{i,l+1} = p_{o,l}, p_{i,1} = p < d/(d-1), p_{i,k} = \infty$ , and  $p_{i,l}, p_{o,l}$  for fixed l satisfies Eq. (10.46). Since we have  $\rho \in L^{\infty,p}_{loc}(E)$ , using Eq. (10.50) iteratively, we have  $\rho \in L^{\infty,p_{o,l}}_{loc}(E)$  for any  $1 \le l \le k$ . As a result, we have  $\rho \in L^{\infty}_{loc}(E)$ . Step 3. Derivatives,  $D\rho$ ,  $D^2\rho$ , and  $D^3\rho$ .

By [LSU88, Chapter IV, section 3, (3.1)], for any function u defined on  $E = (0, \infty) \times \mathbb{R}^d$ , we have

$$\langle \langle G *_2 u \rangle \rangle_{L^p(E)}^{(2m+2)} \le C \langle \langle u \rangle \rangle_{L^p(E)}^{(2m)}, \tag{10.51}$$

where 1 and m is a nonnegative integer.

First, we show the regularity of  $D\rho$ . Note that  $\rho \in L^{\infty}_{loc}(E)$ ,  $\eta \in C^{\infty}_{c}(\mathbb{R}^{d})$ ,  $\nabla \Psi \in L^{\infty}_{loc}(E)$ , we have  $\varphi_{1}, \varphi_{2} \in L^{\infty}(E)$ . Due to Eq. (10.51), we have  $D^{2}\{\varphi_{1}*_{2}G\}$ ,  $D^{2}\{\varphi_{2}*_{2}G\} \in L^{\infty}(E)$ , which also implies  $D\{\varphi_{1}*_{2}G\} \in L^{\infty}_{loc}(E)$ . Hence we have  $D(\rho\eta) = D\{\varphi_{1}*_{2}G\} + D^{2}\{\varphi_{2}*_{2}G\} + D[\psi*G_{\varepsilon}] \in L^{\infty}(S)$ , which gives  $D\rho \in L^{\infty}_{loc}(E)$ .

Then we show the regularity of  $D^2\rho$ . Note that  $\nabla^2\Psi\in L^\infty_{loc}(E)$ , we have  $D\varphi_1, D\varphi_2\in L^\infty(E)$ . Due to Eq. (10.51), we have  $D^3\{\varphi_1*_2G\}, D^3\{\varphi_2*_2G\}\in L^\infty(E)$ , which also implies  $D^2\{\varphi_1*_2G\}\in L^\infty_{loc}(E)$ . Hence we have  $D^2(\rho\eta)=D^2\{\varphi_1*_2G\}+D^3\{\varphi_2*_2G\}+D^2[\psi*G_\varepsilon]\in L^\infty(S)$ , which gives  $D^2\rho\in L^\infty_{loc}(E)$ .

Next we show the regularity of  $D^3\rho$ . Note that  $\nabla^3\Psi\in L^\infty_{loc}(E)$ , we have  $D^2\varphi_1, D^2\varphi_2\in L^\infty(E)$ . Due to Eq. (10.51), we have  $D^4\{\varphi_1*_2G\}, D^4\{\varphi_2*_2G\}\in L^\infty(E)$ , which also implies  $D^3\{\varphi_1*_2G\}\in L^\infty(E)$ 

 $L^{\infty}_{loc}(E)$ . Hence we have  $D^3(\rho\eta) = D^3\{\varphi_1 *_2 G\} + D^4\{\varphi_2 *_2 G\} + D^3[\psi * G] \in L^{\infty}(S)$ , which gives  $D^3\rho \in L^{\infty}_{loc}(E)$ .

# Step 4. Derivatives, $D_t \rho$ , $D_t D \rho$ , and $D_t D^2 \rho$ .

Now we study the regularity of  $D_t \rho$ ,  $D_t D \rho$ ,  $D_t D^2 \rho$ . Note we have  $D_t(\rho \eta) = D_t \{\varphi_1 *_2 G\} - D_t \{D\varphi_1 *_2 G\} + D_t \{\psi *_G G\}$ . Due to Eq. (10.51),  $\varphi_1, D\varphi_2 \in L^{\infty}(E)$  implies that  $D_t \{\varphi_1 *_2 G\}$ ,  $D_t \{D\varphi_1 *_2 G\} \in L^{\infty}(E)$  and hence  $D_t [\rho \eta] \in L^{\infty}(S)$ ,  $D_t \rho \in L^{\infty}_{loc}(E)$ .

Note we have  $D_t D(\rho \eta) = D_t \{D\varphi_1 *_2 G\} + D_t \{D^2\varphi_1 *_2 G\} + D_t \{D\psi * G_{\varepsilon}\}$ . The fact that  $D\varphi_1, D^2\varphi_2 \in L^{\infty}(E)$  implies that  $D_t \{D\varphi_1 *_2 G\}, D_t \{D^2\varphi_1 *_2 G\} \in L^{\infty}(E)$  and hence  $D_t D\rho \in L^{\infty}_{loc}(E)$ .

Note we have  $D_t D^2(\rho \eta) = D_t \{D^2 \varphi_1 *_2 G\} - D_t \{D^3 \varphi_1 *_2 G\} + D_t \{D^2 \psi *_G \varepsilon\}$ . Note that  $\nabla^4 \Psi \in L^{\infty}_{loc}(E)$ , hence  $D^3 \varphi_2 \in L^{\infty}(E)$ . Combining with the fact that  $D^2 \varphi_1 \in L^{\infty}(E)$ , we have  $D_t \{D^2 \varphi_1 *_2 G\}, D_t \{D^3 \varphi_1 *_2 G\} \in L^{\infty}(E)$  and hence  $D_t D^2 \rho \in L^{\infty}_{loc}(E)$ .

## Step 5. Derivatives, $D_t^2 \rho$ .

Finally we show the regularity of  $D_t^2 \rho$ . We have  $D_t^2(\rho \eta) = D_t \{D_t[\varphi_1 *_2 G] - D_t[D\varphi_1 *_2 G] + D_t[\psi * G_{\varepsilon}]\}$ , and

$$D_t[\varphi_1 *_2 G] = [\Delta \varphi_1] *_2 G + \varphi_1(\varepsilon) * G_{\varepsilon}, \tag{10.52}$$

$$D_t[D\varphi_2 *_2 G] = [D\Delta\varphi_2] *_2 G + [D\varphi_2(\varepsilon)] * G_{\varepsilon}. \tag{10.53}$$

Note that  $\nabla^4 \Psi \in L^{\infty}_{loc}(E)$ , we have  $\Delta \varphi_1, D\Delta \varphi_1 \in L^{\infty}_{loc}(E)$ , and  $\varphi_1(\varepsilon), D\varphi_2(\varepsilon) \in L^1(\mathbb{R}^d)$ . Hence according to Eq. (10.51), we have  $D_t\{[\Delta \varphi_1] *_2 G\}, D_t\{[D\Delta \varphi_2] *_2 G\}$ . In addition  $D_t\{\varphi_1(\varepsilon) *_G G_{\varepsilon}\}, D_t\{[D\varphi_2(\varepsilon)] *_G G_{\varepsilon}\} \in L^{\infty}(S)$ . As a result, we have  $D_t^2 \rho \in L^{\infty}_{loc}(E)$ .

### Step 6. Finish the proof.

As a result, we have  $\rho, D\rho, D^2\rho, D^3\rho, D_t\rho, D_tD\rho, D_tD^2\rho, D_t^2\rho \in L^{\infty}_{loc}(E)$ . Sobolev embedding theorem implies that  $\rho, \partial_t \rho, \nabla_{\theta} \rho, \nabla_{\theta}^2 \rho \in C^{0,0}(\mathbb{R}^d)$ . In other words,  $\rho \in C^{1,2}(E)$ , which is the desired result.

**Lemma 10.8.** Assume conditions A1 - A4 hold. Let initialization  $\rho_0 \in \mathcal{K}$  with  $F_{\beta,\lambda}(\rho_0) < \infty$ . Denote the solution of PDE (10.22) to be  $(\rho_t)_{t\geq 0}$ . Then  $\rho_t(\boldsymbol{\theta}) > 0$  for any  $(t,\boldsymbol{\theta}) \in (0,\infty) \times \mathbb{R}^d$ .

*Proof.* Note that  $\rho_t \in C^{1,2}((0,\infty) \times \mathbb{R}^d)$ . By the Harnack's inequality [Eva09], we immediately have  $\rho_t(\boldsymbol{\theta}) > 0$  for any  $(t,\boldsymbol{\theta}) \in (0,\infty) \times \mathbb{R}^d$ .

We say  $\rho_*$  is a fixed point of PDE (10.22), if its solution  $(\rho_t)_{t\geq 0}$  starting from  $\rho_*$  satisfies  $\rho_t \equiv \rho_*$  for any  $t\geq 0$ .

**Lemma 10.9.** Assume conditions A1 - A3 hold. Then any fixed point  $\rho_*$  of PDE (10.22) with  $\rho_* \in \mathcal{K}$  must satisfy the Boltzmann fixed point condition (10.13).

*Proof.* Suppose  $\rho_* \in \mathcal{K}$  is a fixed point of PDE (10.22), taking  $W(\boldsymbol{\theta}) \equiv \Psi_{\lambda}(\boldsymbol{\theta}; \rho_*)$ , then  $\rho_* \in \mathcal{K}$  is a fixed point of the Fokker-Planck equation (10.54).

$$\partial_t \rho_t(\boldsymbol{\theta}) = 2\xi(t)\nabla \cdot (\nabla W(\boldsymbol{\theta})\rho_t(\boldsymbol{\theta})) + 2\xi(t)/\beta \cdot \Delta_{\boldsymbol{\theta}}\rho_t(\boldsymbol{\theta}). \tag{10.54}$$

Since  $\lambda/2 \cdot \|\boldsymbol{\theta}\|_2^2 - 2K_3 \leq \Psi_{\lambda}(\boldsymbol{\theta}; \rho_*) \leq \lambda/2 \cdot \|\boldsymbol{\theta}\|_2^2 + 2K_3$ , the Fokker-Planck equation has a unique fixed point [MV00], which solves

$$\rho_*(\boldsymbol{\theta}) = \frac{1}{Z_{\beta}} \exp\{-\beta W(\boldsymbol{\theta})\}, \quad Z_{\beta} = \int_{\mathbb{R}^d} \exp\{-\beta W(\boldsymbol{\theta})\} d\boldsymbol{\theta}.$$

This is exactly the Boltzmann fixed point condition.

**Lemma 10.10.** Assume conditions A1 - A4 hold. Let  $(\rho_t)_{t\geq 0}$  be the solution of PDE (10.22) for an initialization  $\rho_0 \in \mathcal{K}$ . Then the free energy  $F_{\beta,\lambda}(\rho_t)$  is differentiable with respect to t, with

$$\partial_t F_{\beta,\lambda}(\rho_t) = -2\xi(t) \int_{\mathbb{R}^d} \|\nabla_{\boldsymbol{\theta}}(\Psi_{\lambda}(\boldsymbol{\theta}; \rho_t) + 1/\beta \cdot \log \rho_t(\boldsymbol{\theta}))\|_2^2 \rho_t(\boldsymbol{\theta}) d\boldsymbol{\theta}.$$
 (10.55)

Therefore,  $F_{\beta,\lambda}(\rho_t)$  is non-increasing in t.

*Proof.* Calculate the differential of the free energy along the curve  $\rho_t$ , we have

$$\partial_t F_{\beta,\lambda}(\rho_t) = \int_{\mathbb{R}^d} \Psi_{\lambda}(\boldsymbol{\theta}; \rho_t) \partial_t \rho_t(\boldsymbol{\theta}) d\boldsymbol{\theta} + 1/\beta \cdot \int \log(\rho_t(\boldsymbol{\theta})) \partial_t \rho_t(\boldsymbol{\theta}) d\boldsymbol{\theta}$$
$$= -2\xi(t) \int_{\mathbb{R}^d} \|\nabla_{\boldsymbol{\theta}} (\Psi_{\lambda}(\boldsymbol{\theta}; \rho_t) + 1/\beta \cdot \log \rho_t(\boldsymbol{\theta}))\|_2^2 \rho_t(\boldsymbol{\theta}) d\boldsymbol{\theta}.$$

**Lemma 10.11.** Assume  $K_0 \|\boldsymbol{\theta}\|_2^2 - K_1 \leq \Phi(\boldsymbol{\theta}) \leq K_0 \|\boldsymbol{\theta}\|_2^2 + K_1$  for some positive constant  $K_0, K_1$ . Define

$$\mu_*(\mathrm{d}\boldsymbol{\theta}) = \frac{1}{Z_*} \exp\{-\Phi(\boldsymbol{\theta})\} \mathrm{d}\boldsymbol{\theta}, \quad Z_* = \int_{\mathbb{R}^d} \exp\{-\Phi(\boldsymbol{\theta})\} \mathrm{d}\boldsymbol{\theta}$$
 (10.56)

Let  $\mathcal{D} \equiv \{ f \in L^2(\mathbb{R}^d, \mu_*) \cap C^1(\mathbb{R}^d) : \|\nabla f\|_2 \in L^2(\mathbb{R}^d, \mu_*) \}$ . For any  $f \in \mathcal{D}$ , define

$$I(f) \equiv \int_{\mathbb{R}^d} \|\nabla f(\boldsymbol{\theta})\|_2^2 \cdot \mu_*(\mathrm{d}\boldsymbol{\theta}) < \infty.$$
 (10.57)

Assume  $(f_n)_{n\geq 1}\subseteq \mathcal{D}$ , with  $\lim_{n\to\infty}I(f_n)=0$ , and  $f_n$  converges weakly to  $f_*$  in  $L^2(\mathbb{R}^d,\mu_*)$ . Then  $f_*(\boldsymbol{\theta})\equiv F_*$  for some constant  $F_*$ .

*Proof.* First we show that the measure  $\mu_*$  satisfies the Poincare inequality: for any  $f \in \mathcal{D}$ ,

$$\mu_*((f - \mu_*(f))^2) \le K \cdot I(f),$$
(10.58)

for some constant K.

Let  $\mu$  be the Gaussian distribution  $N(\mathbf{0}, 1/(2K_0) \cdot \mathbf{I}_d)$ . Then for any  $\boldsymbol{\theta} \in \mathbb{R}^d$ ,

$$\mu(\boldsymbol{\theta}) \cdot \exp\{-2K_1\} \le \mu_*(\boldsymbol{\theta}) \le \mu(\boldsymbol{\theta}) \cdot \exp\{2K_1\}. \tag{10.59}$$

Therefore, for any nonnegative measurable function  $f: \mathbb{R}^d \to [0, \infty)$  and  $g: \mathbb{R}^d \times \mathbb{R}^d \to [0, \infty)$ , letting  $(G, G') \sim \mu \times \mu$  and  $(X, X') \sim \mu_* \times \mu_*$ , we have

$$\mathbb{E}[f(G)] \cdot \exp\{-2K_1\} \le \mathbb{E}[f(X)] \le \mathbb{E}[f(G)] \cdot \exp\{2K_1\},$$

$$\mathbb{E}[g(G, G')] \cdot \exp\{-4K_1\} \le \mathbb{E}[g(X, X')] \le \mathbb{E}[g(G, G')] \cdot \exp\{4K_1\}.$$

Note we have the Poincare inequality for the Gaussian distribution  $\mu$ ,

$$Var[f(G)] \le 1/(2K_0) \cdot \mathbb{E}[\|\nabla f(G)\|_2^2]$$
(10.60)

П

for any differentiable f. Therefore, we have

$$\operatorname{Var}[f(X)] = \frac{1}{2} \mathbb{E}[(f(X) - f(X'))^{2}] \leq \frac{1}{2} \exp\{4K_{1}\} \cdot \mathbb{E}[(f(G) - f(G'))^{2}]$$

$$= \exp\{4K_{1}\} \cdot \operatorname{Var}[f(G)] \leq 1/(2K_{0}) \cdot \exp\{4K_{1}\} \cdot \mathbb{E}[\|\nabla f(G)\|_{2}^{2}]$$

$$\leq 1/(2K_{0}) \cdot \exp\{6K_{1}\} \cdot \mathbb{E}[\|\nabla f(X)\|_{2}^{2}].$$

This proves the Poincare inequality (10.58) for  $\mu_*$ .

Since  $\lim_{n\to\infty} I(f_n) = 0$ , due to (10.58), we immediately have  $f_n - \mu_*(f_n)$  converges to 0 in  $L^2(\mathbb{R}^d, \mu_*)$ . Note we assumed  $f_n$  converges weakly to  $f_*$  in  $L^2(\mathbb{R}^d, \mu_*)$ , and  $1 \in L^2(\mathbb{R}^d, \mu_*)$ , we have

$$\lim_{n\to\infty}\mu_*(f_n)=\mu_*(f).$$

Therefore,  $f_n - \mu_*(f_n)$  converges weakly to  $f_* - \mu_*(f_*)$  in  $L^2(\mathbb{R}^d, \mu_*)$ . Hence  $f_*(\theta) \equiv \mu_*(f_*)$ .

**Lemma 10.12.** Assume conditions A1 - A4 hold. Then the solution  $(\rho_t)_{t\geq 0}$  of PDE (10.22) for any initialization  $\rho_0 \in \mathcal{K}$  converges weakly to  $\rho_* \in \mathcal{K}$  as  $t \to \infty$ , where  $\rho_*$  is the unique solution of the Boltzmann fixed point condition, which is the global minimizer of  $F_{\beta,\lambda}$ .

*Proof.* According to Lemma 10.10,  $F_{\beta,\lambda}$  is non-increasing along the solution path. According to Lemma 10.2,  $F_{\beta,\lambda}(\rho_t)$  is lower bounded. Therefore, we have

$$\lim_{t \to \infty} \int_{\mathbb{R}^d} \|\nabla_{\boldsymbol{\theta}}(\Psi_{\lambda}(\boldsymbol{\theta}; \rho_t) + 1/\beta \cdot \log \rho_t(\boldsymbol{\theta}))\|_2^2 \rho_t(\boldsymbol{\theta}) d\boldsymbol{\theta} = 0.$$
 (10.61)

Since  $M(\rho_t)$  is uniformly bounded, by Lemma 10.6,  $(\rho_t)_{t\geq 0}$  as a sequence of probability distribution in  $\mathscr{P}(\mathbb{R}^d)$  is uniformly tight. Hence there exists  $\rho_* \in \mathscr{P}(\mathbb{R}^d)$  and a subsequence  $(\rho_{t_k})_{k\geq 1}$  with  $\lim_{k\to\infty} t_k = \infty$  such that  $(\rho_{t_k})_{k\geq 1}$  converges weakly to  $\rho_*$ . By Lemma 10.6 and Lemma 10.1,  $\{\int \max\{\rho_{t_k}\log\rho_{t_k},0)\}d\theta\}_{k\geq 1}$  is uniformly bounded. Using de la Vallée-Poussin's criteria, we can show that  $(\rho_{t_k})_{k\geq 1}$  is uniformly integrable, and hence  $\rho_*$  is absolute continuous with respect to Lebesgue measure, which means  $\rho_*$  has a density.

Note we have

$$\nabla_{\boldsymbol{\theta}} \Psi_{\lambda}(\boldsymbol{\theta}; \rho_{t}) - \nabla_{\boldsymbol{\theta}} \Psi_{\lambda}(\boldsymbol{\theta}; \rho_{*}) = \int_{\mathbb{R}^{d}} \nabla_{\boldsymbol{\theta}} U(\boldsymbol{\theta}, \boldsymbol{\theta}') (\rho_{t}(\boldsymbol{\theta}') - \rho_{*}(\boldsymbol{\theta}')) d\boldsymbol{\theta}'.$$

According to condition A3,  $\nabla_{\theta}U$  is  $K_3$ -bounded-Lipschitz with respect to  $(\theta, \theta')$ . Therefore,

$$\sup_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\nabla_{\boldsymbol{\theta}} \Psi_{\lambda}(\boldsymbol{\theta}; \rho_t) - \nabla_{\boldsymbol{\theta}} \Psi_{\lambda}(\boldsymbol{\theta}; \rho_*)\|_2 \le K_3 \cdot d_{\text{BL}}(\rho_t, \rho_*) \to 0, \tag{10.62}$$

as  $d_{\rm BL}(\rho_t, \rho_*) \to 0$ . Accordingly, we have

$$\lim_{k \to \infty} \int_{\mathbb{R}^d} \|\nabla_{\boldsymbol{\theta}}(\Psi_{\lambda}(\boldsymbol{\theta}; \rho_{t_k}) - \Psi_{\lambda}(\boldsymbol{\theta}; \rho_*))\|_2^2 \rho_{t_k}(\boldsymbol{\theta}) d\boldsymbol{\theta} \le K_3^2 \cdot \lim_{k \to \infty} d_{\text{BL}}(\rho_{t_k}, \rho_*)^2 = 0.$$
 (10.63)

Combining Eq. (10.63) with Eq. (10.61), we have

$$\lim_{k \to \infty} \int_{\mathbb{R}^d} \|\nabla_{\boldsymbol{\theta}}(\Psi_{\lambda}(\boldsymbol{\theta}; \rho_*) + 1/\beta \cdot \log \rho_{t_k}(\boldsymbol{\theta}))\|_2^2 \rho_{t_k}(\boldsymbol{\theta}) d\boldsymbol{\theta} = 0.$$
 (10.64)

Note we have

$$\int_{\mathbb{R}^{d}} \|\nabla_{\boldsymbol{\theta}}(\Psi_{\lambda}(\boldsymbol{\theta}; \rho_{*}) + 1/\beta \cdot \log \rho_{t_{k}}(\boldsymbol{\theta}))\|_{2}^{2} \rho_{t_{k}}(\boldsymbol{\theta}) d\boldsymbol{\theta}$$

$$= \frac{1}{\beta^{2}} \int_{\mathbb{R}^{d}} \|\nabla_{\boldsymbol{\theta}}(\rho_{t_{k}}(\boldsymbol{\theta}) \exp\{\beta \Psi_{\lambda}(\boldsymbol{\theta}; \rho_{*})\})\|_{2}^{2} \cdot \rho_{t_{k}}(\boldsymbol{\theta})^{-1} \exp\{-2\beta \Psi_{\lambda}(\boldsymbol{\theta}; \rho_{*})\} d\boldsymbol{\theta}$$

$$= \frac{4}{\beta^{2}} \int_{\mathbb{R}^{d}} \|\nabla_{\boldsymbol{\theta}}[(\rho_{t_{k}}(\boldsymbol{\theta}) \exp\{\beta \Psi_{\lambda}(\boldsymbol{\theta}; \rho_{*})\})^{1/2}]\|_{2}^{2} \cdot \exp\{-\beta \Psi_{\lambda}(\boldsymbol{\theta}; \rho_{*})\} d\boldsymbol{\theta}.$$
(10.65)

Define

$$\mu_*(\mathrm{d}\boldsymbol{\theta}) = 1/Z_* \cdot \exp\{-\beta \Psi_{\lambda}(\boldsymbol{\theta}; \rho_*)\} \mu_0(\mathrm{d}\boldsymbol{\theta}), \quad Z_* = \int_{\mathbb{R}^d} \exp\{-\beta \Psi_{\lambda}(\boldsymbol{\theta}; \rho_*)\} \mu_0(\mathrm{d}\boldsymbol{\theta}), \quad (10.66)$$

 $f_k(\boldsymbol{\theta}) = [\exp(\beta \Psi_{\lambda}(\boldsymbol{\theta}; \rho_*)) \rho_{t_k}(\boldsymbol{\theta})]^{1/2} \in \mathcal{D} \equiv \{ f \in L^2(\mathbb{R}^d, \mu_*) \cap C^1(\mathbb{R}^d) : \|\nabla f\|_2 \in L^2(\mathbb{R}^d, \mu_*) \}$   $(f_k \in C^1(\mathbb{R}^d) \text{ because } \rho_t(\boldsymbol{\theta}) > 0 \text{ for any } \boldsymbol{\theta} \in \mathbb{R}^d \text{ and } \rho_t(\boldsymbol{\theta}) \in C^1(\mathbb{R}^d) \text{ for fixed } t), \text{ and } f_*(\boldsymbol{\theta}) = [\exp(\beta \Psi_{\lambda}(\boldsymbol{\theta}; \rho_*)) \rho_*(\boldsymbol{\theta})]^{1/2} \in L^2(\mathbb{R}^d, \mu_*). \text{ Since we have } \rho_{t_k} \text{ converges to } \rho_* \text{ weakly in } L^1(\mathbb{R}^d, \mu_0), \text{ then } f_k \text{ converges weakly to } f_* \text{ in } L^2(\mathbb{R}^d, \mu_*). \text{ Define } I(f) \equiv \int_{\mathbb{R}^d} \|\nabla f(\boldsymbol{\theta})\|_2^2 \cdot \mu_*(\mathrm{d}\boldsymbol{\theta}). \text{ Eq. } (10.64)$ and (10.65) give  $\lim_{k \to \infty} I(f_k) = 0$ . Now we apply Lemma 10.11 with  $\Phi(\boldsymbol{\theta}) = \beta \Psi_{\lambda}(\boldsymbol{\theta}; \rho_*)$ . This  $\Phi$  satisfies  $\beta \lambda/2 \cdot \|\boldsymbol{\theta}\|_2^2 - 2\beta K_2 \leq \Phi(\boldsymbol{\theta}) \leq \beta \lambda/2 \cdot \|\boldsymbol{\theta}\|_2^2 + 2\beta K_2, \text{ where } K_2 \text{ is the constant in Assumption } A2. \text{ Lemma 10.11 implies } f_*(\boldsymbol{\theta}) \equiv F_* \text{ for some constant } F_*.$ 

This proves that  $\rho_*(\boldsymbol{\theta}) = F_* \cdot \exp\{-\beta \Psi_{\lambda}(\boldsymbol{\theta}; \rho_*)\}$ . Combining with the fact that  $\int_{\mathbb{R}^d} \rho_*(\boldsymbol{\theta}) d\boldsymbol{\theta} = 1$ ,  $\rho_*$  satisfies the Boltzmann fixed point condition. According to Lemma 10.4, the Boltzmann fixed point condition has a unique solution  $\rho_*^{\beta,\lambda}$ . Therefore, all the converging weak limit of subsequence of  $\rho_t$  converges to the same point  $\rho_*^{\beta,\lambda}$ . As a result,  $\rho_t$  converges to  $\rho_*^{\beta,\lambda}$  weakly in  $L^1(\mathbb{R}^d)$ .

### 10.3 Proof of Proposition 3, Theorem 4, and Theorem 5

Proposition 3 is given by Lemma 10.6, 10.4, and Lemma 10.9. Theorem 4 is given by Lemma 10.2, 10.4, 10.5, and 10.12.

Now we prove Theorem 5. First, according to Lemma 10.5, for any  $\eta > 0$ , there exists constant K depending on  $\eta, K_0, K_1, K_2, K_3$ , such that as we take  $\beta \geq KD$ , we have

$$R(\rho_*^{\beta,\lambda}) \le \inf_{\rho \in \mathscr{P}(\mathbb{R}^D)} R_{\lambda}(\rho) + \eta/3.$$
 (10.67)

According to Lemma 10.12, we have  $\rho_t$  converges to  $\rho_*^{\beta,\lambda}$  weakly. Therefore, there exists  $T = T(\eta, V, U, \{K_i\}, D, \lambda, \beta) < \infty$ , so that  $d_{\text{BL}}(\rho_t, \rho_*^{\beta,\lambda}) \leq \eta/(3Z)$  for any  $t \geq T$ , where  $Z = Z(\{K_i\})$  is the bounded-Lipschitz constant of R with respect to  $\rho$ . Hence, we have

$$R(\rho_t) \le R(\rho_*^{\beta,\lambda}) + \eta/3 \tag{10.68}$$

for any  $t \geq T$ .

Finally, according to Theorem 3, there exists K' depending on  $K_i$ 's, so that for all  $k \leq 10T/\varepsilon$ , we have

$$|R_N(\boldsymbol{\theta}^k) - R_{\rho_{k\varepsilon}}| \leq K' e^{K'T} \sqrt{1/N \vee \varepsilon} \cdot \left[ \sqrt{D + \log(N(1/\varepsilon \vee 1))} + z \right],$$

with probability at least  $1 - e^{-z^2}$ . Hence there exists  $C_0 = C_0(\eta, \{K_i\}, \delta)$ , so that as  $N, 1/\varepsilon \ge C_0 \exp\{C_0T\}D$  and  $\varepsilon \ge 1/N^{10}$ , we have

$$|R_N(\boldsymbol{\theta}^k) - R(\rho_{k\varepsilon})| \le \eta/3,\tag{10.69}$$

with probability at least  $1 - \delta$ .

Combining Eq. (10.67), (10.68), and (10.69) we get the desired result.

## 10.4 Dependence of convergence time on D and $\eta$

Theorem 5 does not provide any estimate for the dependence of the convergence time on the problem dimensions D and on the accuracy  $\eta$ . However the proof suggests the following heuristic. When  $\rho_t$  is sufficiently close to the minimizer  $\rho_*$ , we heuristically can approximate the free energy dissipation formula (10.2) as

$$\partial_t F_{\beta,\lambda}(\rho_t) \approx -\int_{\mathbb{R}^d} \|\nabla_{\boldsymbol{\theta}}(\Psi_{\lambda}(\boldsymbol{\theta}; \rho_*) + 1/\beta \cdot \log \rho_t(\boldsymbol{\theta}))\|_2^2 \rho_t(\boldsymbol{\theta}) d\boldsymbol{\theta}.$$
 (10.70)

This is the same as the free energy dissipation for the Fokker-Planck equation with potential  $\Psi_{\lambda}(\boldsymbol{\theta}; \rho_*)$ . This suggests that, close to  $\rho_*$ , convergence should be dominated by the speed of convergence in this Fokker-Plank equation, which is controlled by the log-Sobolev constant of the potential  $\Psi_{\lambda}(\boldsymbol{\theta}; \rho_*)$ , to be denote by  $c_*$  [MV00]:

$$F_{\beta,\lambda}(\rho_t) \lesssim F_{\beta,\lambda}(\rho_{t_0}) e^{-c_*(t-t_0)}. \tag{10.71}$$

Note that the log-Sobolev constant can be exponentially small in D. We expect this heuristic to capture the rough dependence of the convergence time T on  $\eta$  and D, hence suggesting  $T = e^{O(D)} \log(1/\eta)$ .

# 11 Numerical Experiments

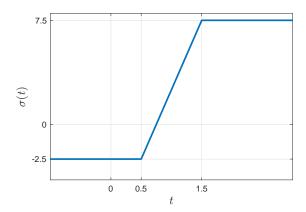
In this section, we discuss numerical experiments whose results were presented in the main text, as well as some additional ones. Some technical details of the figures in the main text are also presented here; in particular, Section 11.1.1 for Figure 1, Section 11.1.2 for Figure 2, Section 11.2 for Figure 3, and Section 11.3 for Figure 4.

### 11.1 Isotropic Gaussians

In this section, we present details of the numerical experiments pertaining to the example of centered isotropic Gaussians:

With probability 1/2: y = +1,  $\boldsymbol{x} \sim N(0, (1 + \Delta)^2 \boldsymbol{I}_d)$ .

With probability 1/2: y = -1,  $\boldsymbol{x} \sim \mathsf{N}(0, (1+\Delta)^2 \boldsymbol{I}_d)$ .



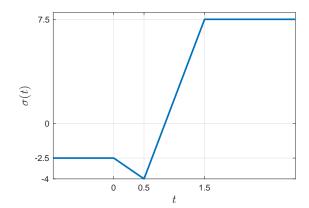


Figure 11.1: The activation functions  $\sigma(t)$  used in Section 11.1 (left plot) and Section 11.3 (right plot).

In all numerical examples in this section, we use the activation  $\sigma_*(\boldsymbol{x};\boldsymbol{\theta}_i) = \sigma(\langle \boldsymbol{w}_i, \boldsymbol{x} \rangle)$ , where  $\sigma(t) = s_1$  if  $t \leq t_1$ ,  $\sigma(t) = s_2$  if  $t \geq t_2$ , and  $\sigma(t)$  interpolated linearly for  $t \in (t_1, t_2)$ . In simulations we use  $t_1 = 0.5$ ,  $t_2 = 1.5$ ,  $s_1 = -2.5$ ,  $s_2 = 7.5$ . This is also used for examples with centered Gaussians in the main text, cf. Figures 1 and 2, and Section 8 in the supplemental information. This activation is plotted in Figure 11.1.

### 11.1.1 Empirical validation of distributional dynamics

Here we discuss empirical validation for the dynamics in the isotropic Gaussian example.

**PDE simulation.** Simulating the PDE (Eq. [13] of the main text) for general d is computationally intensive. In order to simplify the problem, we only consider  $d = \infty$ . In that case, we recall that the risk is given by Eq. (8.10), which we copy here for ease of reference:

$$\overline{R}_{\infty}(\overline{\rho}) = \frac{1}{2} \left( 1 - \int q_{+}(r) \,\overline{\rho}(\mathrm{d}r) \right)^{2} + \frac{1}{2} \left( 1 + \int q_{-}(r) \,\overline{\rho}(\mathrm{d}r) \right)^{2} \,, \tag{11.1}$$

where  $q_{\pm}(t) = \mathbb{E}\{\sigma((1 \pm \Delta)tG)\}, G \sim N(0,1)$ . In addition, from Eq. (8.12),

$$\psi_{\infty}(r;\overline{\rho}) = \frac{1}{2} [\langle q_+, \overline{\rho} \rangle - 1] q_+(r) + \frac{1}{2} [\langle q_-, \overline{\rho} \rangle + 1] q_-(r). \tag{11.2}$$

The PDE is then  $\partial_t \overline{\rho}_t = 2\xi(t)\partial_r [\overline{\rho}_t \partial_r \psi_\infty(r; \overline{\rho}_t)].$ 

The solution to the PDE is approximated, at all time t, by the following multiple-deltas ansatz:

$$\overline{\rho}_t = \frac{1}{J} \sum_{i=1}^J \delta_{r_i(t)} \,, \tag{11.3}$$

where  $J \in \mathbb{N}$  is a pre-chosen parameter. Note that for a fixed J, if the PDE is initialized at  $\overline{\rho}_0$  taking the above form, then for any  $t \geq 0$ ,  $\overline{\rho}_t$  remains in the above form. Then for any smooth test

function  $f: \mathbb{R} \to \mathbb{R}$  with compact support,

$$\frac{1}{J} \sum_{i=1}^{J} f'(r_i(t)) r_i'(t) = \partial_t \langle f, \overline{\rho}_t \rangle = -2\xi(t) \langle f', \overline{\rho}_t \partial_r \psi_\infty(r; \overline{\rho}_t) \rangle$$
 (11.4)

$$= -2\xi(t)\frac{1}{J}\sum_{i=1}^{J} f'(r_i(t))\partial_r\psi_\infty(r_i(t);\overline{\rho}_t). \tag{11.5}$$

Under this ansatz, let us write  $\overline{R}_{\infty}(\overline{\rho}_t) = \overline{R}_{\infty,J}(\mathbf{r}(t))$ , where  $\mathbf{r}(t) = (r_1(t),...,r_J(t))^{\top}$ , and

$$\overline{R}_{\infty,J}(\mathbf{r}) = \frac{1}{2} \left( 1 - \frac{1}{J} \sum_{i=1}^{J} q_{+}(r_{i}) \right)^{2} + \frac{1}{2} \left( 1 + \frac{1}{J} \sum_{i=1}^{J} q_{-}(r_{i}) \right)^{2}.$$
(11.6)

Notice that  $\partial_r \psi_{\infty}(r_i(t); \overline{\rho}_t) = (J/2)(\nabla \overline{R}_{\infty,J}(\mathbf{r}(t)))_i$ . Therefore we obtain

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathbf{r}(t) = -J\xi(t)\nabla\overline{R}_{\infty,J}(\mathbf{r}(t)). \tag{11.7}$$

Hence under the multiple-deltas ansatz, one can simulate numerically the PDE via the above evolution equation of  $\mathbf{r}(t)$ . In particular, given  $\mathbf{r}(t)$ , one approximates  $\mathbf{r}(t + \delta t)$  for some small displacement  $\delta t$  by

$$\mathbf{r}(t+\delta t) \approx \mathbf{r}(t) - J\xi(t)\nabla \overline{R}_{\infty,J}(\mathbf{r}(t))\delta t.$$
 (11.8)

In general, one would want to take a large J to obtain a more accurate approximation. There are certain cases where one can take small J (even J=1). An example of such case is given in the following.

Details of Figure 1 of the main text. For the data generation, we set  $\Delta = 0.8$ . For the SGD simulation, we take d = 40, N = 800, with  $\varepsilon = 10^{-6}$  and  $\xi(t) = 1$ . The weights are initialized as  $(\boldsymbol{w}_i)_{i \leq N} \sim_{iid} \mathsf{N}(0, 0.8^2/d \cdot \mathbf{I}_d)$ . We take a single SGD run. At iteration  $10^3, 4 \times 10^6, 10^7$ , we plot the histogram of  $(\|\boldsymbol{w}_i\|_2)_{i \leq N}$ . This produces the results of the SGD in Figure 1 of the main text.

To obtain results from the PDE, we take J = 400, and generate  $r_i(0) = ||Z_i||_2$ , where  $(Z_i)_{i \leq J} \sim_{iid} \mathsf{N}(0, 0.8^2/d \cdot \mathbf{I}_d)$ . We obtain  $\mathbf{r}(t)$  from t = 0 until  $t = 10^7 \varepsilon$ , by discretizing this interval with  $10^5$  points equally spaced on the  $\log_{10}$  scale and sequentially computing  $\mathbf{r}(t)$  at each point using Eq. (11.8). Note that the SGD result at iteration k corresponds to  $\mathbf{r}(\varepsilon k)$ . We re-simulate the PDE for 100 times, each with an independently generated initialization. The obtained histogram for the PDE, as shown in the figure, is the aggregation of these 100 runs.

Further numerical simulations. Figure 11.2 plots the evolution of  $\overline{\rho}_t$  for  $\Delta = 0.2$ . The setting is identical to the one in Figure 1 of the main text, described in the previous paragraphs.

In Figure 11.3, we plot the evolution of the population risk for the SGD and its PDE prediction counterpart, for  $\Delta=0.2$  and  $\Delta=0.8$ . The setting for the SGD plots is the same as described in the previous paragraphs. We compute the risk attained by the SGD by Monte Carlo averaging over  $10^4$  samples. The setting for the PDE plots tagged "J=400" is almost the same as in the previous paragraphs, except that we take only 1 run. For the PDE plot tagged "J=1", we take J=1 and r(0)=0.8 instead. In the inset plot, we also show the evolution of  $(1/N)\sum_{i=1}^{N} \|\boldsymbol{w}_i\|_2$  of the SGD, and  $(1/J)\sum_{i=1}^{J} r_i(t)$  of the PDE.

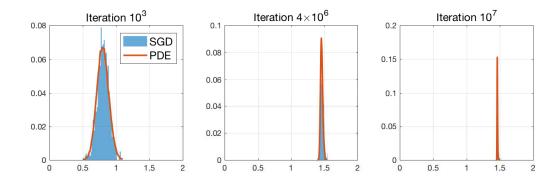


Figure 11.2: Evolution of the reduced distribution  $\overline{\rho}_t$  for  $\Delta = 0.2$ , in the isotropic Gaussians example of Section 11.1.

In Figure 11.4, we plot the function  $\overline{R}_d^{(1)}(r)$ , for d=40 and  $\Delta=0.2$ . (Recall  $\overline{R}_d^{(1)}(r)$  from Eq. [14] of the main text, and see also Section 11.1.3.) On this landscape, we also plot the evolution of the corresponding SGD and PDE, as described in the last paragraph.

Comments. We observe in Figure 11.3 a good match between the SGD and the PDE, even when J=1, for  $\Delta=0.2$ . This can be explained with our theory, which predicts that at  $\Delta=0.2$ , the minimum risk is achieved by the uniform distribution over a sphere of radius  $\|\boldsymbol{w}\|_2 = r_*$  (see also Section 11.1.3). This corresponds to  $\overline{\rho}_t$ , as  $t\to\infty$ , being a delta function and placing probability 1 at  $r_*$ . Furthermore due to the way we initialize the SGD,  $\overline{\rho}_0$  is well concentrated. One can then expect that  $\overline{\rho}_t$  is also well concentrated at all time t, in which case J=1 is sufficient. This claim is reflected in our numerical experiments, shown in Figure 11.2.

We also observe in Figure 11.3 that the case  $\Delta=0.2$  has a rapid transition from a high risk to a lower risk, unlike the case  $\Delta=0.8$ . This is also expected from our theory. As said above,  $\overline{\rho}_t$  is approximately a delta function at all time t, and the position r(t) evolves by gradient flow in the landscape of  $\overline{R}_d^{(1)}(r)$ . This latter claim is well supported by Figure 11.4. As observed in Figure 11.4,  $\overline{R}_d^{(1)}(r)$  is rather benign, and hence the transition of the population risk should be smooth. However the case for  $\Delta=0.8$  is different:  $\overline{\rho}_t$  is not concentrating at large t, as evident in Figure 1 of the main text, even though  $\overline{R}_d^{(1)}(r)$  is generally benign for a vast variety of values of d and  $\Delta$  (see Figure 11.6 and Section 11.1.3).

Note that the computation of the PDE assumes  $d=\infty$ . Furthermore it also requires  $N=\infty$  (recalling Theorem 3 of the main text). The discrepancy to the SGD is due to the fact that d and N are finite in the SGD simulations. Nevertheless in our numerical examples, such discrepancy is insignificant.

#### 11.1.2 Empirical validation of the statics

Here we discuss numerical verification for the statics in the isotropic Gaussian example.

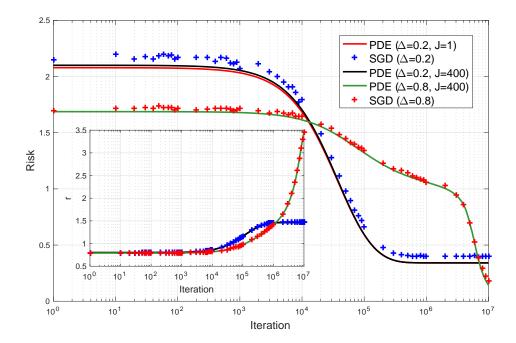


Figure 11.3: The evolution of the population risk and the parameter r of the reduced distribution  $\overline{\rho}_t$ , in the isotropic Gaussians example of Section 11.1.

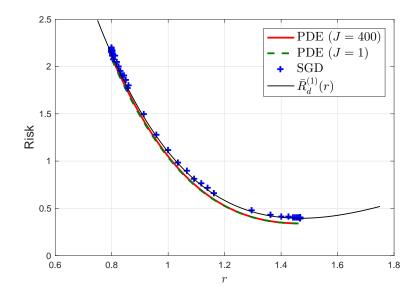


Figure 11.4: The function  $\overline{R}_d^{(1)}(r)$  vs r, as well as the evolution of the SGD and the PDE on this landscape, for  $\Delta=0.2$  and d=40, in the isotropic Gaussians example of Section 11.1. Here the SGD and the PDE evolve from the leftmost point to the rightmost point.

**Optimizing**  $\overline{R}_d(\overline{\rho})$ . For the chosen activation, we have from Eq. (8.8) that

$$\overline{R}_d(\overline{\rho}) = 1 + 2 \int v(r) \,\overline{\rho}(dr) + \int u_d(r_1, r_2) \,\overline{\rho}(dr_1) \,\overline{\rho}(dr_2) \,, \tag{11.9}$$

$$v(r) = -\frac{1}{2}g(0, (1+\Delta)r) + \frac{1}{2}g(0, (1-\Delta)r), \qquad (11.10)$$

$$u_d(r_1, r_2) = \frac{\Gamma(d/2)}{\Gamma(1/2)\Gamma((d-1)/2)} \int_{\theta=0}^{\pi} \hat{u}(r_1, r_2, \theta) \sin^{d-2}\theta d\theta, \qquad (11.11)$$

$$\hat{u}(r_1, r_2, \theta) = \frac{1}{2} f((1 + \Delta)r_1, (1 + \Delta)r_2, \theta) + \frac{1}{2} f((1 - \Delta)r_1, (1 - \Delta)r_2, \theta),$$
(11.12)

$$f(r_1, r_2, \theta) = \int_{x = -\infty}^{+\infty} \sigma(r_1 x) g(r_2 x \cos \theta, r_2 \sin \theta) \phi(x) dx, \qquad (11.13)$$

$$g(a,b) = s_2 + (s_1 - \sigma_{\text{itc}} - \sigma_{\text{sl}}a)\Phi\left(\frac{t_1 - a}{b}\right) + (\sigma_{\text{sl}}a + \sigma_{\text{itc}} - s_2)\Phi\left(\frac{t_2 - a}{b}\right) + \sigma_{\text{sl}}b\left[\phi\left(\frac{t_1 - a}{b}\right) - \phi\left(\frac{t_2 - a}{b}\right)\right].$$

$$(11.14)$$

where  $\sigma_{\rm sl} = (s_2 - s_1)/(t_2 - t_1)$ ,  $\sigma_{\rm itc} = s_1 - \sigma_{\rm sl}t_1$ ,  $\phi(x) = \exp(-x^2/2)/\sqrt{2\pi}$ ,  $\Phi(x) = \int_{-\infty}^x \phi(t) dt$ , and  $\Gamma$  is the Gamma function. To numerically optimize  $\overline{R}_d(\overline{\rho})$ , we perform the following approximation:

$$\inf_{\overline{\rho}} \overline{R}_d(\overline{\rho}) \approx \inf_{p_i \ge 0, \sum_{i=1}^K p_i = 1} \overline{R}_d \left( \sum_{i=1}^K p_i \delta_{o_i} \right). \tag{11.15}$$

Here  $o_i \in \mathbb{R}$ , i = 1, ..., K, are K pre-chosen points. Let  $\mathbf{v} = (v(o_1), ..., v(o_K))^{\top}$  and  $\mathbf{U} = (u_d(o_i, o_j))_{1 \le i,j \le K}$ . Then the approximation becomes

$$\inf_{\overline{\rho}} \overline{R}_d(\overline{\rho}) \approx \inf_{p_i \ge 0, \sum_{i=1}^K p_i = 1} \left\{ 1 + 2\mathbf{v}^\top \mathbf{p} + \mathbf{p}^\top \mathbf{U} \mathbf{p} \right\}, \tag{11.16}$$

which is a quadratic programming problem and can be solved numerically. Here  $\mathbf{v}$  can be computed easily with the explicit formula, and the computation of  $\mathbf{U}$  amounts to numerically evaluating double integrals. In the case  $d = \infty$ , the computation of  $\mathbf{U}$  is much easier, since

$$u_{\infty}(r_1, r_2) = \frac{1}{2}g(0, (1+\Delta)r_1)g(0, (1+\Delta)r_2) + \frac{1}{2}g(0, (1-\Delta)r_1)g(0, (1-\Delta)r_2).$$
 (11.17)

**Details of Figure 2 of the main text.** For the SGD simulation, we take N=800, with  $\varepsilon=3\times 10^{-3}$  and  $\xi(t)=t^{-1/4}$ . The weights are initialized as  $(\boldsymbol{w}_i)_{i\leq N}\sim_{iid} \mathsf{N}(0,0.4^2/d\cdot\mathbf{I}_d)$ . We compute the risk attained by the SGD by Monte Carlo averaging over  $10^4$  samples. We take a single SGD run per  $\Delta$ , per d, and report the risk at iteration  $10^7$ .

For the approximate optimization of  $\overline{R}_d(\overline{\rho})$ , we choose K = 100, and  $o_i$ , i = 1, ..., K, being equally spaced on the interval [0.01, 10].

For the optimization of  $\overline{R}_d^{(1)}(r)$  (recalling Eq. [14] in the main text), we approximate it with  $\min_{i=1,\ldots,K} \overline{R}_d^{(1)}(o_i)$ , for the above chosen  $o_i$  and K.

We find that in general, one needs higher  $\max_{i=1,...,K} o_i$  to produce accurate results for higher  $\Delta$ . For the chosen set of  $o_i$ 's, we choose to plot up until  $\Delta = 0.8$ .

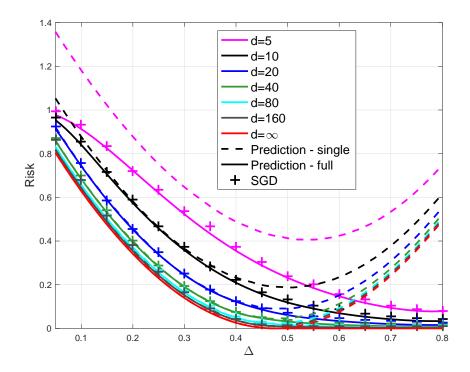


Figure 11.5: The population risk as a function of  $\Delta$  for different values of d, in the isotropic Gaussians example of Section 11.1. Here "Prediction - single" refers to  $\min_{r\geq 0} \overline{R}_d^{(1)}(r)$ , "Prediction - full" refers to the optimized  $R(\rho)$  as described in Section 11.1.2, and "SGD" refers to the risk attained by the SGD.

Further numerical simulations. In Figure 11.5, we extend Figure 2 of the main text to include results for additional values of d. The setting remains the same.

This figure provides further support to the respective discussion in the main text. For the threshold values of  $\Delta$  for which the minimum risk is achieved by a uniform distribution  $\rho_{r_*}^{\text{unif}}$  over a sphere of radius  $\|\boldsymbol{w}\|_2 = r_*$  (see the main text around Eq. [14], and Section 11.1.3).

#### 11.1.3 Checking the condition of Lemma 1 in the main text

We check of the condition of Lemma 1 in the main text. This has two steps: (1) we solve for the minimizer  $r_*$  of  $\overline{R}_d^{(1)}(r) = 1 + 2v(r) + u_d(r,r)$ , where v(r) and  $u_d(r_1,r_2)$  are given by Eq. (11.10) and (11.11) respectively, and (2) we check whether  $v(r) + u_d(r,r_*) \geq v(r_*) + u_d(r_*,r_*)$  for all  $r \geq 0$ . Figure 11.6 suggests that the behavior of  $\overline{R}_d^{(1)}(r)$  is rather benign and hence  $r_*$  can be solved easily by searching for a local minimum. For the second step, we check the condition on a grid of values of r from 0.1 to 10 with a spacing of 0.1, for each value of  $\Delta$  on a grid from 0.01 to 0.99 with a spacing of 0.01. In general, we find that the conditioned is satisfied for  $\Delta \in [\Delta_d^1, \Delta_d^h]$ . Table 1 reports  $\Delta_d^1$  and  $\Delta_d^h$  for a number of values of d for the isotropic Gaussians example with the given activation function.

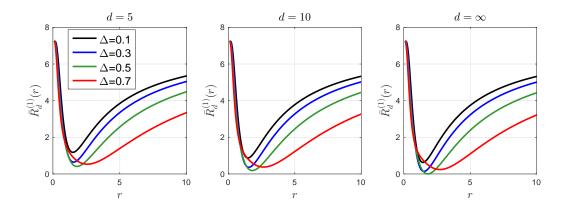


Figure 11.6: The function  $\overline{R}_d^{(1)}(r)$  for different values of d and  $\Delta$ , in the isotropic Gaussians example of Section 11.1.

d	$\Delta_d^{\scriptscriptstyle  ext{I}}$	$\Delta_d^{ ext{n}}$
5	N/A	N/A
10	N/A	N/A
20	0.08	0.38
40	0.03	0.42
80	0.02	0.45
160	0.0	0.46
$\infty$	0	0.47

Table 1:  $\Delta_d^l$  and  $\Delta_d^h$  for different values of d, in the isotropic Gaussians example of Section 11.1. Here "N/A" refers to that no values of  $\Delta$  are found to satisfy the condition of Lemma 1 in the main text. Note that for  $d=\infty$ , the value  $\Delta_\infty^l=0$  is exact, according to Theorem 8.

## 11.2 Centered anisotropic Gaussians with ReLU Activation

In this section, we present details of the numerical experiments pertaining to the example of anisotropic Gaussians with ReLU activation. In particular, we use the activation  $\sigma_*(\boldsymbol{x};\boldsymbol{\theta}) = a \max(\langle \boldsymbol{w}, \boldsymbol{x} \rangle + b, 0)$ , with  $\boldsymbol{\theta} = (\boldsymbol{w}, a, b) \in \mathbb{R}^{d+2}$ . We consider the centered anisotropic Gaussian case:

With probability 1/2: y = +1,  $\boldsymbol{x} \sim N(0, \boldsymbol{\Sigma}_{+})$ .

With probability 1/2: y = -1,  $\boldsymbol{x} \sim N(0, \boldsymbol{\Sigma}_{-})$ .

More specifically, we opt for

$$\Sigma_{+} = \operatorname{Diag}\left(\underbrace{(1+\Delta)^{2}, \dots, (1+\Delta)^{2}}_{s_{0}}, \underbrace{1, \dots, 1}_{d-s_{0}}\right), \tag{11.18}$$

$$\Sigma_{-} = \operatorname{Diag}\left(\underbrace{(1-\Delta)^{2}, \dots, (1-\Delta)^{2}}_{s_{0}}, \underbrace{1, \dots, 1}_{d-s_{0}}\right). \tag{11.19}$$

This setting is used in Figure 3 in the main text.

We consider  $s_0 = \gamma d$  for some  $\gamma \in (0,1)$ . For simplicity, we consider the limit  $d \to \infty$ . For  $\boldsymbol{\theta} \sim \rho$ , let  $\overline{\rho}$  be the joint distribution of four parameters  $\mathbf{r} = (a, b, r_1 = \|\boldsymbol{w}_{1:s_0}\|_2, r_2 = \|\boldsymbol{w}_{(s_0+1):d}\|_2)$ , where  $\boldsymbol{w}_{i:j} = (w_i, ..., w_j)^{\top}$ . Using a similar argument to Section 8, we have, in the limit  $d \to \infty$ , the risk  $R(\rho) = \overline{R}_{\infty}(\overline{\rho})$  for

$$\overline{R}_{\infty}(\overline{\rho}) = \frac{1}{2} \left( 1 - \int aq_{+}(r_{1}, r_{2}, b)\overline{\rho}(d\mathbf{r}) \right)^{2} + \frac{1}{2} \left( 1 + \int aq_{-}(r_{1}, r_{2}, b)\overline{\rho}(d\mathbf{r}) \right)^{2}, \tag{11.20}$$

$$q_{\pm}(r_1, r_2, b) = b \Phi\left(\frac{b}{\sqrt{(1 \pm \Delta)^2 r_1^2 + r_2^2}}\right) + \sqrt{(1 \pm \Delta)^2 r_1^2 + r_2^2} \phi\left(\frac{b}{\sqrt{(1 \pm \Delta)^2 r_1^2 + r_2^2}}\right), \quad (11.21)$$

where  $\phi(x) = \exp(-x^2/2)/\sqrt{2\pi}$  and  $\Phi(x) = \int_{-\infty}^{x} \phi(t) dt$ . Furthermore, letting  $\overline{\rho}_t$  denote the corresponding distribution at time t, the PDE [7] in the main text can be reduced to the following PDE of  $\overline{\rho}_t$ :

$$\partial_{t}\overline{\rho}_{t} = 2\xi(t)\nabla_{\mathbf{r}}\cdot(\overline{\rho}_{t}\nabla_{\mathbf{r}}\psi_{\infty}(\mathbf{r};\overline{\rho}_{t})), \qquad (11.22)$$

$$\psi_{\infty}(\mathbf{r};\overline{\rho}) = \frac{1}{2}\left[\int a'q_{+}(r'_{1},r'_{2},b')d\overline{\rho}(a',b',r'_{1},r'_{2}) - 1\right]aq_{+}(r_{1},r_{2},b)$$

$$+ \frac{1}{2}\left[\int a'q_{-}(r'_{1},r'_{2},b')d\overline{\rho}(a',b',r'_{1},r'_{2}) + 1\right]aq_{-}(r_{1},r_{2},b). \qquad (11.23)$$

**PDE simulation.** As in Section 11.1.1, we posit that the solution to the PDE can be approximated, at all time t, by the multiple-deltas ansatz:

$$\overline{\rho}_t = \frac{1}{J} \sum_{i=1}^{J} \delta_{\mathbf{r}_i(t)} \,, \tag{11.24}$$

where  $J \in \mathbb{N}$  is a pre-chosen parameter, and  $\mathbf{r}_i(t) = (a_i(t), b_i(t), r_{1,i}(t), r_{2,i}(t))$ . Following the same argument as in Section 11.1.1, we obtain the following evolution equation:

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathbf{r}_{i}(t) = -J\xi(t)\nabla_{i}\overline{R}_{\infty,J}(\mathbf{r}_{1}(t),...,\mathbf{r}_{J}(t)), \tag{11.25}$$

for i=1,...,J, where  $\overline{R}_{\infty,J}(\mathbf{r}_1(t),...,\mathbf{r}_J(t))=\overline{R}_{\infty}(\overline{\rho}_t)$  under the ansatz, and  $\nabla_i$  denotes the gradient of  $\overline{R}_{\infty,J}(\mathbf{r}_1,...,\mathbf{r}_J)$  w.r.t.  $\mathbf{r}_i$ . More explicitly,

$$\overline{R}_{\infty,J}(\mathbf{r}_1,\dots,\mathbf{r}_J) = \frac{1}{2} \left( 1 - \frac{1}{J} \sum_{i=1}^J a_i q_+(r_{1,i}, r_{2,i}, b_i) \right)^2 + \frac{1}{2} \left( 1 + \frac{1}{J} \sum_{i=1}^J a_i q_-(r_{1,i}, r_{2,i}, b_i) \right)^2.$$
(11.26)

Again, given  $\mathbf{r}_i(t)$ , one approximates  $\mathbf{r}_i(t+\delta t)$  for some small displacement  $\delta t$  by

$$\mathbf{r}_i(t+\delta t) \approx \mathbf{r}_i(t) - J\xi(t)\nabla_i \overline{R}_{\infty,J}(\mathbf{r}_1,\dots,\mathbf{r}_J)\delta t.$$
 (11.27)

**Details of Figure 3 of the main text.** For the SGD simulation, we take d = 320,  $s_0 = 60$ , N = 800, with  $\varepsilon = 2 \times 10^{-4}$  and  $\xi(t) = t^{-1/4}$ . The weights are initialized as  $(\mathbf{w}_i)_{i \leq N} \sim_{iid} \mathsf{N}(0, 0.8^2/d \cdot \mathbf{I}_d)$ ,  $(a_i)_{i \leq N} = 1$  and  $(b_i)_{i \leq N} = 1$ . We take a single SGD run. We compute the risk attained by the SGD by Monte Carlo averaging over  $10^4$  samples.

To obtain results from the PDE, we take J=400. We initialize  $r_{1,i}(0)=\|Z_{1,i}\|_2$  and  $r_{2,i}(0)=\|Z_{2,i}\|_2$ , where  $(Z_{1,i})_{i\leq N}\sim_{iid}\mathsf{N}(0,0.8^2/d\cdot\mathbf{I}_{s_0})$  and  $(Z_{2,i})_{i\leq N}\sim_{iid}\mathsf{N}(0,0.8^2/d\cdot\mathbf{I}_{d-s_0})$  independently, along with  $a_i(0)=1$ ,  $b_i(0)=1$ . We obtain  $\mathbf{r}_i(t)$  from t=0 until  $t=10^7\varepsilon$ , by discretizing this interval with  $10^5$  points equally spaced on the  $\log_{10}$  scale and sequentially computing  $\mathbf{r}_i(t)$  at each point using Eq. (11.27). Note that the SGD result at iteration  $\ell$  corresponds to  $\mathbf{r}_i(\varepsilon^{4/3}\ell)$ . We take a single run of the PDE.

To produce the inset plot in Figure 3 of the main text, for the "a (mean)" axis, we compute  $\frac{1}{N}\sum_{i=1}^{N}a_{i}$  for the SGD and  $\frac{1}{J}\sum_{i=1}^{J}a_{i}(t)$  for the PDE. Similarly, for the "b (mean)" axis, we compute  $\frac{1}{N}\sum_{i=1}^{N}b_{i}$  for the SGD and  $\frac{1}{J}\sum_{i=1}^{J}b_{i}(t)$  for the PDE, and for the " $r_{1}$  (mean)" axis, we compute  $\frac{1}{N}\sum_{i=1}^{N}\|\boldsymbol{w}_{i,1:s_{0}}\|_{2}$  for the SGD and  $\frac{1}{J}\sum_{i=1}^{J}r_{1,i}(t)$  for the PDE.

Further numerical simulations. In Figure 11.7, we plot the evolution of the four parameters, for the same setting as Figure 3 of the main text. Here "a (mean)", "b (mean)" and " $r_1$  (mean)" hold the same meanings, and " $r_2$  (mean)" refers to  $\frac{1}{N} \sum_{i=1}^{N} \|\boldsymbol{w}_{i,(s_0+1):d}\|_2$  for the SGD and  $\frac{1}{J} \sum_{i=1}^{J} r_{2,i}(t)$  for the PDE.

In Figure 11.8, we plot the population risk's evolution for the same setting as Figure 3 of the main text, apart from that  $\Delta = 0.6$  and  $s_0$  varies.

**Comments.** We observe a good match between the SGD and the PDE in Figure 3 of the main text as well as Figure 11.7, up until iteration  $10^6$ . In general there is less discrepancy with larger  $s_0$ , d and N, recalling that the PDE is computed assuming infinite  $s_0$ , d and N. This is evident from Figure 11.8.

As a note, in Figure 11.8, the PDE evolves differently for different  $s_0$ . This is because the ratio  $s_0/d$  is used to determine the initialization of the PDE.

## 11.3 Isotropic Gaussians: Predictable Failure of SGD

In this section, we consider the isotropic Gaussians example (see Section 11.1 for the setting and notations), with the following activation function:  $\sigma_*(\boldsymbol{x};\boldsymbol{\theta}_i) = \sigma(\langle \boldsymbol{w}_i, \boldsymbol{x} \rangle)$ , where  $\sigma(t) = -2.5$  for  $t \leq 0$ ,  $\sigma(t) = 7.5$  for  $t \geq 1.5$ , and  $\sigma(t)$  linearly interpolates from the knot (0, -2.5) to (0.5, -4), and from (0.5, -4) to (1.5, 7.5). This activation is plotted in Figure 11.1. This corresponds to Section "Predicting failure" and Figure 4 in the main text. The simulation of the PDE can be done in the same way as in Section 11.1.1.

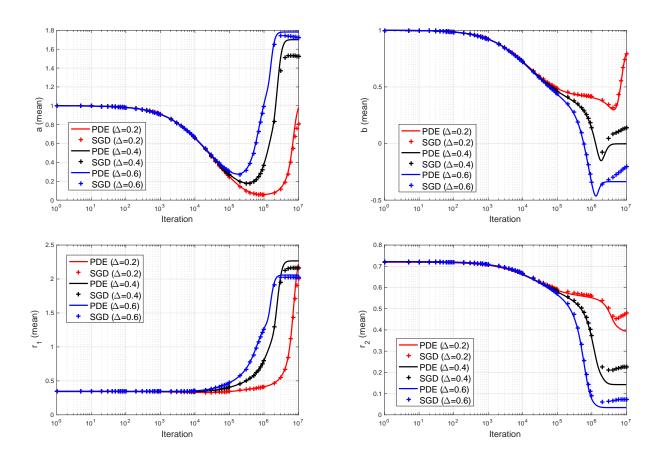


Figure 11.7: The evolution of the four parameters in the anisotropic Gaussians example of Section 11.2.

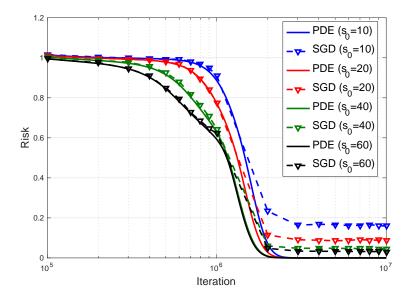


Figure 11.8: The evolution of the population risk for  $\Delta = 0.6$ , d = 320, N = 800 in the anisotropic Gaussians example of Section 11.2.

Rationale of the activation choice. We give an explanation for the choice of the above activation based on our theory. We aim to find an activation  $\sigma_*(x; \theta_i) = \sigma(\langle w_i, x \rangle)$  in which there exists a local minimum that does not generalize well. To simplify the task, we wish for such minimum to be attained at  $\rho_* = \delta_0$ . This minimum does not generalize well, since it implies all the weights are zero and the neuron outputs are constant, rendering the network unable to perform classification. Theorem 6 of the main text suggests taking  $\sigma(t)$  such that

$$\nabla^2 V(\mathbf{0}) + \nabla^2_{1,1} U(\mathbf{0}, \mathbf{0}) \succ 0.$$
 (11.28)

In the isotropic Gaussians case, this becomes

$$\sigma''(0)\left\{ (1-\Delta)^2 - (1+\Delta)^2 + \sigma(0)[(1-\Delta)^2 + (1+\Delta)^2] \right\} > 0.$$
 (11.29)

(Note that the condition  $\nabla V(\mathbf{0}) + \nabla_1 U(\mathbf{0}, \mathbf{0}) = \mathbf{0}$  in Theorem 6 of the main text is trivially satisfied.) Another requirement is that there should still be a minimum whose risk is nearly zero. Hence we do not wish for a dramatic change in the choice of the activation function, as compared to the one used in Section 11.1. That is, we leave  $\sigma(0) < 0$  unchanged. Hence we would want  $\sigma''(0) < 0$ , which is accomplished by our aforementioned choice.

Note that Theorem 6 of the main text also suggests that if the SGD is initialized sufficiently close to this local minimum, the SGD trajectory should converge to it.

**Details of Figure 4 of the main text.** For the data generation, we set  $\Delta = 0.5$ . For the SGD simulation, we take d = 320, N = 800, with  $\varepsilon = 10^{-5}$  and  $\xi(t) = t^{-1/4}$ . We take a single SGD run each for two different initializations: the weights are initialized as  $(\mathbf{w}_i)_{i \leq N} \sim_{iid} \mathsf{N}(0, \kappa^2/d \cdot \mathbf{I}_d)$  for either  $\kappa = 0.1$  or  $\kappa = 0.4$ . We compute the risk attained by the SGD by Monte Carlo averaging over  $10^4$  samples.

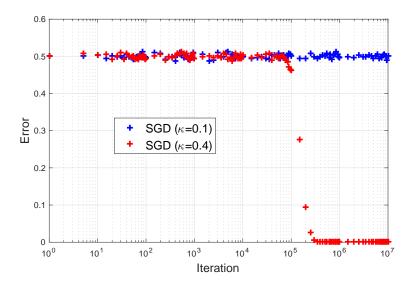


Figure 11.9: The error rate attained by the SGD in the example of Figure 4 of the main text.

To obtain results from the PDE, we take J = 400, and generate  $r_i(0) = ||Z_i||_2$ , where  $(Z_i)_{i \leq N} \sim_{iid} \mathsf{N}(0, \kappa^2/d \cdot \mathbf{I}_d)$ . We obtain  $\boldsymbol{r}(t)$  from t = 0 until  $t = 10^7 \varepsilon$ , by discretizing this interval with  $10^5$  points equally spaced on the  $\log_{10}$  scale and sequentially computing  $\mathbf{r}(t)$  at each point using Eq. (11.8). Note that the SGD result at iteration k corresponds to  $\mathbf{r}(\varepsilon^{4/3}k)$ . We take a single run of the PDE.

To produce the inset plot, we compute  $\frac{1}{N} \sum_{i=1}^{N} \|\boldsymbol{w}_i\|_2$  for the SGD, and  $\frac{1}{J} \sum_{i=1}^{J} r_i(t)$  for the PDE.

As observed from Figure 4 of the main text, the SGD trajectory with  $\kappa = 0.1$  converges to a point where  $\|\boldsymbol{w}_i\|_2$  is nearly zero and the risk is very high, in stark contrast to the SGD trajectory with  $\kappa = 0.4$ , as expected.

**Error plot.** In Figure 11.9, we plot the empirical error rate attained by the SGD in the above example for the two initializations. Here the error rate is defined as the misclassification probability  $\mathbb{P}\{\operatorname{sign}(\hat{y}(\boldsymbol{x};\boldsymbol{\theta})) \neq y\}$ , and is computed with Monte Carlo averaging over  $10^4$  samples. This validates the claim that, in this example, there exists a local minimum which the SGD can converge to, yet has bad generalization (i.e. attains the trivial misclassification rate of 0.5), whereas there is a global minimum which the SGD can also find and yet generalizes well.

# A Concentration inequalities

**Lemma A.1** (Azuma-Höeffding bound). Let  $(\mathbf{X}_k)_{k\geq 0}$ , be a martingale taking values in  $\mathbb{R}^d$  with respect to the filtration  $(\mathcal{F}_k)_{k\geq 0}$ , with  $\mathbf{X}_0 = \mathbf{0}$ . Assume that the following holds almost surely for all  $k\geq 1$ :

$$\mathbb{E}\left\{e^{\langle \boldsymbol{\lambda}, \boldsymbol{X}_k - \boldsymbol{X}_{k-1} \rangle} \middle| \mathcal{F}_{k-1}\right\} \le e^{L^2 \|\boldsymbol{\lambda}\|^2 / 2}. \tag{A.1}$$

Then we have

$$\mathbb{P}\left(\max_{k \le n} \|\boldsymbol{X}_k\|_2 \ge 2L\sqrt{n}(\sqrt{d} + t)\right) \le e^{-t^2}. \tag{A.2}$$

*Proof.* Let  $\mathbf{Z}_k = \mathbf{X}_k - \mathbf{X}_{k-1}$  be the martingale differences. By the subgaussian condition (A.1), we get

$$\mathbb{E}\left\{e^{\langle \boldsymbol{\lambda}, \boldsymbol{X}_n \rangle}\right\} \leq \mathbb{E}\left\{\mathbb{E}\left\{e^{\langle \boldsymbol{\lambda}, \boldsymbol{Z}_n \rangle} | \mathcal{F}_{n-1}\right\} e^{\langle \boldsymbol{\lambda}, \boldsymbol{X}_{n-1} \rangle}\right\}$$
(A.3)

$$\leq e^{L^2 \|\boldsymbol{\lambda}\|^2 / 2} \mathbb{E}\left\{ e^{\langle \boldsymbol{\lambda}, \boldsymbol{X}_{n-1} \rangle} \right\} \leq e^{nL^2 \|\boldsymbol{\lambda}\|_2^2 / 2}. \tag{A.4}$$

Letting  $G \sim N(0, I_d)$  a standard Gaussian vector and  $\xi \geq 0$ ,

$$\mathbb{E}\left\{e^{\xi \|\boldsymbol{X}_n\|_2^2/2}\right\} = \mathbb{E}_{\boldsymbol{G}} \mathbb{E}\left\{e^{\sqrt{\xi}\langle \boldsymbol{G}, \boldsymbol{X}_n \rangle}\right\} \le \mathbb{E}_{\boldsymbol{G}} e^{nL^2 \xi \|\boldsymbol{G}\|_2^2/2} \tag{A.5}$$

$$= \left(1 - nL^2\xi\right)^{-d/2}.\tag{A.6}$$

By Markov inequality, setting  $\xi = 1/(2nL^2)$ , we get, for all  $t \ge 0$ ,

$$\mathbb{P}\Big(\|X_n\|_2 \ge 2L\sqrt{n}(\sqrt{d}+t)\Big) \le e^{d/2 - (\sqrt{d}+t)^2} \le e^{-t^2}.$$
(A.7)

Finally, to upper bound  $\max_{k \leq n} \|\boldsymbol{X}_k\|_2$ , we define the stopping time  $\tau \equiv \min\{k : \|\boldsymbol{X}_k\|_2 \geq 2L\sqrt{n}(\sqrt{d}+t)\}$ , and the martingale  $\overline{\boldsymbol{X}}_k = \boldsymbol{X}_{k \wedge \tau}$ . Since  $\{\max_{k \leq n} \|\boldsymbol{X}_n\|_2 \geq 2L\sqrt{n}(\sqrt{d}+t)\} = \{\|\overline{\boldsymbol{X}}_n\|_2 \geq 2L\sqrt{n}(\sqrt{d}+t)\}$ , the claim follows by applying the previous inequality to  $\overline{\boldsymbol{X}}_n$ .

# B On the generalization to other loss functions

The objective of this section is to show that the framework of this paper can be formally extended to other loss functions  $\ell : \mathbb{R} \times \mathbb{R} \to \mathbb{R}$ . All arguments here will be heuristic, and we defer a rigorous study of this problem to future work.

First of all, we note that the population risk reads

$$R_N(\boldsymbol{\theta}) = \mathbb{E}\left\{\ell\left(y, \frac{1}{N} \sum_{i=1}^N \sigma_*(\boldsymbol{x}; \boldsymbol{\theta}_i)\right)\right\},$$
(B.1)

which naturally leads to the following mean field risk  $R: \mathscr{P}(\mathbb{R}^D) \to \mathbb{R}$ :

$$R(\rho) = \mathbb{E}\left\{\ell\left(y, \int \sigma_*(\boldsymbol{x}; \boldsymbol{\theta}) \, \rho(\mathrm{d}\boldsymbol{\theta})\right)\right\}. \tag{B.2}$$

The corresponding distributional dynamics is formally identical to the one for quadratic loss, cf. Eq. (7.1). The only change is in the definition of  $\Psi(\theta; \rho)$ :

$$\partial_t \rho_t(\boldsymbol{\theta}) = 2\xi(t) \nabla \cdot \left[ \rho_t(\boldsymbol{\theta}) \nabla \Psi(\boldsymbol{\theta}; \rho_t) \right], \tag{B.3}$$

$$\Psi(\boldsymbol{\theta}; \rho) = \frac{\delta R(\rho)}{\delta \rho(\boldsymbol{\theta})} = \mathbb{E} \left\{ \partial_2 \ell \left( y, \int \sigma_*(\boldsymbol{x}; \bar{\boldsymbol{\theta}}) \rho(\mathrm{d}\bar{\boldsymbol{\theta}}) \right) \sigma_*(\boldsymbol{x}; \boldsymbol{\theta}) \right\},$$
(B.4)

where  $\partial_2 \ell$  denotes the derivative of  $\ell$  with respect to its second argument. It is immediate to see that, for the quadratic loss  $\ell(y, \hat{y}) = (y - \hat{y})^2$ , we recover the expressions used in the rest of the paper.