# Dynamical mean-field analysis of adaptive Langevin diffusions: Replica-symmetric fixed point and empirical Bayes

Zhou Fan,[*] Justin Ko,[†] Bruno Loureiro,[‡] Yue M. Lu,[§] Yandi Shen[¶]

## Abstract

In many applications of statistical estimation via sampling, one may wish to sample from a high-dimensional target distribution that is adaptively evolving to the samples already seen. We study an example of such dynamics, given by a Langevin diffusion for posterior sampling in a Bayesian linear regression model with i.i.d. regression design, whose prior continuously adapts to the Langevin trajectory via a maximum marginal-likelihood scheme. Results of dynamical mean-field theory (DMFT) developed in our companion paper establish a precise high-dimensional asymptotic limit for the joint evolution of the prior parameter and law of the Langevin sample. In this work, we carry out an analysis of the equations that describe this DMFT limit, under conditions of approximate time-translation-invariance which include, in particular, settings where the posterior law satisfies a log-Sobolev inequality. In such settings, we show that this adaptive Langevin trajectory converges on a dimension-independent time horizon to an equilibrium state that is characterized by a system of scalar fixed-point equations, and the associated prior parameter converges to a critical point of a replica-symmetric limit for the model free energy. As a by-product of our analyses, we obtain a new dynamical proof that this replica-symmetric limit for the free energy is exact, in models having a possibly misspecified prior and where a log-Sobolev inequality holds for the posterior law.

# Contents

[*]Department of Statistics and Data Science, Yale University
[†]Department of Statistics and Actuarial Science, University of Waterloo
[‡]Departement d'Informatique, École Normale Supérieure, PSL & CNRS
[§]Departments of Electrical Engineering and Applied Mathematics, Harvard University
[¶]Department of Statistics and Data Science, Carnegie Mellon University

# 1 Introduction

Parameter estimation via Monte Carlo sampling is a common paradigm in statistical learning, arising for example in stochastic implementations of Expectation-Maximization estimation in latent variable models [1,2], and contrastive-divergence [3] and diffusion-based learning [4–7] of generative models for data. In these applications, one wishes to learn a parameter using Monte Carlo samples from an associated distribution on a high-dimensional space. Monte Carlo methods whose target distribution continuously adapts to the learned parameter are natural for such tasks, and we refer to [8–12] for several recent proposals of this form.

The goal of our current work is to study the learning dynamics in a particular (classical) instance of this paradigm, namely the estimation of the distribution of regression coefficients in a high-dimensional regression model [13,14]. We will focus on the linear model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta}^* + \boldsymbol{\varepsilon}$$

with a latent and high-dimensional coefficient vector $\boldsymbol{\theta}^* \in \mathbb{R}^d$, whose coordinates have an unknown "prior" distribution $g_*$. Estimation of this prior distribution is a classical example of empirical Bayes inference [15,16], and arises ubiquitously in genetic association analyses where $g_*$ represents the distribution of genetic effect sizes in linear mixed models for complex traits [17–24]. Two recent works [10,25] have established the statistical consistency of nonparametric maximum marginal-likelihood estimators of $g_*$ in settings of high-dimensional regression designs $\mathbf{X} \in \mathbb{R}^{n \times d}$, as $n, d \to \infty$. However, direct computation of this maximum marginal-likelihood estimate is intractable for general regression designs, motivating approaches based on approximate posterior inference schemes.

We will investigate in this work a parametric analogue of a learning procedure proposed in [10], modeling the prior distribution via a parametric model $g(\cdot, \alpha)$ and applying an adaptive diffusion to estimate the parameter $\alpha \in \mathbb{R}^K$. This procedure will take the form of a Langevin diffusion

$$\mathrm{d}\boldsymbol{\theta}^t = \nabla_{\boldsymbol{\theta}} \log \mathsf{P}_{g(\cdot, \widehat{\alpha}^t)}(\boldsymbol{\theta}^t \mid \mathbf{X}, \mathbf{y})\mathrm{d}t + \sqrt{2}\,\mathrm{d}\mathbf{b}^t \tag{1}$$

for sampling from the posterior distribution $\mathsf{P}_{g(\cdot,\widehat{\alpha}^t)}(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y})$ of the regression coefficients, under a prior law $g(\cdot, \widehat{\alpha}^t)$ whose parameter evolves according to a coupled continuous-time dynamics

$$\mathrm{d}\widehat{\alpha}^t = \mathcal{G}\Big(\widehat{\alpha}^t, \frac{1}{d}\sum_{j=1}^d \delta_{\theta_j^t}\Big)\mathrm{d}t. \tag{2}$$

Here $\mathcal{G}(\cdot)$ is a map that implements gradient-based maximum marginal-likelihood learning of $\alpha$ via the empirical distribution of coordinates of $\boldsymbol{\theta}^t$, and we defer a discussion of this motivation to Section 2. The procedure may be understood as an approximation to an idealized dynamics

$$\mathrm{d}\alpha^t = \mathcal{G}(\alpha^t, \mathsf{P}(\theta^t))\mathrm{d}t \tag{3}$$

where $\mathsf{P}(\theta^t)$ denotes the average law of the coordinates $\theta_1^t, \ldots, \theta_d^t$. For these idealized dynamics, the analyses of [10] may be adapted to show that the prior parameter $\alpha^t$ converges to a fixed point of the marginal log-likelihood, under certain conditions for the noise and regression design. Related results have also been shown recently in more general latent variable models in [8, 9, 12, 26], which, in addition, provide convergence guarantees for particle approximations of the McKean-Vlasov type $\mathrm{d}\alpha^t = \mathcal{G}(\alpha^t, \frac{1}{dM}\sum_{j=1}^d \sum_{m=1}^M \delta_{\theta_j^{m,t}})\mathrm{d}t$, having $M$ parallel sampling chains $\{\boldsymbol{\theta}^{1,t}\}_{t\geq 0}, \ldots, \{\boldsymbol{\theta}^{M,t}\}_{t\geq 0}$ for the latent variable $\boldsymbol{\theta} \in \mathbb{R}^d$, in the limit $M \to \infty$.

The aforementioned results are not fully satisfactory in our context of a high-dimensional regression model, and leave open the following two interesting questions about the original dynamics (1–2):

1. *Single chain propagation-of-chaos.* In the limit of increasing dimensions $d \to \infty$, are the idealized dynamics (3) well-approximated by (2) using just a single Langevin chain $\{\boldsymbol{\theta}^t\}_{t\geq 0}$ in $\mathbb{R}^d$?

2. *Characterization of fixed points.* Can the fixed points $\widehat{\alpha}$ of (2) be explicitly characterized? Does (2) exhibit dimension-free convergence to these fixed points, and in what settings is the fixed point representing the maximum marginal-likelihood estimator of $\alpha$ unique?

The purpose of our work is to provide answers to these questions in the context of an i.i.d. regression design. Question 1 is addressed in our companion paper [27], which build upon the recent results of [28, 29] to formalize a dynamical mean-field theory (DMFT) approximation of (2) by (3) over dimension-independent time horizons $t \in [0, T]$, for a general class of such adaptive Langevin dynamics procedures. Our current paper addresses Question 2 by carrying out an analysis of the resulting DMFT system, under an assumption of a uniform log-Sobolev inequality for the posterior law.

## 1.1 Summary of results

Our main results provide an analysis of the DMFT equations that approximate the empirical Bayes Langevin dynamics (1–2) in the high-dimensional limit as $n, d \to \infty$ proportionally. En route to this analysis, we obtain also new results for the DMFT approximation of the standard non-adaptive Langevin diffusion (1) with a fixed prior $g(\cdot) \equiv g(\cdot, \alpha)$. We summarize these results as follows:

1. In the setting of a non-adaptive Langevin diffusion, we formalize a condition of *approximate time-translation-invariance* (TTI) for the DMFT system. We perform an analysis of the dynamical fixed-point equations for the DMFT correlation and response functions under this condition, and show that they recover the static fixed-point equations for the free energy and posterior mean-squared-error predicted by a replica-symmetric ansatz [30, 31].

2. We show that a log-Sobolev inequality (LSI) for the posterior law provides a sufficient condition to guarantee the above approximate-TTI property for the DMFT system, and we discuss several settings of log-concavity, high noise, or large sample size where such an LSI holds. As a consequence, we obtain a new dynamical proof of the validity of the replica-symmetric predictions for the free energy and MSE in the Bayesian linear model with a possibly misspecified prior law, under such an LSI condition.

3. When the LSI holds uniformly over the posterior laws corresponding to the deterministic DMFT trajectory of $\{\alpha^t\}_{t\geq0}$, we show that the empirical Bayes estimate $\widehat{\alpha}^t$ converges on a dimension-free time horizon to a critical point $\alpha^\infty$ of the replica-symmetric limit for the free energy. This is explicitly characterized by a system of scalar fixed-point equations, and we discuss examples of models where this critical point may or may not be unique.

We present and discuss these results and examples in further detail in Section 2.

## 1.2 Further related literature

Approximating the dynamical behavior of many degrees-of-freedom by an effective single-particle problem interacting self-consistently with its environment is an old idea in the statistical physics literature. Relevant to our work is the development of this idea in the context of disordered systems, and in particular the study of high-dimensional Langevin dynamics of soft-spin variants of the Sherrington-Kirkpatrick model [32, 33] and the spherical p-spin model [34–37]. Mathematical proofs of these approximations were first shown for such models in the works of [38–40] using large deviations techniques, and more recently in generalized linear models close to our setting by [28, 29] using different methods around Approximate Message Passing algorithms and iterative Gaussian conditioning. In recent years, DMFT analyses have been applied to study Langevin dynamics and gradient-based optimization in many statistical models and applications, including Gaussian mixture classification [41], matrix and tensor PCA [42–44], phase retrieval and generalized linear models [45,46], and learning in perceptron and neural network models [47–52]. These analyses have uncovered surprising phenomena about the efficacy of gradient-based methods and relationships to landscape complexity for high-dimensional non-convex problems [42].

Understanding the long-time behavior of DMFT systems, in particular in low-temperature regimes characterized by aging or metastability, has been a primary goal in both the physics and mathematics literature since the original inception of these methods (see [53,54] and references within for a review). Mathematically rigorous analyses of long-time dynamics have been obtained previously for spherical 2-spin models in [55] and related statistical models in [44, 56] by leveraging the rotational invariance of these models. However, such analyses of DMFT are (to our knowledge) quite rare in more general settings. Our work takes a step towards filling this gap, by providing a rigorous analysis of the DMFT approximation to Langevin dynamics in a more general model without a rotationally invariant prior, in settings where approximate-TTI holds.

As a by-product of our analyses, we obtain a new proof of a replica formula [57] for the free energy and posterior MSE in the Bayesian linear model. This proof is different from several existing proofs of this result [58–62] and from the Gaussian interpolation methods of Guerra-Talagrand [63,64], and is based instead on deducing a static fixed-point equation from the dynamical fixed-point equations of DMFT. Our current result is specific to a high-temperature regime where a LSI holds for the posterior law, but it applies to models where the prior law is misspecified [30,31,65]. In this misspecified context, the closest mathematical result of which we are aware is [66] which proved the replica-symmetric predictions in a setting where the posterior is log-concave. A complete large deviations analysis of the free energy in a related rank-one matrix estimation model with misspecified prior and noise was carried out in [67], showing that in general the asymptotic free energy is characterized by a Parisi-type variational problem whose solution may not be replica-symmetric. Our results imply for the linear model that this solution must be replica-symmetric under our assumed condition of a LSI for the posterior law.

In the context of adaptive empirical Bayes Langevin dynamics, our results complement the previous analyses of [10] for more general regression designs, and of [8, 12] in general latent variable models. We deduce a dimension-free convergence rate, in contrast to the results of [10] that established convergence (for a nonparametric variant of this algorithm) on a time horizon growing linearly with $n, d$, and without employing a time-dependent and decaying learning rate as in [12]. Under the additional mean-field structure of our current model, we are able to establish convergence of a single-chain implementation of the empirical Bayes Langevin dynamics using (2), rather than for an idealized dynamics as studied in [10, 12] or for an implementation using $M$ parallel chains as studied in [8]. We are also able to give an explicit characterization and analysis of the fixed points to which the dynamics of $\{\widehat{\alpha}^t\}_{t\geq0}$ may converge.

**Notational conventions**

In the context of the posterior law $\mathsf{P}_g(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y})$ for a given prior $g(\cdot)$, we will write

$$\langle f(\boldsymbol{\theta}) \rangle = \mathbb{E}[f(\boldsymbol{\theta}) \mid \mathbf{X}, \mathbf{y}]$$

for the posterior expectation conditioning on the "quenched" variables $\mathbf{X}, \mathbf{y}$. In the context of Langevin dynamics, we will write similarly

$$\langle f(\boldsymbol{\theta}^t) \rangle = \mathbb{E}[f(\boldsymbol{\theta}^t) \mid \mathbf{X}, \mathbf{y}]$$

also for an expectation conditioning on $\mathbf{X}, \mathbf{y}$. In some arguments it is convenient to consider the expectation also conditioned on the initial condition $\boldsymbol{\theta}^0$, and we will denote this by

$$\langle f(\boldsymbol{\theta}^t) \rangle_{\mathbf{x}} = \mathbb{E}[f(\boldsymbol{\theta}^t) \mid \mathbf{X}, \mathbf{y}, \boldsymbol{\theta}^0 = \mathbf{x}].$$

We reserve $\mathbb{E}$ and $\mathbb{P}$ for the full expectation and probability also over $\mathbf{X}, \mathbf{y}, \boldsymbol{\theta}^*, \boldsymbol{\varepsilon}$.

Constants $C, C', c, c' > 0$ throughout are independent of the dimensions $n, d$. For any random variable $\xi$ in a complete and separable normed vector space $(\mathcal{M}, \|\cdot\|)$, we will use $\mathsf{P}(\xi)$ to denote its law. $\mathcal{P}_2(\mathcal{M})$ is the space of probability distributions $\mathsf{P}$ on $(\mathcal{M}, \|\cdot\|)$ such that $\mathbb{E}_{\xi \sim \mathsf{P}} \|\xi\|^2 < \infty$, and $W_1(\cdot)$ and $W_2(\cdot)$ denote the Wasserstein-1 and Wasserstein-2 metrics on $\mathcal{P}_2(\mathcal{M})$.

For $f : \mathbb{R}^d \to \mathbb{R}$, $\nabla f \in \mathbb{R}^d$ is its gradient, $\nabla^2 f \in \mathbb{R}^{d \times d}$ its Hessian, and $\nabla^3 f \in \mathbb{R}^{d \times d \times d}$ the symmetric tensor of its $3^{\mathrm{rd}}$-order partial derivatives. For $f : \mathbb{R} \times \mathbb{R}^K \to \mathbb{R}$, $\partial_\theta f(\theta, \alpha)$ and $\nabla_\alpha f(\theta, \alpha)$ denote its partial derivatives with respect to $\theta \in \mathbb{R}$ and $\alpha \in \mathbb{R}^K$. $\|\cdot\|_2$ is the Euclidean norm for vectors and vectorized Euclidean norm for matrices and tensors. $\mathrm{Tr}\, M$ and $\|M\|_{\mathrm{op}}$ are the matrix trace and Euclidean operator norm. $C([0, T], \mathbb{R}^d)$ is the space of continuous functions $f : [0, T] \to \mathbb{R}^d$ equipped with the norm of uniform convergence $\|f\|_\infty = \sup_{t \in [0,T]} \|f(t)\|_2$. $C^k(\mathbb{R}^d, \mathbb{R}^m)$ is the space of functions $f : \mathbb{R}^d \to \mathbb{R}^m$ that are $k$-times continuously-differentiable. For two probability densities $p, q$ on $\mathbb{R}^d$, $\mathrm{D}_{\mathrm{KL}}(p\|q) = \int q(\log q - \log p)$ is the Kullback-Leibler divergence. For a scalar random variable $X$, $\mathrm{Var}\, X = \mathbb{E}X^2 - (\mathbb{E}X)^2$ and $\mathrm{Ent}\, X = \mathbb{E}X \log X - \mathbb{E}X \log \mathbb{E}X$. For a random vector $X \in \mathbb{R}^k$, $\mathrm{Cov}\, X = \mathbb{E}XX^\top - (\mathbb{E}X)(\mathbb{E}X)^\top \in \mathbb{R}^{k \times k}$.

# 2 Model and main results

## 2.1 Bayesian linear model and adaptive Langevin dynamics

We study a linear model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta}^* + \boldsymbol{\varepsilon} \in \mathbb{R}^n \tag{4}$$

with random effects $\boldsymbol{\theta}^* \in \mathbb{R}^d$. Modeling $\theta_1^*, \ldots, \theta_d^* \overset{iid}{\sim} g$ for a prior density $g(\cdot)$ on the real line and modeling $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$ as Gaussian noise, Bayesian inference for $\boldsymbol{\theta}^*$ is based upon the posterior density

$$\mathsf{P}_g(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y}) = \frac{1}{\mathsf{P}_g(\mathbf{y} \mid \mathbf{X})} \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2\right) \prod_{j=1}^d g(\theta_j). \tag{5}$$

Here $\mathsf{P}_g(\mathbf{y} \mid \mathbf{X})$ is the marginal likelihood of $\mathbf{y}$ (i.e. model evidence or partition function), given by

$$\mathsf{P}_g(\mathbf{y} \mid \mathbf{X}) = \int \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2}\|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2\right) \prod_{j=1}^{d} g(\theta_j)\mathrm{d}\theta_j. \tag{6}$$

We will study (overdamped) Langevin dynamics for sampling from the posterior density (5) in two settings, the first in which the prior law $g(\cdot)$ is fixed but may be misspecified, and the second in which this prior law may adapt to the Langevin trajectory to implement empirical Bayes learning from the observed data $(\mathbf{X}, \mathbf{y})$. In the former setting, we consider the Langevin dynamics

$$\mathrm{d}\boldsymbol{\theta}^t = \nabla_{\boldsymbol{\theta}}\left(-\frac{1}{2\sigma^2}\|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}^t\|_2^2 + \sum_{j=1}^{d} \log g(\theta_j^t)\right)\mathrm{d}t + \sqrt{2}\,\mathrm{d}\mathbf{b}^t \tag{7}$$

where $\{\mathbf{b}^t\}_{t\geq 0}$ is a standard Brownian motion on $\mathbb{R}^d$. In the latter setting, we will model the prior via a parametric model

$$\left\{g(\,\cdot\,,\alpha) : \alpha \in \mathbb{R}^K\right\} \tag{8}$$

and consider the empirical Bayes Langevin dynamics

$$\mathrm{d}\boldsymbol{\theta}^t = \nabla_{\boldsymbol{\theta}}\left(-\frac{1}{2\sigma^2}\|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}^t\|_2^2 + \sum_{j=1}^{d} \log g(\theta_j^t, \widehat{\alpha}^t)\right)\mathrm{d}t + \sqrt{2}\,\mathrm{d}\mathbf{b}^t \tag{9}$$

$$\mathrm{d}\widehat{\alpha}^t = \nabla_{\alpha}\left(\frac{1}{d}\sum_{j=1}^{d} \log g(\theta_j^t, \widehat{\alpha}^t) - R(\widehat{\alpha}^t)\right)\mathrm{d}t. \tag{10}$$

The equation (10) describes a continuous-time evolution of the prior parameter $\alpha \in \mathbb{R}^K$ that is coupled to the Langevin diffusion (9) of the posterior sample, and $R : \mathbb{R}^K \to \mathbb{R}$ is a possible smooth regularizer. (In this work, we will be interested mostly in the behavior of these dynamics when $R(\alpha) \equiv 0$, and we introduce $R(\alpha)$ for theoretical purposes to confine the dynamics of $\widehat{\alpha}^t$ in certain examples.)

To motivate the dynamics (9–10) as a procedure that implements maximum marginal-likelihood learning of $\alpha \in \mathbb{R}^K$, we may consider the free energy (i.e. negative marginal log-likelihood)

$$\widehat{F}(\alpha) = -\frac{1}{d}\log \mathsf{P}_{g(\cdot,\alpha)}(\mathbf{y} \mid \mathbf{X}) \tag{11}$$

as a function of the prior parameter $\alpha \in \mathbb{R}^K$. By the Gibbs variational principle (c.f. [68, Proposition 4.7]),

$$\widehat{F}(\alpha) = \inf_{q \in \mathcal{P}_*(\mathbb{R}^d)} V(q, \alpha) \tag{12}$$

where $\mathcal{P}_*(\mathbb{R}^d)$ is the space of all probability densities on $\mathbb{R}^d$, and

$$V(q, \alpha) = \frac{1}{d}\int \left(\frac{1}{2\sigma^2}\|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|^2 - \sum_{j=1}^{d} \log g(\theta_j, \alpha) + \log q(\boldsymbol{\theta})\right)q(\boldsymbol{\theta})\mathrm{d}\boldsymbol{\theta} + \frac{n}{2d}\log 2\pi\sigma^2 \tag{13}$$

is the Gibbs free energy corresponding to the prior $g(\cdot) = g(\,\cdot\,,\alpha)$. We propose to implement maximum-likelihood learning of $\alpha \in \mathbb{R}^K$ by minimizing the regularized Gibbs free energy $V(q, \alpha) + R(\alpha)$ jointly over $(q, \alpha)$, via a gradient flow in the Wasserstein-2 geometry for $q \in \mathcal{P}_*(\mathbb{R}^d)$ and the standard Euclidean geometry for $\alpha \in \mathbb{R}^K$. The resulting gradient flow equations take the form

$$\frac{\mathrm{d}}{\mathrm{d}t}q_t = -d \cdot \mathrm{grad}_q^{W_2} V(q_t, \alpha^t) := \nabla_{\boldsymbol{\theta}} \cdot \left[q_t(\boldsymbol{\theta})\nabla_{\boldsymbol{\theta}}\left(\frac{1}{2\sigma^2}\|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 - \sum_{j=1}^{d} \log g(\theta_j, \alpha^t)\right)\right] + \mathrm{Tr}\,\nabla_{\boldsymbol{\theta}}^2 q_t(\boldsymbol{\theta}), \quad (14)$$

$$\frac{\mathrm{d}}{\mathrm{d}t}\alpha^t = -\nabla_{\alpha}[V(q_t, \alpha^t) + R(\alpha^t)] = \nabla_{\alpha}\left(\int \frac{1}{d}\sum_{j=1}^{d} \log g(\theta_j, \alpha^t)\,q_t(\boldsymbol{\theta})\mathrm{d}\boldsymbol{\theta} - R(\alpha^t)\right). \tag{15}$$

6

In (14), we identify $\text{grad}_q^{W_2} V(q, \alpha)$ with the Fokker-Planck equation for the density evolution of $\boldsymbol{\theta}^t$ under the Langevin diffusion (7) with prior law $g(\cdot) = g(\cdot, \alpha)$, via its variational interpretation put forth in [69]. Then (9–10) may be understood as a particle implementation of (14–15) that uses a single Langevin trajectory $\{\boldsymbol{\theta}^t\}_{t\geq 0}$ to simulate the dynamics of $q_t$ in (14), and that uses the empirical distribution $\frac{1}{d}\sum_{j=1}^{d}\delta_{\theta_j^t}$ to approximate the expectation over $\boldsymbol{\theta} \sim q_t$ in the dynamics of $\alpha^t$ in (15). An algorithm similar to (9–10) was introduced in [10], with some additional reparametrization ideas to allow for nonparametric modeling of the prior $g(\cdot)$. Here, to simplify technical considerations, we restrict our study to parametric prior models of the form (8).

## 2.2 DMFT equations

The empirical Bayes Langevin diffusion (9–10) is an example of a general class of adaptive Langevin diffusions that we study in our companion work [27]. In particular, the gradient equation (10) for $\widehat{\alpha}^t$ is a function of the empirical distribution of coordinates $\boldsymbol{\theta}^t$,

$$\mathrm{d}\widehat{\alpha}^t = \mathcal{G}\Big(\widehat{\alpha}^t, \frac{1}{d}\sum_{j=1}^{d}\delta_{\theta_j^t}\Big)\mathrm{d}t$$

where $\mathcal{G}: \mathbb{R}^K \times \mathcal{P}_2(\mathbb{R}) \to \mathbb{R}^K$ is the gradient map

$$\mathcal{G}(\alpha, \mathsf{P}) = \mathbb{E}_{\theta\sim\mathsf{P}}[\nabla_\alpha \log g(\theta, \alpha)] - \nabla R(\alpha).$$

Our analyses will rely on a system of dynamical mean-field theory (DMFT) equations, formalized in [27] and building upon the results of [28, 29], that describes a deterministic evolution of a prior parameter $\alpha^t \in \mathbb{R}^K$ and a univariate law $\mathsf{P}(\theta^t) \in \mathcal{P}_2(\mathbb{R})$ that approximate $(\widehat{\alpha}^t, \frac{1}{d}\sum_{j=1}^{d}\delta_{\theta_j^t})$ in the large system limit $n, d \to \infty$. This approximation will hold under the following assumptions, which we assume throughout this work.

**Assumption 2.1** (Linear model and initial conditions).

(a) (Asymptotic scaling) $\lim_{n,d\to\infty} \frac{n}{d} = \delta \in (0, \infty)$.

(b) (Random design) $\mathbf{X} = (x_{ij}) \in \mathbb{R}^{n\times d}$ has independent entries satisfying $\mathbb{E}[x_{ij}] = 0$, $\mathbb{E}[x_{ij}^2] = \frac{1}{d}$, and $\|\sqrt{d}x_{ij}\|_{\psi_2} \leq C$ for a constant $C > 0$, where $\|\cdot\|_{\psi_2}$ is the sub-Gaussian norm.

(c) (Bayesian linear model) $\boldsymbol{\theta}^*, \boldsymbol{\varepsilon}$ are independent of each other and of $\mathbf{X}$, and $\mathbf{y} = \mathbf{X}\boldsymbol{\theta}^* + \boldsymbol{\varepsilon}$. The entries of $\boldsymbol{\theta}^*, \boldsymbol{\varepsilon}$ are distributed as

$$\theta_1^*, \ldots, \theta_d^* \overset{iid}{\sim} g_*, \qquad \varepsilon_1, \ldots, \varepsilon_n \overset{iid}{\sim} \mathcal{N}(0, \sigma^2) \tag{16}$$

for some $\sigma^2 > 0$ and probability density $g_*$ (both fixed and independent of $n, d$), where $g_*$ satisfies the log-Sobolev inequality

$$\text{Ent}_{\theta^*\sim g_*}[f(\theta^*)^2] \leq C_{\text{LSI}}\,\mathbb{E}_{\theta^*\sim g_*}[(f'(\theta^*))^2] \text{ for all } f \in C^1(\mathbb{R}). \tag{17}$$

(d) (Initial conditions) The initialization $\boldsymbol{\theta}^0$ is independent of $\mathbf{X}, \boldsymbol{\theta}^*, \boldsymbol{\varepsilon}$, and

$$\theta_1^0, \ldots, \theta_d^0 \overset{iid}{\sim} g_0 \tag{18}$$

for some probability density $g_0$ (fixed and independent of $n, d$) with finite entropy $\int g_0 \log g_0$ and finite moment-generating-function in a neighborhood of 0. The initialization $\widehat{\alpha}^0$ satisfies $\lim_{n,d\to\infty}\widehat{\alpha}^0 = \alpha^0$ a.s. for a deterministic parameter $\alpha^0 \in \mathbb{R}^K$.

**Assumption 2.2** (Prior model and regularizer).

(a) In the context of a fixed prior, $g(\theta)$ is strictly positive and thrice continuously-differentiable, and $(\log g)'''(\theta)$ is uniformly Hölder continuous over $\theta \in \mathbb{R}$. For a constant $C > 0$ and all $\theta \in \mathbb{R}$,

$$|(\log g)'(\theta)| \leq C(1 + |\theta|), \qquad |(\log g)''(\theta)| \leq C, \qquad |(\log g)'''(\theta)| \leq C,$$

and for some constants $r_0, c_0 > 0$,

$$-(\log g)''(\theta) \geq c_0 \text{ for all } |\theta| \geq r_0.$$

(b) In the context of an adaptive prior, $g(\theta, \alpha)$ is strictly positive, $R(\alpha)$ is nonnegative, and both are thrice continuously-differentiable. For a constant $C > 0$ and all $(\theta, \alpha) \in \mathbb{R} \times \mathbb{R}^K$,

$$\|\nabla_{(\theta,\alpha)} \log g(\theta, \alpha)\|_2 \leq C(1 + |\theta| + \|\alpha\|_2),$$
$$\|\nabla R(\alpha)\|_2 \leq C(1 + \|\alpha\|_2). \tag{19}$$

Furthermore, for each compact subset $S \subset \mathbb{R}^K$, $\theta \mapsto \nabla^3_{(\theta,\alpha)} \log g(\theta, \alpha)$ is Hölder-continuous uniformly over $(\theta, \alpha) \in \mathbb{R} \times S$, and for some constants $C(S), r_0(S), c_0(S) > 0$,

$$\|\nabla^2_{(\theta,\alpha)} \log g(\theta, \alpha)\|_2 \leq C(S), \ \|\nabla^3_{(\theta,\alpha)} \log g(\theta, \alpha)\|_2 \leq C(S) \text{ for all } (\theta, \alpha) \in \mathbb{R} \times S,$$
$$-\partial^2_\theta \log g(\theta, \alpha) \geq c_0(S) \text{ for all } |\theta| \geq r_0(S) \text{ and } \alpha \in S, \tag{20}$$
$$\|\nabla^2 R(\alpha)\|_2 \leq C(S), \ \|\nabla^3 R(\alpha)\|_2 \leq C(S) \text{ for all } \alpha \in S.$$

In particular, Assumption 2.2(b) requires that $g(\cdot) \equiv g(\,\cdot\,, \alpha)$ satisfies Assumption 2.2(a) for each fixed prior parameter $\alpha \in \mathbb{R}^K$. We assume the LSI condition (17) for the true prior $g_*$ to ensure concentration of the free energy (c.f. Proposition 4.11), and we clarify that the conditions of Assumption 2.2(a) imply that the modeled prior $g(\cdot)$ must also satisfy an LSI of the form (17), as reviewed in Lemma C.1.

Under the above Assumptions 2.1 and 2.2(b), the DMFT limit for (9–10) is described by the following construction: Let

$$\theta^* \sim g_*, \qquad \theta^0 \sim g_0, \qquad \varepsilon \sim \mathcal{N}(0, \sigma^2) \tag{21}$$

denote independent scalar variables with the distributions (16) and (18), and let $\delta = \lim \frac{n}{d}$ be as in Assumption 2.1. Let $\{b^t\}_{t \geq 0}$, $\{u^t\}_{t \geq 0}$, and $(w^*, \{w^t\}_{t \geq 0})$ be centered univariate Gaussian processes independent of each other and of $(\theta^*, \theta^0, \varepsilon)$, where $\{b^t\}_{t \geq 0}$ is a standard Brownian motion on $\mathbb{R}$, and $\{u^t\}_{t \geq 0}$ and $(w^*, \{w^t\}_{t \geq 0})$ have covariance kernels

$$\mathbb{E}[u^t u^s] = C_\eta(t, s), \qquad \mathbb{E}[w^t w^s] = C_\theta(t, s), \qquad \mathbb{E}[w^t w^*] = C_\theta(t, *), \qquad \mathbb{E}[(w^*)^2] = C_\theta(*, *) \tag{22}$$

defined self-consistently in (28) below. We consider a system of stochastic differential equations

$$\mathrm{d}\theta^t = \left[ -\frac{\delta}{\sigma^2}(\theta^t - \theta^*) + \partial_\theta \log g(\theta^t, \alpha^t) + \int_0^t R_\eta(t, s)(\theta^s - \theta^*)\mathrm{d}s + u^t \right]\mathrm{d}t + \sqrt{2}\,\mathrm{d}b^t \tag{23}$$

$$\mathrm{d}\left(\frac{\partial \theta^t}{\partial u^s}\right) = \left[ -\left(\frac{\delta}{\sigma^2} - \partial^2_\theta \log g(\theta^t, \alpha^t)\right)\frac{\partial \theta^t}{\partial u^s} + \int_s^t R_\eta(t, s')\frac{\partial \theta^{s'}}{\partial u^s}\mathrm{d}s' \right]\mathrm{d}t \tag{24}$$

for univariate processes $\{\theta^t\}_{t \geq 0}$ and $\{\frac{\partial \theta^t}{\partial u^s}\}_{t \geq s \geq 0}$ adapted to the filtration $\mathcal{F}^\theta_t := \mathcal{F}(\{b^s\}_{s \leq t}, \{u^s\}_{s \leq t}, \theta^*, \theta^0)$, with the initial conditions

$$\theta^t|_{t=0} = \theta^0, \qquad \left.\frac{\partial \theta^t}{\partial u^s}\right|_{t=s} = 1.$$

These are driven by a deterministic and continuous $\mathbb{R}^K$-valued process $\{\alpha^t\}_{t \geq 0}$ representing the asymptotic limit of $\{\widehat{\alpha}^t\}_{t \geq 0}$. We consider likewise univariate processes $\{\eta^t\}_{t \geq 0}$ and $\{\frac{\partial \eta^t}{\partial w^s}\}_{t \geq s \geq 0}$ defined by

$$\eta^t = -\frac{1}{\sigma^2}\int_0^t R_\theta(t, s)\big(\eta^s + w^* - \varepsilon\big)\mathrm{d}s - w^t \tag{25}$$

$$\frac{\partial \eta^t}{\partial w^s} = -\frac{1}{\sigma^2}\left[\int_s^t R_\theta(t, s')\frac{\partial \eta^{s'}}{\partial w^s}\mathrm{d}s' - R_\theta(t, s)\right] \tag{26}$$

adapted to the filtration $\mathcal{F}^\eta_t := \mathcal{F}(\{w^s\}_{s \leq t}, w^*, \varepsilon)$. The deterministic process $\{\alpha^t\}_{t \geq 0}$ above is defined self-consistently by

$$\frac{\mathrm{d}}{\mathrm{d}t}\alpha^t = \mathcal{G}(\alpha^t, \mathsf{P}(\theta^t)), \qquad \mathcal{G}(\alpha, \mathsf{P}) = \mathbb{E}_{\theta \sim \mathsf{P}}[\nabla_\alpha \log g(\theta, \alpha)] - \nabla R(\alpha) \tag{27}$$

with initial condition $\alpha^t|_{t=0} = \alpha^0$ given in Assumption 2.2, where $\mathsf{P}(\theta^t)$ is the law of $\theta^t$ in (23). The covariance and response functions $C_\theta, C_\eta, R_\theta, R_\eta$ are also defined for all $t \geq s \geq 0$ self-consistently via the above processes by

$$C_\theta(t,s) = \mathbb{E}[\theta^t \theta^s], \quad C_\theta(t,*) = \mathbb{E}[\theta^t \theta^*], \quad C_\theta(*,*) = \mathbb{E}[(\theta^*)^2],$$

$$C_\eta(t,s) = \frac{\delta}{\sigma^4} \mathbb{E}[(\eta^t + w^* - \varepsilon)(\eta^s + w^* - \varepsilon)] \tag{28}$$

$$R_\theta(t,s) = \mathbb{E}\Big[\frac{\partial \theta^t}{\partial u^s}\Big], \quad R_\eta(t,s) = \frac{\delta}{\sigma^2} \mathbb{E}\Big[\frac{\partial \eta^t}{\partial w^s}\Big].$$

This DMFT system (22–28) describes the $n, d \to \infty$ limit of the empirical Bayes Langevin dynamics (9–10). In the setting of a fixed prior $g(\cdot) \equiv g(\,\cdot\,, \alpha^0)$, the DMFT limit for the standard Langevin diffusion (7) is the same, upon replacing $\mathcal{G}(\alpha, \mathsf{P})$ in (27) by $\mathcal{G}(\alpha, \mathsf{P}) = 0$ so that $\alpha^t = \alpha^0$ for all $t \geq 0$.

Fixing any time horizon $T > 0$, let us set

$$\eta^* = -w^*$$

and denote by

$$\mathsf{P}(\theta^*, \{\theta^t\}_{t \in [0,T]}) \in \mathcal{P}_2(\mathbb{R} \times C([0,T], \mathbb{R})), \quad \mathsf{P}(\eta^*, \varepsilon, \{\eta^t\}_{t \in [0,T]}) \in \mathcal{P}_2(\mathbb{R} \times \mathbb{R} \times C([0,T], \mathbb{R}))$$

the joint laws of sample paths $(\theta^*, \{\theta^t\}_{t \in [0,T]})$ and $(\eta^*, \varepsilon, \{\eta^t\}_{t \in [0,T]})$ in this DMFT system. We write $\theta_j^*, \theta_j^t, \varepsilon_i, \eta_i^*, \eta_i^t$ for the coordinates of $\boldsymbol{\theta}^*, \boldsymbol{\theta}^t, \boldsymbol{\varepsilon}, \boldsymbol{\eta}^* = \mathbf{X}\boldsymbol{\theta}^*, \boldsymbol{\eta}^t = \mathbf{X}\boldsymbol{\theta}^t$, and $\xrightarrow{W_2}$ for Wasserstein-2 convergence in the spaces $\mathcal{P}_2(\mathbb{R} \times C([0,T], \mathbb{R}))$ and $\mathcal{P}_2(\mathbb{R} \times \mathbb{R} \times C([0,T], \mathbb{R}))$ as $n, d \to \infty$. The main result we will use from our companion work [27] is summarized in the following theorem.

**Theorem 2.3.** *(a) Suppose Assumptions 2.1 and 2.2(a) hold, and identify $g(\cdot) \equiv g(\,\cdot\,, \alpha^0)$. Let $\{\boldsymbol{\theta}^t\}_{t \geq 0}$ be the solution to the dynamics (7) with fixed prior $g(\cdot)$, and denote $\boldsymbol{\eta}^* = \mathbf{X}\boldsymbol{\theta}^*$ and $\boldsymbol{\eta}^t = \mathbf{X}\boldsymbol{\theta}^t$. Then for each fixed $T \geq 0$, there exists a solution up to time $T$ of the DMFT system (22–28) with (27) replaced by $\mathcal{G}(\alpha, \mathsf{P}) = 0$, such that almost surely as $n, d \to \infty$,*

$$\frac{1}{d} \sum_{j=1}^d \delta_{\theta_j^*, \{\theta_j^t\}_{t \in [0,T]}} \xrightarrow{W_2} \mathsf{P}(\theta^*, \{\theta^t\}_{t \in [0,T]}), \quad \frac{1}{n} \sum_{i=1}^n \delta_{\eta_i^*, \varepsilon_i, \{\eta_i^t\}_{t \in [0,T]}} \xrightarrow{W_2} \mathsf{P}(\eta^*, \varepsilon, \{\eta^t\}_{t \in [0,T]}). \tag{29}$$

*(b) Suppose Assumptions 2.1 and 2.2(b) hold, and let $\{\boldsymbol{\theta}^t, \widehat{\alpha}^t\}_{t \geq 0}$ be the solution to the empirical Bayes Langevin dynamics (9–10). Then for each fixed $T > 0$, there exists a solution up to time $T$ of the DMFT system (22–28) such that almost surely as $n, d \to \infty$, (29) holds and also*

$$\{\widehat{\alpha}^t\}_{t \in [0,T]} \to \{\alpha^t\}_{t \in [0,T]} \text{ in } C([0,T], \mathbb{R}^K).$$

Both parts of this theorem follow from [27, Theorem 2.5], and we explain the details of the reduction to [27, Theorem 2.5] in Appendix A. The above solutions to the DMFT systems are unique in certain domains of exponential growth for $\{\alpha^t\}$ and for the correlation and response functions, and we refer readers to [27, Theorem 2.4] for details of this uniqueness claim.

For our analyses of dynamics with fixed prior $g(\cdot)$ in the setting of Theorem 2.3(a), we will require a second result from [27] that gives an interpretation for the DMFT response functions $R_\theta(t,s)$ and $R_\eta(t,s)$ as coordinate averages of single-particle responses in the Langevin diffusion (7). We defer a statement of this result to Section 4.1.

## 2.3 Replica-symmetric characterization of equilibrium for a fixed prior

This and the next section describe the main results of our current paper. We discuss results pertaining to the dynamics (7) with a fixed prior $g(\cdot)$ in this section, and results pertaining to the empirical Bayes dynamics (9–10) in Section 2.4 to follow.

### 2.3.1 Approximately-TTI DMFT systems

We first introduce a set of conditions for the correlation and response functions of the DMFT system that characterize an approximate time-translation-invariance (TTI) property. Under these conditions, we establish convergence of the joint law of $(\theta^*, \theta^t)$ in the DMFT equations to a replica-symmetric fixed point as $t \to \infty$.

**Definition 2.4.** In the setting of a fixed prior $g(\cdot)$ [i.e. with $\mathcal{G}(\alpha, \mathsf{P}) = 0$ in (27)], the solution of the DMFT system (22–28) is *approximately-TTI* if it satisfies the following conditions:

1. There exists a scalar value $c_\theta(*) \in \mathbb{R}$ and functions $c_\theta^{\mathrm{tti}}, c_\eta^{\mathrm{tti}} : [0, \infty) \to \mathbb{R}$ such that, for some $\varepsilon : [0, \infty) \to [0, \infty)$ satisfying $\lim_{s \to \infty} \varepsilon(s) = 0$ and for all $t \geq s \geq 0$,

$$|C_\theta(t, s) - c_\theta^{\mathrm{tti}}(t - s)| \leq \varepsilon(s), \tag{30}$$

$$|C_\eta(t, s) - c_\eta^{\mathrm{tti}}(t - s)| \leq \varepsilon(s), \tag{31}$$

$$|C_\theta(s, *) - c_\theta(*)| \leq \varepsilon(s). \tag{32}$$

Furthermore, there exist values $c_\theta^{\mathrm{tti}}(\infty), c_\eta^{\mathrm{tti}}(\infty) \geq 0$ and finite positive measures $\mu_\theta, \mu_\eta$ supported on $[\iota, \infty)$ for some $\iota > 0$ (strictly) such that

$$c_\theta^{\mathrm{tti}}(\tau) = c_\theta^{\mathrm{tti}}(\infty) + \int_\iota^\infty e^{-a\tau} \mathrm{d}\mu_\theta(a), \qquad c_\eta^{\mathrm{tti}}(\tau) = c_\eta^{\mathrm{tti}}(\infty) + \int_\iota^\infty e^{-a\tau} \mathrm{d}\mu_\eta(a). \tag{33}$$

2. There exist functions $r_\theta^{\mathrm{tti}}, r_\eta^{\mathrm{tti}} : [0, \infty) \to \mathbb{R}$ such that, for some $\varepsilon : [0, \infty) \to [0, \infty)$ satisfying $\lim_{t \to \infty} \varepsilon(t) = 0$,

$$\int_0^t |R_\theta(t, s) - r_\theta^{\mathrm{tti}}(t - s)| \mathrm{d}s \leq \varepsilon(t), \tag{34}$$

$$\int_0^t |R_\eta(t, s) - r_\eta^{\mathrm{tti}}(t - s)| \mathrm{d}s \leq \varepsilon(t). \tag{35}$$

Furthermore, $r_\theta^{\mathrm{tti}}, r_\eta^{\mathrm{tti}}, c_\theta^{\mathrm{tti}}, c_\eta^{\mathrm{tti}}$ satisfy the fluctuation-dissipation relations

$$r_\theta^{\mathrm{tti}}(\tau) = -c_\theta^{\mathrm{tti}\prime}(\tau), \qquad r_\eta^{\mathrm{tti}}(\tau) = -c_\eta^{\mathrm{tti}\prime}(\tau). \tag{36}$$

We show that if the DMFT system is approximately-TTI in the above sense, then its $t \to \infty$ limit is characterized by a system of "static" scalar fixed-point equations. To describe this characterization, consider a scalar Gaussian convolution model

$$y = \theta^* + z \in \mathbb{R}. \tag{37}$$

Let

$$\mathsf{P}_{g,\omega}(\theta \mid y) = \frac{1}{\mathsf{P}_{g,\omega}(y)} \sqrt{\frac{\omega}{2\pi}} \exp\left(-\frac{\omega}{2}(y - \theta)^2\right) g(\theta) \tag{38}$$

be the posterior distribution of $\theta$ in this model, assuming a prior law $\theta \sim g(\cdot)$ and independent Gaussian noise $z \sim \mathcal{N}(0, \omega^{-1})$, where

$$\mathsf{P}_{g,\omega}(y) = \int \sqrt{\frac{\omega}{2\pi}} \exp\left(-\frac{\omega}{2}(y - \theta)^2\right) g(\theta) \mathrm{d}\theta \tag{39}$$

denotes the marginal density of $y$ under these assumptions. Let the true model be $y = \theta^* + z$ with $\theta^* \sim g_*$ and independent noise $z \sim \mathcal{N}(0, \omega_*^{-1})$, and denote by

$$\mathsf{P}_{g_*,\omega_*;g,\omega}(\theta^*, \theta) \tag{40}$$

the joint law of the true parameter $\theta^*$ and a posterior sample $\theta$ under the generating process

$$\theta^* \sim g_*, \ z \sim \mathcal{N}(0, \omega_*^{-1}) \ (\text{independent}) \quad \Rightarrow \quad y = \theta^* + z \quad \Rightarrow \quad \theta \mid y \sim \mathsf{P}_{g,\omega}(\cdot \mid y) \tag{41}$$

(where $\theta \mid y$ is defined with misspecified prior law $g(\cdot)$ and misspecified noise variance $\omega^{-1}$). We write $\langle f(\theta) \rangle_{g,\omega}$ for the posterior average with respect to $\mathsf{P}_{g,\omega}(\cdot \mid y)$ depending implicitly on $y$, and $\mathbb{E}_{g_*,\omega_*} f(y)$ for the expectation under the true model $y = \theta^* + z$. Thus, an expectation over the joint law $\mathsf{P}_{g_*,\omega_*;g,\omega}$ in (40) takes the form

$$\mathbb{E}_{(\theta^*,\theta) \sim \mathsf{P}_{g_*,\omega_*;g,\omega}} f(\theta^*,\theta) = \mathbb{E}_{g_*,\omega_*} \langle f(\theta^*,\theta) \rangle_{g,\omega}.$$

**Theorem 2.5.** *Suppose Assumptions 2.1 and 2.2(a) hold. Consider the Langevin diffusion (7) with a fixed prior $g(\cdot)$, and suppose that the corresponding solution of the DMFT system in Theorem 2.3(a) is approximately-TTI. Define, from the quantities of Definition 2.4,*

$$\mathrm{mse} = c_\theta^{\mathrm{tti}}(0) - c_\theta^{\mathrm{tti}}(\infty), \qquad \mathrm{mse}_* = \mathbb{E}[\theta^{*2}] - 2c_\theta(*) + c_\theta^{\mathrm{tti}}(\infty),$$

$$\mathrm{ymse} = \frac{\sigma^4}{\delta}\big(c_\eta^{\mathrm{tti}}(0) - c_\eta^{\mathrm{tti}}(\infty)\big), \qquad \mathrm{ymse}_* = \frac{\sigma^4}{\delta}\big(2c_\eta^{\mathrm{tti}}(0) - c_\eta^{\mathrm{tti}}(\infty)\big) - \sigma^2. \tag{42}$$

*Then there are unique values $\omega, \omega_* > 0$ (given $\mathrm{mse}, \mathrm{mse}_*$) for which $\mathrm{mse}, \mathrm{mse}_*, \omega, \omega_*$ satisfy the fixed-point equations*

$$\omega = \delta(\sigma^2 + \mathrm{mse})^{-1}, \qquad \omega_* = \delta(\sigma^2 + \mathrm{mse}_*)^{-1},$$

$$\mathrm{mse} = \mathbb{E}_{g_*,\omega_*}[\langle (\theta - \langle\theta\rangle_{g,\omega})^2 \rangle_{g,\omega}], \qquad \mathrm{mse}_* = \mathbb{E}_{g_*,\omega_*}[(\theta^* - \langle\theta\rangle_{g,\omega})^2]. \tag{43}$$

*The quantities $\mathrm{ymse}, \mathrm{ymse}_*$ are related to these fixed points by*

$$\mathrm{ymse} = \sigma^2\Big(1 - \frac{\omega\sigma^2}{\delta}\Big), \qquad \mathrm{ymse}_* = \sigma^2 + \frac{\omega\sigma^4}{\delta}\Big(\frac{\omega}{\omega_*} - 2\Big). \tag{44}$$

*Furthermore, letting $\mathsf{P}(\theta^*,\theta^t)$ be the joint law of $(\theta^*,\theta^t)$ in the DMFT system, as $t \to \infty$,*

$$\mathsf{P}(\theta^*,\theta^t) \overset{W_2}{\to} \mathsf{P}_{g_*,\omega_*;g,\omega}. \tag{45}$$

**Remark 2.6.** Let $\langle f(\boldsymbol{\theta}) \rangle$ and $\langle f(\boldsymbol{\theta},\boldsymbol{\theta}') \rangle$ denote the expectation over independent samples $\boldsymbol{\theta}, \boldsymbol{\theta}' \sim \mathsf{P}_g(\cdot \mid \mathbf{X}, \mathbf{y})$ from the posterior law (5) with a fixed prior $g(\cdot)$. Then the asymptotic overlaps

$$\lim_{n,d\to\infty} d^{-1}\langle \boldsymbol{\theta}^\top \boldsymbol{\theta} \rangle, \qquad \lim_{n,d\to\infty} d^{-1}\langle \boldsymbol{\theta}^\top \boldsymbol{\theta}' \rangle, \qquad \lim_{n,d\to\infty} d^{-1}\langle \boldsymbol{\theta}^\top \boldsymbol{\theta}^* \rangle$$

are predicted in the DMFT system, respectively, by

$$c_\theta^{\mathrm{tti}}(0) = \lim_{t\to\infty} C_\theta(t,t), \qquad c_\theta^{\mathrm{tti}}(\infty) = \lim_{t,\tau\to\infty} C_\theta(t,t+\tau), \qquad c_\theta(*) = \lim_{t\to\infty} C_\theta(t,*).$$

Thus $\mathrm{mse}$ and $\mathrm{mse}_*$ as defined in (42) represent the DMFT predictions for

$$\lim_{n,d\to\infty} d^{-1}\langle \|\boldsymbol{\theta} - \langle\boldsymbol{\theta}\rangle\|_2^2 \rangle, \qquad \lim_{n,d\to\infty} d^{-1}\|\boldsymbol{\theta}^* - \langle\boldsymbol{\theta}\rangle\|_2^2.$$

Similarly, one may check that $\mathrm{ymse}$ and $\mathrm{ymse}_*$ as defined in (42) represent the DMFT predictions for

$$\lim_{n,d\to\infty} n^{-1}\langle \|\mathbf{X}\boldsymbol{\theta} - \mathbf{X}\langle\boldsymbol{\theta}\rangle\|_2^2 \rangle, \qquad \lim_{n,d\to\infty} n^{-1}\|\mathbf{X}\boldsymbol{\theta}^* - \mathbf{X}\langle\boldsymbol{\theta}\rangle\|_2^2.$$

These fixed-point equations (43) that characterize $\mathrm{mse}$ and $\mathrm{mse}_*$ coincide with those derived via the replica method (with misspecified prior) under a replica-symmetric ansatz, c.f. [30, 31, 65].

We clarify that Theorem 2.5 does not claim that the joint solution $(\mathrm{mse}, \mathrm{mse}_*, \omega, \omega_*)$ of the fixed-point equations (43) is unique. In settings with multiple such fixed points, the theorem pertains to the specific choice of this fixed point that arises from the $t \to \infty$ limit of the DMFT dynamics.

### 2.3.2 Asymptotic MSE and free energy under a posterior LSI

To motivate Definition 2.4, it is illustrative to consider the example of a fixed Gaussian prior $g(\cdot)$, where the Langevin diffusion for $\boldsymbol{\theta}^t$ is a linear Ornstein-Uhlenbeck process. Then $C_\theta, C_\eta, R_\theta, R_\eta$ of the DMFT system may be computed explicitly, as we show in Appendix B, and it is directly checked from their explicit forms that the DMFT system is indeed approximately-TTI.

Generalizing this Gaussian prior example, we consider a setting where the posterior distribution (5) satisfies a log-Sobolev inequality.

**Assumption 2.7.** There exists a constant $C_{\mathrm{LSI}} > 0$ and a $\mathbf{X}$-dependent event $\mathcal{E}(\mathbf{X})$ holding almost surely for all large $n, d$, for which

(a) (LSI for posterior) On $\mathcal{E}(\mathbf{X})$, for all $\mathbf{y} \in \mathbb{R}^n$, the posterior distribution $\mathsf{P}_g(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y})$ satisfies

$$\mathrm{Ent}[f(\boldsymbol{\theta})^2 \mid \mathbf{X}, \mathbf{y}] \leq C_{\mathrm{LSI}} \, \mathbb{E}[\|\nabla f(\boldsymbol{\theta})\|_2^2 \mid \mathbf{X}, \mathbf{y}] \text{ for all } f \in C^1(\mathbb{R}^d). \tag{46}$$

(b) (LSI for larger noise) On $\mathcal{E}(\mathbf{X})$, for every noise variance $\tilde{\sigma}^2 \in [\sigma^2, \infty)$, (46) holds also for the posterior law $\mathsf{P}_g(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y})$ defined with $\tilde{\sigma}^2$ in place of $\sigma^2$ (with a uniform constant $C_{\mathrm{LSI}} > 0$ for all $\tilde{\sigma}^2 \geq \sigma^2$).

For clarity of interpretation, we list in the following proposition three concrete settings in which these LSI conditions hold by currently known techniques. A proof of Proposition 2.8 is given in Appendix C.

**Proposition 2.8.** *Suppose $\mathbf{X}$ satisfies Assumption 2.1(a–b), and $g(\cdot)$ satisfies Assumption 2.2(a). Let $C, r_0, c_0 > 0$ be the constants of Assumption 2.2(a), and define*

$$C_0 = \frac{2.01}{c_0} \exp\Big(\frac{8r_0^2(c_0 + C)^2}{\pi c_0}\Big).$$

*Suppose, in addition, that at least one of the following conditions hold:*

*(a) (global log-concavity) $-(\log g)''(\theta) \geq c_0$ for all $\theta \in \mathbb{R}$, or*

*(b) (high noise) $\sigma^2 > C_0(4\sqrt{\delta}\,\mathbf{1}\{\delta > 1\} + (\sqrt{\delta} + 1)^2\mathbf{1}\{\delta \leq 1\})$, or*

*(c) (large sample size) $\delta > 1$ and $(\sqrt{\delta} - 1)^2 > 4C_0 C\sqrt{\delta}$.*

*Then Assumption 2.7 holds for a constant $C_{\mathrm{LSI}} > 0$ depending only on $\delta, C, r_0, c_0$.*

Under the posterior LSI condition of Assumption 2.7(a), we verify that the solution of the DMFT system must be approximately-TTI in the sense of Definition 2.4.

**Theorem 2.9.** *Consider the dynamics (7) with a fixed prior $g(\cdot)$, and suppose Assumptions 2.1, 2.2(a), and 2.7(a) hold. Then the DMFT system given by Theorem 2.3(a) is approximately-TTI, where the statements of Definition 2.4 hold with $\varepsilon(t) = Ce^{-ct}$ and some constants $C, c > 0$.*

As a consequence, we obtain the following corollary showing that the asymptotic free energy and mean-squared-errors associated to the posterior distribution $\mathsf{P}_g(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y})$ in the linear model (with a possibly misspecified prior) are given by their replica-symmetric predictions, and furthermore the joint empirical distribution of coordinates of $\boldsymbol{\theta}^*$ and a posterior sample $\boldsymbol{\theta} \sim \mathsf{P}_g(\cdot \mid \mathbf{X}, \mathbf{y})$ converges to the preceding law $\mathsf{P}_{g_*,\omega_*;g,\omega}$ in the scalar Gaussian convolution model. (Our analysis for the free energy uses an I-MMSE relation, for which we require the posterior LSI condition of Assumption 2.7(b) for an extended range of noise variances.)

**Corollary 2.10.** *Suppose Assumptions 2.1, 2.2(a), and 2.7(a) hold for dynamics (7) with a fixed prior $g(\cdot)$. Let $\mathsf{P}_g(\mathbf{y} \mid \mathbf{X})$ be the marginal likelihood of $\mathbf{y}$ in (6), let $\langle f(\boldsymbol{\theta}) \rangle$ denote the posterior expectation under $\mathsf{P}_g(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y})$, and define*

$$\begin{aligned}
\mathrm{MSE} &= d^{-1}\langle\|\boldsymbol{\theta} - \langle\boldsymbol{\theta}\rangle\|_2^2\rangle, & \mathrm{MSE}_* &= d^{-1}\|\boldsymbol{\theta}^* - \langle\boldsymbol{\theta}\rangle\|_2^2 \\
\mathrm{YMSE} &= n^{-1}\langle\|\mathbf{X}\boldsymbol{\theta} - \langle\mathbf{X}\boldsymbol{\theta}\rangle\|_2^2\rangle, & \mathrm{YMSE}_* &= n^{-1}\|\mathbf{X}\boldsymbol{\theta}^* - \langle\mathbf{X}\boldsymbol{\theta}\rangle\|_2^2.
\end{aligned} \tag{47}$$

*Let $\mathrm{mse}, \mathrm{mse}_*, \omega, \omega_*, \mathrm{ymse}, \mathrm{ymse}_*$ be as defined by (42–43) for the corresponding (approximately-TTI) DMFT system, let $\mathsf{P}_{g,\omega}(y)$ be the marginal density of $y$ in (39), and let $\mathbb{E}_{g_*,\omega_*}$ denote the expectation over $y = \theta^* + z$ in (37) with $\theta^* \sim g_*$ and $z \sim \mathcal{N}(0, \omega_*^{-1})$.*

*(a) Almost surely,*

$$\lim_{n,d\to\infty} \mathrm{MSE} = \mathrm{mse}, \qquad \lim_{n,d\to\infty} \mathrm{MSE}_* = \mathrm{mse}_*,$$

$$\lim_{n,d\to\infty} \mathrm{YMSE} = \mathrm{ymse}, \qquad \lim_{n,d\to\infty} \mathrm{YMSE}_* = \mathrm{ymse}_*,$$

$$\lim_{n,d\to\infty} \left\langle W_2\left(\frac{1}{d}\sum_{j=1}^d \delta_{(\theta_j^*,\theta_j)}, \mathsf{P}_{g_*,\omega_*;g,\omega}\right)^2 \right\rangle = 0. \tag{48}$$

*(b) If furthermore Assumption 2.7(b) holds, then almost surely,*

$$\lim_{n,d\to\infty} \frac{1}{d}\log \mathsf{P}_g(\mathbf{y}\mid\mathbf{X}) = \mathbb{E}_{g_*,\omega_*}\log\mathsf{P}_{g,\omega}(y) + \frac{1}{2}\left(\delta + \log\frac{2\pi}{\omega} - \delta\log\frac{2\pi\delta}{\omega} + (1-\delta)\frac{\omega}{\omega_*} + \omega\sigma^2\left(\frac{\omega}{\omega_*} - 2\right)\right).$$

As discussed in Remark 2.6, the fixed-point equations characterizing these limits of the mean-squared-error quantities $\mathrm{MSE}, \mathrm{MSE}_*, \mathrm{YMSE}, \mathrm{YMSE}_*$ are those derived via the replica method under an assumption of replica symmetry. One may verify that the limit of the free energy in part (b) agrees also with the replica prediction that was computed in [31, Eq. (20)].

The proof of Theorem 2.5 is given in Section 3, and the proofs of Theorem 2.9 and Corollary 2.10 are given in Section 4.

## 2.4 Convergence of empirical Bayes Langevin dynamics

We now discuss results pertaining to the empirical Bayes Langevin dynamics (9–10) with a data-adaptive evolution of the prior law.

### 2.4.1 A general condition for dimension-free convergence

We impose the following strengthening of Assumption 2.7, ensuring that $\{\alpha^t\}_{t\geq 0}$ of the DMFT solution remains confined to a bounded domain where the posterior log-Sobolev conditions of Assumption 2.7 hold uniformly.

**Assumption 2.11.** Let $\{\alpha^t\}_{t\geq 0}$ be the $\alpha$-component of the DMFT system. There exists a compact subset $S \subset \mathbb{R}^K$ such that

$$\alpha^t \in S \text{ for all } t \geq 0.$$

Furthermore, there exists a (bounded) open neighborhood $O \supset S$ and an $\mathbf{X}$-dependent event $\mathcal{E}(\mathbf{X})$ on which the statements of Assumption 2.7(a–b) hold with a uniform constant $C_{\mathrm{LSI}} > 0$ for every prior $g \in \{g(\cdot,\alpha) : \alpha \in O\}$.

Under this condition, we will show dimension-free convergence of the prior parameter $\{\alpha^t\}_{t\geq 0}$ to a fixed point of the replica-symmetric free energy. To state this result, let us recall the free energy

$$\widehat{F}(\alpha) = -\frac{1}{d}\log\mathsf{P}_{g(\cdot,\alpha)}(\mathbf{y}\mid\mathbf{X})$$

of the linear model from (11), and denote by

$$F(\alpha) = -\mathbb{E}_{g_*,\omega_*}\log\mathsf{P}_{g(\cdot,\alpha),\omega}(Y) - \frac{1}{2}\left(\delta + \log\frac{2\pi}{\omega} - \delta\log\frac{2\pi\delta}{\omega} + (1-\delta)\frac{\omega}{\omega_*} + \omega\sigma^2\left(\frac{\omega}{\omega_*} - 2\right)\right) \tag{49}$$

its asymptotic limit prescribed by Corollary 2.10, both viewed as a function of $\alpha \in O \subset \mathbb{R}^K$. Here, the fixed points $(\omega,\omega_*) \equiv (\omega(\alpha),\omega_*(\alpha))$ implicitly depend on $\alpha$ and are well-defined by Theorem 2.9 for all $\alpha \in O$. Recalling the law $\mathsf{P}_{g_*,\omega_*;g,\omega}(\theta^*,\theta)$ from (40), let us abbreviate this law with $g \equiv g(\cdot,\alpha)$ and fixed points $(\omega(\alpha),\omega_*(\alpha))$ as

$$\mathsf{P}_\alpha \equiv \mathsf{P}_{g_*,\omega_*(\alpha);g(\cdot,\alpha),\omega(\alpha)}. \tag{50}$$

We write $\theta \sim \mathsf{P}_\alpha$ as shorthand for the $\theta$-marginal of $(\theta^*,\theta) \sim \mathsf{P}_\alpha$. We write also $\langle\cdot\rangle_\alpha$ for the expectation under the posterior law $\mathsf{P}_{g(\cdot,\alpha)}(\theta\mid\mathbf{X},\mathbf{y})$ in the linear model. The following lemma strengthens Corollary 2.10(b) to convergence of $F(\alpha)$ and its gradient, uniformly over the compact subset $S \subset O$ containing $\{\alpha^t\}_{t\geq 0}$, and shows also that a true prior parameter $\alpha^* \in O$ must be a global minimizer of $F(\alpha)$.

13

**Lemma 2.12.** *Suppose Assumptions 2.1, 2.2(b), and 2.11 hold, and let $S \subset O \subset \mathbb{R}^K$ be the domains of Assumption 2.11. Then*

*(a) $\widehat{F}(\alpha)$ and $F(\alpha)$ are continuously differentiable on $O$ with gradients*

$$\nabla \widehat{F}(\alpha) = -\left\langle \frac{1}{d} \sum_{j=1}^{d} \nabla_\alpha \log g(\theta_j, \alpha) \right\rangle_\alpha, \qquad \nabla F(\alpha) = -\mathbb{E}_{\theta \sim \mathsf{P}_\alpha} [\nabla_\alpha \log g(\theta, \alpha)]. \tag{51}$$

*(b) Almost surely*

$$\lim_{n,d \to \infty} \sup_{\alpha \in S} |\widehat{F}(\alpha) - F(\alpha)| = 0, \qquad \lim_{n,d \to \infty} \sup_{\alpha \in S} \|\nabla \widehat{F}(\alpha) - \nabla F(\alpha)\|_2 = 0.$$

*(c) If $g_*(\cdot) = g(\cdot, \alpha^*)$ for some $\alpha^* \in O$, then $F(\alpha^*) = \inf_{\alpha \in O} F(\alpha)$.*

We now show that under the uniform LSI condition of Assumption 2.11, the DMFT solution $\{\alpha^t\}_{t \geq 0}$ converges as $t \to \infty$ to a critical point $\alpha^\infty$ of the asymptotic free energy $F(\alpha)$ (with possible additional regularization by $R(\alpha)$). Consequently, for a dimension-independent time horizon $T > 0$ and large system sizes $n, d$, the learned prior parameter $\widehat{\alpha}^T$ will be close to $\alpha^\infty$, and the Langevin sample $\widehat{\boldsymbol{\theta}}^T$ will have entrywise statistics close to those in the scalar Gaussian convolution model described by Theorem 2.5 for the limiting prior $g(\cdot) = g(\cdot, \alpha^\infty)$.

**Theorem 2.13.** *Suppose Assumptions 2.1, 2.2(b), and 2.11 hold. Let $O \subset \mathbb{R}^K$ be as in Assumption 2.11, define $F(\alpha)$ for $\alpha \in O$ by (49), and denote*

$$\mathrm{Crit} = \{\alpha \in S : \nabla F(\alpha) + \nabla R(\alpha) = 0\}.$$

*Consider the empirical Bayes Langevin dynamics (9–10), and let $\{\alpha^t\}_{t \geq 0}$ be the deterministic approximation of $\{\widehat{\alpha}^t\}_{t \geq 0}$ in the solution of the DMFT system in Theorem 2.3(b). Then $\{\alpha^t\}_{t \geq 0}$ satisfies*

$$\lim_{t \to \infty} \mathrm{dist}(\alpha^t, \mathrm{Crit}) = 0.$$

*In particular, if all points of* Crit *are isolated, then there exists a limit*

$$\alpha^\infty = \lim_{t \to \infty} \alpha^t \in \mathrm{Crit}. \tag{52}$$

*Consequently, for any $\varepsilon > 0$, there exists a time horizon $T := T(\varepsilon) > 0$ independent of $n, d$ such that for any fixed $t > T(\varepsilon)$, the solution $\{(\boldsymbol{\theta}^t, \widehat{\alpha}^t)\}_{t \geq 0}$ of (9–10) satisfies almost surely*

$$\limsup_{n,d \to \infty} \|\widehat{\alpha}^t - \alpha^\infty\|_2 < \varepsilon, \qquad \limsup_{n,d \to \infty} W_2\left(\frac{1}{d} \sum_{j=1}^{d} \delta_{(\theta_j^*, \theta_j^t)}, \mathsf{P}_{\alpha^\infty}\right) < \varepsilon. \tag{53}$$

The proof of Theorem 2.13 is given in Section 5.

Supposing that $g_*(\cdot) = g(\cdot, \alpha^*)$ for a true prior parameter $\alpha^* \in O$, in settings where $R(\alpha) = 0$ and the critical point $\alpha^\infty \in \mathrm{Crit}$ of $F(\alpha)$ is unique, Lemma 2.12(c) ensures that $\alpha^\infty = \alpha^*$, and Theorem 2.13 then provides a guarantee for estimation of this true prior parameter as $n, d \to \infty$. In general, $F(\alpha)$ may have multiple critical points. Theorem 2.13 ensures convergence to a point $\alpha^\infty \in \mathrm{Crit}$ that is specified deterministically by the initial conditions of Assumption 2.1(d), and successful learning of $\alpha^*$ may require multiple initializations from different starting values of $\widehat{\alpha}^0$. We discuss both types of settings in the following examples.

### 2.4.2 Examples

We develop some further implications of Theorem 2.13 in a few specific examples of parametric models for $g(\cdot, \alpha)$. We explore also via numerical simulation the convergence of $(\boldsymbol{\theta}^t, \widehat{\alpha}^t)$, the landscape of the replica-symmetric free energy $F(\alpha)$, and the nature of its critical point set Crit in a few settings where a posterior log-Sobolev inequality may not hold.

**Example 2.14.** Consider the Gaussian prior

$$g(\theta, \alpha) = \sqrt{\frac{\omega_0}{2\pi}} \exp\left(-\frac{\omega_0}{2}(\theta - \alpha)^2\right)$$

with varying mean $\alpha \in \mathbb{R}$ and a fixed and known prior variance $\omega_0^{-1}$, and suppose $g_*(\theta) = g(\theta, \alpha^*)$. Consider the empirical Bayes dynamics driven by

$$\mathcal{G}(\alpha, \mathsf{P}) = \mathbb{E}_{\theta \sim \mathsf{P}}[\partial_\alpha \log g(\theta, \alpha)]$$

in (27), with no regularizer (i.e. $R(\alpha) \equiv 0$).

We verify in Section 5.2 that Assumptions 2.2(b) and 2.11 hold for this example, for a subset $O \subset \mathbb{R}$ containing $\alpha^*$. The posterior mean in the Gaussian convolution model (37) is given explicitly by

$$\langle \theta \rangle_{g(\cdot, \alpha), \omega} = \frac{\omega_0}{\omega_0 + \omega}\alpha + \frac{\omega}{\omega_0 + \omega}y.$$

Then the condition $\alpha \in \mathrm{Crit}$ is $0 = \nabla F(\alpha) = \mathbb{E}_{\theta \sim \mathsf{P}_\alpha}[\omega_0(\alpha - \theta)]$, i.e.

$$\alpha = \mathbb{E}_{\theta \sim \mathsf{P}_\alpha}[\theta] = \mathbb{E}_{g_*, \omega_*}[\langle \theta \rangle_{g(\cdot, \alpha), \omega}] = \mathbb{E}_{g_*, \omega}\left[\frac{\omega_0}{\omega_0 + \omega}\alpha + \frac{\omega}{\omega_0 + \omega}y\right] = \frac{\omega_0}{\omega_0 + \omega}\alpha + \frac{\omega}{\omega_0 + \omega}\alpha^*,$$

so Crit consists of the unique critical point $\alpha^*$. Theorem 2.13 then holds with $\alpha^\infty = \alpha^*$, i.e. over a dimension-independent time horizon, $\widehat{\alpha}^t$ converges to $\alpha^*$ (in the limit $n, d \to \infty$ followed by $t \to \infty$ as described in Theorem 2.13), and the empirical distribution of coordinates of the Langevin sample $\boldsymbol{\theta}^t$ converges to that of the posterior distribution for the true prior $\mathcal{N}(\alpha^*, \omega_0^{-1})$. □

**Example 2.15.** Consider more generally a log-concave location prior

$$g(\theta, \alpha) = \exp\left(-f(\theta - \alpha)\right)$$

where $\alpha \in \mathbb{R}$ and $f : \mathbb{R} \to \mathbb{R}$ is a fixed strongly convex function, such that $f$ is thrice continuously-differentiable with Hölder-continuous third derivative, and

$$f'(0) = 0, \qquad C \geq f''(x) \geq c_0, \qquad |f'''(x)| \leq C$$

for some constants $C, c_0 > 0$ and all $x \in \mathbb{R}$. Suppose again $g_*(\theta) = g(\theta, \alpha^*)$, and consider the empirical Bayes dynamics driven by

$$\mathcal{G}(\alpha, \mathsf{P}) = \mathbb{E}_{\theta \sim \mathsf{P}}[\partial_\alpha \log g(\theta, \alpha)]$$

with no regularizer.

We verify in Section 5.2 that Assumptions 2.2(b) and 2.11 hold for this example, for a subset $O \subset \mathbb{R}$ containing $\alpha^*$. Furthermore, we show in Section 5.2 via an adaptation of the Brascamp-Lieb argument of [8, Theorem 3] that $F(\alpha)$ must be strongly convex on $O$. Hence, Crit consists again of the unique critical point $\alpha = \alpha^*$, and Theorem 2.13 holds for $\alpha^\infty = \alpha^*$. □

We next consider two canonical examples where the prior $g(\theta, \alpha)$ is a Gaussian mixture model that is not log-concave in $\theta$, and where the landscape of $F(\alpha)$ is also not necessarily convex in $\alpha$. We will check the uniform log-Sobolev condition of Assumption 2.11 and also characterize analytically the landscape of the free energy $F(\alpha)$ for sufficiently large $\delta = \lim \frac{n}{d}$, and explore by simulation the learning dynamics and free energy landscape for some smaller values of $\delta$.

The sub-level sets of $F(\alpha)$ may not be bounded in these examples. To confine $\{\alpha^t\}_{t \geq 0}$ to a bounded subset of $\mathbb{R}^K$, we introduce an additional regularizer: Fix a radius $D > 0$, and let $\mathcal{B}(D) = \{\alpha \in \mathbb{R}^K : \|\alpha\|_2 < D\}$ be the open ball of radius $D$. For a smooth function $r : [0, \infty) \to [0, \infty)$ having bounded derivatives of all orders and satisfying

$$r(x) = 0 \text{ for all } x \in [0, D], \quad r(x) \geq x - D \text{ for all } x \geq D + 1, \quad r'(x) > 0 \text{ for all } x > D, \tag{54}$$

we fix the regularizer $R : \mathbb{R}^K \to \mathbb{R}$ as

$$R(\alpha) = r(\|\alpha\|_2). \tag{55}$$

15

Note that $R(\alpha) = 0$ for all $\alpha \in \mathcal{B}(D)$, so adding such a regularizer does not change the critical points $\alpha \in \mathrm{Crit} \cap \mathcal{B}(D)$. We show in Proposition 5.2 of Section 5.2 that adding such a regularizer indeed confines the dynamics of $\{\alpha^t\}_{t \geq 0}$ to a bounded domain.

We will study analytically a large-$\delta$ limit under a reparametrization of the noise variance $\sigma^2$ by $s^2 = \sigma^2/\delta$, corresponding to a rescaling of the regression design $\mathbf{X}$ to have entries of variance $1/n$ and a rescaling of the noise $\boldsymbol{\varepsilon}$ to have entries $\mathcal{N}(0, s^2)$. The setting $\delta \to \infty$ with fixed $s^2 > 0$ is a limiting regime in which each coordinate of the posterior distribution of $\boldsymbol{\theta}$ does not contract around its mode, the Bayes-optimal mean-squared-error for estimating $\boldsymbol{\theta}$ remains bounded away from 0, and the landscape of $F(\alpha)$ approaches (up to an additive constant) the log-likelihood landscape in the scalar Gaussian convolution model $y = \theta + z$ where $\theta \sim g(\cdot, \alpha)$ and $z \sim \mathcal{N}(0, s^2)$. We denote by

$$G_{s^2}(\alpha) = -\mathbb{E}_{g_*, s^{-2}}[\log \mathsf{P}_{g(\cdot, \alpha), s^{-2}}(y)] \tag{56}$$

the negative population log-likelihood in this model as a function of the prior parameter $\alpha$, when the true distribution of $y$ is given by $y = \theta^* + z$ with $\theta^* \sim g_*$.

**Proposition 2.16.** *Suppose Assumptions 2.1 and 2.2(b) hold, and the regularizer $R(\alpha)$ is given by (54–55) with $\alpha^0 \in \mathcal{B}(D)$. Fix $s^2 = \sigma^2/\delta$, and define*

$$\mathrm{Crit}_G = \{\alpha \in \mathcal{B}(D) : \nabla G_{s^2}(\alpha) = 0\}.$$

*Then, for any $s^2 > 0$, there exists a constant $\delta_0 := \delta_0(s^2) > 0$ and a function $\iota : [\delta_0, \infty) \to (0, \infty)$ with $\lim_{\delta \to \infty} \iota(\delta) = 0$ such that if $\delta > \delta_0$, then Assumption 2.11 holds. Furthermore,*

1. *Each point of $\mathrm{Crit} \cap \mathcal{B}(D)$ belongs to a ball of radius $\iota(\delta)$ around some point of $\mathrm{Crit}_G$.*

2. *For each point $\alpha \in \mathrm{Crit}_G$ where $\nabla^2 G_{s^2}(\alpha)$ is non-singular, there is exactly one point of $\mathrm{Crit}$ in the ball of radius $\iota(\delta)$ around $\alpha$.*

*In particular, if $g_* = g(\cdot, \alpha^*)$ for some $\alpha^* \in \mathcal{B}(D)$, and if $\alpha^*$ is the unique point of $\mathrm{Crit}_G$ and $\nabla^2 G_{s^2}(\alpha^*)$ is non-singular, then $\alpha^*$ is also the unique point of $\mathrm{Crit} \cap \mathcal{B}(D)$.*

**Example 2.17.** Consider a $K$-component Gaussian mixture prior

$$g(\theta, \alpha) = \sum_{k=1}^{K} p_k \sqrt{\frac{\omega_k}{2\pi}} \exp\left(-\frac{\omega_k}{2}(\theta - \alpha_k)^2\right)$$

with fixed mixture weights $p_1, \ldots, p_K$ and variances $\omega_1^{-1}, \ldots, \omega_K^{-1}$, parametrized by the mixture means $\alpha \in \mathbb{R}^K$. Let us suppose that $g_*(\theta) = g(\theta, \alpha^*)$ for some $\alpha^* \in \mathbb{R}^K$, and the variances $\omega_1^{-1}, \ldots, \omega_K^{-1}$ are distinct. We consider the empirical Bayes dynamics driven by

$$\mathcal{G}(\alpha, \mathsf{P}) = \mathbb{E}_{\theta \sim \mathsf{P}}[\nabla_\alpha \log g(\theta, \alpha)] - \nabla R(\alpha),$$

where $R(\alpha)$ is a regularizer of the form (54–55) for which $\alpha^0, \alpha^* \in \mathcal{B}(D)$.

We verify in Section 5.2 that Assumption 2.2(b) holds. Then, for fixed $s^2 > 0$ and all sufficiently large $\delta$, Proposition 2.16 ensures that the confinement and log-Sobolev conditions of Assumption 2.11 also hold, and the proposition further establishes a 1-to-1 correspondence between the critical points of $F$ and the (non-singular) critical points of $G_{s^2}(\alpha)$ in $\mathcal{B}(D)$. We note that, here, $G_{s^2}(\alpha)$ is the negative population log-likelihood in the Gaussian mixture model

$$\mathsf{P}_{g(\cdot, \alpha), s^{-2}}(y) = \sum_{k=1}^{K} p_k \cdot \frac{1}{\sqrt{2\pi(\omega_k^{-1} + s^2)}} \exp\left(-\frac{1}{2(\omega_k^{-1} + s^2)}(y - \alpha_k)^2\right) \tag{57}$$

having the same mixture means $\alpha \in \mathbb{R}^K$ as the prior, and elevated mixture variances $\omega_k^{-1} + s^2$. The optimization landscape of $G_{s^2}(\alpha)$ is well-studied in the literature, see e.g. [70–73], and in general $G_{s^2}(\alpha)$
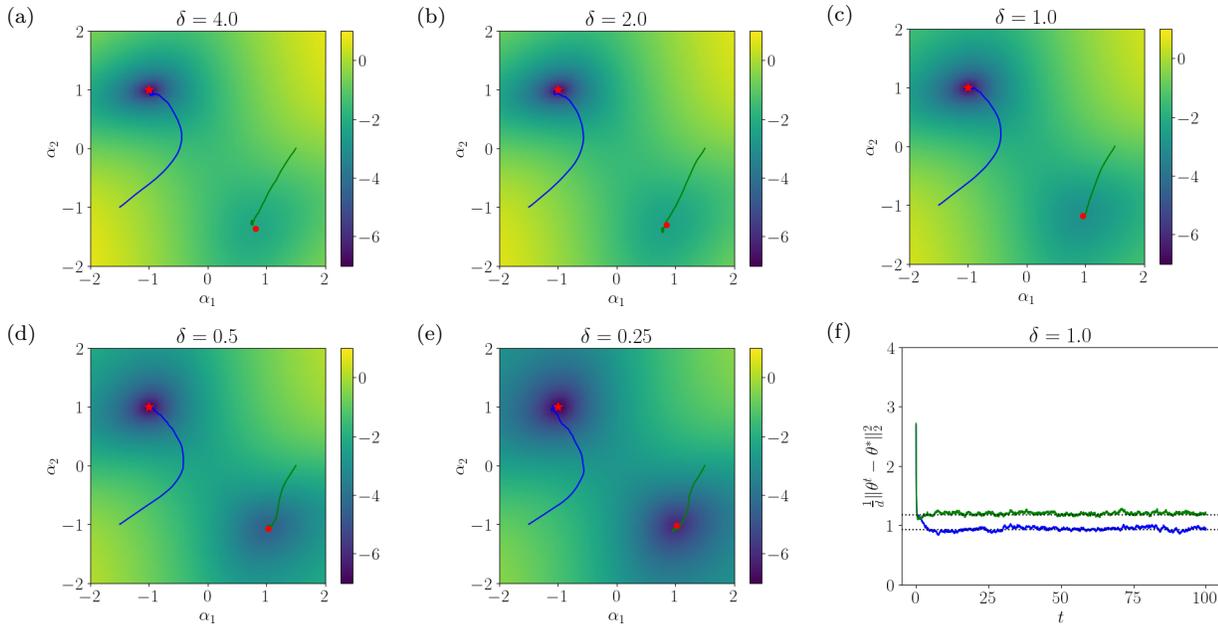
16

Figure 1: Simulations for the Gaussian mixture prior model $\frac{1}{2}\mathcal{N}(\alpha_1, 1) + \frac{1}{2}\mathcal{N}(\alpha_2, 0.25)$ of Example 2.17, with true mixture means $\alpha^* = (1, -1)$ and linear model noise variance $\sigma^2 = \delta s^2$ for $s = 0.5$. Empirical Bayes Langevin dynamics is run for a single instance $(\mathbf{X}, \mathbf{y})$ with $\max(n, d) = 5000$, initialization $\theta_j^0 \overset{iid}{\sim} \mathcal{N}(0, 1)$, and an Euler-Maruyama discretization of the dynamics. (a–e) Landscape of the replica-symmetric free energy $F(\alpha)$ plotted (for visual clarity) as $\log(F(\alpha) - F(\alpha^*) + 10^{-3})$, for $\delta \in \{4, 2, 1, 0.5, 0.25\}$. Two stable fixed points of $0 = \nabla F(\alpha)$ are depicted in red, with star indicating the true parameter $\alpha^* = (-1, 1)$ and circle indicating a second fixed point $\alpha^\dagger$ near $(1, -1)$. Sample paths $\{\widehat{\alpha}^t\}_{t \geq 0}$ from two different initial states $\widehat{\alpha}^0$ are shown in blue and green. (f) Mean-squared-error $\frac{1}{d}\|\boldsymbol{\theta}^t - \boldsymbol{\theta}^*\|_2^2$ across iterations for these same two initial states, at $\delta = 1$. The predicted value for a posterior sample $\boldsymbol{\theta} \sim \mathsf{P}_{g(\cdot, \alpha)}(\cdot \mid \mathbf{X}, \mathbf{y})$ is $\frac{1}{d}\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2^2 \approx \mathrm{mse}(\alpha) + \mathrm{mse}_*(\alpha)$, depicted by dashed lines for $\alpha \in \{\alpha^\dagger, \alpha^*\}$.

may have local minimizers in $\mathcal{B}(D)$ that are different from $\alpha^*$. In such settings, Proposition 2.16 implies that Crit must also have critical points different from $\alpha^*$ for large $\delta$.

We depict in Figure 1 a simulation of the landscape of $F(\alpha)$ and the dynamics (9–10) across a range of values $\delta \in [0.25, 4]$, in a simple setting of $\frac{1}{2}\mathcal{N}(\alpha_1, 1) + \frac{1}{2}\mathcal{N}(\alpha_2, 0.25)$ with $K = 2$ mixture components and true mixture means $\alpha^* = (-1, 1)$. The Almeida-Thouless condition for stability of the replica-symmetric phase was computed in [31, Eq. (25)] to be (in our notation)

$$1 - \frac{\omega^2}{\delta}\mathbb{E}_{g_*, \omega_*}[\mathrm{Var}_{g, \omega}[\theta]^2] \geq 0 \tag{58}$$

where $g(\cdot) = g(\cdot, \alpha)$ and $(\omega, \omega_*) = (\omega(\alpha), \omega_*(\alpha))$. We have verified that this condition holds at each tested $\delta > 0$ and parameter $\alpha \in \mathbb{R}^K$ depicted in Figure 1, and thus we conjecture that the depicted replica-symmetric free energy function $F(\alpha)$ is indeed the correct asymptotic limit of $-\frac{1}{d}\log \mathsf{P}_{g(\cdot, \alpha)}(\mathbf{y} \mid \mathbf{X})$ as $n, d \to \infty$ (even in settings where our assumption of a log-Sobolev inequality for the posterior law may not hold). We observe, not only for large $\delta$ but across a range of values $\delta \in [0.25, 4]$, that the landscape $F(\alpha)$ has two local minimizers, one fixed at the true parameter $\alpha^* = (-1, 1)$ and a second minimizer $\alpha^\dagger$ whose location depends on $\delta$. As $\delta$ decreases, this second minimizer approaches $(1, -1)$ — characterizing a prior law with mixture means matching those of $g_* = g(\cdot, \alpha^*)$ but with the mixture variances reversed — and the free energy difference $F(\alpha^\dagger) - F(\alpha^*)$ approaches 0, indicating that it becomes increasingly difficult to distinguish $\alpha^\dagger$ from the true parameter $\alpha^*$. The dynamics $\{\widehat{\alpha}^t\}_{t \geq 0}$ follow a smooth trajectory to one of $\alpha^\dagger$ or $\alpha^*$, depending on the initial state $\widehat{\alpha}^0$.
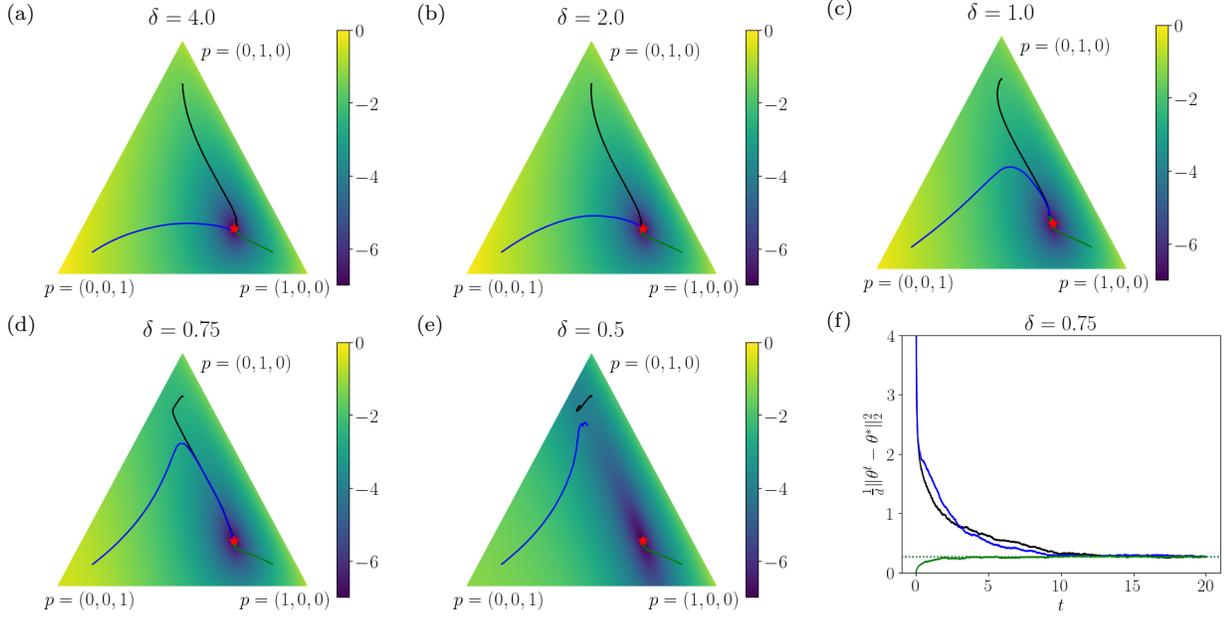
Figure 2: Simulations for the Gaussian mixture prior model $p_1(\alpha)\mathcal{N}(0, 0.04) + p_2(\alpha)\mathcal{N}(0, 1) + p_3(\alpha)\mathcal{N}(0, 25)$ of Example 2.18, with true weights $p(\alpha^*) = (0.6, 0.2, 0.2)$ and linear model noise variance $\sigma^2 = \delta s^2$ for $s = 0.2$. Empirical Bayes Langevin dynamics are run for two initializations $\widehat{\alpha}^0$ with random $\theta_j^0 \overset{iid}{\sim} \mathcal{N}(0, 1)$ (black and blue), and an initialization $\widehat{\alpha}^0$ near $\alpha^*$ with ground truth $\theta_j^0 = \theta_j^*$ (green). The remaining setup is the same as in Figure 1. (a–e) Landscape of the replica-symmetric free energy $F(\alpha)$ for $\delta \in \{4, 2, 1, 0.75, 0.5\}$, plotted as $\log(F(\alpha) - F(\alpha^*) + 10^{-3})$ in the coordinates $p(\alpha)$ on the simplex. The unique critical point $p(\alpha^*)$ is depicted as the red star. Sample paths of $\{p(\widehat{\alpha}^t)\}_{t \geq 0}$ are shown in green, black, and blue. (f) Mean-squared-error $\frac{1}{d}\|\boldsymbol{\theta}^t - \boldsymbol{\theta}^*\|_2^2$ across iterations for these same three initial states, at $\delta = 0.75$. The predicted value of $\mathrm{mse}(\alpha^*) + \mathrm{mse}_*(\alpha^*)$ for a posterior sample is depicted by the dashed line.

**Example 2.18.** Consider a $K + 1$-component Gaussian mixture prior

$$g(\theta, \alpha) = \sum_{k=0}^{K} p_k(\alpha) \sqrt{\frac{\omega_k}{2\pi}} \exp\left(-\frac{\omega_k}{2}(\theta - \mu_k)^2\right), \qquad p_k(\alpha) = \frac{e^{\alpha_k}}{e^{\alpha_0} + \ldots + e^{\alpha_K}}$$

with fixed means $\mu_0, \ldots, \mu_K$ and variances $\omega_0^{-1}, \ldots, \omega_K^{-1}$, parametrized instead by the mixture weights $p_k(\alpha) = e^{\alpha_k}/(e^{\alpha_0} + \ldots + e^{\alpha_K})$. Let us suppose that $g_*(\theta) = g(\theta, \alpha^*)$ for some $\alpha^* \in \mathbb{R}^{K+1}$, and the parameter pairs $(\mu_0, \omega_0), \ldots, (\mu_K, \omega_K)$ are distinct. We again consider the dynamics driven by

$$\mathcal{G}(\alpha, \mathsf{P}) = \mathbb{E}_{\theta \sim \mathsf{P}}[\nabla_\alpha \log g(\theta, \alpha)] - \nabla R(\alpha),$$

where $R(\alpha)$ is a regularizer of the form (54–55) such that $\alpha^0, \alpha^* \in \mathcal{B}(D)$. This parametrization is over-parametrized by a single parameter — however, defining the $K$-dimensional linear subspace $E = \{\alpha \in \mathbb{R}^{K+1} : \alpha_0 + \ldots + \alpha_K = 0\}$, a direct calculation (c.f. Section 5.2) verifies that $\nabla_\alpha \log g(\theta, \alpha) \in E$ and $\nabla R(\alpha) \in E$ if $\alpha \in E$. Thus, initializing $\widehat{\alpha}^0 \in E$ ensures $\widehat{\alpha}^t \in E$ for all $t \geq 0$, and we may apply our preceding results upon identifying $E$ isometrically with $\mathbb{R}^K$.

We verify in Section 5.2 that Assumption 2.2(b) holds. Then again for fixed $s^2 > 0$ and all large $\delta$, Proposition 2.16 ensures that Assumption 2.11 also holds, and there is a 1-to-1 correspondence between the critical points of $F$ and $G_{s^2}(\alpha)$ on $\mathcal{B}(D)$. Here, $G_{s^2}(\alpha)$ is the negative population log-likelihood of the Gaussian mixture model

$$\mathsf{P}_{g(\cdot, \alpha), s^{-2}}(y) = \sum_{k=0}^{K} p_k(\alpha) \sqrt{\frac{1}{2\pi(\omega_k^{-1} + s^2)}} \exp\left(-\frac{1}{2(\omega_k^{-1} + s^2)}(y - \mu_k)^2\right). \tag{59}$$

18

Letting $S = \{(p_0, \ldots, p_K) : p_0 + \ldots + p_K = 1, \, p_0, \ldots, p_K > 0\}$ be the open probability simplex, the mapping $\alpha \in E \mapsto p(\alpha) \in S$ is a 1-to-1 smooth parametrization with smooth inverse, and the function $G_{s^2}$ is strictly convex in the parametrization by $(p_0, \ldots, p_K) \in S$. Thus $p^* = p(\alpha^*) \in S$ is the unique critical point where $\nabla_p G_{s^2} = 0$, and the Hessian $\nabla_p^2 G_{s^2}$ is nonsingular at $p^*$. This implies that $\alpha^* \in E$ is also the unique critical point where $\nabla_\alpha G_{s^2} = 0$, and $\nabla_\alpha^2 G_{s^2}$ is also non-singular at $\alpha^*$. So for large $\delta$, Proposition 2.16 ensures that $\alpha^*$ must be the unique point of $\mathrm{Crit} \cap \mathcal{B}(D)$.

Figure 2 depicts the simulated landscape of $F(\alpha)$ and dynamics (9–10) in a scaled-mixture-of-normals model $p_1(\alpha)\mathcal{N}(0, 0.04) + p_2(\alpha)\mathcal{N}(0, 1) + p_3(\alpha)\mathcal{N}(0, 25)$ with all components having mean 0, across a range of values $\delta \in [0.5, 4]$, and with true weights $p(\alpha^*) = (0.6, 0.2, 0.2)$. (We have again verified that the Almeida-Thouless stability condition (58) holds at each depicted $\delta > 0$ and parameter value $\alpha \in \mathbb{R}^K$ in these figures.) We observe for all tested values $\delta \in [0.5, 4]$ that $\alpha^*$ is the unique local minimizer and critical point of $F(\alpha)$. However, as $\delta$ decreases, the landscape of $F(\alpha)$ flattens around $\alpha^*$ along a direction representing a family of priors $g(\,\cdot\,, \alpha)$ having the same first two moments as $g(\,\cdot\,, \alpha^*)$, reflecting that the problem of learning $g(\,\cdot\,, \alpha^*)$ beyond its second moment becomes increasingly ill-conditioned. The learned parameter $\{\widehat{\alpha}^t\}_{t \geq 0}$ successfully converges to $\alpha^*$ from several different initial states $\widehat{\alpha}^0$ when $\delta \geq 0.75$, with mixing of Langevin dynamics becoming increasingly slower as $\delta$ decreases. For $\delta = 0.5$, the learned parameter $\{\widehat{\alpha}^t\}_{t \geq 0}$ fails to converge to $\alpha^*$ under the tested time horizon from random initializations of $\boldsymbol{\theta}^0$, but does converge to $\alpha^*$ under a ground-truth initialization $\boldsymbol{\theta}^0 = \boldsymbol{\theta}^*$ and $\widehat{\alpha}^0$ close to $\alpha^*$.

# 3 Analysis of approximately-TTI DMFT systems

In this section, we prove Theorem 2.5 on the equilibrium properties of the solution to the DMFT equations under an assumption of approximate time-translation-invariance (from an out-of-equilibrium initialization). We assume throughout this section that Assumptions 2.1 and 2.2(a) hold, and that the solution to the DMFT system in Theorem 2.3(a) approximating the dynamics (7) with the fixed prior $g(\cdot)$ is approximately-TTI in the sense of Definition 2.4. We denote by $\{\theta^t\}_{t \geq 0}$, $\{\eta^t\}_{t \geq 0}$, and $C_\theta, C_\eta, R_\theta, R_\eta$ the components of this DMFT solution.

## 3.1 Analysis of $\theta$-equation

We first derive, from analysis of the evolution (23) for $\{\theta^t\}_{t \geq 0}$, a representation of $c_\theta^{\mathrm{tti}}(0), c_\theta^{\mathrm{tti}}(\infty), c_\theta(*)$ in terms of $c_\eta^{\mathrm{tti}}(0), c_\eta^{\mathrm{tti}}(\infty)$, assuming a condition $c_\eta^{\mathrm{tti}}(0) - c_\eta^{\mathrm{tti}}(\infty) < \delta/\sigma^2$ which ensures long-time stability of $\{\theta^t\}_{t \geq 0}$ under (23). This condition will be checked in our subsequent analysis of the evolution of $\{\eta^t\}_{t \geq 0}$.

**Lemma 3.1.** *Suppose $c_\eta^{\mathrm{tti}}(0) - c_\eta^{\mathrm{tti}}(\infty) < \delta/\sigma^2$. Set $\omega = \delta/\sigma^2 - (c_\eta^{\mathrm{tti}}(0) - c_\eta^{\mathrm{tti}}(\infty))$ and $\omega_* = \omega^2/c_\eta^{\mathrm{tti}}(\infty)$. Then*

$$c_\theta^{\mathrm{tti}}(0) = \mathbb{E}_{g_*, \omega_*} \langle \theta^2 \rangle_{g, \omega}, \quad c_\theta^{\mathrm{tti}}(\infty) = \mathbb{E}_{g_*, \omega_*} \langle \theta \rangle_{g, \omega}^2, \quad c_\theta(*) = \mathbb{E}_{g_*, \omega_*}[\langle \theta \rangle_{g, \omega} \theta^*]. \tag{60}$$

The main idea of the proof is to apply the explicit form of $c_\theta^{\mathrm{tti}}$ in (33) together with its fluctuation dissipation relation with $r_\theta^{\mathrm{tti}}$ in (36) to approximate $C_\theta, R_\theta$ at large times by correlation and response functions $C_\theta^{(M)}, R_\theta^{(M)}$ that admit an interpretation as the effect of marginalization over auxiliary variables $(x_1^t, \ldots, x_M^t)$ in a *Markovian* joint evolution of $(\theta^t, x_1^t, \ldots, x_M^t)$ conditional on $\theta^*$. In contrast to the original high-dimensional dynamics of $\{\boldsymbol{\theta}^t\}_{t \geq 0}$ in $\mathbb{R}^d$, here $M$ does not depend on $(n, d)$, and the dynamics of $\{x_1^t, \ldots, x_M^t\}$ will be decoupled given $\{\theta^t\}_{t \geq 0}$. This decoupling allows us to provide a simple explicit form for the $\theta$-marginal of the stationary distribution of $(\theta, x_1, \ldots, x_M)$ conditional on $\theta^*$, which in the limit $M \to \infty$ will match the conditional distribution $\theta \mid \theta^*$ under the limit law $\mathsf{P}_{g_*, \omega_*; g, \omega}(\theta, \theta^*)$.

To implement this idea, we will exhibit a coupling of the processes $\{\theta^t\}_{t \geq 0}$ driven by $C_\theta, R_\theta$ and $\{\theta_{M, T_0}^t\}_{t \geq 0}$ driven by $C_\theta^{(M)}, R_\theta^{(M)}$ from time $T_0$ onwards, and then analyze the convergence of $\{\theta_{M, T_0}^t\}_{t \geq 0}$ under the equivalent Markovian representation of its dynamics. The main technical challenge is to ensure either that the discretization error $\varepsilon(M)$ obtained by approximating $C_\theta, R_\theta$ by $C_\theta^{(M)}, R_\theta^{(M)}$ does not compound exponentially over time, or that the convergence time of $\{\theta_{M, T_0}^t\}_{t \geq 0}$ in the equivalent Markovian dynamics is independent of the approximation dimension $M$. We will take the first approach here, by adapting ideas around sticky and reflection couplings developed in [74, 75] to a setting of non-Markovian DMFT dynamics for $\{\theta^t\}_{t \geq 0}$ and $\{\theta_{M, T_0}^t\}_{t \geq 0}$.

### 3.1.1 Comparison with an auxiliary process

Let us fix a positive integer $M$ and define two sequences $\{a_m\}_{m=0}^M$ and $\{c_m\}_{m=1}^M$ by

$$a_m = \iota + \frac{m}{\sqrt{M}} \text{ for } m = 0, \ldots, M, \qquad \frac{c_m^2}{a_m} = \mu_\eta([a_{m-1}, a_m)) \text{ for } m = 1, \ldots, M \tag{61}$$

where $\mu_\eta$ is given in Definition 2.4. We set

$$R_\eta^{(M)}(\tau) = \sum_{m=1}^M c_m^2 e^{-a_m \tau}, \qquad C_\eta^{(M)}(t, s) = \sum_{m=1}^M \frac{c_m^2}{a_m}(e^{-a_m|t-s|} - e^{-a_m(t+s)}) + c_\eta^{\mathrm{tti}}(\infty).$$

A direct calculation of the covariance shows that

$$C_\eta^{(M)}(t, s) = \mathbb{E}[u_M^t u_M^s] \quad \text{for} \quad u_M^t = z + \sum_{m=1}^M c_m \int_0^t e^{-a_m(t-s)} \sqrt{2}\, \mathrm{d}b_m^s, \tag{62}$$

where $z \sim \mathcal{N}(0, c_\eta^{\mathrm{tti}}(\infty))$ and $\{b_1^t\}_{t\geq 0}, \ldots, \{b_M^t\}_{t\geq 0}$ are standard Brownian motions independent of each other and of $z$. In particular, $C_\eta^{(M)}(t, s)$ is a positive-semidefinite covariance kernel on $[0, \infty)$.

For convenience, let us set

$$U(\theta, \theta^*) = -\frac{\delta}{\sigma^2}(\theta - \theta^*) + (\log g)'(\theta),$$

so the DMFT equation (23) reads

$$\mathrm{d}\theta^t = \left[U(\theta^t, \theta^*) + \int_0^t R_\eta(t, s)(\theta^s - \theta^*)\mathrm{d}s + u^t\right]\mathrm{d}t + \sqrt{2}\, \mathrm{d}b^t. \tag{63}$$

Let $\{u_M^t\}_{t\geq 0}$ be a centered Gaussian process with covariance kernel $C_\eta^{(M)}$, defined in the probability space of $\{u^t, \theta^t\}_{t\geq 0}$ and independent of $\theta^*$. Fixing a time $T_0 > 0$, let $\{\tilde{b}^t\}_{t\geq T_0}$ be a standard Brownian motion initialized at $\tilde{b}^{T_0} = 0$, independent of $\{u^t\}_{t\geq 0}$, $\theta^*$, and $\{\theta^t\}_{t\in[0,T]}$. We consider an auxiliary process $\{\theta_{M,T_0}^t\}_{t\geq 0}$ defined by

$$\theta_{M,T_0}^t = \theta^t \text{ for } t \in [0, T_0],$$
$$\mathrm{d}\theta_{M,T_0}^t = \left[U(\theta_{M,T_0}^t, \theta^*) + \int_0^t R_\eta^{(M)}(t - s)(\theta_{M,T_0}^s - \theta^*)\mathrm{d}s + u_M^t\right]\mathrm{d}t + \sqrt{2}\, \mathrm{d}\tilde{b}^t \text{ for } t \geq T_0. \tag{64}$$

We proceed to construct a coupling of $\{u_M^t\}_{t\geq 0}$ with $\{u^t\}_{t\geq 0}$ and of $\{\tilde{b}^t\}_{t\geq T_0}$ with $\{b^t - b^{T_0}\}_{t\geq T_0}$ defining the DMFT solution $\{\theta^t\}_{t\geq 0}$, to yield a coupling of $\{\theta_{M,T_0}^t\}_{t\geq T_0}$ with $\{\theta^t\}_{t\geq T_0}$.

**Lemma 3.2.** *For any $M, T_0, T > 0$, there exists a coupling of $\{u_M^t\}_{t\geq 0}$ and $\{u^t\}_{t\geq 0}$ such that*

$$\sup_{t\in[T_0, T_0+T]} \mathbb{E}(u_M^t - u^t)^2 \leq \varepsilon(M) + \sqrt{T}\, \varepsilon(T_0),$$

*where $\varepsilon(M)$ does not depend on $T_0, T$ and $\varepsilon(T_0)$ does not depend on $M, T$, and $\lim_{M\to\infty} \varepsilon(M) = 0$ and $\lim_{T_0\to\infty} \varepsilon(T_0) = 0$.*

*Proof.* Define the covariance kernel $C_\eta^{(\infty)}(t, s) = \int_\iota^\infty \left(e^{-a|t-s|} - e^{-a(t+s)}\right)\mu_\eta(\mathrm{d}a) + c_\eta^{\mathrm{tti}}(\infty)$ representing the $M \to \infty$ limit of (62). We will couple Gaussian processes with covariance kernels $(C_\eta^{(M)}, C_\eta^{(\infty)})$ and with $(C_\eta^{(\infty)}, C_\eta)$ respectively.

*Coupling of $(C_\eta^{(M)}, C_\eta^{(\infty)})$.* Let $M' > M$ be any positive integer for which $\sqrt{M'}$ is an integer multiple of $\sqrt{M}$, and let $\{\tilde{a}_m\}_{m=0}^{M'}$ and $\{\tilde{c}_m\}_{m=1}^{M'}$ be the sequences as defined above with $M'$ in place of $M$. Note then that the grid points $\{a_j\}_{j=0}^M$ are a subset of the grid points $\{\tilde{a}_i\}_{i=0}^{M'}$. Let

$$u_{M'}^t = z + \sum_{i=1}^{M'} \tilde{c}_i \int_0^t e^{-\tilde{a}_i(t-s)} \sqrt{2}\, \mathrm{d}\tilde{b}_i^s \tag{65}$$

20

where $z \sim \mathcal{N}(0, c_\eta^{\mathrm{tti}}(\infty))$ and $\{\tilde{b}_1^t\}_{t\geq 0}, \dots \{\tilde{b}_{M'}^t\}_{t\geq 0}$ are standard Brownian motions independent of each other and of $z$. Then (62) shows that $\{u_{M'}^t\}_{t\geq 0}$ has covariance $C_\eta^{(M')}$. Now, for each $j = 1, \dots, M$, let

$$I_j = \{i : a_{j-1} < \tilde{a}_i \leq a_j\}, \qquad b_j^t = \sum_{i \in I_j} \tilde{c}_i \tilde{b}_i^t \Big/ \sqrt{\sum_{i \in I_j} \tilde{c}_i^2}$$

and set

$$u_M^t = z + \sum_{j=1}^M c_j \int_0^t e^{-a_j(t-s)} \sqrt{2}\, \mathrm{d}b_j^s.$$

Then $\{b_1^t\}_{t\geq 0}, \dots, \{b_M^t\}_{t\geq 0}$ are standard Brownian motions independent of each other and of $z$, so (62) shows also that $\{u^t\}_{t\geq 0}$ is a Gaussian process with covariance $C_\eta^{(M)}$.

We may now bound

$$\mathbb{E}[(u_M^t - u_{M'}^t)^2] \leq 4\, \mathbb{E}\Big[\Big(\sum_{i:\tilde{a}_i > a_M} \tilde{c}_i \int_0^t e^{-\tilde{a}_i(t-s)} \mathrm{d}\tilde{b}_i^s\Big)^2\Big]$$

$$+ 4\, \mathbb{E}\Big[\Big(\sum_{j=1}^M \sum_{i \in I_j} \tilde{c}_i \int_0^t e^{-\tilde{a}_i(t-s)} \mathrm{d}\tilde{b}_s^i - \sum_{j=1}^M c_j \int_0^t e^{-a_j(t-s)} \mathrm{d}b_j^s\Big)^2\Big].$$

Since $a_M = \iota + \sqrt{M}$, the first term equals $\sum_{i:\tilde{a}_i > \iota + \sqrt{M}} \tilde{c}_i^2/\tilde{a}_i = \sum_{i:\tilde{a}_i > \iota + \sqrt{M}} \mu_\eta([\tilde{a}_{i-1}, \tilde{a}_i))$, which is at most some $\varepsilon_1(M)$ satisfying $\lim_{M\to\infty} \varepsilon_1(M) = 0$, by finiteness of the measure $\mu_\eta$. The second term is bounded as

$$\mathbb{E}\Big[\Big(\sum_{j=1}^M \sum_{i \in I_j} \tilde{c}_i \int_0^t e^{-\tilde{a}_i(t-s)} \mathrm{d}\tilde{b}_i^s - \sum_{j=1}^M c_j \int_0^t e^{-a_j(t-s)} \mathrm{d}b_j^s\Big)^2\Big]$$

$$= \mathbb{E}\Big[\Big(\sum_{j=1}^M \sum_{i \in I_j} \int_0^t \Big(\tilde{c}_i e^{-\tilde{a}_i(t-s)} - \frac{c_j \tilde{c}_i}{\sqrt{\sum_{\ell \in I_j} \tilde{c}_\ell^2}} e^{-a_j(t-s)}\Big) \mathrm{d}\tilde{b}_i^s\Big)^2\Big]$$

$$= \sum_{j=1}^M \sum_{i \in I_j} \int_0^t \Big(\tilde{c}_i e^{-\tilde{a}_i s} - \frac{c_j \tilde{c}_i}{\sqrt{\sum_{\ell \in I_j} \tilde{c}_\ell^2}} e^{-a_j s}\Big)^2 \mathrm{d}s \leq 2(\mathrm{I} + \mathrm{II}),$$

where

$$\mathrm{I} = \sum_{j=1}^M \sum_{i \in I_j} \int_0^t \tilde{c}_i^2 (e^{-\tilde{a}_i s} - e^{-a_j s})^2 \mathrm{d}s, \qquad \mathrm{II} = \sum_{j=1}^M \sum_{i \in I_j} \int_0^t \tilde{c}_i^2 \Big(1 - \frac{c_j}{\sqrt{\sum_{\ell \in I_j} \tilde{c}_\ell^2}}\Big)^2 e^{-2a_j s} \mathrm{d}s.$$

Let $\Delta = 1/\sqrt{M}$ be the spacing of $\{a_j\}_{j=0}^M$. Then, since $|\tilde{a}_i - a_j| \leq \Delta$ and $\tilde{a}_i \leq a_j$ for all $i \in I_j$,

$$\mathrm{I} \leq \sum_{j=1}^M \sum_{i \in I_j} \int_0^t \tilde{c}_i^2 e^{-2a_{j-1}s} s^2 \Delta^2 \mathrm{d}s \leq \sum_{j=1}^M \int_0^t \Big(\sum_{i \in I_j} \frac{\tilde{c}_i^2}{\tilde{a}_i}\Big) a_j e^{-2a_{j-1}s} s^2 \Delta^2 \mathrm{d}s \stackrel{(*)}{=} \sum_{j=1}^M c_j^2 \int_0^t e^{-2a_{j-1}s} s^2 \Delta^2 \mathrm{d}s$$

where we use $\sum_{i \in I_j} \tilde{c}_i^2/\tilde{a}_i = \sum_{i \in I_j} \mu_\eta([\tilde{a}_{i-1}, \tilde{a}_i)) = \mu_\eta([a_{j-1}, a_j)) = c_j^2/a_j$ in $(*)$. Evaluating this integral, for an absolute constant $C > 0$,

$$\mathrm{I} \leq C\Delta^2 \sum_{j=1}^M \frac{c_j^2}{a_{j-1}^3} \leq \frac{C\Delta^2}{\iota^2} \sum_{j=1}^M \frac{c_j^2}{a_j} \leq \frac{C\Delta^2}{\iota^2} \mu_\eta([\iota, \infty)) \leq \varepsilon_2(M)$$

where $\lim_{M\to\infty} \varepsilon_2(M) = 0$. For II, since $c_j^2 = \sum_{\ell \in I_j} \frac{a_j}{\tilde{a}_\ell} \tilde{c}_\ell^2$, we have $|c_j^2 - \sum_{\ell \in I_j} \tilde{c}_\ell^2| \leq \Delta \sum_{\ell \in I_j} \tilde{c}_\ell^2/\tilde{a}_\ell = \Delta c_j^2/a_j$, and hence

$$\mathrm{II} \leq \sum_{j=1}^M \sum_{i \in I_j} \frac{\tilde{c}_i^2}{2a_j} \frac{(c_j - \sqrt{\sum_{\ell \in I_j} \tilde{c}_\ell^2})^2}{\sum_{\ell \in I_j} \tilde{c}_\ell^2} \leq \sum_{j=1}^M \frac{(c_j^2 - \sum_{\ell \in I_j} \tilde{c}_\ell^2)^2}{2a_j c_j^2} \leq \frac{\Delta^2}{2} \sum_{j=1}^M \frac{c_j^2}{a_j^3} \leq \frac{\Delta^2}{2\iota^2} \sum_{j=1}^M \frac{c_j^2}{a_j} \leq \varepsilon_3(M)$$

21

where $\lim_{M\to\infty}\varepsilon_3(M)=0$. In summary, we have shown that $\sup_{t\geq 0}\mathbb{E}[(u_M^t-u_{M'}^t)^2]\leq\varepsilon(M)$ for some $\varepsilon(M)\to 0$ as $M\to\infty$.

Now note that for any fixed $T_0$ and $T$, $\{u_{M'}^t\}_{t\in[T_0,T_0+T]}$ has covariance kernel $C_\eta^{(M')}$ that converges uniformly to $C_\eta^{(\infty)}$ over $[T_0,T_0+T]$ as $M'\to\infty$. It is direct to check from its definitions that $C_\eta^{(\infty)}$ satisfies the condition (235) of Lemma D.1. So by Lemma D.1, there exists a coupling of $\{u_{M'}^t\}_{t\in[T_0,T_0+T]}$ and a Gaussian process $\{u_\infty^t\}_{t\in[T_0,T_0+T]}$ with covariance $\{C_\eta^{(\infty)}(t,s)\}_{s,t\in[T_0,T_0+T]}$ such that $\sup_{t\in[T_0,T_0+T]}\mathbb{E}(u_{M'}^t-u_\infty^t)^2\to 0$ as $M'\to\infty$. Combining this with the above bound $\sup_{t\geq 0}\mathbb{E}[(u_M^t-u_{M'}^t)^2]\leq\varepsilon(M)$ and taking $M'\to\infty$ shows $\sup_{t\in[T_0,T_0+T]}\mathbb{E}(u_M^t-u_\infty^t)^2\leq\varepsilon(M)$.

*Coupling of $(C_\eta^{(\infty)},C_\eta)$.* By the approximation (31) for $C_\eta$ in Definition 2.4, we have for any $t\geq s\geq 0$ that

$$|C_\eta(t,s)-C_\eta^{(\infty)}(t,s)|\leq\varepsilon(s)+\int_\iota^\infty e^{-a(t+s)}\mathrm{d}\mu_\eta(a),$$

so there exists a (different) function $\varepsilon(T_0)$ with $\lim_{T_0\to\infty}\varepsilon(T_0)=0$ such that

$$\sup_{s,t\in[T_0,T_0+T]}|C_\eta(t,s)-C_\eta^{(\infty)}(t,s)|\leq\varepsilon(T_0).$$

Here $C_\eta^{(\infty)}$ satisfies (235) for a constant $C_0>0$ depending only on $\mu_\eta$, so by Lemma D.1, there exists a coupling of $\{u_\infty^t\}_{t\in[T_0,T_0+T]}$ with $\{u^t\}_{t\in[T_0,T_0+T]}$, the latter having covariance $\{C_\eta(t,s)\}_{s,t\in[T_0,T_0+T]}$, for which $\sup_{t\in[T_0,T_0+T]}\mathbb{E}(u_\infty^t-u^t)^2\leq C(\sqrt{T\,\varepsilon(T_0)}+\varepsilon(T_0))$ for a constant $C>0$.

Combining these two couplings yields a coupling of $\{u_M^t\}_{t\in[T_0,T_0+T]}$ with $\{u^t\}_{t\in[T_0,T_0+T]}$ such that $\sup_{t\in[T_0,T_0+T]}\mathbb{E}(u_M^t-u^t)^2\leq\varepsilon(M)+C(\sqrt{T\,\varepsilon(T_0)}+\varepsilon(T_0))$, and extending this arbitrarily to a full coupling of $\{u_M^t\}_{t\geq 0}$ and $\{u^t\}_{t\geq 0}$ and adjusting the value of $\varepsilon(T_0)$ shows the lemma. $\qquad\square$

**Lemma 3.3.** *Suppose $c_\eta^{\mathrm{tti}}(0)-c_\eta^{\mathrm{tti}}(\infty)<\delta/\sigma^2$. Then for any $M,T_0,T>0$, there exists a coupling of the processes $\{\theta^t\}_{t\geq 0}$ and $\{\theta_{M,T_0}^t\}_{t\geq 0}$ defined by (63) and (64) such that*

$$\sup_{t\in[0,T_0+T]}\mathbb{E}|\theta^t-\theta_{M,T_0}^t|\leq\varepsilon(M)+\sqrt{T}\,\varepsilon(T_0),$$

*where $\varepsilon(M)$ does not depend on $T_0,T$ and $\varepsilon(T_0)$ does not depend on $M,T$, and $\lim_{M\to\infty}\varepsilon(M)=0$ and $\lim_{T_0\to\infty}\varepsilon(T_0)=0$.*

*Proof.* To ease notation, let us write $\tilde{\theta}^t=\theta_{M,T_0}^t$ and $\tilde{u}^t=u_M^t$. We couple $\{u^t\}_{t\geq 0}$ and $\{\tilde{u}^t\}_{t\geq 0}$ according to Lemma 3.2. By definition, $\{\theta^t\}_{t\in[0,T_0]}$ and $\{\tilde{\theta}^t\}_{t\in[0,T_0]}$ coincide up to time $T_0$.

To construct the coupling of $\theta^t$ and $\tilde{\theta}^t$ for times $t\in[T_0,T_0+T]$, we adapt the ideas of [74,75]: Fix some $\varepsilon>0$, and let $h:[0,\infty)\to[0,1]$ be a function such that $h(0)=0$, $h(x)>0$ for $x>0$, $h(x)=1$ for $x\geq\varepsilon$, and both $x\mapsto h(x)$ and $x\mapsto\sqrt{1-h(x)^2}$ are Lipschitz. Let $\{b^t\}_{t\geq T_0}$ and $\{\tilde{b}^t\}_{t\geq T_0}$ be two standard Brownian motions initialized at $b^{T_0}=\tilde{b}^{T_0}=0$, independent of each other and of $\{u^t\}_{t\geq 0}$, $\{\tilde{u}^t\}_{t\geq 0}$, $\theta^*$, and $\{\theta^t\}_{t\in[0,T_0]}$. We define a coupling of $\{\theta^t\}_{t\geq T_0}$ and $\{\tilde{\theta}^t\}_{t\geq T_0}$ by the joint evolutions, for $t\geq T_0$,

$$\mathrm{d}\theta^t=\Big[U(\theta^t,\theta^*)+\int_0^t R_\eta(t,s)(\theta^s-\theta^*)\mathrm{d}s+u^t\Big]\mathrm{d}t+h(|\theta^t-\tilde{\theta}^t|)\sqrt{2}\,\mathrm{d}b^t+\sqrt{2(1-h(|\theta^t-\tilde{\theta}^t|)^2)}\,\mathrm{d}\tilde{b}^t,$$

$$\mathrm{d}\tilde{\theta}^t=\Big[U(\tilde{\theta}^t,\theta^*)+\int_0^t R_\eta^{(M)}(t-s)(\tilde{\theta}^s-\theta^*)\mathrm{d}s+\tilde{u}^t\Big]\mathrm{d}t-h(|\theta^t-\tilde{\theta}^t|)\sqrt{2}\,\mathrm{d}b^t+\sqrt{2(1-h(|\theta^t-\tilde{\theta}^t|)^2)}\,\mathrm{d}\tilde{b}^t.$$

Thus the coupling of the Brownian motions defining these processes is by reflection at times $t\geq T_0$ where $|\theta^t-\tilde{\theta}^t|\geq\varepsilon$, and it transitions to a synchronous coupling as $|\theta^t-\tilde{\theta}^t|\to 0$. Lévy's characterization of Brownian motion shows that the resulting marginal laws of $\{\theta^t\}_{t\geq T_0}$ and $\{\tilde{\theta}^t\}_{t\geq T_0}$ indeed coincide with those of (63) and (64).

22

Let us write as shorthand

$$\xi^t = \theta^t - \tilde{\theta}^t$$

$$v^t = U(\theta^t, \theta^*) + \int_0^t R_\eta(t, s)(\theta^s - \theta^*)\mathrm{d}s + u^t$$

$$\tilde{v}^t = U(\tilde{\theta}^t, \theta^*) + \int_0^t R_\eta^{(M)}(t - s)(\tilde{\theta}^s - \theta^*)\mathrm{d}s + \tilde{u}^t.$$

We derive a SDE for $|\xi^t|$ that is analogous to [75, Eq. (66)]: For any $t \geq T_0$, since $\mathrm{d}\xi^t = (v^t - \tilde{v}^t)\mathrm{d}t + 2\sqrt{2}h(|\xi^t|)\mathrm{d}b^t$, Itô's formula yields

$$\mathrm{d}(\xi^t)^2 = 2\xi^t[(v^t - \tilde{v}^t)\mathrm{d}t + 2\sqrt{2}h(|\xi^t|)\mathrm{d}b^t] + 8h(|\xi^t|)^2\mathrm{d}t.$$

For a small constant $\beta > 0$, let $S_\beta : [0, \infty) \to [0, \infty)$ be a twice continuously-differentiable approximation to the square root, satisfying $S_\beta(x) = \sqrt{x}$ for $x \geq \beta$, $\sup_{0 \leq x \leq \beta} |S_\beta(x)| \leq C$, $\sup_{0 \leq x \leq \beta} |S'_\beta(x)| \leq C\beta^{-1/2}$, and $\sup_{0 \leq x \leq \beta} |S''_\beta(x)| \leq C\beta^{-3/2}$ for a universal constant $C > 0$. (A specific construction is given in [75, Eq. (68)].) Then again by Itô's formula, for any $t \geq T_0$,

$$\mathrm{d}S_\beta((\xi^t)^2) = S'_\beta((\xi^t)^2)\Big[2\xi^t(v^t - \tilde{v}^t)\mathrm{d}t + 4\sqrt{2}\,\xi^t h(|\xi^t|)\mathrm{d}b^t + 8h(|\xi^t|)^2\mathrm{d}t\Big] + 16S''_\beta((\xi^t)^2)(\xi^t)^2 h(|\xi^t|)^2\,\mathrm{d}t. \quad (66)$$

We may take the limit $\beta \to 0$ via a dominated convergence argument: Applying $S'_\beta(x) = x^{-1/2}/2$ for $x \geq \beta$ and the bound $|S'_\beta(x)| \leq C\beta^{-1/2}$ for $x < \beta$, we have $|S'_\beta((\xi^t)^2)\xi^t(v^t - \tilde{v}^t)| \leq \max(C, 1/2)|v^t - \tilde{v}^t|$. Since $v^t - \tilde{v}^t$ is continuous and hence integrable over $[T_0, t]$, by dominated convergence

$$\lim_{\beta \to 0} \int_{T_0}^t S'_\beta((\xi^t)^2)\xi^t(v^t - \tilde{v}^t)\mathrm{d}t = \int_{T_0}^t \lim_{\beta \to 0} S'_\beta((\xi^t)^2)\xi^t(v^t - \tilde{v}^t)\mathrm{d}t = \int_{T_0}^t \frac{\mathrm{sign}(\xi^t)}{2}(v^t - \tilde{v}^t)\mathrm{d}t.$$

Applying the Lipschitz bound $h(|\xi^t|) \leq h(0) + C|\xi^t| = C|\xi^t|$ and a similar dominated convergence argument for the other terms of (66), we obtain in the limit $\beta \to 0$ that for $t \geq T_0$,

$$\mathrm{d}|\xi^t| = \mathrm{sign}(\xi^t)(v^t - \tilde{v}^t)\mathrm{d}t + 2\sqrt{2}\,\mathrm{sign}(\xi^t)h(|\xi^t|)\mathrm{d}b^t$$

which is the analogue of [75, Eq. (66)]. (There is no term corresponding to a local time of $\xi^t$ at 0 that would instead arise under a pure reflection coupling.)

Now let $A : [0, \infty) \to [0, \infty)$ be any continuously-differentiable function, and let $f : [0, \infty) \to [0, \infty)$ be any continuously-differentiable function with absolutely continuous first derivative (for which Itô's formula applies, c.f. [76, Theorem 71]), and satisfying $f'(r) \in [0, 1]$ and $f''(r) \leq 0$ for all $r \geq 0$. Set

$$r^t = |\xi^t| + \int_0^t A(t - s)|\xi^s|\mathrm{d}s.$$

Then $\mathrm{d}r^t = \mathrm{d}|\xi^t| + [A(0)|\xi^t| + \int_0^t A'(t - s)|\xi^s|\mathrm{d}s]\,\mathrm{d}t$. Applying Itô's formula and taking expectations gives, for $t \geq T_0$,

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathbb{E}f(r^t) = \mathbb{E}\Big[f'(r^t)\Big(\mathrm{sign}(\xi^t)(v^t - \tilde{v}^t) + A(0)|\xi^t| + \int_0^t A'(t - s)|\xi^s|\mathrm{d}s\Big) + 4f''(r^t)h(|\xi^t|)^2\Big]. \quad (67)$$

Let us define $\kappa : [0, \infty) \to \mathbb{R}$ by

$$\kappa(r) = \inf\left\{\frac{-(\log g)'(x) + (\log g)'(y)}{x - y} : |x - y| = r\right\} \quad (68)$$

so that $[-(\log g)'(\theta^t) + (\log g')(\tilde{\theta}^t)]/\xi^t \geq \kappa(|\xi^t|)$. Let us set also

$$\Delta_t = \int_0^t \Big(|R_\eta(t, s) - R_\eta^{(M)}(t - s)| \cdot \mathbb{E}|\theta^s - \theta^*|\Big)\mathrm{d}s + \mathbb{E}|u^t - \tilde{u}^t|. \quad (69)$$

23

Then, under our assumption $f'(r) \in [0, 1]$, we have the bound

$$\mathbb{E}\Big[f'(r^t)\operatorname{sign}(\xi^t)(v^t - \tilde{v}^t)\Big]$$

$$= \mathbb{E}\Big[f'(r^t)\operatorname{sign}(\xi^t)\Big(-\frac{\delta}{\sigma^2}\xi^t - \big(-(\log g)'(\theta^t) + (\log g)'(\tilde{\theta}^t)\big)$$

$$+ \int_0^t R_\eta^{(M)}(t-s)(\theta^s - \tilde{\theta}^s)\mathrm{d}s + \int_0^t \big(R_\eta(t,s) - R_\eta^{(M)}(t-s)\big)(\theta^s - \theta^*)\mathrm{d}s + (u^t - \tilde{u}^t)\Big)\Big]$$

$$\leq -\frac{\delta}{\sigma^2}\mathbb{E}[f'(r^t)|\xi^t|] - \mathbb{E}[f'(r^t)\kappa(|\xi^t|)|\xi^t|] + \mathbb{E}\Big[f'(r^t)\int_0^t R_\eta^{(M)}(t-s)|\xi^s|\mathrm{d}s\Big] + \Delta_t.$$

Applying this to (67), for all $t \geq T_0$,

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathbb{E}f(r^t) \leq \mathbb{E}\Big[-\Big(\frac{\delta}{\sigma^2} - A(0)\Big)f'(r^t)|\xi^t| - f'(r^t)\kappa(|\xi^t|)|\xi^t| + 4f''(r^t)h(|\xi^t|)^2\Big]$$

$$+ \mathbb{E}\Big[f'(r^t)\int_0^t \big(A'(t-s) + R_\eta^{(M)}(t-s)\big)|\xi^s|\mathrm{d}s\Big] + \Delta_t. \tag{70}$$

Let us now choose the functions $A(\cdot)$ and $f(\cdot)$. For some small enough $c_0 \in (0, \iota)$, let

$$A(0) = \frac{\delta}{\sigma^2} - c_0, \qquad A(\tau) = A(0)e^{-c_0\tau} - \int_0^\tau e^{-c_0(\tau-s)}\Big(\sum_{m=1}^M c_m^2 e^{-a_m s}\Big)\mathrm{d}s.$$

This choice of $A(\tau)$ satisfies $A'(\tau) = -c_0 A(\tau) - \sum_{m=1}^M c_m^2 e^{-a_m\tau}$, i.e.

$$A'(\tau) + R_\eta^{(M)}(\tau) = -c_0 A(\tau). \tag{71}$$

We will require that $A(\tau) \geq 0$ for all $\tau \geq 0$. To check this condition, observe that explicitly evaluating the integral defining $A(\tau)$ yields

$$e^{c_0\tau}A(\tau) = A(0) - \sum_{m=1}^M \frac{c_m^2}{a_m - c_0}\Big(1 - e^{-(a_m-c_0)\tau}\Big) \geq \frac{\delta}{\sigma^2} - c_0 - \sum_{m=1}^M \frac{c_m^2}{a_m - c_0} \geq \frac{\delta}{\sigma^2} - c_0 - \sum_{m=1}^M \frac{c_m^2}{a_m} \cdot \frac{\iota}{\iota - c_0},$$

the last inequality using $a_m \geq \iota \geq c_0$. Further bounding $\sum_{m=1}^M c_m^2/a_m = \sum_{m=1}^M \mu_\eta([a_{m-1}, a_m)) \leq \mu_\eta([\iota, \infty)) = c_\eta^{\mathrm{tti}}(0) - c_\eta^{\mathrm{tti}}(\infty)$, this shows $e^{c_0\tau}A(\tau) \geq \frac{\delta}{\sigma^2} - c_0 - \frac{\iota}{\iota - c_0}(c_\eta^{\mathrm{tti}}(0) - c_\eta^{\mathrm{tti}}(\infty))$. Then by the given assumption that $c_\eta^{\mathrm{tti}}(0) - c_\eta^{\mathrm{tti}}(\infty) < \delta/\sigma^2$, we obtain $A(\tau) \geq 0$ for a sufficiently small choice of $c_0 \in (0, \iota)$ and all $\tau \geq 0$, as desired. Applying (71) and $A(0) = \delta/\sigma^2 - c_0$ into (70), and recalling the definition $r^t = |\xi^t| + \int_0^t A(t-s)|\xi^s|\mathrm{d}s$, we get for all $t \geq T_0$ that

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathbb{E}f(r^t) \leq \mathbb{E}\Big[\underbrace{-c_0 f'(r^t)r^t - f'(r^t)\kappa(|\xi^t|)|\xi^t| + 4f''(r^t)h(|\xi^t|)^2}_{:=F(r^t, \xi^t)}\Big] + \Delta_t. \tag{72}$$

We next proceed to bound the above quantity $\mathbb{E}[F(r^t, \xi^t)]$. Observe that under the convexity-at-infinity condition for $-\log g(\theta)$ in Assumption 2.2(a), there must exist constants $R_0, \kappa_0 > 0$ for which

$$\kappa(r) \geq -\kappa_0 \text{ for all } r \geq 0, \qquad \kappa(r) \geq 0 \text{ for all } r \geq R_0. \tag{73}$$

Let us denote $\kappa_-(r) = \max(-\kappa(r), 0)$. Then $-\kappa(r)r \leq \kappa_-(r)r$ and $\kappa_-(r)r \in [0, \kappa_0 R_0]$ for all $r \geq 0$. Recall the constant $\varepsilon > 0$ for which $h(x) = 1$ when $x \geq \varepsilon$, and define $K : (\varepsilon, \infty) \to [0, \kappa_0 R_0]$ by

$$K(r) = \sup_{t \geq T_0}\mathbb{E}\Big[\kappa_-(|\xi^t|)|\xi^t| \,\Big|\, r^t = r, |\xi^t| > \varepsilon\Big].$$

Then define $f : \mathbb{R} \to \mathbb{R}$ by

$$f(0) = 0, \qquad f'(r) = \exp\Big(-\frac{1}{4}\int_0^{\max(r, 2\kappa_0 R_0/c_0)} K(s)ds\Big) \text{ for } r \geq 0.$$

24

Note that $f'(r)$ is absolutely continuous as required, with $f'(r) \in [c_1, 1]$ for all $r \geq 0$ where $c_1 = \exp(-\frac{\kappa_0^2 R_0^2}{2c_0})$, and $f''(r) = -\frac{1}{4} K(r) f'(r) \mathbf{1}\{r < 2\kappa_0 R_0/c_0\} \leq 0$. By these definitions, for any $r > \varepsilon$ and $t \geq T_0$, we have

$$\mathbb{E}[F(r^t, \xi^t) \mid r^t = r, |\xi^t| > \varepsilon] \leq \mathbb{E}\Big[-c_0 f'(r^t) r^t + f'(r^t) \kappa_-(|\xi^t|) |\xi^t| + 4 f''(r^t) h(|\xi^t|)^2 \mid r^t = r, |\xi^t| > \varepsilon\Big]$$

$$\leq \mathbb{E}\Big[-c_0 f'(r) r + f'(r) K(r) + 4 f''(r) \mid r^t = r, |\xi^t| > \varepsilon\Big].$$

When $r \geq 2\kappa_0 R_0/c_0$ we may apply $f''(r) = 0$ and $K(r) \leq \kappa_0 R_0 \leq c_0 r/2$, to bound this above by $-(c_0/2) f'(r) r$. When $r \in (\varepsilon, 2\kappa_0 R_0/c_0)$ we may instead apply $f''(r) = -\frac{1}{4} K(r) f'(r)$ to see that this equals $-c_0 f'(r) r$. Thus for all $r > \varepsilon$ and $t \geq T_0$,

$$\mathbb{E}[F(r^t, \xi^t) \mid r^t = r, |\xi^t| > \varepsilon] \leq -(c_0/2) f'(r) r.$$

For any $r \geq 0$, on the event $|\xi^t| \leq \varepsilon$ (which occurs with probability 1 when $r^t \leq \varepsilon$ since $A(t) \geq 0$), let us use $f''(r^t) h(|\xi^t|)^2 \leq 0$ and $-f'(r^t) \kappa_-(|\xi^t|) |\xi^t| \leq \varepsilon \kappa_0$ to bound

$$\mathbb{E}[F(r^t, \xi^t) \mid r^t = r, |\xi^t| \leq \varepsilon] \leq -c_0 f'(r) r + \varepsilon \kappa_0.$$

Combining these cases and taking the full expectation over $\mathbf{1}\{|\xi^t| > \varepsilon\}$ and over $r^t$, we get for all $t \geq T_0$ that

$$\mathbb{E}[F(r^t, \xi^t)] \leq -(c_0/2) \mathbb{E}[f'(r^t) r^t] + \varepsilon \kappa_0.$$

Applying $f'(r^t) \geq c_1$ and $r^t \geq f(r^t)$ and putting this bound into (72), for all $t \in [T_0, T_0 + T]$,

$$\frac{\mathrm{d}}{\mathrm{d}t} \mathbb{E} f(r^t) \leq -(c_0 c_1/2) \mathbb{E} f(r^t) + \varepsilon \kappa_0 + \max_{t \in [T_0, T_0 + T]} \Delta_t.$$

Since $f(r^{T_0}) = f(0) = 0$, this differential inequality yields for all $t \in [T_0, T_0 + T]$,

$$\mathbb{E} f(r^t) \leq \Big(\varepsilon \kappa_0 + \max_{t \in [T_0, T_0 + T]} \Delta_t\Big) \frac{1 - e^{-(c_0 c_1/2)(t - T_0)}}{c_0 c_1/2} \leq \frac{2}{c_0 c_1} \Big(\varepsilon \kappa_0 + \max_{t \in [T_0, T_0 + T]} \Delta_t\Big).$$

Since also $r^t \leq c_1 f(r^t)$ from the lower bound $f'(r) \geq c_1$, this gives $\mathbb{E} r^t \leq (2/c_0)(\varepsilon \kappa_0 + \max_{t \in [T_0, T_0 + T]} \Delta_t)$. Applying that $A(t) \geq 0$ for all $t \geq 0$, we have $|\xi^t| \leq r^t$, so this gives finally

$$\max_{t \in [T_0, T_0 + T]} \mathbb{E}|\theta^t - \tilde{\theta}^t| = \max_{t \in [T_0, T_0 + T]} \mathbb{E}|\xi^t| \leq (2/c_0)\Big(\varepsilon \kappa_0 + \max_{t \in [T_0, T_0 + T]} \Delta_t\Big).$$

We may choose $\varepsilon$ such that $\varepsilon \kappa_0 \leq \max_{t \in [T_0, T_0 + T]} \Delta_t$. Thus, to conclude the proof, it remains to show

$$\sup_{t \in [T_0, T_0 + T]} \Delta_t \leq \varepsilon(M) + \sqrt{T} \varepsilon(T_0). \tag{74}$$

We note that under Definition 2.4, $\mathbb{E}|\theta^t - \theta^*| \leq [\mathbb{E}(\theta^t - \theta^*)^2]^{1/2} = (C_\theta(t,t) - 2C_\theta(t,*) + \mathbb{E}\theta^{*2})^{1/2} \leq C$ for a constant $C > 0$ and all $t \geq 0$. Then for the first term of $\Delta_t$, by property (35) of Definition 2.4,

$$\int_0^t |R_\eta(t,s) - R_\eta^{(M)}(t-s)| \cdot \mathbb{E}|\theta^s - \theta^*| \, \mathrm{d}s$$

$$\leq C \int_0^t |R_\eta(t,s) - r_\eta^{\mathrm{tti}}(t-s)| \mathrm{d}s + C \int_0^t |r_\eta^{\mathrm{tti}}(t-s) - R_\eta^{(M)}(t-s)| \mathrm{d}s$$

$$\leq C \varepsilon(t) + \int_0^t |r_\eta^{\mathrm{tti}}(t-s) - R_\eta^{(M)}(t-s)| \mathrm{d}s$$

where here $\lim_{t \to \infty} \varepsilon(t) = 0$. Recalling the sequences $\{a_m\}_{m=0}^M, \{c_m\}_{m=1}^M$ defining $R_\eta^{(M)}$,

$$r_\eta^{\mathrm{tti}}(\tau) - R_\eta^{(M)}(\tau) = \int_\iota^\infty a e^{-a\tau} \mu_\eta(\mathrm{d}a) - \sum_{m=1}^M c_m^2 e^{-a_m \tau}$$

$$= \sum_{m=1}^M \int_{a_{m-1}}^{a_m} (a e^{-a\tau} - a_m e^{-a_m \tau}) \mu_\eta(\mathrm{d}a) + \int_{a:a > a_M} a e^{-a\tau} \mu_\eta(\mathrm{d}a),$$

hence using the fact that $h(a) = ae^{-a\tau}$ satisfies $|h'(a)| \leq 2e^{-a\tau/2}$ and $|a_m - a_{m-1}| = 1/\sqrt{M}$,

$$
\begin{aligned}
\int_0^t |r_\eta^{\text{tti}}(\tau) - R_\eta^{(M)}(\tau)| \mathrm{d}\tau &\leq \sum_{m=1}^M \int_{a_{m-1}}^{a_m} \left( \int_0^t |ae^{-a\tau} - a_m e^{-a_m\tau}| \mathrm{d}\tau \right) \mu_\eta(\mathrm{d}a) + \int_{a:a>a_M} \left( \int_0^t ae^{-a\tau} \mathrm{d}\tau \right) \mu_\eta(\mathrm{d}a) \\
&\leq \sum_{m=1}^M \int_{a_{m-1}}^{a_m} \left( \int_0^t (2/\sqrt{M})e^{-a_{m-1}\tau/2} \mathrm{d}\tau \right) \mu_\eta(\mathrm{d}a) + \mu_\eta([a_M, \infty)) \\
&\leq \frac{4}{\sqrt{M}} \sum_{m=1}^M \int_{a_{m-1}}^{a_m} \frac{1}{a_{m-1}} \mu_\eta(\mathrm{d}a) + \mu_\eta([a_M, \infty)) \\
&\leq \frac{4}{\iota\sqrt{M}} \mu_\eta([\iota, \sqrt{M})) + \mu_\eta([\sqrt{M}, \infty)) \leq \varepsilon(M),
\end{aligned}
$$

where $\lim_{M\to\infty} \varepsilon(M) = 0$. This bounds the first term of $\Delta_t$ by $\varepsilon(M) + Ct \cdot \varepsilon(t)$. Bounding also the second term $\mathbb{E}|u^t - \tilde{u}^t|$ of $\Delta_t$ by Lemma 3.2, we have

$$
\sup_{t\in[T_0,T_0+T]} \Delta_t \leq \varepsilon(M) + \sqrt{T}\,\varepsilon(T_0) + \sup_{t\in[T_0,T_0+T]} Ct \cdot \varepsilon(t)
$$

which implies (74) upon adjusting $\varepsilon(T_0)$. This completes the proof. $\qquad\square$

Adapting part of the previous argument, we record here a uniform bound on $\mathbb{E}(\theta^t)^4$ for the solution $\{\theta^t\}_{t\geq 0}$ of the DMFT equation.

**Lemma 3.4.** *Suppose $c_\eta^{\text{tti}}(0) - c_\eta^{\text{tti}}(\infty) < \delta/\sigma^2$. Then $\sup_{t\geq 0} \mathbb{E}(\theta^t)^4 \leq C$ and $\sup_{t\geq 0} \mathbb{E}(\theta_{M,T_0}^t)^4 \leq C$ for a constant $C > 0$ and all $M, T_0 > 0$.*

*Proof.* We prove the statement for $\{\theta^t\}_{t\geq 0}$. Let $A : [0, \infty) \to [0, \infty)$ be defined by

$$
A(0) = \frac{\delta}{\sigma^2} - c_0, \qquad A(\tau) = A(0)e^{-c_0\tau} - \int_0^t e^{-c_0(\tau-s)} r_\eta^{\text{tti}}(s) \mathrm{d}s
$$

for a small enough constant $c_0 \in (0, \iota)$. Here, by the conditions of Definition 2.4, $r_\eta^{\text{tti}}(s) = -c_\eta^{\text{tti}\prime}(s) = \int_\iota^\infty ae^{-as} \mathrm{d}\mu_\eta(a)$, and the same argument as in the preceding proof verifies that $\inf_{\tau\in[0,\infty)} A(\tau)$ is bounded below by a positive constant for a sufficiently small choice of $c_0 \in (0, \iota)$ and all $\tau \geq 0$.

Let $f : \mathbb{R} \to [0, \infty)$ be a smooth approximation to the absolute value, satisfying $f(x) = |x|$ for all $|x| \geq 1$, and $1 + f'(x) \cdot x \geq f(x) \geq |x|$, and $|f'(x)| \leq 1$, and $|f''(x)| \leq C$ for all $x \in \mathbb{R}$ and an absolute constant $C > 0$. Let $\bar{\theta}^t = \theta^t - \theta^*$, and set $r^t = f(\bar{\theta}^t) + \int_0^t A(t-s)f(\bar{\theta}^s) \mathrm{d}s$. Then by the DMFT equation (23) and Itô's formula,

$$
\mathrm{d}\bar{\theta}^t = \left[ -\frac{\delta}{\sigma^2} \bar{\theta}^t + (\log g)'(\theta^t) + \int_0^t R_\eta(t,s)\bar{\theta}^s \mathrm{d}s + u^t \right] \mathrm{d}t + \sqrt{2}\,\mathrm{d}b^t,
$$

$$
\mathrm{d}r^t = f'(\bar{\theta}^t)\mathrm{d}\bar{\theta}^t + f''(\bar{\theta}^t)\mathrm{d}t + \left[ A(0)f(\bar{\theta}^t) + \int_0^t A'(t-s)f(\bar{\theta}^s)\mathrm{d}s \right]\mathrm{d}t,
$$

$$
\mathrm{d}(r^t)^4 = 4(r^t)^3\mathrm{d}r^t + 12(r^t)^2 f'(\bar{\theta}^t)^2 \mathrm{d}t.
$$

Applying $r^t \geq 0$ and the bounds $f'(\bar{\theta}^t)\bar{\theta}^t \geq f(\bar{\theta}^t) - 1$, $|f'(\bar{\theta}^t)| \leq 1$, $|\bar{\theta}^s| \leq f(\bar{\theta}^s)$, and $|f''(\bar{\theta}^t)| \leq C$ from the definition of $f(\cdot)$, this gives

$$
\begin{aligned}
\frac{\mathrm{d}}{\mathrm{d}t} \mathbb{E}(r^t)^4 \leq \mathbb{E}\bigg[ 4(r^t)^3 \bigg( &-\frac{\delta}{\sigma^2}[f(\bar{\theta}^t) - 1] + f'(\bar{\theta}^t)(\log g)'(\theta^t) + \int_0^t |R_\eta(t,s)|f(\bar{\theta}^s)\mathrm{d}s + |u^t| + C \\
&+ A(0)f(\bar{\theta}^t) + \int_0^t A'(t-s)f(\bar{\theta}^s)\mathrm{d}s \bigg) + 12(r^t)^2 \bigg].
\end{aligned}
$$

Then, using $A(0) = \delta/\sigma^2 - c_0$ and $A'(t-s) + r_\eta^{\mathrm{tti}}(t-s) = -c_0 A(t-s)$ from the definition of $A(\cdot)$,

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathbb{E}(r^t)^4 \le \mathbb{E}\left[4(r^t)^3\left(-c_0 r^t + \frac{\delta}{\sigma^2} + f'(\bar{\theta}^t)(\log g)'(\theta^t) + \int_0^t |R_\eta(t,s) - r_\eta^{\mathrm{tti}}(t-s)|f(\bar{\theta}^s)\mathrm{d}s + |u^t| + C\right) + 12(r^t)^2\right].$$

When $|\bar{\theta}^t| \ge 1$, we must have $f'(\bar{\theta}^t) = \mathrm{sign}(\bar{\theta}^t) = |\bar{\theta}^t|/(\theta^t - \theta^*)$. Recalling the function $\kappa(r)$ from (68), let us bound in this case

$$f'(\bar{\theta}^t)[(\log g)'(\theta^t) - (\log g)'(\theta^*)] = -|\bar{\theta}^t| \cdot \frac{-(\log g)'(\theta^t) + (\log g)'(\theta^*)}{\theta^t - \theta^*} \le -\kappa(|\bar{\theta}^t|)|\bar{\theta}^t| \le \kappa_0 R_0,$$

where $\kappa_0 R_0$ is the deterministic upper bound for $-\kappa(r)r$. For $|\bar{\theta}^t| \le 1$, let us apply instead the Lipschitz bound $|f'(\bar{\theta}^t)[(\log g)'(\theta^t) - (\log g)'(\theta^*)]| \le L$ where $L$ is the Lipschitz constant of $(\log g)'$ under Assumption 2.2(a). We also apply $|R_\eta(t,s) - r_\eta^{\mathrm{tti}}(t-s)| \le \varepsilon(t)$ from Definition 2.4, and $\int_0^t f(\bar{\theta}^s)\mathrm{d}s \le r^t/\inf_{\tau \in [0,\infty)} A(\tau)$ by the definition of $r^t$, where we recall that $\inf_{\tau \in [0,\infty)} A(\tau)$ is bounded below by a positive constant. Thus, for some constant $C' > 0$, this yields

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathbb{E}(r^t)^4 \le -4c_0\mathbb{E}(r^t)^4 + C'\,\mathbb{E}\left[\varepsilon(t)(r^t)^4 + (r^t)^3(1 + |(\log g)'(\theta^*)| + |u^t|) + (r^t)^2\right].$$

Since $u^t$ is a centered Gaussian variable, we note that $\mathbb{E}(u^t)^4 = 3[\mathbb{E}(u^t)^2]^2 = 3C(t,t)^2$ which is bounded uniformly for all $t \ge 0$ under Definition 2.4. Also $\mathbb{E}[|(\log g)'(\theta^*)|^4]$ is finite by the Lipschitz continuity of $(\log g)'$ and finiteness of moments of $\theta^*$ under Assumption 2.1. Then by Hölder's inequality, $\mathbb{E}[(r^t)^3(1 + |(\log g)'(\theta^*)| + |u^t|) + (r^t)^2] \le C(\mathbb{E}[(r^t)^4]^{3/4} + \mathbb{E}[(r^t)^4]^{1/2})$ for some $C > 0$. Thus, for some $C, T, R > 0$ sufficiently large depending on $C', c_0$, the above implies

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathbb{E}(r^t)^4 \le C\,\mathbb{E}(r^t)^4, \qquad \frac{\mathrm{d}}{\mathrm{d}t}\mathbb{E}(r^t)^4 \le -4c_0\mathbb{E}(r^t)^4 + c_0\mathbb{E}(r^t)^4 < 0 \text{ whenever } t \ge T \text{ and } \mathbb{E}(r^t)^4 \ge R.$$

This implies that $\sup_{t \ge 0} \mathbb{E}(r^t)^4$ is bounded by a constant depending only on $C, T, R$. Then $\sup_{t \ge 0} \mathbb{E}(\theta^t)^4$ is also bounded since $\theta^t = \bar{\theta}^t + \theta^*$ and $|\bar{\theta}^t| \le f(\bar{\theta}^t) \le r^t$.

The argument to bound $\mathbb{E}(\theta_{M,T_0})^4$ is the same upon replacing $R_\eta(t,s)$ and $u^t$ by $R_\eta^{(M)}(t,s)$ and $u_M^t$ for $s, t \ge T_0$, and we omit the details for brevity. $\qquad\square$

### 3.1.2 Convergence of the auxiliary process

Extending the definition (40) of $\mathsf{P}_{g_*,\omega_*;g,\omega}$, let $\mathsf{P}_{g_*,\omega_*;g,\omega}^{\otimes 2}$ denote the law of a triple $(\theta^*, \theta, \theta')$ where $\theta^*, \theta$ are generated according to (41) defining $\mathsf{P}_{g_*,\omega_*;g,\omega}$ and $\theta'$ is a second independent copy of $\theta$ drawn from the posterior measure conditional on $y$.

**Lemma 3.5.** *Suppose $c_\eta^{\mathrm{tti}}(0) - c_\eta^{\mathrm{tti}}(\infty) < \delta/\sigma^2$. Fix any $M, T_0 > 0$, set $\omega^{(M)} = \delta/\sigma^2 - \sum_{m=1}^M c_m^2/a_m$ and $\omega_*^{(M)} = (\omega^{(M)})^2/c_\eta^{\mathrm{tti}}(\infty)$, and let $\{\theta_{M,T_0}^t\}_{t \ge 0}$ be the process (64). Then for any $T, T' > 0$,*

$$W_1\left(\mathsf{P}(\theta^*, \theta_{M,T_0}^{T_0+T}, \theta_{M,T_0}^{T_0+T+T'}), \mathsf{P}_{g_*,\omega_*^{(M)};g,\omega^{(M)}}^{\otimes 2}\right) \le C\sqrt{M}(e^{-cT} + e^{-cT'})$$

*for some constants $C, c > 0$ not depending on $M, T_0, T, T'$.*

*Proof.* Let $z \sim \mathcal{N}(0, c_\eta^{\mathrm{tti}}(\infty))$ and let $\{\tilde{b}^t\}_{t \ge T_0}$ and $\{b_m^t\}_{t \ge 0}$ for $m = 1, \ldots, M$ be $M + 1$ standard Brownian motions. These are all independent of each other, of $\theta^*$, and of $\{\theta^t\}_{t \in [0,T]}$. We note that the law of $\{\theta_{M,T_0}^t\}_{t \ge 0}$ defined by (64) coincides with the marginal law of $\{\theta_{M,T_0}^t\}_{t \ge 0}$ in the joint process

$$\theta_{M,T_0}^t = \theta^t \text{ for } t \in [0, T_0],$$

$$\mathrm{d}\theta_{M,T_0}^t = \left[-\frac{\delta}{\sigma^2}(\theta_{M,T_0}^t - \theta^*) + (\log g)'(\theta_{M,T_0}^t) + z + \sum_{m=1}^M c_m x_m^t\right]\mathrm{d}t + \sqrt{2}\,\mathrm{d}\tilde{b}^t \text{ for } t > T_0, \tag{75}$$

$$\mathrm{d}x_m^t = [-a_m x_m^t + c_m(\theta_{M,T_0}^t - \theta^*)]\mathrm{d}t + \sqrt{2}\,\mathrm{d}b_m^t \text{ for } 1 \le m \le M, t \ge 0 \tag{76}$$

27

with initial conditions $x_1^0 = \ldots = x_M^0 = 0$. Indeed, given $\{\theta_{M,T_0}^t\}_{t \geq 0}$, the equations (76) for $\{x_m^t\}_{t \geq 0}$ are linear and have the explicit solutions

$$x_m^t = c_m \int_0^t e^{-a_m(t-s)}(\theta_{M,T_0}^s - \theta^*)\mathrm{d}s + \int_0^t e^{-a_m(t-s)}\sqrt{2}\,\mathrm{d}b_m^s. \tag{77}$$

Substituting these solutions into (75) gives (64), upon identifying $u_M^t = z + \sum_{m=1}^M c_m \int_0^t e^{-a_m(t-s)}\sqrt{2}\,\mathrm{d}b_m^s$. Here $\{u_M^t\}_{t \geq 0}$ is a centered Gaussian process independent of $\theta^*$ and $\{\tilde{b}^t\}_{t \geq T_0}$, with covariance kernel exactly $C_\eta^{(M)}(t,s)$ by (62). Thus the marginal law of $\{\theta_{M,T_0}^t\}_{t \geq 0}$ coincides with the definition in (64).

For $t \geq T_0$, let $x^t = (\theta_{M,T_0}^t, x_1^t, \ldots, x_M^t)$ and $b^t = (\tilde{b}^t, b_1^t, \ldots, b_M^t)$. Conditional on $\theta^*$ and $z$, the evolution of $\{x^t\}_{t \geq T_0}$ defined by (75–76) is a standard (Markovian) Langevin diffusion given by

$$\mathrm{d}x^t = -\nabla H(x^t \mid \theta^*, z)\mathrm{d}t + \sqrt{2}\,\mathrm{d}b^t$$

with Hamiltonian

$$H(x \mid \theta^*, z) = H(\theta, x_1, \ldots, x_M \mid \theta^*, z) \tag{78}$$

$$= \frac{1}{2}\underbrace{\left(\frac{\delta}{\sigma^2} - \sum_{m=1}^M \frac{c_m^2}{a_m}\right)}_{=\omega^{(M)}}(\theta - \theta^*)^2 - \log g(\theta) - z\,\theta + \sum_{m=1}^M \frac{a_m}{2}\left(\frac{c_m}{a_m}(\theta - \theta^*) - x_m\right)^2$$

$$= \frac{\omega^{(M)}}{2}(\theta - \theta^* - z/\omega^{(M)})^2 - \log g(\theta) + \sum_{m=1}^M \frac{a_m}{2}\left(\frac{c_m}{a_m}(\theta - \theta^*) - x_m\right)^2 + \text{const.}, \tag{79}$$

for an additive constant not depending on $x = (\theta, x_1, \ldots, x_M)$. Note that the given condition $c_\eta^{\mathrm{tti}}(0) - c_\eta^{\mathrm{tti}}(\infty) < \delta/\sigma^2$ implies that $\omega^{(M)} > 0$ strictly.

Convergence of $\{x^t\}_{t \geq T_0}$ in Wasserstein-1 to a stationary law then follows from the results of [74]: For any $x = (\theta, x_1, \ldots, x_M)$ and $x' = (\theta', x_1', \ldots, x_M')$, we have

$$(x - x')^\top(\nabla H(x \mid \theta^*, z) - \nabla H(x' \mid \theta^*, z)) = (\theta - \theta')(-(\log g)'(\theta) + (\log g)'(\theta')) + (x - x')^\top L(x - x')$$

where

$$L = \begin{pmatrix} \frac{\delta}{\sigma^2} & -c_1 & \cdots & -c_M \\ -c_1 & a_1 & & \\ \vdots & & \ddots & \\ -c_M & & & a_M \end{pmatrix}.$$

By the positivity of the Schur complement $\omega^{(M)} = \delta/\sigma^2 - \sum_{m=1}^M c_m^2/a_m \geq \delta/\sigma^2 - (c_\eta^{\mathrm{tti}}(0) - c_\eta^{\mathrm{tti}}(\infty))$ and of $a_m \geq \iota$, this matrix $L$ is strictly positive-definite, with smallest eigenvalue bounded away from 0 independently of $M$. Then, recalling the function $\kappa(r)$ from (68),

$$(x - x')^\top(\nabla H(x \mid \theta^*, z) - \nabla H(x' \mid \theta^*, z)) \geq (\theta - \theta')^2\kappa(|\theta - \theta'|) + c\|x - x'\|_2^2$$

for a constant $c > 0$. Recalling also that $\kappa(r)$ is positive for all $r > R_0$ and some $R_0 > 0$, and considering separately the cases where $|\theta - \theta'| \leq R_0$ and $|\theta - \theta'| > R_0$, this verifies that

$$\inf_{\|x-x'\|_2 = r} \frac{(x - x')^\top(\nabla H(x \mid \theta^*, z) - \nabla H(x' \mid \theta^*, z))}{\|x - x'\|_2^2} > c'$$

for all $r > R_0'$ and some constants $c', R_0' > 0$. Then by [74, Corollary 3], the Langevin diffusion $\{x^t\}_{t \geq T_0}$ has the unique stationary law

$$\mathsf{P}^\infty(x) \propto \exp(-H(x \mid \theta^*, z)). \tag{80}$$

Let us write $x^\infty \sim \mathsf{P}^\infty$ and $\langle f(x^\infty)\rangle$ for a sample and Gibbs average under this stationary law. Let us write also $W_1(\cdot)$ for the Wasserstein-1 distance conditional on $\theta^*, z$, and $\mathsf{P}(x^{T_0+T} \mid x^{T_0} = x)$ for the conditional

28

law of $x^{T_0+T}$ given $\theta^*, z$ and the initial condition $x^{T_0} = x$. Then also by [74, Corollary 2 and 3], there exist constants $C, c > 0$ such that for any $T > 0$,

$$W_1(\mathsf{P}(x^{T_0+T} \mid x^{T_0} = x), \mathsf{P}^\infty) \le Ce^{-cT}W_1(\delta_x, \mathsf{P}^\infty) \le Ce^{-cT}(\|x\|_2 + \langle \|x^\infty\|_2 \rangle).$$

Similarly, for any $T' > 0$,

$$W_1(\mathsf{P}(x^{T_0+T+T'} \mid x^{T_0+T} = x), \mathsf{P}^\infty) \le Ce^{-cT'}(\|x\|_2 + \langle \|x^\infty\|_2 \rangle).$$

Combining the two conditional couplings that attain these Wasserstein-1 bounds, and taking the average over the sample path $\{x^t\}_{t \ge T_0}$ (which we denote by $\langle f(x^t) \rangle$, still conditional on $\theta^*, z$),

$$W_1(\mathsf{P}(x^{T_0+T}, x^{T_0+T+T'}), (\mathsf{P}^\infty)^{\otimes 2}) \le C(e^{-cT} + e^{-cT'})(\langle \|x^{T_0}\|_2 \rangle + \langle \|x^{T_0+T}\|_2 \rangle + \langle \|x^\infty\|_2 \rangle)$$

where $(\mathsf{P}^\infty)^{\otimes 2}$ is the law of two independent samples from $\mathsf{P}^\infty$. The explicit form (77) for each $\{x_m^t\}_{t \ge 0}$ implies that $\langle |x_m^t| \rangle \le c_m \int_0^t e^{-\iota(t-s)}(\langle |\theta_{M,T_0}^s| \rangle + |\theta^*|)\mathrm{d}s + (a_m)^{-1/2}$, and hence

$$\langle \|x^t\|_2 \rangle \le C\sqrt{M}\Big(1 + |\theta^*| + \int_0^t e^{-\iota(t-s)}\langle |\theta_{M,T_0}^s| \rangle \mathrm{d}s\Big)$$

for a constant $C > 0$. Then, taking the full expectation over $\theta^*, z$ and applying $\mathbb{E}\langle |\theta_{M,T_0}^t| \rangle \le C$ by Lemma 3.4, we get $\mathbb{E}\langle \|x^t\|_2 \rangle \le C'\sqrt{M}$ for a constant $C' > 0$ and all $t \ge 0$. Then, applying this above gives

$$\mathbb{E}W_1(\mathsf{P}(x^{T_0+T}, x^{T_0+T+T'}), (\mathsf{P}^\infty)^{\otimes 2}) \le C\sqrt{M}(e^{-cT} + e^{-cT'}) \tag{81}$$

for some (different) constants $C, c > 0$.

Finally, note that the stationary law $\mathsf{P}^\infty(x)$ defined by (80) with Hamiltonian (79) describes a joint law (conditional on $\theta^*, z$) of $(\theta, x_1, \ldots, x_M)$ where $x_m \mid \theta$ is Gaussian and independent across $m = 1, \ldots, M$, and $\theta$ has marginal law given exactly by $\mathsf{P}(\theta \mid y)$ in the Gaussian convolution model (37) with observation $y = \theta^* + z/\omega^{(M)}$. Here, the noise variable $z/\omega^{(M)}$ is Gaussian with variance $c_\eta^{\mathrm{tti}}(\infty)/(\omega^{(M)})^2 = (\omega_*^{(M)})^{-1}$, so the joint law of $\theta^*$ and the $(\theta, \theta')$-marginals of the conditional law $(\mathsf{P}^\infty)^{\otimes 2}$ given $(\theta^*, z)$ is precisely $\mathsf{P}_{g_*, \omega_*^{(M)}; g, \omega^{(M)}}^{\otimes 2}$. Then, taking the $(\theta, \theta')$-marginals of the coupling (conditional on $\theta^*, z$) that attains $W_1(\mathsf{P}(x^{T_0+T}, x^{T_0+T+T'}), (\mathsf{P}^\infty)^{\otimes 2})$ and combining with the identity coupling of $\theta^*$, we have

$$W_1(\mathsf{P}(\theta^*, \theta_{M,T_0}^{T_0+T}, \theta_{M,T_0}^{T_0+T+T'}), \mathsf{P}_{g_*, \omega_*^{(M)}; g, \omega^{(M)}}^{\otimes 2}) \le W_1(\mathsf{P}(x^{T_0+T}, x^{T_0+T+T'}), (\mathsf{P}^\infty)^{\otimes 2}).$$

Taking the full expectation over $\theta^*, z$ on both sides and applying the bound (81) shows the lemma. $\qquad \square$

**Lemma 3.6.** *Suppose $c_\eta^{\mathrm{tti}}(0) - c_\eta^{\mathrm{tti}}(\infty) < \delta/\sigma^2$. For any $M > 0$, let*

$$\omega^{(M)} = \delta/\sigma^2 - \sum_{m=1}^M c_m^2/a_m, \quad \omega_*^{(M)} = (\omega^{(M)})^2/c_\eta^{\mathrm{tti}}(\infty),$$

$$\omega = \delta/\sigma^2 - (c_\eta^{\mathrm{tti}}(0) - c_\eta^{\mathrm{tti}}(\infty)), \quad \omega_* = \omega^2/c_\eta^{\mathrm{tti}}(\infty).$$

*Then $\lim_{M \to \infty} W_1(\mathsf{P}_{g_*, \omega_*^{(M)}; g, \omega^{(M)}}^{\otimes 2}, \mathsf{P}_{g_*, \omega_*; g, \omega}^{\otimes 2}) = 0$.*

*Proof.* Let $(\theta^*, \theta, \theta') \sim \mathsf{P}_{g_*, \omega_*; g, \omega}^{\otimes 2}$, i.e. $\theta, \theta'$ are two independent draws from the posterior law $\mathsf{P}_{g,\omega}(\theta \mid y)$ in the scalar Gaussian convolution model (37) where $y = \theta^* + \omega^{*-1/2}z$ and $z \sim \mathcal{N}(0, 1)$. Let $\langle \cdot \rangle_{g,\omega}$ be average over $\theta, \theta'$ conditional on $\theta^*, z$, and let $\mathcal{F}$ be the class of 1-Lipschitz functions $f(\theta^*, \theta, \theta')$. Then, for any $f \in \mathcal{F}$,

$$\mathbb{E}_{g_*, \omega_*}\langle f(\theta^*, \theta, \theta') \rangle_{g,\omega} = \mathbb{E}\frac{\int f(\theta^*, \theta, \theta') \exp(-\frac{\omega}{2}[(\theta^* + \omega^{*-1/2}z - \theta)^2 + (\theta^* + \omega^{*-1/2}z - \theta')^2])g(\theta)g(\theta')\mathrm{d}(\theta, \theta')}{\int \exp(-\frac{\omega}{2}[(\theta^* + \omega^{*-1/2}z - \theta)^2 + (\theta^* + \omega^{*-1/2}z - \theta')^2])g(\theta)g(\theta')\mathrm{d}(\theta, \theta')}$$

where $\mathbb{E}$ on the right side is over $\theta^* \sim g_*$ and $z \sim \mathcal{N}(0,1)$. Writing $\langle \cdot \rangle$ for $\langle \cdot \rangle_{g,\omega}$ and $\kappa_2$ for its associated posterior covariance, the above is continuously-differentiable in $(\omega, \omega^*)$ with

$$\partial_\omega \mathbb{E}\langle f(\theta^*, \theta, \theta') \rangle = \mathbb{E}\Big[\kappa_2\Big(f(\theta^*, \theta, \theta'), -\tfrac{1}{2}[(\theta^* + \omega^{*-1/2}z - \theta)^2 + (\theta^* + \omega^{*-1/2}z - \theta')^2]\Big)\Big]$$

$$\partial_{\omega^*} \mathbb{E}\langle f(\theta^*, \theta, \theta') \rangle = \mathbb{E}\Big[\kappa_2\Big(f(\theta^*, \theta, \theta'), \tfrac{\omega z}{\omega_*^{3/2}}[(\theta^* + \omega^{*-1/2}z - \theta) + (\theta^* + \omega^{*-1/2}z - \theta')]\Big)\Big]$$

By the 1-Lipschitz bound for $f$ and the identity $\operatorname{Var} X = \tfrac{1}{2}\mathbb{E}[(X - X')^2]$ where $X'$ is an independent copy of $X$, we have $\kappa_2(f(\theta^*, \theta, \theta'), f(\theta^*, \theta, \theta')) \leq C(\kappa_2(\theta, \theta) + \kappa_2(\theta', \theta'))$ for an absolute constant $C > 0$. Then, applying Cauchy-Schwarz to $\kappa_2(\cdot)$ above, we get that $(\omega, \omega_*) \mapsto \mathbb{E}_{g_*,\omega_*}\langle f(\theta^*, \theta, \theta')\rangle_{g,\omega}$ is locally Lipschitz-continuous uniformly over $f \in \mathcal{F}$. Since $\lim_{M\to\infty}\sum_{m=1}^M c_m^2/a_m = \mu_\eta([\iota, \infty)) = c_\eta^{\mathrm{tti}}(0) - c_\eta^{\mathrm{tti}}(\infty)$, we have $\lim_{M\to\infty}(\omega^{(M)}, \omega_*^{(M)}) = (\omega, \omega_*)$. Then this local Lipschitz continuity implies as desired

$$\lim_{M\to\infty} W_1(\mathsf{P}^{\otimes 2}_{g_*,\omega_*^{(M)};g,\omega^{(M)}}, \mathsf{P}^{\otimes 2}_{g_*,\omega_*;g,\omega}) = \lim_{M\to\infty} \sup_{f\in\mathcal{F}}\Big|\mathbb{E}_{g_*,\omega_*}\langle f(\theta^*, \theta, \theta')\rangle_{g,\omega} - \mathbb{E}_{g_*,\omega_*^{(M)}}\langle f(\theta^*, \theta, \theta')\rangle_{g,\omega^{(M)}}\Big| = 0.$$

$\square$

We now complete the proof of Lemma 3.1.

*Proof of Lemma 3.1.* By Lemmas 3.3, 3.5, and 3.6, for any $M, T_0, T, T' > 0$,

$$W_1(\mathsf{P}(\theta^*, \theta^{T_0+T}, \theta^{T_0+T+T'}), \mathsf{P}^{\otimes 2}_{g_*,\omega_*;g,\omega}) \leq \varepsilon(M) + 2\sqrt{T + T'}\,\varepsilon(T_0) + C\sqrt{M}(e^{-cT} + e^{-cT'}).$$

Setting $T = T' = t$, choosing $T_0 \equiv T_0(t)$ so that $\lim_{t\to\infty} T_0(t) = \infty$ and $\lim_{t\to\infty}\sqrt{2t}\,\varepsilon(T_0(t)) = 0$, and taking $t \to \infty$ followed by $M \to \infty$, this shows

$$\lim_{t\to\infty} W_1(\mathsf{P}(\theta^*, \theta^{T_0(t)+t}, \theta^{T_0(t)+2t}), \mathsf{P}^{\otimes 2}_{g_*,\omega_*;g,\omega}) = 0.$$

In particular, we have the weak convergence in distribution of $(\theta^*, \theta^{T_0(t)+t}, \theta^{T_0(t)+2t})$ to $\mathsf{P}^{\otimes 2}_{g_*,\omega_*;g,\omega}$. Lemma 3.4 implies that $(\theta^*, \theta^{T_0(t)+t}, \theta^{T_0(t)+2t})$ is uniformly bounded in $L^4$ and hence uniformly integrable in $L^2$, so this implies

$$\lim_{t\to\infty} W_2(\mathsf{P}(\theta^*, \theta^{T_0(t)+t}, \theta^{T_0(t)+2t}), \mathsf{P}^{\otimes 2}_{g_*,\omega_*;g,\omega}) = 0. \tag{82}$$

Then, under Definition 2.4 and by definition of the law $\mathsf{P}^{\otimes 2}_{g_*,\omega_*;g,\omega}$, we have as desired

$$c_\theta^{\mathrm{tti}}(0) = \lim_{t\to\infty} C_\theta(T_0(t) + t, T_0(t) + t) = \lim_{t\to\infty}\mathbb{E}[(\theta^{T_0(t)+t})^2] = \mathbb{E}_{g_*,\omega_*}\langle\theta^2\rangle_{g,\omega},$$

$$c_\theta^{\mathrm{tti}}(\infty) = \lim_{t\to\infty} C_\theta(T_0(t) + t, T_0(t) + 2t) = \lim_{t\to\infty}\mathbb{E}[\theta^{T_0(t)+t}\theta^{T_0(t)+2t}] = \mathbb{E}_{g_*,\omega_*}\langle\theta\rangle_{g,\omega}^2,$$

$$c_\theta(*) = \lim_{t\to\infty} C_\theta(T_0(t) + t, *) = \lim_{t\to\infty}\mathbb{E}[\theta^{T_0(t)+t}\theta^*] = \mathbb{E}_{g_*,\omega_*}[\langle\theta\rangle_{g,\omega}\theta^*].$$

$\square$

## 3.2 Analysis of $\eta$-equation

We next derive from an analysis of the evolution (25) for $\{\eta^t\}_{t\geq 0}$ a representation of $c_\eta^{\mathrm{tti}}(0), c_\eta^{\mathrm{tti}}(\infty)$ in terms of $c_\theta^{\mathrm{tti}}(0), c_\theta^{\mathrm{tti}}(\infty), c_\theta(*)$.

**Lemma 3.7.** *It holds that*

$$c_\eta^{\mathrm{tti}}(0) = \frac{\delta}{\sigma^4}\left[\frac{\mathbb{E}\theta^{*2} + \sigma^2 + c_\theta^{\mathrm{tti}}(\infty) - 2c_\theta(*)}{\big(1 + \sigma^{-2}(c_\theta^{\mathrm{tti}}(0) - c_\theta^{\mathrm{tti}}(\infty))\big)^2} + \frac{c_\theta^{\mathrm{tti}}(0) - c_\theta^{\mathrm{tti}}(\infty)}{1 + \sigma^{-2}(c_\theta^{\mathrm{tti}}(0) - c_\theta^{\mathrm{tti}}(\infty))}\right], \tag{83}$$

$$c_\eta^{\mathrm{tti}}(\infty) = \frac{\delta}{\sigma^4}\frac{\mathbb{E}\theta^{*2} + \sigma^2 + c_\theta^{\mathrm{tti}}(\infty) - 2c_\theta(*)}{(1 + \sigma^{-2}(c_\theta^{\mathrm{tti}}(0) - c_\theta^{\mathrm{tti}}(\infty)))^2}, \tag{84}$$

*and in particular* $c_\eta^{\mathrm{tti}}(0) - c_\eta^{\mathrm{tti}}(\infty) < \delta/\sigma^2$.

The argument is similar to the analysis of $\{\theta^t\}_{t\geq 0}$, where we may approximate the dynamics of $\{\eta^t\}_{t\geq 0}$ at large times by a Markovian joint evolution of a system $(\eta^t, x_1^t, \ldots, x_M^t)$. Our argument here is simpler than before, as the dynamics of $(\eta^t, x_1^t, \ldots, x_M^t)$ will be linear, from which we may explicitly analyze the convergence of $\eta^t$ and show that it is independent of $M$; thus we will apply a simple Gronwall argument to bound the propagation of the discretization error $\varepsilon(M)$ over time.

### 3.2.1 Comparison with an auxiliary process

We again fix a positive integer $M$, and define $\{a_m\}_{m=0}^M$ and $\{c_m\}_{m=1}^M$ by

$$a_m = \iota + \frac{m}{\sqrt{M}} \text{ for } m = 0, \ldots, M, \qquad \frac{c_m^2}{a_m} = \mu_\theta([a_{m-1}, a_m))$$

with $\mu_\theta$ now instead of $\mu_\eta$. For convenience, let us introduce $\xi^t = \eta^t + w^* - \varepsilon$ and $v^t = -w^t + w^* - \varepsilon$, so the DMFT equation (25) for $\{\eta^t\}_{t\geq 0}$ is equivalently

$$\xi^t = -\frac{1}{\sigma^2}\int_0^t R_\theta(t,s)\xi^s \mathrm{d}s + v^t. \tag{85}$$

Here, $\{v^t\}_{t\geq 0}$ is a centered Gaussian process with covariance $\mathbb{E}[v^t v^s] = C_\theta(t,s) - C_\theta(t,*) - C_\theta(s,*) + \mathbb{E}(\theta^*)^2 + \sigma^2$. We set

$$R_\theta^{(M)}(\tau) = \sum_{m=1}^M c_m^2 e^{-a_m \tau}, \qquad C_\theta^{(M)}(t,s) = \sum_{m=1}^M \frac{c_m^2}{a_m}(e^{-a_m|t-s|} - e^{-a_m(t+s)}) + c_\theta^{\text{tti}}(\infty)$$

and define an auxiliary process $\{\xi_{M,T_0}^t\}_{t\geq 0}$ by

$$\xi_{M,T_0}^t = \xi^t \text{ for } t \in [0, T_0)$$

$$\xi_{M,T_0}^t = -\frac{1}{\sigma^2}\int_0^t R_\theta^{(M)}(t-s)\xi_{M,T_0}^s \mathrm{d}s + v_M^t \text{ for } t \geq T_0 \tag{86}$$

where $\{v_M^t\}_{t\geq 0}$ is a centered Gaussian process with covariance $\mathbb{E}[v_M^t v_M^s] = C_\theta^{(M)}(t,s) - 2c_\theta(*) + \mathbb{E}(\theta^*)^2 + \sigma^2$, defined in the probability space of $\{\xi^t\}_{t\geq 0}$. (We check in the proof of Lemma 3.10 below that this is indeed a positive-semidefinite covariance kernel.) We note that the process $\{\xi_{M,T_0}^t\}_{t\geq 0}$ may be discontinuous at $T_0$; this is inconsequential for our subsequent analysis.

**Lemma 3.8.** *For any $M, T_0, T > 0$, there exists a coupling of $\{\xi^t\}_{t\geq 0}$ and $\{\xi_{M,T_0}^t\}_{t\geq 0}$ such that*

$$\sup_{t \in [0, T_0 + T]} \mathbb{E}(\xi^t - \xi_{M,T_0}^t)^2 \leq Ce^{CT}(\varepsilon(M) + \sqrt{T}\,\varepsilon(T_0))$$

*where $\varepsilon(M)$ does not depend on $T_0, T$ and $\varepsilon(T_0)$ does not depend on $M, T$, and $\lim_{M \to \infty} \varepsilon(M) = 0$ and $\lim_{T_0 \to \infty} \varepsilon(T_0) = 0$.*

*Proof.* Applying the approximation (30) and arguments analogous to Lemma 3.2, we have that

$$\sup_{s,t \in [T_0, T_0 + T]} |\mathbb{E}[v_M^t v_M^s] - \mathbb{E}[v^t v^s]|$$

$$\leq \sup_{s,t \in [T_0, T_0 + T]} |C_\theta^{(M)}(t,s) - C_\theta(t,s)| + |c_\theta(*) - C_\theta(t,*)| + |c_\theta(*) - C_\theta(s,*)| \leq \varepsilon(M) + \varepsilon(T_0),$$

and hence there exists a coupling of $\{v_{M,T_0}^t\}_{t\geq 0}$ and $\{v^t\}_{t\geq 0}$ such that

$$\sup_{t \in [T_0, T_0 + T]} \mathbb{E}(v^t - v_M^t)^2 \leq \varepsilon(M) + \sqrt{T}\,\varepsilon(T_0).$$

We bound $\xi^t - \xi^t_{M,T_0}$ under this coupling of $\{v^t\}_{t \geq 0}$ with $\{v^t_M\}_{t \geq 0}$: Let us write $\tilde{\xi}^t = \xi^t_{M,T_0}$. We have $\xi^t = \tilde{\xi}^t$ for $t \in [0, T_0)$, while for $t \in [T_0, T_0 + T]$,

$$\mathbb{E}(\xi^t - \tilde{\xi}^t)^2 \leq 3 \Big[ \mathbb{E} \Big( \int_0^t R_\theta^{(M)}(t-s)|\xi^s - \tilde{\xi}^s| \mathrm{d}s \Big)^2 + \mathbb{E} \Big( \int_0^t |R_\theta(t,s) - R_\theta^{(M)}(t-s)||\xi^s| \mathrm{d}s \Big)^2 + \mathbb{E}(v^t - v^t_M)^2 \Big]. \quad (87)$$

From the explicit definition of $R_\theta^{(M)}(t-s)$, the first term of (87) satisfies

$$\mathbb{E} \Big( \int_0^t |R_\theta^{(M)}(t-s)||\xi^s - \tilde{\xi}^s| \mathrm{d}s \Big)^2 = \mathbb{E} \Big( \int_{T_0}^t |R_\theta^{(M)}(t-s)||\xi^s - \tilde{\xi}^s| \mathrm{d}s \Big)^2 \leq C \int_{T_0}^t \mathbb{E}(\xi^s - \tilde{\xi}^s)^2 \mathrm{d}s$$

for a constant $C > 0$. Following the argument used to bound (74), the second term of (87) is bounded by $\varepsilon(M) + C\varepsilon(t)^2$ where $\varepsilon(t) \to 0$ as $t \to \infty$, while the third term is bounded by $\varepsilon(M) + \sqrt{T}\varepsilon(T_0)$ under the above coupling. Then by Gronwall's inequality,

$$\sup_{t \in [T_0, T_0+T]} \mathbb{E}(\xi^t - \tilde{\xi}^t)^2 \leq Ce^{CT} \Big( \varepsilon(M) + \sup_{t \in [T_0, T_0+T]} \varepsilon(t)^2 + \sqrt{T}\varepsilon(T_0) \Big),$$

which implies the lemma upon adjusting $\varepsilon(T_0)$. $\qquad \square$

### 3.2.2 Convergence of the auxiliary process

**Lemma 3.9.** *The value $\sigma_Z^2 = \mathbb{E}\theta^{*2} + c_\theta^{\mathrm{tti}}(\infty) - 2c_\theta(*) + \sigma^2$ is positive.*

*Proof.* Let $\{\boldsymbol{\theta}^t\}_{t \geq 0}$ be the Langevin diffusion (7) for which the DMFT system of Theorem 2.5 is the large-$(n, d)$ limit. By Theorem 4.3 to follow,

$$C_\theta(t, s) - C_\theta(t, *) - C_\theta(s, *) + \mathbb{E}(\theta^*)^2 = \lim_{n,d \to \infty} \mathbb{E} \Big[ \frac{1}{d} \sum_{i=1}^d (\theta_i^t - \theta_i^*)(\theta_i^s - \theta_i^*) \Big]$$

Since $\{\boldsymbol{\theta}^t\}_{t \geq 0}$ is Markovian (conditional on $\mathbf{X}, \mathbf{y}, \boldsymbol{\theta}^*$), we have for all $t \geq s$ that

$$\mathbb{E} \Big[ \frac{1}{d} \sum_{i=1}^d (\theta_i^t - \theta_i^*)(\theta_i^s - \theta_i^*) \Big] = \mathbb{E} \Big[ \mathbb{E} \Big[ \frac{1}{d} \sum_{i=1}^d (\theta_i^t - \theta_i^*)(\theta_i^s - \theta_i^*) \Big| \boldsymbol{\theta}^s, \mathbf{X}, \mathbf{y}, \boldsymbol{\theta}^* \Big] \Big] = \mathbb{E} \Big[ \frac{1}{d} \sum_{i=1}^d (\theta_i^s - \theta_i^*)^2 \Big] \geq 0,$$

hence $C_\theta(t, s) - C_\theta(t, *) - C_\theta(s, *) + \mathbb{E}(\theta^*)^2 \geq 0$. Setting $s = t/2$ and taking the limit $t \to \infty$ under Definition 2.4 shows $c_\theta^{\mathrm{tti}}(\infty) - 2c_\theta(*) + \mathbb{E}\theta^{*2} \geq 0$, and the lemma follows. $\qquad \square$

**Lemma 3.10.** *Let $c = (c_1, \ldots, c_M)$, $A = \mathrm{diag}(a_1, \ldots, a_M)$, $\Lambda = A + cc^\top/\sigma^2$, and consider the 2-dimensional Gaussian law $\mathcal{N}(0, \Sigma_M)$ with*

$$\Sigma_M = \begin{pmatrix} \rho_M^2 & \kappa_M \\ \kappa_M & \rho_M^2 \end{pmatrix}, \qquad \kappa_M = \sigma_Z^2 \cdot \Big[ 1 - c^\top \Lambda^{-1} c/\sigma^2 \Big]^2, \qquad \rho_M^2 = \kappa_M + c^\top \Lambda^{-1} c,$$

*where $\sigma_Z^2 = \mathbb{E}(\theta^*)^2 + c_\theta^{\mathrm{tti}}(\infty) - 2c_\theta(*) + \sigma^2$. Then there exists an error $\varepsilon(T)$ not depending on $T_0, M$ and satisfying $\lim_{T \to \infty} \varepsilon(T) = 0$, such that for any $M, T_0, T, T' > 0$,*

$$W_2(\mathsf{P}(\xi_{M,T_0}^{T_0+T}, \xi_{M,T_0}^{T_0+T+T'}), \mathcal{N}(0, \Sigma_M)) \leq \varepsilon(T) + \varepsilon(T').$$

*Proof.* Let $z \sim \mathcal{N}(0, \sigma_Z^2)$, where $\sigma_Z^2 > 0$ by Lemma 3.9, and let $\{b_m^t\}_{t \geq 0}$ for $m = 1, \ldots, M$ be standard Brownian motions. We assume these are independent of each other and of $\{\xi^t\}_{t \in [0,T]}$. Then the law of $\{\xi_{M,T_0}^t\}_{t \geq 0}$ coincides with the marginal law of $\{\xi_{M,T_0}^t\}_{t \geq 0}$ in the joint process

$$\xi_{M,T_0}^t = \xi^t \text{ for } t \in [0, T_0)$$

$$\xi_{M,T_0}^t = \sum_{m=1}^M c_m x_m^t + z \text{ for } t \geq T_0 \qquad (88)$$

$$\mathrm{d}x_m^t = -[a_m x_m^t + c_m \xi_{M,T_0}^t/\sigma^2]\mathrm{d}t + \sqrt{2}\,\mathrm{d}b_m^t \text{ for } 1 \leq m \leq M, t \geq 0 \qquad (89)$$

with initial conditions $x_1^0 = \ldots = x_M^0 = 0$. Indeed, given $\{\xi_{M,T_0}^t\}_{t \geq 0}$, the equations (89) for $\{x_m^t\}_{t \geq 0}$ have the explicit solutions

$$x_m^t = -\frac{1}{\sigma^2} \int_0^t c_m e^{-a_m(t-s)} \xi_{M,T_0}^s \mathrm{d}s + \int_0^t e^{-a_m(t-s)} \sqrt{2}\, \mathrm{d}b_m^s,$$

and substituting this into (88) gives (86) upon identifying $v_M^t = z + \int_0^t \sum_{m=1}^M c_m e^{-a_m(t-s)} \sqrt{2}\, \mathrm{d}b_m^s$. It is direct to check that $\{v_M^t\}_{t \geq 0}$ thus defined has covariance $C_\theta^{(M)}(t,s) - 2c_\theta(*) + \mathbb{E}(\theta^*)^2 + \sigma^2$, so this coincides with the law of $\{\xi_{M,T_0}^t\}_{t \geq 0}$ defined by (86).

Let us denote $\tilde{\xi}^t = \xi_{M,T_0}^t$, $x^t = (x_1^t, \ldots, x_M^t)$, and $b^t = (b_1^t, \ldots, b_M^t)$. For $t \geq T_0$, the evolution of $(\tilde{\xi}^t, x^t) \in \mathbb{R}^{M+1}$ is a (Markovian) Ornstein-Uhlenbeck process. Substituting (88) into (89), we have

$$\mathrm{d}x^t = -[\Lambda x^t + cz/\sigma^2]\mathrm{d}t + \sqrt{2}\, \mathrm{d}b^t \text{ for } t \geq T_0$$

where $c = (c_1, \ldots, c_M)$ and $\Lambda = A + cc^\top/\sigma^2$ with $A = \mathrm{diag}(a_1, \ldots, a_M)$. This has the solution, for $t \geq T_0$,

$$x^t = e^{-\Lambda(t-T_0)} x^{T_0} + \frac{z}{\sigma^2} \Lambda^{-1}(e^{-\Lambda(t-T_0)} - I)c + \int_{T_0}^t e^{-\Lambda(t-s)} \sqrt{2}\, \mathrm{d}b^s.$$

Substituting back into (88),

$$\tilde{\xi}^t = c^\top e^{-\Lambda(t-T_0)} x^{T_0} + z\left[1 + \frac{1}{\sigma^2} c^\top \Lambda^{-1}(e^{-\Lambda(t-T_0)} - I)c\right] + \int_{T_0}^t c^\top e^{-\Lambda(t-s)} \sqrt{2}\, \mathrm{d}b^s \text{ for } t \geq T_0. \quad (90)$$

Here, we note that the equation (85) implies that $\{\xi^t\}_{t \geq 0}$ is itself a Gaussian process (given by a linear functional of $\{v^t\}_{t \geq 0}$), so $x^{T_0}$ with coordinates

$$x_m^{T_0} = \underbrace{-\frac{c_m}{\sigma^2} \int_0^{T_0} e^{-a_m(T_0-s)} \xi^s \mathrm{d}s}_{=U_m} + \underbrace{\int_0^{T_0} e^{-a_m(T_0-s)} \sqrt{2}\, \mathrm{d}b_m^s}_{=V_m} \quad (91)$$

is a Gaussian vector. Consequently, the form (90) shows that for any $T, T' > 0$, $(\tilde{\xi}^{T_0+T}, \tilde{\xi}^{T_0+T+T'})$ has a centered bivariate Gaussian law. To conclude the proof of the lemma, it suffices to show

$$|\mathbb{E}[(\tilde{\xi}^{T_0+T})^2] - \rho_M^2|,\ |\mathbb{E}[(\tilde{\xi}^{T_0+T+T'})^2] - \rho_M^2| \leq \varepsilon(T) + \varepsilon(T') \quad (92)$$

$$|\mathbb{E}\tilde{\xi}^{T_0+T}\tilde{\xi}^{T_0+T+T'}] - \kappa_M| \leq \varepsilon(T) + \varepsilon(T') \quad (93)$$

for some errors $\varepsilon(T), \varepsilon(T')$ that hold uniformly over all $M, T_0 > 0$.

For (92), we may compute from the solution (90) that

$$\mathbb{E}[(\tilde{\xi}^{T_0+T})^2] = \underbrace{c^\top e^{-\Lambda T} \mathbb{E}[x^{T_0}(x^{T_0})^\top] e^{-\Lambda T} c}_{=\mathrm{I}} + \underbrace{\sigma_Z^2 \cdot \left[1 + \frac{1}{\sigma^2} c^\top \Lambda^{-1}(e^{-\Lambda T} - I)c\right]^2 + c^\top \Lambda^{-1}(I - e^{-2\Lambda T})c}_{=\mathrm{II}}.$$

Observe that $\|\Lambda^{-1/2}c\|_2^2 \leq \sum_{m=1}^M c_m^2/a_m = \sum_{m=1}^M \mu_\theta([a_{m-1}, a_m)) \leq \mu_\theta([\iota, \infty))$. Hence $\|\Lambda^{-1/2}c\|_2 \leq C$ for a constant $C > 0$ not depending on $M$. Since also $\lambda_{\min}(\Lambda) \geq \iota > 0$, we have $\|e^{-\Lambda T}\|_{\mathrm{op}} \leq e^{-\iota T}$, so

$$|\mathrm{II} - \rho_M^2| \leq \varepsilon(T)$$

for an error $\varepsilon(T)$ not depending on $M$. To bound I, write $x^{T_0} = U + V$ where $U, V \in \mathbb{R}^M$ have the coordinates $U_m, V_m$ in (91). Then, from the bound $\mathbb{E}[(u^\top x^{T_0})^2] \leq 2\mathbb{E}[(u^\top U)^2] + 2\mathbb{E}[(u^\top V)^2]$ for each unit vector $u \in \mathbb{R}^M$, we have

$$\|\mathbb{E}[x^{T_0}(x^{T_0})^\top]\|_{\mathrm{op}} \leq 2\|\mathbb{E}[UU^\top]\|_{\mathrm{op}} + 2\|\mathbb{E}[VV^\top]\|_{\mathrm{op}}.$$

For the second term, $\|\mathbb{E}[VV^\top]\|_{\mathrm{op}} = \|\operatorname{diag}(a_m^{-1}(1 - e^{-a_m T_0}))\|_{\mathrm{op}} \le \iota^{-1}$. For the first term,

$$\|\mathbb{E}[UU^\top]\|_{\mathrm{op}} \le \mathbb{E}\|U\|_2^2 = \mathbb{E}\sum_{m=1}^M \frac{c_m^2}{\sigma^4}\left(\int_0^{T_0} e^{-a_m(T_0-s)}\xi^s \mathrm{d}s\right)^2$$

$$\le \sum_{m=1}^M \frac{c_m^2}{\sigma^4}\int_0^{T_0} e^{-a_m(T_0-s)}\mathrm{d}s \cdot \int_0^{T_0} e^{-a_m(T_0-s)}\mathbb{E}(\xi^s)^2 \mathrm{d}s.$$

Noting that $\mathbb{E}(\xi^t)^2 = (\sigma^4/\delta)C_\eta(t,t) \le C$ for all $t \ge 0$ under Definition 2.4, this gives $\|\mathbb{E}[UU^\top]\|_{\mathrm{op}} \le C'\sum_{m=1}^M c_m^2/a_m^2 \le C'\mu_\theta([\iota,\infty))/\iota$. Combining these bounds shows $\|\mathbb{E}[x^{T_0}(x^{T_0})^\top]\|_{\mathrm{op}} \le C$ for a constant $C > 0$ not depending on $M, T_0$. Then, combining with the previous bounds $\|\Lambda^{-1/2}u\|_2 \le C$ and $\lambda_{\min}(\Lambda) \ge \iota$, this shows $|\mathrm{I}| \le \varepsilon(T)$, so $|\mathbb{E}[(\tilde\xi^{T_0+T})^2] - \rho_M^2| \le \varepsilon(T)$. The bound for $\mathbb{E}[(\tilde\xi^{T_0+T+T'})^2]$ in (92) holds similarly.

For (93), we may compute similarly from (90)

$$\mathbb{E}[(\tilde\xi^{T_0+T})\tilde\xi^{T_0+T+T'}] = u^\top e^{-\Lambda T}\mathbb{E}[x^{T_0}(x^{T_0})^\top]e^{-\Lambda(T+T')}u$$

$$+ \sigma_Z^2 \cdot \left[1 + \frac{1}{\sigma^2}u^\top\Lambda^{-1}(e^{-\Lambda(T+T')} - I)u\right]\left[1 + \frac{1}{\sigma^2}u^\top\Lambda^{-1}(e^{-\Lambda T} - I)u\right]$$

$$+ u^\top\Lambda^{-1}(e^{-\Lambda T'} - e^{-\Lambda(2T+T')})u,$$

and the arguments to show (93) from this form are the same as above. $\qquad\square$

**Lemma 3.11.** *Consider the 2-dimensional Gaussian law $\mathcal{N}(0,\Sigma_\infty)$ with*

$$\Sigma_\infty = \begin{pmatrix} \rho_\infty^2 & \kappa_\infty \\ \kappa_\infty & \rho_\infty^2 \end{pmatrix}, \quad \kappa_\infty = \frac{\mathbb{E}\theta^{*2} + \sigma^2 + c_\theta^{\mathrm{tti}}(\infty) - 2c_\theta(*)}{(1 + \sigma^{-2}(c_\theta^{\mathrm{tti}}(0) - c_\theta^{\mathrm{tti}}(\infty))^2}, \quad \rho_\infty^2 = \kappa_\infty + \frac{c_\theta^{\mathrm{tti}}(0) - c_\theta^{\mathrm{tti}}(\infty)}{1 + \sigma^{-2}(c_\theta^{\mathrm{tti}}(0) - c_\theta^{\mathrm{tti}}(\infty))}.$$

*Then $\lim_{M\to\infty}\|\Sigma_M - \Sigma_\infty\|_{\mathrm{op}} = 0$.*

*Proof.* This follows from noting that $c^\top\Lambda^{-1}c = \frac{\sum_{m=1}^M c_m^2/a_m}{1 + \sigma^{-2}\sum_{m=1}^M c_m^2/a_m}$ via the Sherman-Morrison identity, and $\sum_{m=1}^M c_m^2/a_m \to \int_\iota^\infty \mu_\theta(\mathrm{d}a) = c_\theta^{\mathrm{tti}}(0) - c_\theta^{\mathrm{tti}}(\infty)$ as $M \to \infty$. $\qquad\square$

We now complete the proof of Lemma 3.7.

*Proof of Lemma 3.7.* By Lemmas 3.8, 3.10, and 3.11, it holds that

$$W_2(\mathsf{P}(\xi^{T_0+T},\xi^{T_0+T+T'}),\mathcal{N}(0,\Sigma_\infty)) \le Ce^{C(T+T')}(\varepsilon(M) + \sqrt{T+T'}\,\varepsilon(T_0)) + \varepsilon(T) + \varepsilon(T') + \varepsilon(M).$$

Taking first the limit $M \to \infty$, then choosing $T = T' = t$ and $T_0 \equiv T_0(t)$ such that $\lim_{t\to\infty} T_0(t) = \infty$ and $\lim_{t\to\infty} e^{2Ct}\sqrt{2t}\,\varepsilon(T_0(t)) = 0$ and taking $t \to \infty$, this shows $W_2(\mathsf{P}(\xi^{T_0(t)+t},\xi^{T_0(t)+2t}),\mathcal{N}(0,\Sigma_\infty)) \to 0$ as $t \to \infty$. Under Definition 2.4, this implies

$$\frac{\sigma^4}{\delta}c_\eta^{\mathrm{tti}}(0) = \lim_{t\to\infty}\frac{\sigma^4}{\delta}C_\eta(T_0(t)+t, T_0(t)+t) = \lim_{t\to\infty}\mathbb{E}[(\xi^{T_0(t)+t})^2] = \rho_\infty^2,$$

$$\frac{\sigma^4}{\delta}c_\eta^{\mathrm{tti}}(\infty) = \lim_{t\to\infty}\frac{\sigma^4}{\delta}C_\eta(T_0(t)+t, T_0(t)+2t) = \lim_{t\to\infty}\mathbb{E}[\xi^{T_0(t)+t}\xi^{T_0(t)+2t}] = \kappa_\infty.$$

This shows the desired forms of $c_\eta^{\mathrm{tti}}(0)$ and $c_\eta^{\mathrm{tti}}(\infty)$, and we have also from these forms that

$$c_\eta^{\mathrm{tti}}(0) - c_\eta^{\mathrm{tti}}(\infty) = \frac{\delta}{\sigma^2}\left[\frac{\sigma^{-2}(c_\theta^{\mathrm{tti}}(0) - c_\theta^{\mathrm{tti}}(\infty))}{1 + \sigma^{-2}(c_\theta^{\mathrm{tti}}(0) - c_\theta^{\mathrm{tti}}(\infty))}\right] < \frac{\delta}{\sigma^2}.$$

$\qquad\square$

34

## 3.3 Completing the proof

*Proof of Theorem 2.5.* By Lemmas 3.1 and 3.7, we have five equations (60), (83), (84) for the five variables $c_\theta^{\mathrm{tti}}(0), c_\theta^{\mathrm{tti}}(\infty), c_\theta(*), c_\eta^{\mathrm{tti}}(0), c_\eta^{\mathrm{tti}}(\infty)$. Defining $\mathrm{mse}, \mathrm{mse}_*$ by (42), these equations show

$$\omega = \frac{\delta}{\sigma^2} - (c_\eta^{\mathrm{tti}}(0) - c_\eta^{\mathrm{tti}}(\infty)) = \frac{\delta}{\sigma^2 + (c_\theta^{\mathrm{tti}}(0) - c_\theta^{\mathrm{tti}}(\infty))} = \frac{\delta}{\sigma^2 + \mathrm{mse}},$$

$$\omega_* = \frac{\omega^2}{c_\eta^{\mathrm{tti}}(\infty)} = \frac{\delta}{\mathbb{E}\theta^{*2} + \sigma^2 + c_\theta(\infty) - 2c_\theta(*)} = \frac{\delta}{\sigma^2 + \mathrm{mse}_*},$$

as well as

$$\mathrm{mse} = c_\theta^{\mathrm{tti}}(0) - c_\theta^{\mathrm{tti}}(\infty) = \mathbb{E}_{g_*,\omega_*}[\langle\theta^2\rangle_{g,\omega} - \langle\theta\rangle_{g,\omega}^2] = \mathbb{E}_{g_*,\omega_*}\langle(\theta - \langle\theta\rangle_{g,\omega})^2\rangle_{g,\omega},$$

$$\mathrm{mse}_* = \mathbb{E}\theta^{*2} - 2\mathbb{E}_{g_*,\omega_*}[\theta^*\langle\theta\rangle_{g,\omega}] + \mathbb{E}_{g_*,\omega_*}\langle\theta\rangle_{g,\omega}^2 = \mathbb{E}_{g_*,\omega_*}(\theta^* - \langle\theta\rangle_{g,\omega})^2.$$

This verifies that the fixed-point equations (43) hold, where it is clear that $\omega, \omega_*$ are uniquely defined from $\mathrm{mse}, \mathrm{mse}_*$ via (43). Defining $\mathrm{ymse}, \mathrm{ymse}_*$ by (42), we have also from the above forms of $\omega, \omega_*$ that

$$\mathrm{ymse} = \frac{\sigma^4}{\delta}(c_\eta^{\mathrm{tti}}(0) - c_\eta^{\mathrm{tti}}(\infty)) = \sigma^2\Big(1 - \frac{\omega\sigma^2}{\delta}\Big),$$

$$\mathrm{ymse}^* = \frac{\sigma^4}{\delta}(2c_\eta^{\mathrm{tti}}(0) - c_\eta^{\mathrm{tti}}(\infty)) - \sigma^2 = \sigma^2 + \frac{\omega\sigma^4}{\delta}\Big(\frac{\omega}{\omega_*} - 2\Big),$$

verifying (44). Finally, the statement (45) is a consequence of (82) shown in the proof of Lemma 3.1. □

# 4 Analysis of fixed-prior Langevin dynamics under LSI

In this section, we prove Theorem 2.9 and Corollary 2.10 that verify Definition 2.4 and deduce the replica-symmetric limits for the Bayes-optimal mean-squared-errors and free energy, under Assumption 2.7 of a log-Sobolev inequality (LSI) for the posterior law.

## 4.1 Preliminaries

### 4.1.1 Properties of Langevin dynamics

We review in this section two general results on a Langevin diffusion of the form

$$\mathrm{d}\boldsymbol{\theta}^t = \nabla U(\boldsymbol{\theta}^t)\mathrm{d}t + \sqrt{2}\,\mathrm{d}\mathbf{b}^t \tag{94}$$

with an equilibrium measure $e^{U(\boldsymbol{\theta})}$. The first is a fluctuation-dissipation relation for its correlation and response functions at equilibrium, and the second is a Bismut-Elworthy-Li representation for the spatial derivative of its Markov semigroup. For bounded observables, similar fluctuation-dissipation theorems have been stated and shown in [77, 78] and Bismut-Elworthy-Li formulae in [79, 80]. We give versions of these results here for a class of unbounded observables which may have linear growth

$$\mathcal{A} = \{f \in C^2(\mathbb{R}^d, \mathbb{R}) : \nabla f, \nabla^2 f \text{ are globally bounded}\},$$

and a class of drift coefficients

$$\mathcal{B} = \{U \in C^3(\mathbb{R}^d, \mathbb{R}) : \nabla^2 U, \nabla^3 U \text{ are globally bounded and Hölder continuous}\}, \tag{95}$$

drawing upon some analyses of our companion work [27, Appendix A].

We write

$$P_t f(\boldsymbol{\theta}) = \mathbb{E}[f(\boldsymbol{\theta}^t) \mid \boldsymbol{\theta}^0 = \boldsymbol{\theta}], \qquad \mathrm{L}f(\boldsymbol{\theta}) = \nabla U^\top \nabla f(\boldsymbol{\theta}) + \mathrm{Tr}\,\nabla^2 f(\boldsymbol{\theta})$$

for the Markov semigroup and infinitesimal generator associated to (94). It is shown in [27, Proposition A.2] that

$$f \in \mathcal{A}, U \in \mathcal{B} \quad \Rightarrow \quad \nabla P_t f(\boldsymbol{\theta}), \nabla^2 P_t f(\boldsymbol{\theta}) \text{ are uniformly bounded over } t \in [0, T], \boldsymbol{\theta} \in \mathbb{R}^d \tag{96}$$

for any fixed $T > 0$. In particular, $P_t f \in \mathcal{A}$ for each fixed $t > 0$.

**Lemma 4.1.** *Suppose $U \in \mathcal{B}$, and (94) has the unique stationary distribution $q(\boldsymbol{\theta}) = e^{U(\boldsymbol{\theta})}$ with finite third moments. Let $\{\boldsymbol{\theta}^t\}_{t \geq 0}$ be the solution to (94) with initial condition $\boldsymbol{\theta}^0 = \mathbf{x}$, and let $A \in \mathcal{A}$ and $B \in \mathcal{B}$.*

(a) *Define the response function $R^{\mathbf{x}}_{AB}(t,s) = P_s(\nabla B^\top \nabla P_{t-s}A)(\mathbf{x})$. Then $R^{\mathbf{x}}_{AB}(t,s)$ satisfies the following condition: Fix any continuous bounded function $h : [0, \infty) \to \mathbb{R}$. For each $\varepsilon > 0$, let $\{\boldsymbol{\theta}^{t,\varepsilon}\}_{t \geq 0}$ denote the solution of the perturbed dynamics*

$$d\boldsymbol{\theta}^{t,\varepsilon} = \nabla[U(\boldsymbol{\theta}^{t,\varepsilon}) + \varepsilon h(t)B(\boldsymbol{\theta}^{t,\varepsilon})]dt + \sqrt{2}\, d\mathbf{b}^t$$

*with the same initial condition $\boldsymbol{\theta}^{0,\varepsilon} = \mathbf{x}$. Then for any $t > 0$,*

$$\lim_{\varepsilon \to 0} \frac{1}{\varepsilon}\left(\mathbb{E}[A(\boldsymbol{\theta}^{t,\varepsilon}) \mid \boldsymbol{\theta}^{0,\varepsilon} = \mathbf{x}] - \mathbb{E}[A(\boldsymbol{\theta}^t) \mid \boldsymbol{\theta}^0 = \mathbf{x}]\right) = \int_0^t R^{\mathbf{x}}_{AB}(t,s)h(s)ds.$$

(b) *Define the correlation function $C^{\mathbf{x}}_{AB}(t,s) = \mathbb{E}[A(\boldsymbol{\theta}^t)B(\boldsymbol{\theta}^s) \mid \boldsymbol{\theta}^0 = \mathbf{x}]$. Then for any $t \geq s \geq 0$, averaging over an initial condition $\mathbf{x} \sim q$ drawn from the stationary distribution,*

$$\partial_t \mathbb{E}_{\mathbf{x} \sim q} C^{\mathbf{x}}_{AB}(t,s) = -\mathbb{E}_{\mathbf{x} \sim q} R^{\mathbf{x}}_{AB}(t,s).$$

*Proof.* Part (a) is an application of [27, Proposition A.1] of our companion paper (specialized to this setting of dynamics with a fixed and non-adaptive prior).

For part (b), we will use also from [27, Proposition A.2] that for $A \in \mathcal{A}$, we have $\partial_t P_t A = L P_t A$. Since the dynamics (103) are Markovian with stationary distribution $q(\boldsymbol{\theta})$, we have

$$\mathbb{E}_{\mathbf{x} \sim q} C^{\mathbf{x}}_{AB}(t,s) = \mathbb{E}_{\mathbf{x} \sim q}[\mathbb{E}[A(\boldsymbol{\theta}^{t-s}) \mid \boldsymbol{\theta}^0 = \mathbf{x}]B(\mathbf{x})] = \mathbb{E}_{\mathbf{x} \sim q}(B \cdot P_{t-s}A)[\mathbf{x}].$$

To differentiate under the integral in $t$, note that $\partial_t(B \cdot P_{t-s}A) = B \cdot L P_{t-s}A$. By the uniform boundedness of $\nabla P_t A, \nabla^2 P_t A$ over $t \in [0,T]$, the Lipschitz-continuity of $\nabla B, \nabla U$, and finiteness of third moments of $q$, we have that $(B \cdot L P_t A)[\mathbf{x}]$ is uniformly integrable with respect to $\mathbf{x} \sim q$ over $t \in [0,T]$. Thus dominated convergence applies to show

$$\partial_t \mathbb{E}_{\mathbf{x} \sim q} C^{\mathbf{x}}_{AB}(t,s) = \partial_t \mathbb{E}_{\mathbf{x} \sim q}(B \cdot P_{t-s}A)[\mathbf{x}] = \mathbb{E}_{\mathbf{x} \sim q}(B \cdot L P_{t-s}A)[\mathbf{x}].$$

On the other hand, using also that both $\nabla B^\top \nabla P_t A$ and $B \cdot L P_t A$ are integrable with respect to $\mathbf{x} \sim q$, we have via integration-by-parts

$$\mathbb{E}_{\mathbf{x} \sim q} R^{\mathbf{x}}_{AB}(t,s) = \mathbb{E}_{\mathbf{x} \sim q}(\nabla B^\top \nabla P_{t-s}A)[\mathbf{x}] = \int q(\boldsymbol{\theta})(\nabla B^\top \nabla P_{t-s}A)[\boldsymbol{\theta}]d\boldsymbol{\theta}$$

$$= -\int B(\boldsymbol{\theta})\sum_{j=1}^d \partial_j[q\, \partial_j(P_{t-s}A)](\boldsymbol{\theta})d\boldsymbol{\theta}$$

$$= -\int B(\boldsymbol{\theta})\left[q\, \mathrm{Tr}\, \nabla^2(P_{t-s}A) + \nabla(P_{t-s}A)^\top \nabla q\right](\boldsymbol{\theta})d\boldsymbol{\theta}$$

$$= -\int q(\boldsymbol{\theta})B(\boldsymbol{\theta})\left[\mathrm{Tr}\, \nabla^2(P_{t-s}A) + \nabla(P_{t-s}A)^\top \nabla \log q\right](\boldsymbol{\theta})d\boldsymbol{\theta}$$

$$= -\mathbb{E}_{\mathbf{x} \sim q}(B \cdot L P_{t-s}A)[\mathbf{x}].$$

$\square$

**Lemma 4.2.** *Suppose $U \in \mathcal{B}$, and consider the solution $(\boldsymbol{\theta}^t, \mathbf{V}^t) \in \mathbb{R}^d \times \mathbb{R}^{d \times d}$ to*

$$d\boldsymbol{\theta}^t = \nabla U(\boldsymbol{\theta}^t)dt + \sqrt{2}\, d\mathbf{b}^t, \qquad d\mathbf{V}^t = [\nabla^2 U(\boldsymbol{\theta}^t)]\mathbf{V}^t dt \qquad (97)$$

*with initial condition $(\boldsymbol{\theta}^0, \mathbf{V}^0) = (\mathbf{x}, \mathbf{I})$, adapted to the canonical filtration of the Brownian motion $\{\mathbf{b}^t\}_{t \geq 0}$. Then for any $f \in \mathcal{A}$ and any $t > 0$,*

$$\nabla P_t f(\mathbf{x}) = \mathbb{E}[\mathbf{V}^{t\top} \nabla f(\boldsymbol{\theta}^t) \mid (\boldsymbol{\theta}^0, \mathbf{V}^0) = (\mathbf{x}, \mathbf{I})] \qquad (98)$$

$$= \frac{1}{t\sqrt{2}}\mathbb{E}\left[f(\boldsymbol{\theta}^t)\int_0^t \mathbf{V}^{s\top}d\mathbf{b}^s \,\middle|\, (\boldsymbol{\theta}^0, \mathbf{V}^0) = (\mathbf{x}, \mathbf{I})\right] \qquad (99)$$

*Proof.* The first identity (98) is the statement of [27, Eq. (184)] (again specialized to this setting of dynamics with a fixed and non-adaptive prior).

For the second identity (99), we use from [27, Proposition A.2] that for $f \in \mathcal{A}$ and any fixed $t \geq 0$, $(s, \boldsymbol{\theta}) \mapsto P_{t-s}f(\boldsymbol{\theta})$ is $C^1$ in $s \in [0, t]$ and $C^2$ in $\boldsymbol{\theta}$, with $\partial_s P_{t-s}f(\boldsymbol{\theta}) = -\mathrm{L}P_{t-s}f(\boldsymbol{\theta})$. Then Itô's formula applied to $g(s, \boldsymbol{\theta}) = P_{t-s}f(\boldsymbol{\theta})$ gives

$$
\begin{aligned}
f(\boldsymbol{\theta}^t) = g(t, \boldsymbol{\theta}^t) &= g(0, \boldsymbol{\theta}^0) + \int_0^t \partial_s g(s, \boldsymbol{\theta}^s) \mathrm{d}s + \int_0^t \nabla_{\boldsymbol{\theta}} g(s, \boldsymbol{\theta}^s)^\top \mathrm{d}\boldsymbol{\theta}^s + \int_0^t \mathrm{Tr}\, \nabla_{\boldsymbol{\theta}}^2 g(s, \boldsymbol{\theta}^s) \mathrm{d}s \\
&= P_t f(\mathbf{x}) + \int_0^t (\partial_s + \mathrm{L}) P_{t-s}f(\boldsymbol{\theta}^s) \mathrm{d}s + \sqrt{2} \int_0^t \nabla P_{t-s}f(\boldsymbol{\theta}^s)^\top \mathrm{d}\mathbf{b}^s \\
&= P_t f(\mathbf{x}) + \sqrt{2} \int_0^t \nabla P_{t-s}f(\boldsymbol{\theta}^s)^\top \mathrm{d}\mathbf{b}^s.
\end{aligned}
$$

Since $\nabla^2 U$ is bounded, $\{\mathbf{V}^t\}_{t \in [0,T]}$ is bounded over finite time horizons, so $\int_0^t \mathbf{V}^{s\top} \mathrm{d}\mathbf{b}^s$ is a martingale. Multiplying both sides by this martingale and taking expectations gives

$$
\mathbb{E}\left[ f(\boldsymbol{\theta}^t) \int_0^t \mathbf{V}^{s\top} \mathrm{d}\mathbf{b}^s \,\Big|\, (\boldsymbol{\theta}^0, \mathbf{V}^0) = (\mathbf{x}, \mathbf{I}) \right] = \sqrt{2} \int_0^t \mathbb{E}\left[ \mathbf{V}^{s\top} \nabla P_{t-s}f(\boldsymbol{\theta}^s) \mid (\boldsymbol{\theta}^0, \mathbf{V}^0) = (\mathbf{x}, \mathbf{I}) \right] \mathrm{d}s.
$$

Since $P_t f \in \mathcal{A}$, we may apply (98) with $P_{t-s}f$ in place of $f$ to get

$$
\int_0^t \mathbb{E}\left[ \mathbf{V}^{s\top} \nabla P_{t-s}f(\boldsymbol{\theta}^s) \mid (\boldsymbol{\theta}^0, \mathbf{V}^0) = (\mathbf{x}, \mathbf{I}) \right] \mathrm{d}s = \int_0^t \nabla P_s (P_{t-s}f)(\mathbf{x}) \mathrm{d}s = t \cdot \nabla P_t f(\mathbf{x}).
$$

Substituting above and rearranging shows (99). $\qquad\square$

### 4.1.2 Interpretation of the DMFT correlation and response

We remark that under Assumption 2.2(a), the log-posterior density $\log \mathsf{P}_g(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y})$ belongs to the function class $\mathcal{B}$, and $\mathsf{P}_g(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y})$ is the unique stationary distribution of (7). Fixing $\mathbf{X}, \mathbf{y}, \boldsymbol{\theta}^*$, consider the coordinate functions

$$
e_j(\boldsymbol{\theta}) = \theta_j, \qquad e_j^*(\boldsymbol{\theta}) = \theta_j^*, \qquad x_i(\boldsymbol{\theta}) = \frac{\sqrt{\delta}}{\sigma^2}([\mathbf{X}\boldsymbol{\theta}]_i - y_i). \tag{100}
$$

(Here, $e_j^*$ is a constant function not depending on $\boldsymbol{\theta}$.) We define their associated correlation and response matrices

$$
\begin{aligned}
\mathbf{C}_\theta(t, s) = (C_{e_j e_k}^{\boldsymbol{\theta}_0}(t, s))_{j,k=1}^d, \quad \mathbf{C}_\theta(t, *) = (C_{e_j e_k^*}^{\boldsymbol{\theta}_0}(t, 0))_{j,k=1}^d, \quad \mathbf{R}_\theta(t, s) = (R_{e_j e_k}^{\boldsymbol{\theta}_0}(t, s))_{j,k=1}^d \\
\mathbf{C}_\eta(t, s) = (C_{x_j x_k}^{\boldsymbol{\theta}_0}(t, s))_{j,k=1}^n, \qquad \mathbf{R}_\eta(t, s) = (R_{x_j x_k}^{\boldsymbol{\theta}_0}(t, s))_{j,k=1}^n
\end{aligned} \tag{101}
$$

where $C_{AB}^{\boldsymbol{\theta}_0}(t, s)$ and $R_{AB}^{\boldsymbol{\theta}_0}(t, s)$ are the correlation and response functions as defined in Lemma 4.1 for these coordinate functions, under the dynamics (7) with fixed prior $g(\cdot)$ and the given initial condition $\boldsymbol{\theta}^0$ of Assumption 2.1.

The following result is a direct application of [27, Theorem 2.8].

**Theorem 4.3** ( [27]). *Suppose Assumptions 2.1 and 2.2(a) hold, and let $C_\theta, C_\eta, R_\theta, R_\eta$ be the correlation and response functions of the DMFT system in Theorem 2.3(a) approximating the dynamics (7). Then almost surely as $n, d \to \infty$,*

$$
d^{-1} \mathrm{Tr}\, \mathbf{C}_\theta(t, s) \to C_\theta(t, s), \qquad d^{-1} \mathrm{Tr}\, \mathbf{C}_\theta(t, *) \to C_\theta(t, *), \qquad n^{-1} \mathrm{Tr}\, \mathbf{C}_\eta(t, s) \to C_\eta(t, s)
$$

$$
d^{-1} \mathrm{Tr}\, \mathbf{R}_\theta(t, s) \to R_\theta(t, s), \qquad n^{-1} \mathrm{Tr}\, \mathbf{R}_\eta(t, s) \to R_\eta(t, s).
$$

37

## 4.2 Posterior bounds and Wasserstein-2 convergence

Fixing the prior $g(\cdot)$ and the data $(\mathbf{X}, \mathbf{y})$, let us write for convenience

$$q(\boldsymbol{\theta}) = \mathsf{P}_g(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y}) \propto \exp\left(-\frac{1}{2\sigma^2}\|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 + \sum_{j=1}^d \log g(\theta_j)\right) \tag{102}$$

for the posterior density. The Langevin diffusion (7) with this fixed prior is then

$$d\boldsymbol{\theta}^t = \nabla \log q(\boldsymbol{\theta}^t)dt + \sqrt{2}\,d\mathbf{b}^t. \tag{103}$$

We will use the notations

$$\langle f(\boldsymbol{\theta})\rangle = \mathbb{E}_{\boldsymbol{\theta}\sim q}[f(\boldsymbol{\theta})], \qquad P_t f(\mathbf{x}) = \langle f(\boldsymbol{\theta}^t)\rangle_{\mathbf{x}} = \mathbb{E}[f(\boldsymbol{\theta}^t) \mid \boldsymbol{\theta}^0 = \mathbf{x}]$$

where the former is an average under the posterior law $q(\cdot)$ conditional on $\mathbf{X}, \mathbf{y}$, and the latter is an average over $\{\boldsymbol{\theta}^t\}_{t\geq 0}$ solving (103) conditional on $\mathbf{X}, \mathbf{y}$ and also the initial condition $\boldsymbol{\theta}^0 = \mathbf{x}$. We write as shorthand

$$P_t(\mathbf{x}) = \langle \boldsymbol{\theta}^t\rangle_{\mathbf{x}} = (P_t e_1, \ldots, P_t e_d)[\mathbf{x}] \in \mathbb{R}^d.$$

We reserve $\langle f(\boldsymbol{\theta}^t)\rangle$ for the expectation conditional on $\mathbf{X}, \mathbf{y}$ but averaging also over $\boldsymbol{\theta}^0$.

For constants $C_0, C_{\mathrm{LSI}} > 0$, define the $(\mathbf{X}, \boldsymbol{\theta}^*, \boldsymbol{\varepsilon})$-dependent event

$$\mathcal{E}(C_0, C_{\mathrm{LSI}}) = \left\{\|\mathbf{X}\|_{\mathrm{op}} \leq C_0, \ \|\boldsymbol{\theta}^*\|_2^2, \|\boldsymbol{\varepsilon}\|_2^2 \leq C_0 d, \text{ the LSI (46) holds for } q(\boldsymbol{\theta})\right\}. \tag{104}$$

Note that under Assumptions 2.1 and 2.7(a), this event holds almost surely for all large $n, d$ for some sufficiently large choices of constants $C_0, C_{\mathrm{LSI}} > 0$. All subsequent constants $C, C', c, c' > 0$ in this section may change from instance to instance, and are dimension-free and depend only on

$$C_0, C_{\mathrm{LSI}} \text{ above, } \delta, \sigma^2, g_* \text{ of Assumption 2.1, } C, c_0, r_0 \text{ of Assumption 2.2(a), and } \log g(0). \tag{105}$$

We record the following elementary bounds for the posterior expectation $\langle f(\boldsymbol{\theta})\rangle = \mathbb{E}_{\boldsymbol{\theta}\sim q} f(\boldsymbol{\theta})$.

**Lemma 4.4.** *Suppose Assumption 2.2(a) holds. Then on the event where $\|\mathbf{X}\|_{\mathrm{op}} \leq C_0$, there exists a constant $C > 0$ for which*

$$\langle \|\boldsymbol{\theta}\|_2^2\rangle \leq C(d + \|\mathbf{y}\|_2^2), \tag{106}$$

$$\langle \|\nabla \log q(\boldsymbol{\theta})\|_2^2\rangle \leq C(d + \|\mathbf{y}\|_2^2), \tag{107}$$

$$\|\nabla^2 \log q(\boldsymbol{\theta})\|_{\mathrm{op}} \leq C. \tag{108}$$

*In particular, on $\mathcal{E}(C_0, C_{\mathrm{LSI}})$, for a constant $C' > 0$ we have $\langle \|\boldsymbol{\theta}\|_2^2\rangle \leq C'd$ and $\langle \|\nabla \log q(\boldsymbol{\theta})\|_2^2\rangle \leq C'd$.*

*Proof.* (108) is immediate from the form of $\log q(\boldsymbol{\theta})$, the bound $\|\mathbf{X}\|_{\mathrm{op}} \leq C_0$, and Assumption 2.2(a).

For (106), write $\mathbb{E}_g, \mathbb{P}_g$ for the expectation and probability over the prior $\theta \sim g$ and $\theta_j \overset{iid}{\sim} g$. We note that under Assumption 2.2(a), we have

$$\log g(\theta) = \log g(0) + \theta(\log g)'(0) + \int_0^\theta \int_0^x (\log g)''(u)du\,dx \leq C(1 + |\theta|) - (c_0/2)(|\theta| - r_0)^2 \leq C' - c'\theta^2$$

for some constants $C, C', c' > 0$ depending only on the constants of Assumption 2.2(a) and on $\log g(0)$. Then $g$ is subgaussian, and for some constants $C, c > 0$ (c.f. [81, Eq. (3.1)])

$$\mathbb{E}_g\|\boldsymbol{\theta}\|_2^2 \leq Cd, \qquad \mathbb{P}_g[\|\boldsymbol{\theta}\|_2^2 - \mathbb{E}_g\|\boldsymbol{\theta}\|_2^2 \geq du] \leq Ce^{-cdu} \text{ for all } u \geq 1. \tag{109}$$

Write

$$q(\boldsymbol{\theta}) = \frac{1}{Z}\exp\left(-\frac{\|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2}{2\sigma^2}\right)\prod_{j=1}^d g(\theta_j), \qquad Z = \mathbb{E}_g\left[\exp\left(-\frac{\|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2}{2\sigma^2}\right)\right].$$

38

We have by Jensen's inequality $-\log Z \le \mathbb{E}_g[\|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2/2\sigma^2] \le C(d + \|\mathbf{y}\|_2^2 + \mathbb{E}_g\|\boldsymbol{\theta}\|_2^2) \le C'(d + \|\mathbf{y}\|_2^2)$. Then for any $M > 0$, also bounding the exponential from above by 1,

$$\left\langle \|\boldsymbol{\theta}\|_2^2 \mathbf{1}\{\|\boldsymbol{\theta}\|_2^2 \ge M\}\right\rangle \le \frac{1}{Z}\,\mathbb{E}_g\left[\|\boldsymbol{\theta}\|_2^2 \mathbf{1}\{\|\boldsymbol{\theta}\|_2^2 \ge M\}\right] \le e^{C'(d+\|\mathbf{y}\|_2^2)}\mathbb{E}_g\left[\|\boldsymbol{\theta}\|_2^2\mathbf{1}\{\|\boldsymbol{\theta}\|_2^2 \ge M\}\right].$$

Integrating the tail bound (109) shows that this is less than $d + \|\mathbf{y}\|_2^2$ for $M = C(d + \|\mathbf{y}\|_2^2)$ and a sufficiently large choice of constant $C > 0$. Thus

$$\langle\|\boldsymbol{\theta}\|_2^2\rangle \le M + \left\langle\|\boldsymbol{\theta}\|_2^2\mathbf{1}\{\|\boldsymbol{\theta}\|_2^2 \ge M\}\right\rangle \le C'(d + \|\mathbf{y}\|_2^2).$$

This shows (106). Since $\nabla \log q(\boldsymbol{\theta})$ is $C$-Lipschitz by (108), and $\|\nabla \log q(0)\|_2^2 \le 2\|\mathbf{X}^\top\mathbf{y}/\sigma^2\|_2^2 + 2d \cdot (\log g)'(0)^2 \le C(d + \|\mathbf{y}\|_2^2)$, the statement (107) follows from (106). $\qquad\square$

**Remark 4.5.** In a later proof, we will require that (106) holds in a form

$$\langle\|\boldsymbol{\theta}\|_2^2\rangle \le Cd + (C/\sigma^2)\|\mathbf{y}\|_2^2 \tag{110}$$

for all large noise variances $\sigma^2 > 0$, where $C > 0$ is a constant not depending on $\sigma^2$. This may be seen from the above arguments: Writing now $C, C' > 0$ for constants not depending on $\sigma^2$, the above shows $-\log Z \le (C'/\sigma^2)(d + \|\mathbf{y}\|_2^2)$, and hence $\langle\|\boldsymbol{\theta}\|_2^2\mathbf{1}\{\|\boldsymbol{\theta}\|_2^2 \ge M\}\rangle \le d + \|\mathbf{y}\|_2^2/\sigma^2$ for $M = Cd + (C/\sigma^2)\|\mathbf{y}\|_2^2$ with a sufficiently large choice of constant $C > 0$.

**Lemma 4.6.** *Suppose Assumption 2.2(a) holds. Let $\{\boldsymbol{\theta}^t\}_{t\ge 0}$ be the solution to (9) with initial condition $\boldsymbol{\theta}^0 \sim q_0$, let $q_t(\boldsymbol{\theta}^t)$ be the law of $\boldsymbol{\theta}^t$, and let $W_2(\cdot)$ the Wasserstein-2 distance, all conditional on $\mathbf{X}, \boldsymbol{\theta}^*, \boldsymbol{\varepsilon}$ (and averaging over $\boldsymbol{\theta}^0$). Then on the event $\mathcal{E}(C_0, C_{\mathrm{LSI}})$, there exists a constant $C > 0$ such that*

$$W_2(q_t, q) \le Ce^{-(2/C_{\mathrm{LSI}})t}\, W_2(q_0, q) \text{ for all } t \ge 0. \tag{111}$$

*Proof.* For $t \in [0, 1]$ we may apply a simple synchronous coupling and Grönwall argument: Let $\{\boldsymbol{\theta}^t\}_{t\ge 0}$ and $\{\tilde{\boldsymbol{\theta}}^t\}_{t\ge 0}$ be the solutions of (9) with initial conditions $\boldsymbol{\theta}^0 \sim q_0$ and $\tilde{\boldsymbol{\theta}}^0 \sim q$, coupled by the same Brownian motion. Then $\frac{\mathrm{d}}{\mathrm{d}t}\|(\boldsymbol{\theta}^t - \tilde{\boldsymbol{\theta}}^t)\|_2 \le \|\frac{\mathrm{d}}{\mathrm{d}t}(\boldsymbol{\theta}^t - \tilde{\boldsymbol{\theta}}^t)\|_2 = \|\nabla\log q(\boldsymbol{\theta}^t) - \nabla\log q(\tilde{\boldsymbol{\theta}}^t)\|_2 \le C\|\boldsymbol{\theta}^t - \tilde{\boldsymbol{\theta}}^t\|_2$ by definition of the Langevin equation (9) and by (108). Hence

$$\|\boldsymbol{\theta}^t - \tilde{\boldsymbol{\theta}}^t\|_2 \le e^{Ct}\|\boldsymbol{\theta}^0 - \tilde{\boldsymbol{\theta}}^0\|_2. \tag{112}$$

Letting $(\boldsymbol{\theta}^0, \tilde{\boldsymbol{\theta}}^0)$ be the coupling of $(q_0, q)$ for which $\langle\|\boldsymbol{\theta}^0 - \tilde{\boldsymbol{\theta}}^0\|_2^2\rangle = W_2(q_0, q)^2$, we have that $(\boldsymbol{\theta}^t, \tilde{\boldsymbol{\theta}}^t)$ is a coupling of $(q_t, q)$, so

$$W_2(q_t, q)^2 \le \langle\|\boldsymbol{\theta}^t - \tilde{\boldsymbol{\theta}}^t\|_2^2\rangle \le e^{2Ct}\langle\|\boldsymbol{\theta}^0 - \tilde{\boldsymbol{\theta}}^0\|_2^2\rangle = e^{2Ct}W_2(q_0, q)^2.$$

Thus $W_2(q_t, q) \le C'W_2(q_0, q)$ for all $t \in [0, 1]$, which implies (111) for $t \in [0, 1]$ and some $C > 0$.

For $t \ge 1$, under the curvature-dimension lower bound $-\nabla^2 \log q(\boldsymbol{\theta}) \succeq -L\,\mathrm{Id}$ for a constant $L > 0$ that is implied by (108), we apply from [82, Lemma 4.2] that

$$\mathrm{D}_{\mathrm{KL}}(q_1\|q) \le \left(\frac{1}{4\alpha} + \frac{L}{2}\right)W_2(q_0, q)^2, \qquad \alpha = \frac{e^{2L} - 1}{2L}. \tag{113}$$

Under the LSI condition of $\mathcal{E}(C_0, C_{\mathrm{LSI}})$, we have the exponential contraction of relative entropy (c.f. [83, Theorem 5.2.1])

$$\mathrm{D}_{\mathrm{KL}}(q_t\|q) \le e^{-2(t-1)/C_{\mathrm{LSI}}}\,\mathrm{D}_{\mathrm{KL}}(q_1\|q) \text{ for all } t \ge 1. \tag{114}$$

We have also the $T_2$-transportation inequality (c.f. [83, Theorem 9.6.1])

$$W_2(q_t, q)^2 \le C_{\mathrm{LSI}}\,\mathrm{D}_{\mathrm{KL}}(q_t\|q), \tag{115}$$

and (111) for $t \ge 1$ follows follows from combining (113), (114), and (115). $\qquad\square$

## 4.3 Properties of the correlation and response

In this section, on the event $\mathcal{E}(C_0, C_{\mathrm{LSI}})$, we now show approximate time-translation-invariance at large times for the correlation and response matrices $\mathbf{C}_\theta, \mathbf{C}_\eta, \mathbf{R}_\theta, \mathbf{R}_\eta$ defined in Section 4.1.2. We may write these using our Markov semigroup notation as

$$\mathbf{C}_\theta(t,s) = \Big(\langle e_k(\boldsymbol{\theta}^s) P_{t-s} e_j(\boldsymbol{\theta}^s)\rangle_{\boldsymbol{\theta}^0}\Big)_{j,k=1}^d, \qquad \mathbf{C}_\eta(t,s) = \Big(\langle x_k(\boldsymbol{\theta}^s) P_{t-s} x_j(\boldsymbol{\theta}^s)\rangle_{\boldsymbol{\theta}^0}\Big)_{j,k=1}^n,$$

$$\mathbf{R}_\theta(t,s) = \Big(\langle \nabla e_k(\boldsymbol{\theta}^s)^\top \nabla P_{t-s} e_j(\boldsymbol{\theta}^s)\rangle_{\boldsymbol{\theta}^0}\Big)_{j,k=1}^d, \qquad \mathbf{R}_\eta(t,s) = \Big(\langle \nabla x_k(\boldsymbol{\theta}^s)^\top \nabla P_{t-s} x_j(\boldsymbol{\theta}^s)\rangle_{\boldsymbol{\theta}^0}\Big)_{j,k=1}^n.$$

**Lemma 4.7.** *Suppose Assumption 2.2(a) holds. Let $\mathbf{C}_\theta, \mathbf{C}_\eta, \mathbf{R}_\theta, \mathbf{R}_\eta$ be defined for the dynamics (7), and set*

$$\mathbf{C}_\theta^\infty(\tau) = \Big(\langle e_k(\boldsymbol{\theta}) P_\tau e_j(\boldsymbol{\theta})\rangle\Big)_{j,k=1}^d, \qquad \mathbf{C}_\eta^\infty(\tau) = \Big(\langle x_k(\boldsymbol{\theta}) P_\tau x_j(\boldsymbol{\theta})\rangle\Big)_{j,k=1}^n,$$

$$\mathbf{R}_\theta^\infty(\tau) = \Big(\langle \nabla e_k(\boldsymbol{\theta})^\top \nabla P_\tau e_j(\boldsymbol{\theta})\rangle\Big)_{j,k=1}^d, \qquad \mathbf{R}_\eta^\infty(\tau) = \Big(\langle \nabla x_k(\boldsymbol{\theta})^\top \nabla P_\tau x_j(\boldsymbol{\theta})\rangle\Big)_{j,k=1}^n$$

*where $\langle\cdot\rangle$ is expectation under the posterior law $q(\cdot)$. Then on $\mathcal{E}(C_0, C_{\mathrm{LSI}}) \cap \{\|\boldsymbol{\theta}^0\|_2^2 \le C_0 d\}$, there exist constants $C, c > 0$ such that for all $t \ge s \ge 0$,*

$$|\operatorname{Tr} \mathbf{C}_\theta(t,s) - \operatorname{Tr} \mathbf{C}_\theta^\infty(t-s)| \le Cde^{-cs} \tag{116}$$

$$|\operatorname{Tr} \mathbf{R}_\theta(t,s) - \operatorname{Tr} \mathbf{R}_\theta^\infty(t-s)| \le Cde^{-cs} \tag{117}$$

$$|\operatorname{Tr} \mathbf{C}_\eta(t,s) - \operatorname{Tr} \mathbf{C}_\eta^\infty(t-s)| \le Cde^{-cs} \tag{118}$$

$$|\operatorname{Tr} \mathbf{R}_\eta(t,s) - \operatorname{Tr} \mathbf{R}_\eta^\infty(t-s)| \le Cde^{-cs} \tag{119}$$

*Proof.* Momentarily let $q_t$ be the law of $\boldsymbol{\theta}^t$ conditional on $(\mathbf{X}, \mathbf{y})$ and also on a fixed initial condition $\boldsymbol{\theta}^0 = \mathbf{x}$. For any fixed $t \ge 0$, denote by $\boldsymbol{\varphi}^t \sim q$ a random vector such that $(\boldsymbol{\theta}^t, \boldsymbol{\varphi}^t)$ is a coupling of $(q_t, q)$ for which $\langle\|\boldsymbol{\theta}^t - \boldsymbol{\varphi}^t\|_2^2\rangle_\mathbf{x} = W_2(q_t, q)^2$, where $W_2(\cdot)$ is the Wasserstein-2 distance conditional on $\mathbf{X}, \mathbf{y}$ and $\boldsymbol{\theta}^0 = \mathbf{x}$. Then observe that for any $M$-Lipschitz function $f$, we have

$$\langle\|f(\boldsymbol{\theta}^t) - f(\boldsymbol{\varphi}^t)\|_2^2\rangle_\mathbf{x} \le M^2 \langle\|\boldsymbol{\theta}^t - \boldsymbol{\varphi}^t\|_2^2\rangle_\mathbf{x} = M^2 W_2(q_t, q)^2. \tag{120}$$

Furthermore $W_2(q_t, q)^2 \le Ce^{-ct} W_2(\delta_\mathbf{x}, q)^2 \le 2Ce^{-ct}(\|\mathbf{x}\|_2^2 + \langle\|\boldsymbol{\theta}\|_2^2\rangle)$ for all $t \ge 0$ by Lemma 4.6. Then, applying (120) with $f(\mathbf{x}) = \mathbf{x}$ and $f(\mathbf{x}) = \nabla\log q(\mathbf{x})$, and applying (106–108) on the event $\mathcal{E}(C_0, C_{\mathrm{LSI}})$, we have the basic estimates

$$\langle\|\boldsymbol{\theta}^t - \boldsymbol{\varphi}^t\|_2^2\rangle_\mathbf{x} \le Ce^{-ct}(\|\mathbf{x}\|_2^2 + d), \quad \langle\|\nabla\log q(\boldsymbol{\theta}^t) - \nabla\log q(\boldsymbol{\varphi}^t)\|_2^2\rangle_\mathbf{x} \le Ce^{-ct}(\|\mathbf{x}\|_2^2 + d),$$
$$\langle\|\boldsymbol{\theta}^t\|_2^2\rangle_\mathbf{x} \le C(\|\mathbf{x}\|_2^2 + d), \quad \langle\|\nabla\log q(\boldsymbol{\theta}^t)\|_2^2\rangle_\mathbf{x} \le C(\|\mathbf{x}\|_2^2 + d) \tag{121}$$
$$\langle\|\boldsymbol{\varphi}^t\|_2^2\rangle = \langle\|\boldsymbol{\theta}\|_2^2\rangle \le Cd, \quad \langle\|\nabla\log q(\boldsymbol{\varphi}^t)\|_2^2\rangle = \langle\|\nabla\log q(\boldsymbol{\theta})\|_2^2\rangle \le Cd.$$

We note that also

$$\|P_t(\mathbf{x}) - P_t(\tilde{\mathbf{x}})\|_2^2 \le e^{Ct}\|\mathbf{x} - \tilde{\mathbf{x}}\|_2^2, \tag{122}$$

$$\|P_t(\mathbf{x}) - P_t(\tilde{\mathbf{x}})\|_2^2 \le Ce^{-ct}(\|\mathbf{x}\|_2^2 + \|\tilde{\mathbf{x}}\|_2^2 + d). \tag{123}$$

Indeed, (122) follows from (112) and Jensen's inequality. Also by Jensen's inequality and (121),

$$\|P_t(\mathbf{x}) - \langle\boldsymbol{\theta}\rangle\|_2^2 = \|\langle\boldsymbol{\theta}^t - \boldsymbol{\varphi}^t\rangle_\mathbf{x}\|_2^2 \le Ce^{-ct}(\|\mathbf{x}\|_2^2 + d), \tag{124}$$

and applying this bound for both $P_t(\mathbf{x})$ and $P_t(\tilde{\mathbf{x}})$ yields (123).

For (116), note that for any $s, \tau \ge 0$,

$$\operatorname{Tr} \mathbf{C}_\theta(s+\tau, s) = \sum_{j=1}^d \langle e_j(\boldsymbol{\theta}^s) P_\tau e_j(\boldsymbol{\theta}^s)\rangle_{\boldsymbol{\theta}^0}. \tag{125}$$

Now let $q_s$ be the law of $\boldsymbol{\theta}^s$ conditional on $(\mathbf{X}, \mathbf{y})$ and the given initial condition $\boldsymbol{\theta}^0$ of Assumption 2.2, and let $(\boldsymbol{\theta}^s, \boldsymbol{\varphi}^s)$ be the optimal Wasserstein-2 coupling of $(q_s, q)$ as above. Then

$$
\begin{aligned}
|\operatorname{Tr}\mathbf{C}_\theta(s+\tau, s) - \operatorname{Tr}\mathbf{C}_\theta^\infty(\tau)| &\le \sum_{j=1}^d \left\langle \left| e_j(\boldsymbol{\theta}^s) P_\tau e_j(\boldsymbol{\theta}^s) - e_j(\boldsymbol{\varphi}^s) P_\tau e_j(\boldsymbol{\varphi}^s) \right| \right\rangle_{\boldsymbol{\theta}^0} \\
&\le \sum_{j=1}^d \left\langle \left| (\theta_j^s - \varphi_j^s) P_\tau e_j(\boldsymbol{\theta}^s) \right| \right\rangle_{\boldsymbol{\theta}^0} + \left\langle \left| \varphi_j^s (P_\tau e_j(\boldsymbol{\theta}^s) - P_\tau e_j(\boldsymbol{\varphi}^s)) \right| \right\rangle_{\boldsymbol{\theta}^0} \\
&\le \underbrace{\left\langle \|\boldsymbol{\theta}^s - \boldsymbol{\varphi}^s\|_2^2 \right\rangle_{\boldsymbol{\theta}^0}^{1/2} \left\langle \|P_\tau(\boldsymbol{\theta}^s)\|_2^2 \right\rangle_{\boldsymbol{\theta}^0}^{1/2}}_{\text{I}} + \underbrace{\left\langle \|\boldsymbol{\varphi}^s\|_2^2 \right\rangle^{1/2} \left\langle \|P_\tau(\boldsymbol{\theta}^s) - P_\tau(\boldsymbol{\varphi}^s)\|_2^2 \right\rangle_{\boldsymbol{\theta}^0}^{1/2}}_{\text{II}}
\end{aligned}
\tag{126}
$$

where we recall our shorthand $P_t(\mathbf{x}) = \langle \boldsymbol{\theta}^t \rangle_{\mathbf{x}} \in \mathbb{R}^d$.

We have $\text{I} \le Cd e^{-cs}$ for all $s \ge 0$ by (121) and $\langle \|P_\tau(\boldsymbol{\theta}^s)\|_2^2 \rangle_{\boldsymbol{\theta}^0} \le \langle \|\boldsymbol{\theta}^{s+\tau}\|_2^2 \rangle_{\boldsymbol{\theta}^0}$ which follows from Jensen's inequality. For II, by (122) and (121), we have

$$
\left\langle \|P_\tau(\boldsymbol{\theta}^s) - P_\tau(\boldsymbol{\varphi}^s)\|_2^2 \right\rangle_{\boldsymbol{\theta}^0} \le e^{C\tau} \left\langle \|\boldsymbol{\theta}^s - \boldsymbol{\varphi}^s\|_2^2 \right\rangle_{\boldsymbol{\theta}^0} \le e^{C\tau} \cdot Cd e^{-cs}.
$$

Choosing a large enough constant $s_0 > 0$, for $\tau \le s/s_0$, this gives $\langle \|P_\tau(\boldsymbol{\theta}^s) - P_\tau(\boldsymbol{\varphi}^s)\|_2^2 \rangle_{\boldsymbol{\theta}^0} \le C' d e^{-c's}$. For $\tau > s/s_0$, applying instead (123) and (121), we have $\langle \|P_\tau(\boldsymbol{\theta}^s) - P_\tau(\boldsymbol{\varphi}^s)\|_2^2 \rangle_{\boldsymbol{\theta}^0} \le Ce^{-c\tau}(\langle \|\boldsymbol{\theta}^s\|_2^2 \rangle_{\boldsymbol{\theta}^0} + \langle \|\boldsymbol{\varphi}^s\|_2^2 \rangle + d) \le C'd e^{-c's}$. Thus, for some $C, c > 0$,

$$
\left\langle \|P_\tau(\boldsymbol{\theta}^s) - P_\tau(\boldsymbol{\varphi}^s)\|_2^2 \right\rangle_{\boldsymbol{\theta}^0} \le Cd e^{-cs} \quad \text{for all } s, \tau \ge 0.
$$

Thus also $\text{II} \le Cd e^{-cs}$, and applying these bounds for I and II to (126) shows (116).

For (117), note that

$$
\operatorname{Tr}\mathbf{R}_\theta(s+\tau, s) = \sum_{j=1}^d \langle (\partial_j P_\tau e_j)[\boldsymbol{\theta}^s] \rangle_{\boldsymbol{\theta}^0}, \qquad \operatorname{Tr}\mathbf{R}_\theta^\infty(\tau) = \sum_{j=1}^d \langle (\partial_j P_\tau e_j)[\boldsymbol{\theta}] \rangle.
\tag{127}
$$

Let $dP_t(\mathbf{x}) \in \mathbb{R}^{d \times d}$ be the Jacobian of the vector map $\mathbf{x} \mapsto P_t(\mathbf{x})$. By (98) of Lemma 4.2 applied with $f = e_j$ for each $j = 1, \ldots, d$, we have

$$
dP_t(\mathbf{x}) = \langle \mathbf{V}^t \rangle_{\mathbf{x}}
\tag{128}
$$

where (with slight extension of the notation) we write $\langle \cdot \rangle_{\mathbf{x}}$ for the average over $\{\boldsymbol{\theta}^t, \mathbf{V}^t\}_{t \ge 0}$ solving (97) with initial condition $(\boldsymbol{\theta}^0, \mathbf{V}^0) = (\mathbf{x}, \mathbf{I})$. For $t \ge 1$, let us write also $\nabla P_t e_j(\mathbf{x}) = \nabla P_1 f(\mathbf{x})$ with $f = P_{t-1} e_j$. Noting that $f \in \mathcal{A}$ by (96), we may apply (99) of Lemma 4.2 with this $f$. Doing so for each $j = 1, \ldots, d$ gives

$$
dP_t(\mathbf{x}) = \frac{1}{\sqrt{2}} \left\langle P_{t-1}(\boldsymbol{\theta}^1) \left( \int_0^1 (\mathbf{V}^s)^\top d\mathbf{b}^s \right)^\top \right\rangle_{\mathbf{x}} \quad \text{for } t \ge 1.
\tag{129}
$$

In particular,

$$
\sum_{j=1}^d (\partial_j P_\tau e_j)[\mathbf{x}] = \langle \operatorname{Tr}\mathbf{V}^\tau \rangle_{\mathbf{x}} = \frac{1}{\sqrt{2}} \left\langle P_{\tau-1}(\boldsymbol{\theta}^1)^\top \int_0^1 (\mathbf{V}^s)^\top d\mathbf{b}^s \right\rangle_{\mathbf{x}}
\tag{130}
$$

with the second equality holding for $\tau \ge 1$.

Now let $\{\boldsymbol{\theta}^t, \mathbf{V}^t\}_{t \ge 0}$ and $\{\tilde{\boldsymbol{\theta}}^t, \tilde{\mathbf{V}}^t\}_{t \ge 0}$ be the solutions to (97) with initial conditions $(\boldsymbol{\theta}^0, \mathbf{V}^0) = (\mathbf{x}, \mathbf{I})$ and $(\tilde{\boldsymbol{\theta}}^0, \tilde{\mathbf{V}}^0) = (\tilde{\mathbf{x}}, \mathbf{I})$, coupled by the same Brownian motion $\{\mathbf{b}^t\}_{t \ge 0}$, and write $\langle \cdot \rangle_{\mathbf{x}, \tilde{\mathbf{x}}}$ for the associated average over $\{\boldsymbol{\theta}^t, \mathbf{V}^t, \tilde{\boldsymbol{\theta}}^t, \tilde{\mathbf{V}}^t\}_{t \ge 0}$ conditional on these initial conditions. By the form of (97) and by (108), $\frac{d}{dt}\|\mathbf{V}^t\|_{\text{op}} \le \|\nabla^2 \log q(\boldsymbol{\theta}^t) \cdot \mathbf{V}^t\|_{\text{op}} \le C\|\mathbf{V}^t\|_{\text{op}}$, so

$$
\|\mathbf{V}^t\|_{\text{op}} \le e^{Ct} \|\mathbf{V}^0\|_{\text{op}} = e^{Ct}.
\tag{131}
$$

Then also

$$\frac{\mathrm{d}}{\mathrm{d}t}\|\mathbf{V}^t - \tilde{\mathbf{V}}^t\|_F \leq \|[\nabla^2 \log q(\boldsymbol{\theta}^t) - \nabla^2 \log q(\tilde{\boldsymbol{\theta}}^t)]\mathbf{V}^t\|_F + \|[\nabla^2 \log q(\tilde{\boldsymbol{\theta}}^t)](\mathbf{V}^t - \tilde{\mathbf{V}}^t)\|_F$$

$$\leq \|\nabla^2 \log q(\boldsymbol{\theta}^t) - \nabla^2 \log q(\tilde{\boldsymbol{\theta}}^t)\|_F \|\mathbf{V}^t\|_{\mathrm{op}} + \|\nabla^2 \log q(\tilde{\boldsymbol{\theta}}^t)\|_{\mathrm{op}} \|\mathbf{V}^t - \tilde{\mathbf{V}}^t\|_F.$$

Applying $\nabla^2 \log q(\boldsymbol{\theta}^t) - \nabla^2 \log q(\tilde{\boldsymbol{\theta}}^t) = \mathrm{diag}((\log q)''(\theta_j^t) - (\log q)''(\tilde{\theta}_j^t))$, the bound $\|\nabla^2 \log q(\boldsymbol{\theta})\|_{\mathrm{op}} \leq C$ from (108), $|(\log g)'''(\theta)| \leq C$ under Assumption 2.2(a), and (112),

$$\frac{\mathrm{d}}{\mathrm{d}t}\|\mathbf{V}^t - \tilde{\mathbf{V}}^t\|_F \leq C\|\boldsymbol{\theta}^t - \tilde{\boldsymbol{\theta}}^t\|_2 \|\mathbf{V}^t\|_{\mathrm{op}} + C\|\mathbf{V}^t - \tilde{\mathbf{V}}^t\|_F \leq Ce^{Ct}\|\mathbf{x} - \tilde{\mathbf{x}}\|_2 \cdot e^{Ct} + C\|\mathbf{V}^t - \tilde{\mathbf{V}}^t\|_F.$$

Integrating this bound,

$$\|\mathbf{V}^t - \tilde{\mathbf{V}}^t\|_F \leq C\|\mathbf{x} - \tilde{\mathbf{x}}\|_2 \text{ for all } t \in [0, 1].$$

So it follows from the first equality of (130) that for $\tau \in [0, 1]$,

$$\left|\sum_{j=1}^d \partial_j P_\tau e_j(\mathbf{x}) - \partial_j P_\tau e_j(\tilde{\mathbf{x}})\right| = \left|\langle \mathrm{Tr}(\mathbf{V}^\tau - \tilde{\mathbf{V}}^\tau)\rangle_{\mathbf{x},\tilde{\mathbf{x}}}\right| \leq \sqrt{d}\langle\|\mathbf{V}^\tau - \tilde{\mathbf{V}}^\tau\|_F\rangle_{\mathbf{x},\tilde{\mathbf{x}}} \leq C\sqrt{d}\|\mathbf{x} - \tilde{\mathbf{x}}\|_2.$$

Hence by (127) and (121), for $\tau \in [0, 1]$,

$$|\mathrm{Tr}\,\mathbf{R}_\theta(s + \tau, s) - \mathrm{Tr}\,\mathbf{R}_\theta^\infty(\tau)| \leq C\sqrt{d}\langle\|\boldsymbol{\theta}^s - \boldsymbol{\varphi}^s\|_2\rangle_{\boldsymbol{\theta}^0} \leq C'de^{-cs}. \tag{132}$$

For $\tau \geq 1$, we apply instead the second equality of (130) and Cauchy-Schwarz to obtain

$$\sqrt{2}\left|\sum_{j=1}^d \partial_j P_\tau e_j(\mathbf{x}) - \partial_j P_\tau e_j(\tilde{\mathbf{x}})\right|$$

$$\leq \left|\left\langle P_{\tau-1}(\boldsymbol{\theta}^1)^\top \int_0^1 (\mathbf{V}^s)^\top \mathrm{d}\mathbf{b}^s - P_{\tau-1}(\tilde{\boldsymbol{\theta}}^1)^\top \int_0^1 (\tilde{\mathbf{V}}^s)^\top \mathrm{d}\mathbf{b}^s\right\rangle_{\mathbf{x},\tilde{\mathbf{x}}}\right|$$

$$\leq \left\langle\left\|P_{\tau-1}(\boldsymbol{\theta}^1) - P_{\tau-1}(\tilde{\boldsymbol{\theta}}^1)\right\|_2^2\right\rangle_{\mathbf{x},\tilde{\mathbf{x}}}^{1/2}\left\langle\left\|\int_0^1 (\mathbf{V}^s)^\top \mathrm{d}\mathbf{b}^s\right\|_2^2\right\rangle_{\mathbf{x}}^{1/2} + \left\langle\left\|P_{\tau-1}(\tilde{\boldsymbol{\theta}}^1)\right\|_2^2\right\rangle_{\tilde{\mathbf{x}}}^{1/2}\left\langle\left\|\int_0^1 (\mathbf{V}^s - \tilde{\mathbf{V}}^s)^\top \mathrm{d}\mathbf{b}^s\right\|_2^2\right\rangle_{\mathbf{x},\tilde{\mathbf{x}}}^{1/2}$$

$$= \left\langle\left\|P_{\tau-1}(\boldsymbol{\theta}^1) - P_{\tau-1}(\tilde{\boldsymbol{\theta}}^1)\right\|_2^2\right\rangle_{\mathbf{x},\tilde{\mathbf{x}}}^{1/2}\left\langle\int_0^1 \|\mathbf{V}^s\|_F^2 \mathrm{d}s\right\rangle_{\mathbf{x}}^{1/2} + \left\langle\left\|P_{\tau-1}(\tilde{\boldsymbol{\theta}}^1)\right\|_2^2\right\rangle_{\tilde{\mathbf{x}}}^{1/2}\left\langle\int_0^1 \|\mathbf{V}^s - \tilde{\mathbf{V}}^s\|_F^2 \mathrm{d}s\right\rangle_{\mathbf{x},\tilde{\mathbf{x}}}^{1/2}$$

$$\leq C\sqrt{d}\left\langle\left\|P_{\tau-1}(\boldsymbol{\theta}^1) - P_{\tau-1}(\tilde{\boldsymbol{\theta}}^1)\right\|_2^2\right\rangle_{\mathbf{x},\tilde{\mathbf{x}}}^{1/2} + C\|\mathbf{x} - \tilde{\mathbf{x}}\|_2\left\langle\left\|P_{\tau-1}(\tilde{\boldsymbol{\theta}}^1)\right\|_2^2\right\rangle_{\tilde{\mathbf{x}}}^{1/2}. \tag{133}$$

Note that $\langle\|P_{\tau-1}(\tilde{\boldsymbol{\theta}}^1)\|_2^2\rangle_{\tilde{\mathbf{x}}} \leq \langle\|\tilde{\boldsymbol{\theta}}^\tau\|_2^2\rangle_{\tilde{\mathbf{x}}} \leq C(\|\tilde{\mathbf{x}}\|_2^2 + d)$ by (121). Applying (122–123), (112), and (121),

$$\left\langle\left\|P_{\tau-1}(\boldsymbol{\theta}^1) - P_{\tau-1}(\tilde{\boldsymbol{\theta}}^1)\right\|_2^2\right\rangle_{\mathbf{x},\tilde{\mathbf{x}}} \leq e^{2C(\tau-1)}\langle\|\boldsymbol{\theta}^1 - \tilde{\boldsymbol{\theta}}^1\|_2^2\rangle_{\mathbf{x},\tilde{\mathbf{x}}} \leq Ce^{2C\tau}\|\mathbf{x} - \tilde{\mathbf{x}}\|_2^2 \text{ for all } \tau \geq 1,$$

$$\left\langle\left\|P_{\tau-1}(\boldsymbol{\theta}^1) - P_{\tau-1}(\tilde{\boldsymbol{\theta}}^1)\right\|_2^2\right\rangle_{\mathbf{x},\tilde{\mathbf{x}}} \leq Ce^{-c(\tau-1)}(\langle\|\boldsymbol{\theta}^1\|_2^2\rangle_{\mathbf{x}} + \langle\|\tilde{\boldsymbol{\theta}}^1\|_2^2\rangle_{\tilde{\mathbf{x}}} + d)$$

$$\leq C'e^{-c\tau}(\|\mathbf{x}\|_2^2 + \|\tilde{\mathbf{x}}\|_2^2 + d) \text{ for all } \tau \geq 2.$$

Choosing a large enough constant $s_0 > 0$, if $\tau \in [1, s/s_0]$, then we may apply the former bound, (121), and Cauchy-Schwarz to (133) to get

$$|\mathrm{Tr}\,\mathbf{R}_\theta(s + \tau, s) - \mathrm{Tr}\,\mathbf{R}_\theta^\infty(\tau)| \leq \left\langle\left|\sum_{j=1}^d \partial_j P_\tau e_j(\boldsymbol{\theta}^s) - \partial_j P_\tau e_j(\boldsymbol{\varphi}^s)\right|\right\rangle_{\boldsymbol{\theta}^0}$$

$$\leq Ce^{C\tau}\sqrt{d}\langle\|\boldsymbol{\theta}^s - \boldsymbol{\varphi}^s\|_2\rangle_{\boldsymbol{\theta}^0} + C\left\langle\|\boldsymbol{\theta}^s - \boldsymbol{\varphi}^s\|_2(\|\boldsymbol{\varphi}^s\|_2 + \sqrt{d})\right\rangle_{\boldsymbol{\theta}^0}$$

$$\leq C'd(e^{C\tau} + 1)e^{-cs} \leq C''de^{-c's}. \tag{134}$$

If $\tau \geq s/s_0$, applying instead the latter bound to (133),

$$|\operatorname{Tr}\mathbf{R}_\theta(s+\tau,s) - \operatorname{Tr}\mathbf{R}_\theta^\infty(\tau)| \leq Ce^{-c\tau}\sqrt{d}\left(\langle\|\boldsymbol{\theta}^s\|_2^2\rangle_{\boldsymbol{\theta}^0} + \langle\|\boldsymbol{\varphi}^s\|_2^2\rangle + d\right)^{1/2} + C\left\langle\|\boldsymbol{\theta}^s - \boldsymbol{\varphi}^s\|_2(\|\boldsymbol{\varphi}^s\|_2 + \sqrt{d})\right\rangle_{\boldsymbol{\theta}^0}$$

$$\leq C'de^{-c's}. \tag{135}$$

Combining these bounds for $\tau \in [0,1]$, $\tau \in [1, s/s_0]$, and $\tau \geq s/s_0$ in (132), (134), and (135) shows (117).

The arguments for $\mathbf{C}_\eta$ and $\mathbf{R}_\eta$ in (118–119) are similar: For (118), recall the definitions (100) and note that

$$\left|\operatorname{Tr}\mathbf{C}_\eta(s+\tau,s) - \operatorname{Tr}\mathbf{C}_\eta^\infty(\tau)\right|$$

$$\leq \sum_{i=1}^n \left\langle\left|x_i(\boldsymbol{\theta}^s)P_\tau x_i(\boldsymbol{\theta}^s) - x_i(\boldsymbol{\varphi}^s)P_\tau x_i(\boldsymbol{\varphi}^s)\right|\right\rangle_{\boldsymbol{\theta}^0}$$

$$\leq \sum_{i=1}^n \left\langle\left|(x_i(\boldsymbol{\theta}^s) - x_i(\boldsymbol{\varphi}^s))P_\tau x_i(\boldsymbol{\theta}^s)\right|\right\rangle_{\boldsymbol{\theta}^0} + \left\langle\left|x_i(\boldsymbol{\varphi}^s)(P_\tau x_i(\boldsymbol{\theta}^s) - P_\tau x_i(\boldsymbol{\varphi}^s))\right|\right\rangle_{\boldsymbol{\theta}^0}$$

$$\leq \frac{\delta}{\sigma^4}\left\langle\|\mathbf{X}(\boldsymbol{\theta}^s - \boldsymbol{\varphi}^s)\|_2^2\right\rangle_{\boldsymbol{\theta}^0}^{1/2}\left\langle\|\mathbf{X}P_\tau(\boldsymbol{\theta}^s) - \mathbf{y}\|_2^2\right\rangle_{\boldsymbol{\theta}^0}^{1/2} + \frac{\delta}{\sigma^4}\left\langle\|\mathbf{X}\boldsymbol{\varphi}^s - \mathbf{y}\|_2^2\right\rangle^{1/2}\left\langle\|\mathbf{X}P_\tau(\boldsymbol{\theta}^s) - \mathbf{X}P_\tau(\boldsymbol{\varphi}^s)\|_2^2\right\rangle_{\boldsymbol{\theta}^0}^{1/2}.$$

The desired result (118) follows from the conditions $\|\mathbf{X}\|_{\mathrm{op}} \leq C_0$, $\|\mathbf{y}\|_2^2 \leq C_0d$, and the preceding bounds for (126).

For (119), note that

$$\operatorname{Tr}\mathbf{R}_\eta(s+\tau,s) = \frac{\delta}{\sigma^4}\left\langle\operatorname{Tr}\mathbf{X}[\mathrm{d}P_\tau(\boldsymbol{\theta}^s)]\mathbf{X}^\top\right\rangle_{\boldsymbol{\theta}^0}, \qquad \operatorname{Tr}\mathbf{R}_\eta^\infty(\tau) = \frac{\delta}{\sigma^4}\left\langle\operatorname{Tr}\mathbf{X}[\mathrm{d}P_\tau(\boldsymbol{\theta})]\mathbf{X}^\top\right\rangle. \tag{136}$$

Let $\{\boldsymbol{\theta}^t, \mathbf{V}^t\}_{t\geq 0}$ and $\{\tilde{\boldsymbol{\theta}}^t, \tilde{\mathbf{V}}^t\}_{t\geq 0}$ be the solutions of (97) with initial conditions $(\mathbf{x}, \mathbf{I})$ and $(\tilde{\mathbf{x}}, \mathbf{I})$. If $\tau \in [0,1]$, we apply (128) to obtain

$$|\operatorname{Tr}\mathbf{X}[\mathrm{d}P_\tau(\mathbf{x})]\mathbf{X}^\top - \operatorname{Tr}\mathbf{X}[\mathrm{d}P_\tau(\tilde{\mathbf{x}})]\mathbf{X}^\top| \leq \left\langle\|\mathbf{V}^\tau - \tilde{\mathbf{V}}^\tau\|_F\right\rangle_{\mathbf{x},\tilde{\mathbf{x}}} \cdot \|\mathbf{X}^\top\mathbf{X}\|_F \leq \sqrt{d}\|\mathbf{X}\|_{\mathrm{op}}^2 \cdot \left\langle\|\mathbf{V}^\tau - \tilde{\mathbf{V}}^\tau\|_F\right\rangle_{\mathbf{x},\tilde{\mathbf{x}}},$$

which leads to the bound (132) up to a different constant depending on the bound $C_0$ for $\|\mathbf{X}\|_{\mathrm{op}}$. If $\tau \geq 1$, we apply (129) to obtain

$$\sqrt{2}\left|\operatorname{Tr}\mathbf{X}[\mathrm{d}P_\tau(\mathbf{x})]\mathbf{X}^\top - \operatorname{Tr}\mathbf{X}[\mathrm{d}P_\tau(\tilde{\mathbf{x}})]\mathbf{X}^\top\right|$$

$$\leq \left\langle\left|\left(P_{\tau-1}(\boldsymbol{\theta}^1) - P_{\tau-1}(\tilde{\boldsymbol{\theta}}^1)\right)^\top\mathbf{X}^\top\mathbf{X}\left(\int_0^1 \mathbf{V}^{s\top}\mathrm{db}^s\right)\right|\right\rangle_{\mathbf{x},\tilde{\mathbf{x}}} + \left\langle\left|P_{\tau-1}(\tilde{\boldsymbol{\theta}}^1)^\top\mathbf{X}^\top\mathbf{X}\left(\int_0^1 (\mathbf{V}^s - \tilde{\mathbf{V}}^s)^\top\mathrm{db}^s\right)\right|\right\rangle_{\mathbf{x},\tilde{\mathbf{x}}}$$

$$\leq \|\mathbf{X}\|_{\mathrm{op}}^2\left[\left\langle\|P_{\tau-1}(\boldsymbol{\theta}^1) - P_{\tau-1}(\tilde{\boldsymbol{\theta}}^1)\|^2\right\rangle_{\mathbf{x},\tilde{\mathbf{x}}}^{1/2}\left\langle\int_0^1\|\mathbf{V}^s\|_F^2\mathrm{d}s\right\rangle_{\mathbf{x}}^{1/2} + \left\langle\|P_{\tau-1}(\tilde{\boldsymbol{\theta}}^1)\|^2\right\rangle_{\tilde{\mathbf{x}}}^{1/2}\left\langle\int_0^1\|\mathbf{V}^s - \tilde{\mathbf{V}}^s\|_F^2\mathrm{d}s\right\rangle_{\mathbf{x},\tilde{\mathbf{x}}}^{1/2}\right].$$

This can be bounded in the same way as (133), (134), and (135) up to different constants depending on the bound $C_0$ for $\|\mathbf{X}\|_{\mathrm{op}}$. This shows (119). $\qquad\square$

**Lemma 4.8.** *Suppose Assumption 2.2(a) holds. Let $\{\boldsymbol{\theta}^t\}_{t\geq 0}$ be the solution to (7). Then on the event $\mathcal{E}(C_0, C_{\mathrm{LSI}}) \cap \{\|\boldsymbol{\theta}^0\|_2^2 \leq C_0d\}$, there exist constants $C, c > 0$ such that for all $t \geq s \geq 0$,*

$$\left|\operatorname{Tr}\mathbf{C}_\theta(t,s) - P_s(\boldsymbol{\theta}^0)^\top\langle\boldsymbol{\theta}\rangle\right| \leq Cde^{-c(t-s)} \tag{137}$$

$$\left|\operatorname{Tr}\mathbf{R}_\theta(t,s)\right| \leq Cde^{-c(t-s)} \tag{138}$$

$$\left|\operatorname{Tr}\mathbf{C}_\eta(t,s) - \frac{\delta}{\sigma^4}(\mathbf{X}P_s(\boldsymbol{\theta}^0) - \mathbf{y})^\top(\mathbf{X}\langle\boldsymbol{\theta}\rangle - \mathbf{y})\right| \leq Cde^{-c(t-s)} \tag{139}$$

$$\left|\operatorname{Tr}\mathbf{R}_\eta(t,s)\right| \leq Cde^{-c(t-s)} \tag{140}$$

*and furthermore*

$$\left|P_s(\boldsymbol{\theta}^0)^\top\langle\boldsymbol{\theta}\rangle - \|\langle\boldsymbol{\theta}\rangle\|_2^2\right| \leq Cde^{-cs} \tag{141}$$

$$\left|(\mathbf{X}P_s(\boldsymbol{\theta}^0) - \mathbf{y})^\top(\mathbf{X}\langle\boldsymbol{\theta}\rangle - \mathbf{y}) - \|\mathbf{X}\langle\boldsymbol{\theta}\rangle - \mathbf{y}\|_2^2\right| \leq Cde^{-cs} \tag{142}$$

*Proof.* For (137) and (141), observe that

$$\big|\operatorname{Tr}\mathbf{C}_\theta(s+\tau,s)-P_s(\boldsymbol{\theta}^0)^\top\langle\boldsymbol{\theta}\rangle\big|=\big|\langle\boldsymbol{\theta}^{s\top}(P_\tau\boldsymbol{\theta}^s-\langle\boldsymbol{\theta}\rangle)\rangle_{\boldsymbol{\theta}^0}\big|$$

$$\leq\langle\|\boldsymbol{\theta}^s\|_2^2\rangle_{\boldsymbol{\theta}^0}^{1/2}\langle\|P_\tau\boldsymbol{\theta}^s-\langle\boldsymbol{\theta}\rangle\|_2^2\rangle_{\boldsymbol{\theta}^0}^{1/2}\leq Cde^{-c\tau},$$

the last inequality applying (121) and (124). Similarly

$$\big|P_s(\boldsymbol{\theta}^0)^\top\langle\boldsymbol{\theta}\rangle-\|\langle\boldsymbol{\theta}\rangle\|_2^2\big|=\big|(P_s\boldsymbol{\theta}^0-\langle\boldsymbol{\theta}\rangle)^\top\langle\boldsymbol{\theta}\rangle\big|\leq\|\langle\boldsymbol{\theta}\rangle\|_2\cdot\|P_s\boldsymbol{\theta}^0-\langle\boldsymbol{\theta}\rangle\|_2\leq Cde^{-cs}.$$

For (138), recall from (127) that $\operatorname{Tr}\mathbf{R}_\theta(s+\tau,s)=\sum_{j=1}^d\langle(\partial_j P_\tau e_j)[\boldsymbol{\theta}^s]\rangle_{\boldsymbol{\theta}^0}$. Then by the first equality of (130) and (131), we have $|\operatorname{Tr}\mathbf{R}_\theta(s+\tau,s)|\leq Cd$ for any $\tau\in[0,1]$. For $\tau\geq1$, we apply instead the second equality of (130) where $\int_0^t\mathbf{V}^{s\top}\mathrm{d}\mathbf{b}^s$ is a martingale. Then $\langle\langle\boldsymbol{\theta}\rangle^\top\int_0^1\mathbf{V}^{s\top}\mathrm{d}\mathbf{b}^s\rangle_\mathbf{x}=0$ for any initial condition $\mathbf{x}\in\mathbb{R}^d$, so for any $\tau\geq1$,

$$\left|\sum_{j=1}^d\partial_j P_\tau e_j(\mathbf{x})\right|=\left|\left\langle\left(P_{\tau-1}(\boldsymbol{\theta}^1)-\langle\boldsymbol{\theta}\rangle\right)^\top\int_0^1\mathbf{V}^{s\top}\mathrm{d}\mathbf{b}^s\right\rangle_\mathbf{x}\right|$$

$$\leq\langle\|P_{\tau-1}(\boldsymbol{\theta}^1)-\langle\boldsymbol{\theta}\rangle\|^2\rangle_\mathbf{x}^{1/2}\left\langle\int_0^1\|\mathbf{V}^s\|_F^2\mathrm{d}s\right\rangle_\mathbf{x}^{1/2}\leq Ce^{-c\tau}\sqrt{d}(\|\mathbf{x}\|_2+\sqrt{d}),$$

the last inequality using the estimates (124) and (131). Then by (121) and Jensen's inequality,

$$|\operatorname{Tr}\mathbf{R}_\theta(s+\tau,s)|\leq\left\langle\left|\sum_{j=1}^d(\partial_j P_\tau e_j)[\boldsymbol{\theta}^s]\right|\right\rangle_{\boldsymbol{\theta}^0}\leq C'de^{-c'\tau}.$$

Combining these cases $\tau\in[0,1]$ and $\tau\geq1$ gives (138).

The arguments for (139), (140), and (142), are analogous to the above, and we omit these for brevity. $\square$

## 4.4  The DMFT system is approximately-TTI

We now prove Theorem 2.9, that under the log-Sobolev condition of Assumption 2.7(a), the DMFT system of Theorem 2.3(a) is approximately-TTI in the sense of Definition 2.4.

**Lemma 4.9.** *Under Assumptions 2.1, 2.2(a), and 2.7(a), the DMFT system prescribed by Theorem 2.3(a) satisfies the conditions of Definition 2.4(1) with $\varepsilon(t)=Ce^{-ct}$ for some constants $C,c>0$.*

*Proof.* We restrict to the almost sure event where the convergence statements of Theorem 4.3 hold, and where $\mathcal{E}(C_0,C_{\mathrm{LSI}})\cap\{\|\boldsymbol{\theta}^0\|_2^2\leq C_0d\}$ holds for all large $n,d$.

Consider first the statements for $C_\theta(t,s)$. Applying $\|\boldsymbol{\theta}^0\|_2^2\leq C_0d$ and (137) of Lemma 4.8, for some constants $C,c>0$,

$$\limsup_{n,d\to\infty}\big|d^{-1}\operatorname{Tr}\mathbf{C}_\theta(t,s)-d^{-1}P_s(\boldsymbol{\theta}^0)^\top\langle\boldsymbol{\theta}\rangle\big|\leq Ce^{-ct}\quad\text{for all }s\leq t/2.\tag{143}$$

By Theorem 4.3, $\lim_{n,d\to\infty}d^{-1}\operatorname{Tr}\mathbf{C}_\theta(t,s)=C_\theta(t,s)$ for all $t\geq s\geq0$. Then, for each $s\geq0$ and $t\geq2s$,

$$\limsup_{n,d\to\infty}d^{-1}P_s(\boldsymbol{\theta}^0)^\top\langle\boldsymbol{\theta}\rangle\leq C_\theta(t,s)+Ce^{-ct},\qquad\liminf_{n,d\to\infty}d^{-1}P_s(\boldsymbol{\theta}^0)^\top\langle\boldsymbol{\theta}\rangle\geq C_\theta(t,s)-Ce^{-ct}.$$

Taking $t\to\infty$ on the right side of both statements shows that for each $s\geq0$, there exists a limit

$$\tilde{c}_\theta(s):=\lim_{n,d\to\infty}d^{-1}P_s(\boldsymbol{\theta}^0)^\top\langle\boldsymbol{\theta}\rangle=\lim_{t\to\infty}C_\theta(t,s).\tag{144}$$

Next, (141) of Lemma 4.8 implies for some $C,c>0$,

$$\limsup_{n,d\to\infty}\big|d^{-1}P_s(\boldsymbol{\theta}^0)^\top\langle\boldsymbol{\theta}\rangle-d^{-1}\|\langle\boldsymbol{\theta}\rangle\|_2^2\big|\leq Ce^{-cs}\quad\text{for all }s\geq0.\tag{145}$$

44

Then

$$\limsup_{n,d\to\infty} d^{-1}\|\langle\boldsymbol{\theta}\rangle\|_2^2 \le \tilde{c}_\theta(s) + Ce^{-cs}, \qquad \liminf_{n,d\to\infty} d^{-1}\|\langle\boldsymbol{\theta}\rangle\|_2^2 \ge \tilde{c}_\theta(s) - Ce^{-cs}.$$

Taking $s\to\infty$ on the right side of both statements shows that there exists a limit

$$c_\theta^{\mathrm{tti}}(\infty) := \lim_{n,d\to\infty} d^{-1}\|\langle\boldsymbol{\theta}\rangle\|_2^2 = \lim_{s\to\infty}\tilde{c}_\theta(s). \tag{146}$$

Now consider $\mathbf{C}_\theta^\infty(\tau)$ as defined in Lemma 4.7. Let $-\mathrm{L} = \int_0^\infty a\,\mathrm{d}E_a$ be the spectral decomposition of $-\mathrm{L}$ as a positive, self-adjoint operator on $L^2(q)$ (c.f. [83, Theorem A.4.2]), where $\{E_a\}_{a\ge0}$ is a family of orthogonal projections onto an increasing family of closed linear subspaces of $L^2(q)$. In particular, $E_0 f = \langle f(\boldsymbol{\theta})\rangle$ is the projection onto the constant functions. For each $\tau\ge0$ and all $f,g\in L^2(q)$, we then have

$$\langle f(\boldsymbol{\theta})P_\tau g(\boldsymbol{\theta})\rangle = \int_0^\infty e^{-a\tau}\mathrm{d}\langle f(\boldsymbol{\theta})E_a g(\boldsymbol{\theta})\rangle \tag{147}$$

understood as a Stieltjes integral with respect to the bounded-variation function $a\mapsto\langle f(\boldsymbol{\theta})E_a g(\boldsymbol{\theta})\rangle$ (c.f. [83, Proposition 3.1.6(iii)]). The LSI on the event $\mathcal{E}(C_0,C_{\mathrm{LSI}})$ implies a spectral gap, i.e. the spectrum of $-\mathrm{L}$ is included in $\{0\}\cup[1/C_{\mathrm{LSI}},\infty)$. Thus, fixing any constant $\iota\in(0,1/C_{\mathrm{LSI}})$, we have

$$d^{-1}\operatorname{Tr}\mathbf{C}_\theta^\infty(\tau) = d^{-1}\sum_{j=1}^d\langle e_j(\boldsymbol{\theta})P_\tau e_j(\boldsymbol{\theta})\rangle = d^{-1}\sum_{j=1}^d\langle e_j(\boldsymbol{\theta})E_0 e_j(\boldsymbol{\theta})\rangle + d^{-1}\sum_{j=1}^d\int_\iota^\infty e^{-a\tau}\mathrm{d}\langle e_j(\boldsymbol{\theta})E_a e_j(\boldsymbol{\theta})\rangle$$

$$= d^{-1}\|\langle\boldsymbol{\theta}\rangle\|_2^2 + d^{-1}\sum_{j=1}^d\int_\iota^\infty e^{-a\tau}\mathrm{d}\langle e_j(\boldsymbol{\theta})E_a e_j(\boldsymbol{\theta})\rangle,$$

the first equality applying (147), and the second equality applying $\sum_j\langle e_j(\boldsymbol{\theta})E_0 e_j(\boldsymbol{\theta})\rangle = \sum_j\langle\theta_j\langle\theta_j\rangle\rangle = \|\langle\boldsymbol{\theta}\rangle\|^2$. Define $(n,d,\mathbf{X},\mathbf{y})$-dependent scalars $c_{\theta,d},m_{\theta,d}>0$ and a positive measure $\mu_{\theta,d}$ on $[\iota,\infty)$ by

$$c_{\theta,d} = d^{-1}\|\langle\boldsymbol{\theta}\rangle\|_2^2, \quad m_{\theta,d} = d^{-1}\sum_{j=1}^d\int_\iota^\infty\mathrm{d}\langle e_j(\boldsymbol{\theta})E_a e_j(\boldsymbol{\theta})\rangle, \quad \mu_{\theta,d}(S) = d^{-1}\sum_{j=1}^d\int_S\mathrm{d}\langle e_j(\boldsymbol{\theta})E_a e_j(\boldsymbol{\theta})\rangle, \tag{148}$$

noting that $a\mapsto d^{-1}\sum_{j=1}^d\langle e_j(\boldsymbol{\theta})E_a e_j(\boldsymbol{\theta})\rangle$ is nondecreasing and hence defines a valid distribution function for $\mu_{\theta,d}$. Then

$$d^{-1}\operatorname{Tr}\mathbf{C}_\theta^\infty(\tau) = c_{\theta,d} + \int_\iota^\infty e^{-a\tau}\mu_{\theta,d}(\mathrm{d}a), \qquad m_{\theta,d} = \mu_{\theta,d}([\iota,\infty)). \tag{149}$$

Applying (149) with $\tau=0$,

$$c_{\theta,d} + m_{\theta,d} = d^{-1}\operatorname{Tr}\mathbf{C}_\theta^\infty(0) = d^{-1}\langle\|\boldsymbol{\theta}\|_2^2\rangle \le C. \tag{150}$$

In particular, $\mu_{\theta,d}$ is finite and uniformly bounded in total variation norm for all $(n,d)$. We claim that $\tau\mapsto d^{-1}\operatorname{Tr}\mathbf{C}_\theta^\infty(\tau)$ is uniformly equicontinuous over all $(n,d)$: Observe that

$$\left|d^{-1}\operatorname{Tr}\mathbf{C}_\theta^\infty(\tau) - d^{-1}\operatorname{Tr}\mathbf{C}_\theta^\infty(\tau')\right| = \left|d^{-1}\langle\boldsymbol{\theta}^\top[P_\tau - P_{\tau'}](\boldsymbol{\theta})\rangle\right|$$

$$\le d^{-1}\langle\|\boldsymbol{\theta}\|_2^2\rangle^{1/2}\cdot\langle\|[P_\tau - P_{\tau'}](\boldsymbol{\theta})\|_2^2\rangle^{1/2}$$

$$= d^{-1}\langle\|\boldsymbol{\theta}\|_2^2\rangle^{1/2}\cdot\langle\|\boldsymbol{\theta} - P_{|\tau-\tau'|}(\boldsymbol{\theta})\|_2^2\rangle^{1/2}. \tag{151}$$

[84, Theorem II.2.1] implies $\|P_t(\mathbf{x}) - \mathbf{x}\|_2^2 \le C(1+\|\mathbf{x}\|_2^2)t$ for all $t\in[0,1]$ and a constant $C>0$. This and (121) imply that the right side of (151) is at most $C'|\tau-\tau'|$ for a constant $C'>0$ and all $|\tau-\tau'|\le1$, so $\tau\mapsto d^{-1}\operatorname{Tr}\mathbf{C}_\theta^\infty(\tau)$ is uniformly equicontinuous as claimed. We note that for any $M>0$, by the relation (149), $c_{\theta,d} + \mu_{\theta,d}([0,M)) + e^{-M\tau}\mu_{\theta,d}([M,\infty)) \ge d^{-1}\operatorname{Tr}\mathbf{C}_\theta^\infty(\tau)$. Then setting $\tau=1/M$ and rearranging yields

$$(1-e^{-1})\mu_{\theta,d}([M,\infty)) \le c_{\theta,d} + m_{\theta,d} - d^{-1}\operatorname{Tr}\mathbf{C}_\theta^\infty(1/M) = d^{-1}\operatorname{Tr}\mathbf{C}_\theta^\infty(0) - d^{-1}\operatorname{Tr}\mathbf{C}_\theta^\infty(1/M).$$

45

So this uniform equicontinuity implies that the measures $\mu_{\theta,d}$ are uniformly tight.

Then, there exists a subsequence $\{(n_k, d_k)\}_{k\geq 1}$ of $(n,d)$ along which $\mu_{\theta,d} \Rightarrow \mu_\theta$ weakly, for some finite positive measure $\mu_\theta$ on $[\iota,\infty)$. Recalling also that $c_{\theta,d} = d^{-1}\|\langle\boldsymbol{\theta}\rangle\|_2^2 \to c_\theta^{\mathrm{tti}}(\infty)$ as $n, d \to \infty$ by the definition (146), and setting

$$c_\theta^{\mathrm{tti}}(\tau) = c_\theta^{\mathrm{tti}}(\infty) + \int_\iota^\infty e^{-a\tau}\mu_\theta(\mathrm{d}a), \tag{152}$$

this weak convergence applied to (149) implies $\lim_{k\to\infty} d_k^{-1} \operatorname{Tr} \mathbf{C}_\theta^\infty(\tau) = c_\theta^{\mathrm{tti}}(\tau)$. Combining this with the convergence $\lim_{k\to\infty} d_k^{-1} \operatorname{Tr} \mathbf{C}_\theta(s+\tau,s) = C_\theta(s+\tau,s)$ by Theorem 4.3, for any $s, \tau \geq 0$ we have

$$\left| C_\theta(s+\tau,s) - c_\theta^{\mathrm{tti}}(\tau) \right| \leq \limsup_{k\to\infty} \left| d_k^{-1} \operatorname{Tr} \mathbf{C}_\theta(s+\tau,s) - d_k^{-1} \operatorname{Tr} \mathbf{C}_\theta^\infty(\tau) \right| \leq Ce^{-cs}, \tag{153}$$

where the last inequality holds by (116). Since $C_\theta(t,s)$ is non-random, this implies that $c_\theta^{\mathrm{tti}}(\tau)$ is also non-random for every $\tau \geq 0$, and thus also the measure $\mu_\theta$ is non-random. This shows (30).

The statement (31) follows analogously: By arguments parallel to (144) and (146), applying Theorem 4.3 and (139) and (142) shows that there exist limits

$$\tilde{c}_\eta(s) := \lim_{n,d\to\infty} \frac{\delta}{n\sigma^4}(\mathbf{X}P_s(\boldsymbol{\theta}^0) - \mathbf{y})^\top(\mathbf{X}\langle\boldsymbol{\theta}\rangle - \mathbf{y}) = \lim_{t\to\infty} C_\eta(t,s), \tag{154}$$

$$c_\eta^{\mathrm{tti}}(\infty) := \lim_{n,d\to\infty} \frac{\delta}{n\sigma^4}\|\mathbf{X}\langle\boldsymbol{\theta}\rangle - \mathbf{y}\|_2^2 = \lim_{s\to\infty} \tilde{c}_\eta(s). \tag{155}$$

Note that

$$n^{-1}\operatorname{Tr}\mathbf{C}_\eta^\infty(\tau) = n^{-1}\sum_{i=1}^n \langle x_i(\boldsymbol{\theta})P_\tau x_i(\boldsymbol{\theta})\rangle = \frac{\delta}{n\sigma^4}\|\mathbf{X}\langle\boldsymbol{\theta}\rangle - \mathbf{y}\|_2^2 + \frac{1}{n}\sum_{i=1}^n \int_\iota^\infty e^{-a\tau}\mathrm{d}\langle x_i(\boldsymbol{\theta})E_a x_i(\boldsymbol{\theta})\rangle.$$

Defining

$$c_{\eta,n} = \frac{\delta}{n\sigma^4}\|\mathbf{X}\langle\boldsymbol{\theta}\rangle - \mathbf{y}\|_2^2, \quad \mu_{\eta,n}(S) = \frac{1}{n}\sum_{i=1}^n \int_S \mathrm{d}\langle x_i(\boldsymbol{\theta})E_a x_i(\boldsymbol{\theta})\rangle, \quad m_{\eta,n} = \mu_{\eta,n}([\iota,\infty)), \tag{156}$$

we have

$$c_{\eta,n} + m_{\eta,n} = n^{-1}\operatorname{Tr}\mathbf{C}_\eta^\infty(0) = \frac{\delta}{n\sigma^4}\langle\|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2\rangle \leq C. \tag{157}$$

So along some subsequence $\{(n_k, d_k)\}_{k\geq 1}$, we have $c_{\eta,n} \to c_\eta^{\mathrm{tti}}(\infty)$, $\mu_{\eta,n} \Rightarrow \mu_\eta$ weakly for a finite positive measure $\mu_\eta$ on $[\iota,\infty)$, and $\lim_{k\to\infty} n_k^{-1}\operatorname{Tr}\mathbf{C}_\eta^\infty(\tau) = c_\eta^{\mathrm{tti}}(\tau)$ for the quantity

$$c_\eta^{\mathrm{tti}}(\tau) = c_\eta^{\mathrm{tti}}(\infty) + \int_\iota^\infty e^{-a\tau}\mu_\eta(\mathrm{d}a).$$

By an argument parallel to (153) using Theorem 4.3 and (118), this shows $|C_\eta(s+\tau,s) - c_\eta^{\mathrm{tti}}(\tau)| \leq Ce^{-cs}$, establishing (31).

Finally, for (32), observe that by Theorem 4.3, $\lim_{n,d\to\infty} d^{-1}P_s(\boldsymbol{\theta}^0)^\top\boldsymbol{\theta}^* = C_\theta(s,*)$. Noting that

$$\limsup_{n,d\to\infty} d^{-1}\left|(P_s(\boldsymbol{\theta}^0) - \langle\boldsymbol{\theta}\rangle)^\top\boldsymbol{\theta}^*\right| \leq \limsup_{n,d\to\infty} d^{-1}\|P_s\boldsymbol{\theta}^0 - \langle\boldsymbol{\theta}\rangle\|_2 \cdot \|\boldsymbol{\theta}^*\|_2 \leq Ce^{-cs}$$

by (124), this implies the existence of the limit

$$c_\theta(*) := \lim_{n,d\to\infty} d^{-1}\langle\boldsymbol{\theta}\rangle^\top\boldsymbol{\theta}^* = \lim_{s\to\infty} C_\theta(s,*), \tag{158}$$

which satisfies $|C_\theta(s,*) - c_\theta(*)| \leq Ce^{-cs}$. This shows (32). $\qquad\square$

**Lemma 4.10.** *Under Assumptions 2.1, 2.2(a), and 2.7(a), the DMFT system prescribed by Theorem 2.3(a) satisfies the conditions of Definition 2.4(2) with $\varepsilon(t) = Ce^{-ct}$ for some constants $C, c > 0$.*

*Proof.* We again restrict to the almost sure event where the convergence statements of Theorem 4.3 hold, and where $\mathcal{E}(C_0, C_{\text{LSI}}) \cap \{\|\boldsymbol{\theta}^0\|_2^2 \leq C_0 d\}$ holds for all large $n, d$.

Consider first $R_\theta(t, s)$. By (138) of Lemma 4.8 and the convergence $R_\theta(t, s) = \lim_{n,d \to \infty} d^{-1} \mathbf{R}_\theta(t, s)$ of Theorem 4.3,

$$|R_\theta(t, s)| \leq C e^{-ct} \text{ for all } s \leq t/2. \tag{159}$$

For $s \geq t/2$, note that the forms of (149) and (152) imply that both $d^{-1} \operatorname{Tr} \mathbf{C}_\theta^\infty(\tau)$ and $c_\theta^{\text{tti}}(\tau)$ are convex and differentiable in $\tau \geq 0$. Then, along the subsequence $\{(n_k, d_k)\}_{k \geq 1}$ of the preceding proof, the pointwise convergence $\lim_{k \to \infty} d_k^{-1} \operatorname{Tr} \mathbf{C}_\theta^\infty(\tau) = c_\theta^{\text{tti}}(\tau)$ implies also $\lim_{k \to \infty} d_k^{-1} \partial_\tau \operatorname{Tr} \mathbf{C}_\theta^\infty(\tau) = \partial_\tau c_\theta^{\text{tti}}(\tau)$ for each $\tau \geq 0$ (c.f. [85, Theorem 25.7]). By the fluctuation-dissipation relation of Lemma 4.1 applied with $A = B = e_j$ for each $j = 1, \ldots, d$, we have $\partial_\tau \operatorname{Tr} \mathbf{C}_\theta^\infty(\tau) = -\operatorname{Tr} \mathbf{R}_\theta^\infty(\tau)$. Then, defining $r_\theta^{\text{tti}}(\tau) = -\partial_\tau c_\theta^{\text{tti}}(\tau)$, this shows $\lim_{k \to \infty} d_k^{-1} \operatorname{Tr} \mathbf{R}_\theta^\infty(\tau) = r_\theta^{\text{tti}}(\tau)$. Combining with $\lim_{k \to \infty} d_k^{-1} \operatorname{Tr} \mathbf{R}_\theta(s + \tau, s) = R_\theta(s + \tau, s)$ from Theorem 4.3, for any $s, \tau \geq 0$ we have that

$$\left| R_\theta(s + \tau, s) - r_\theta^{\text{tti}}(\tau) \right| \leq \limsup_{k \to \infty} \left| d_k^{-1} \operatorname{Tr} \mathbf{R}_\theta(s + \tau, s) - d_k^{-1} \operatorname{Tr} \mathbf{R}_\theta^\infty(\tau) \right| \leq C e^{-cs},$$

where the last inequality applies (117). In particular, for any $t \geq 0$,

$$|R_\theta(t, s) - r_\theta^{\text{tti}}(t - s)| \leq C e^{-c't} \text{ for all } s \in [t/2, t]. \tag{160}$$

Together, (159) and (160) imply (34). The statement (35) follows analogously, and we omit this for brevity. $\square$

*Proof of Theorem 2.9.* This follows from Lemmas 4.9 and 4.10. $\square$

## 4.5 Limit MSE and free energy

We now show Corollary 2.10 on the asymptotic values of the mean-squared-errors and the free energy.

**Proposition 4.11.** *Suppose Assumptions 2.1, 2.2(a), and 2.7(a) hold. Let* $\text{YMSE}_*$ *and the marginal likelihood* $\mathsf{P}_g(\mathbf{y} \mid \mathbf{X})$ *be as defined in Corollary 2.10. Let* $\mathbb{E}[\cdot \mid \mathbf{X}]$ *denote the expectation with respect to* $\theta_j^* \overset{iid}{\sim} g_*$ *and* $\varepsilon_i \overset{iid}{\sim} \mathcal{N}(0, \sigma^2)$ *conditioning on* $\mathbf{X}$*. Then almost surely,*

$$\lim_{n,d \to \infty} d^{-1} \log \mathsf{P}_g(\mathbf{y} \mid \mathbf{X}) - d^{-1} \mathbb{E}[\log \mathsf{P}_g(\mathbf{y} \mid \mathbf{X}) \mid \mathbf{X}] = 0, \qquad \lim_{n,d \to \infty} \text{YMSE}_* - \mathbb{E}[\text{YMSE}_* \mid \mathbf{X}] = 0. \tag{161}$$

*Proof.* We condition on $\mathbf{X}$ throughout, and restrict to the $\mathbf{X}$-dependent event

$$\{\|\mathbf{X}\|_{\text{op}} \leq C_0 \text{ and (46) holds}\}.$$

Note that by assumption, this event holds a.s. for all large $n, d$ and does not depend on $\boldsymbol{\theta}^*, \boldsymbol{\varepsilon}$.

For the first statement, let us consider

$$Z(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon}) = \log \int \exp \left( -\frac{1}{2\sigma^2} \|\mathbf{X}\boldsymbol{\theta}^* + \boldsymbol{\varepsilon} - \mathbf{X}\boldsymbol{\theta}\|_2^2 + \sum_{j=1}^d \log g(\theta_j) \right) d\boldsymbol{\theta}$$

(which coincides with $\log \mathsf{P}_g(\mathbf{y} \mid \mathbf{X})$ up to an additive constant) as a function of $(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon})$. Then

$$\nabla_{\boldsymbol{\theta}^*} Z(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon}) = -\frac{1}{\sigma^2} \mathbf{X}^\top (\mathbf{X}\boldsymbol{\theta}^* + \boldsymbol{\varepsilon} - \mathbf{X}\langle\boldsymbol{\theta}\rangle), \qquad \nabla_{\boldsymbol{\varepsilon}} Z(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon}) = -\frac{1}{\sigma^2} (\mathbf{X}\boldsymbol{\theta}^* + \boldsymbol{\varepsilon} - \mathbf{X}\langle\boldsymbol{\theta}\rangle).$$

Under Assumption 2.1, note that $\boldsymbol{\theta}^*$ and $\boldsymbol{\varepsilon}$ have independent subgaussian entries, so there are constants $C_1, c > 0$ such that (c.f. [81, Eq. (3.1)])

$$\mathbb{P}[\|\boldsymbol{\theta}^*\|_2^2 + \|\boldsymbol{\varepsilon}\|_2^2 > C_1 d] \leq e^{-cd}. \tag{162}$$

When $\|\boldsymbol{\theta}^*\|_2^2 + \|\boldsymbol{\varepsilon}\|_2^2 \leq C_1 d$, we have the bound $\langle \|\boldsymbol{\theta}\|_2^2 \rangle \leq C d$ from (106). Applying this and $\|\mathbf{X}\|_{\text{op}} \leq C_0$,

$$\|\nabla_{(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon})} Z(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon})\|_2 \mathbf{1}\{\|\boldsymbol{\theta}^*\|_2^2 + \|\boldsymbol{\varepsilon}\|_2^2 \leq C_1 d\} \leq L\sqrt{d}$$

47

for a constant $L > 0$. Thus $Z(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon})$ is $L\sqrt{d}$-Lipschitz on $\{\|\boldsymbol{\theta}^*\|_2^2 + \|\boldsymbol{\varepsilon}\|_2^2 \leq C_1 d\}$, so its Lipschitz extension

$$\tilde{Z}(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon}) = \inf_{\mathbf{x} \in \mathbb{R}^{d+n} : \|\mathbf{x}\|_2^2 \leq C_1 d} Z(\mathbf{x}) + L\sqrt{d}\|\mathbf{x} - (\boldsymbol{\theta}^*, \boldsymbol{\varepsilon})\|_2$$

is globally $L\sqrt{d}$-Lipschitz on $\mathbb{R}^{d+n}$ and $\tilde{Z}(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon}) = Z(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon})$ over $\{\|\boldsymbol{\theta}^*\|_2^2 + \|\boldsymbol{\varepsilon}\|_2^2 \leq C_1 d\}$. Under Assumption 2.1, the joint distribution of $(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon})$ satisfies a log-Sobolev inequality by tensorization, implying the Lipschitz concentration

$$\mathbb{P}[|\tilde{Z}(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon}) - \mathbb{E}[\tilde{Z}(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon}) \mid \mathbf{X}]| \geq td \mid \mathbf{X}] \leq 2e^{-t^2 d/(2L^2)}. \tag{163}$$

We may bound

$$\begin{aligned}
&|\mathbb{E}[\tilde{Z}(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon}) \mid \mathbf{X}] - \mathbb{E}[Z(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon}) \mid \mathbf{X}]| \\
&\leq \mathbb{E}\Big[\mathbf{1}\{\|\boldsymbol{\theta}^*\|_2^2 + \|\boldsymbol{\varepsilon}\|_2^2 \geq C_1 d\}\Big(|\tilde{Z}(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon})| + |Z(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon})|\Big) \; \Big| \; \mathbf{X}\Big] \\
&\leq \mathbb{P}[\|\boldsymbol{\theta}^*\|_2^2 + \|\boldsymbol{\varepsilon}\|_2^2 \geq C_1 d \mid \mathbf{X}]^{1/2}\Big((\mathbb{E}[\tilde{Z}(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon}) \mid \mathbf{X}]^2)^{1/2} + (\mathbb{E}[Z(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon}) \mid \mathbf{X}]^2)^{1/2}\Big)
\end{aligned}$$

Applying the upper bound $Z(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon}) \leq \log \int \exp(\sum_{j=1}^d \log g(\theta_j)) \mathrm{d}\boldsymbol{\theta} = 0$, Jensen's inequality lower bound $Z(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon}) \geq \mathbb{E}_g[-\frac{1}{2\sigma^2}\|\mathbf{X}\boldsymbol{\theta}^* + \boldsymbol{\varepsilon} - \mathbf{X}\boldsymbol{\theta}\|_2^2]$ where $\mathbb{E}_g[\cdot]$ is the expectation over $\theta_j \overset{iid}{\sim} g$, and $|\tilde{Z}(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon}) - Z(0)| = |\tilde{Z}(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon}) - \tilde{Z}(0)| \leq L\sqrt{d}(\|\boldsymbol{\theta}^*\|_2^2 + \|\boldsymbol{\varepsilon}\|_2^2)^{1/2}$, we obtain

$$|\mathbb{E}[\tilde{Z}(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon}) \mid \mathbf{X}] - \mathbb{E}[Z(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon}) \mid \mathbf{X}]| \leq \mathbb{P}[\|\boldsymbol{\theta}^*\|_2^2 + \|\boldsymbol{\varepsilon}\|_2^2 \geq C_1 d \mid \mathbf{X}]^{1/2} \cdot Cd \leq e^{-c'd}$$

for all large $n, d$, the last inequality applying (162). Thus (163) and (162) imply

$$\mathbb{P}[|Z(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon}) - \mathbb{E}[Z(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon}) \mid \mathbf{X}]| \geq td + e^{-c'd} \mid \mathbf{X}] \leq 2e^{-t^2 d/(2L^2)} + e^{-cd},$$

implying the first statement of (161) by the Borel-Cantelli lemma.

For the second statement, let us write

$$n\,\mathrm{YMSE}_*(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon}) = \|\mathbf{X}\boldsymbol{\theta}^* + \boldsymbol{\varepsilon} - \mathbf{X}\langle\boldsymbol{\theta}\rangle\|_2^2$$

viewed also as a function of $(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon})$. Writing $\kappa_2(\cdot)$ for the covariance associated to the posterior mean $\langle\cdot\rangle$, differentiating in $(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon})$ gives, for any unit vectors $\mathbf{u} \in \mathbb{R}^d$ and $\mathbf{v} \in \mathbb{R}^n$,

$$\begin{aligned}
\mathbf{u}^\top \nabla_{\boldsymbol{\theta}^*}[n\,\mathrm{YMSE}_*] &= 2(\mathbf{X}\boldsymbol{\theta}^* + \boldsymbol{\varepsilon} - \mathbf{X}\langle\boldsymbol{\theta}\rangle)^\top \mathbf{X}\mathbf{u} - \frac{2}{\sigma^2}\kappa_2\Big(\boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{X}\mathbf{u}, (\mathbf{X}\boldsymbol{\theta}^* + \boldsymbol{\varepsilon} - \mathbf{X}\langle\boldsymbol{\theta}\rangle)^\top \mathbf{X}\boldsymbol{\theta}\Big), \\
\mathbf{v}^\top \nabla_{\boldsymbol{\varepsilon}}[n\,\mathrm{YMSE}_*] &= 2(\mathbf{X}\boldsymbol{\theta}^* + \boldsymbol{\varepsilon} - \mathbf{X}\langle\boldsymbol{\theta}\rangle)^\top \mathbf{v} - \frac{2}{\sigma^2}\kappa_2\Big(\boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{v}, (\mathbf{X}\boldsymbol{\theta}^* + \boldsymbol{\varepsilon} - \mathbf{X}\langle\boldsymbol{\theta}\rangle)^\top \mathbf{X}\boldsymbol{\theta}\Big).
\end{aligned}$$

The Poincaré inequality implied by the assumed LSI for $\mathsf{P}_g(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y})$ shows, for any vector $\mathbf{x} \in \mathbb{R}^d$,

$$\kappa_2(\mathbf{x}^\top \boldsymbol{\theta}, \mathbf{u}^\top \boldsymbol{\theta}) \leq C\|\mathbf{x}\|_2^2.$$

On the event $\{\|\boldsymbol{\theta}^*\|_2^2 + \|\boldsymbol{\varepsilon}\|_2^2 \leq C_1 d\}$, applying this Poincaré bound, Cauchy-Schwarz for $\kappa_2(\cdot)$, and $\|\mathbf{X}\|_{\mathrm{op}} \leq C_0$ and $\langle\|\boldsymbol{\theta}\|_2^2\rangle \leq Cd$ from (106), we obtain $|\mathbf{u}^\top \nabla_{\boldsymbol{\theta}^*}[n\,\mathrm{YMSE}_*]| \leq C\sqrt{d}$ and $|\mathbf{v}^\top \nabla_{\boldsymbol{\varepsilon}}[n\,\mathrm{YMSE}_*]| \leq C\sqrt{d}$ for any unit vectors $\mathbf{u}, \mathbf{v}$, and hence

$$\|\nabla_{\boldsymbol{\theta}^*, \boldsymbol{\varepsilon}}[n\,\mathrm{YMSE}_*]\|_2 \mathbf{1}\{\|\boldsymbol{\theta}^*\|_2^2 + \|\boldsymbol{\varepsilon}\|_2^2 \leq C_1 d\} \leq L\sqrt{d}$$

for some constant $L > 0$. So $n\,\mathrm{YMSE}_*$ is $L\sqrt{d}$-Lipschitz in $(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon})$ on $\{\|\boldsymbol{\theta}^*\|_2^2 + \|\boldsymbol{\varepsilon}\|_2^2 \leq C_1 d\}$. For any $(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon})$, we also have the bound $|n\,\mathrm{YMSE}_*(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon})| \leq C(\|\boldsymbol{\theta}^*\|_2^2 + \|\boldsymbol{\varepsilon}\|_2^2)^{1/2}$ by (106), so the second statement of (161) follows from the same Lipschitz extension and concentration argument as above. $\quad\square$

*Proof of Corollary 2.10(a).* We restrict to the almost sure event where $\mathcal{E}(C_0, C_{\mathrm{LSI}})$ holds for all large $n, d$. Observe that by (148) and (150),

$$\mathrm{MSE} = d^{-1}\langle\|\boldsymbol{\theta} - \langle\boldsymbol{\theta}\rangle\|_2^2\rangle = d^{-1}\big(\langle\|\boldsymbol{\theta}\|_2^2\rangle - \|\langle\boldsymbol{\theta}\rangle\|_2^2\big) = m_d = \mu_d([\iota, \infty)),$$

so $\lim_{n,d\to\infty} \mathrm{MSE} = \mu_\theta([\iota,\infty)) = c_\theta^{\mathrm{tti}}(0) - c_\theta^{\mathrm{tti}}(\infty)$ by (152). Also

$$\mathrm{MSE}_* = d^{-1}\|\boldsymbol{\theta}^* - \langle\boldsymbol{\theta}\rangle\|_2^2 = d^{-1}\big(\|\boldsymbol{\theta}^*\|_2^2 - 2\langle\boldsymbol{\theta}\rangle^\top\boldsymbol{\theta}^* + \|\langle\boldsymbol{\theta}\rangle\|_2^2\big),$$

so $\lim_{n,d\to\infty} \mathrm{MSE}_* = \mathbb{E}[\theta^{*2}] - 2c_\theta(*) + c_\theta^{\mathrm{tti}}(\infty)$ by Assumption 2.1 and the definitions (146) and (158). Thus $\mathrm{MSE} \to \mathrm{mse}$ and $\mathrm{MSE}_* \to \mathrm{mse}_*$ for the quantities $\mathrm{mse}, \mathrm{mse}_*$ defined in (42).

Similarly

$$\mathrm{YMSE} = n^{-1}\langle\|\mathbf{X}\boldsymbol{\theta} - \mathbf{X}\langle\boldsymbol{\theta}\rangle\|_2^2\rangle = n^{-1}\big(\langle\|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2\rangle - \|\mathbf{X}\langle\boldsymbol{\theta}\rangle - \mathbf{y}\|_2^2\big).$$

Then by (156) and (157), $\lim_{n,d\to\infty} n^{-1}\|\mathbf{X}\langle\boldsymbol{\theta}\rangle - \mathbf{y}\|_2^2 = \frac{\sigma^4}{\delta}c_\eta^{\mathrm{tti}}(\infty)$ and $\mathrm{YMSE} = \frac{\sigma^4}{\delta}\mu_{\eta,n}([\iota,\infty)) \to \frac{\sigma^4}{\delta}(c_\eta^{\mathrm{tti}}(0) - c_\eta^{\mathrm{tti}}(\infty)) = \mathrm{ymse}$ as defined in (42). For $\mathrm{YMSE}_*$, writing $\mathbb{E}[\cdot \mid \mathbf{X}]$ for the expectation over $(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon})$ as in Proposition 4.11, observe first that

$$n^{-1}\mathbb{E}[\|\mathbf{X}\langle\boldsymbol{\theta}\rangle - \mathbf{y}\|_2^2 \mid \mathbf{X}] = n^{-1}\mathbb{E}[\|\mathbf{X}\langle\boldsymbol{\theta}\rangle - \mathbf{X}\boldsymbol{\theta}^*\|_2^2 \mid \mathbf{X}] - 2n^{-1}\mathbb{E}[\boldsymbol{\varepsilon}^\top(\mathbf{X}\langle\boldsymbol{\theta}\rangle - \mathbf{X}\boldsymbol{\theta}^*) + \sigma^2 \mid \mathbf{X}],$$

and Gaussian integration-by-parts gives

$$\mathbb{E}[\boldsymbol{\varepsilon}^\top(\mathbf{X}\langle\boldsymbol{\theta}\rangle - \mathbf{X}\boldsymbol{\theta}^*) \mid \mathbf{X}] = \mathbb{E}[\boldsymbol{\varepsilon}^\top\mathbf{X}\langle\boldsymbol{\theta}\rangle \mid \mathbf{X}] = \mathbb{E}[\langle\|\mathbf{X}\boldsymbol{\theta} - \mathbf{X}\boldsymbol{\theta}^*\|_2^2\rangle \mid \mathbf{X}] - \mathbb{E}[\|\mathbf{X}\langle\boldsymbol{\theta}\rangle - \mathbf{X}\boldsymbol{\theta}^*\|_2^2 \mid \mathbf{X}]$$
$$= \mathbb{E}[\langle\|\mathbf{X}\boldsymbol{\theta} - \mathbf{X}\langle\boldsymbol{\theta}\rangle\|_2^2\rangle \mid \mathbf{X}] = n\,\mathbb{E}[\mathrm{YMSE} \mid \mathbf{X}].$$

Thus

$$\mathbb{E}[\mathrm{YMSE}_* \mid \mathbf{X}] = n^{-1}\mathbb{E}[\|\mathbf{X}\langle\boldsymbol{\theta}\rangle - \mathbf{X}\boldsymbol{\theta}^*\|_2^2 \mid \mathbf{X}] = n^{-1}\mathbb{E}[\|\mathbf{X}\langle\boldsymbol{\theta}\rangle - \mathbf{y}\|_2^2 \mid \mathbf{X}] + 2\mathbb{E}[\mathrm{YMSE} \mid \mathbf{X}] - \sigma^2. \qquad (164)$$

We remark that $n^{-1}\|\mathbf{X}\langle\boldsymbol{\theta}\rangle - \mathbf{y}\|_2^2$ and YMSE are bounded for all large $n, d$ on the event $\mathcal{E}(C_0, C_{\mathrm{LSI}})$, by the bound for $\langle\|\boldsymbol{\theta}\|_2^2\rangle \le C$ from (106). Thus, applying $\mathrm{YMSE} \to \mathrm{ymse}$ and $n^{-1}\|\mathbf{X}\langle\boldsymbol{\theta}\rangle - \mathbf{y}\|_2^2 \to \frac{\sigma^4}{\delta}c_\eta^{\mathrm{tti}}(\infty)$ as argued above and dominated convergence, the right side of (164) converges to $\mathrm{ymse}_* = \frac{\sigma^4}{\delta}(2c_\eta^{\mathrm{tti}}(0) - c_\eta^{\mathrm{tti}}(\infty)) - \sigma^2$ as defined in (42). Then the concentration of $\mathrm{YMSE}_*$ established in Proposition 4.11 combined with (164) show $\lim_{n,d\to\infty} \mathrm{YMSE}_* = \mathrm{ymse}_*$.

To show the last statement (48), conditional on $\mathbf{X}, \boldsymbol{\theta}^*, \boldsymbol{\varepsilon}$ and averaging over the initial condition $\boldsymbol{\theta}^0 \sim q_0 = g_0^{\otimes d}$, let $q_t$ be the conditional law of $\boldsymbol{\theta}^t$. Consider a coupling of a posterior sample $\boldsymbol{\theta} \sim q$ with $\boldsymbol{\theta}^t \sim q_t$ such that $\langle\|\boldsymbol{\theta}^t - \boldsymbol{\theta}\|_2^2\rangle = W_2(q_t, q)^2$, where $\langle\cdot\rangle$ denotes the expectation under this coupling and $W_2(\cdot)$ is the Wasserstein-2 distance, both conditional on $\mathbf{X}, \boldsymbol{\theta}^*, \boldsymbol{\varepsilon}$. For a given realization of $(\boldsymbol{\theta}^t, \boldsymbol{\theta})$ from this coupling, considering the coordinatewise coupling of $\frac{1}{d}\sum_{j=1}^d \delta_{(\theta_j^*, \theta_j^t)}$ with $\frac{1}{d}\sum_{j=1}^d \delta_{(\theta_j^*, \theta_j)}$ shows

$$W_2\left(\frac{1}{d}\sum_{j=1}^d \delta_{(\theta_j^*, \theta_j^t)}, \frac{1}{d}\sum_{j=1}^d \delta_{(\theta_j^*, \theta_j)}\right)^2 \le \frac{1}{d}\sum_{j=1}^d (\theta_j^t - \theta_j)^2 = \frac{1}{d}\|\boldsymbol{\theta}^t - \boldsymbol{\theta}\|_2^2.$$

Then

$$\left\langle W_2\left(\frac{1}{d}\sum_{j=1}^d \delta_{(\theta_j^*, \theta_j^t)}, \frac{1}{d}\sum_{j=1}^d \delta_{(\theta_j^*, \theta_j)}\right)^2 \right\rangle \le \frac{1}{d}\langle\|\boldsymbol{\theta}^t - \boldsymbol{\theta}\|_2^2\rangle = \frac{1}{d}W_2(q_t, q)^2.$$

Applying Lemmas 4.4 and 4.6, $W_2(q_t, q)^2 \le Ce^{-ct}(\langle\|\boldsymbol{\theta}\|_2^2\rangle + \langle\|\boldsymbol{\theta}^0\|_2^2\rangle) \le C'de^{-ct}$ on the event $\mathcal{E}(C_0, C_{\mathrm{LSI}})$, for some constants $C, C', c > 0$. So on this event,

$$\limsup_{n,d\to\infty} \left\langle W_2\left(\frac{1}{d}\sum_{j=1}^d \delta_{(\theta_j^*, \theta_j^t)}, \frac{1}{d}\sum_{j=1}^d \delta_{(\theta_j^*, \theta_j)}\right)^2 \right\rangle \le C'e^{-ct}. \qquad (165)$$

Now by Theorem 2.3(a), for each fixed $t \ge 0$, almost surely with respect to the randomness of both $\mathbf{X}, \boldsymbol{\theta}^*, \boldsymbol{\varepsilon}$ and $\boldsymbol{\theta}^0, \{\mathbf{b}^t\}_{t\ge0}$ defining $\{\boldsymbol{\theta}^t\}_{t\ge0}$, we have

$$\lim_{n,d\to\infty} W_2\left(\frac{1}{d}\sum_{j=1}^d \delta_{(\theta_j^*, \theta_j^t)}, \mathsf{P}(\theta^*, \theta^t)\right)^2 = 0 \qquad (166)$$

where $\mathsf{P}(\theta^*, \theta^t)$ here is the law of $(\theta^*, \theta^t)$ in the DMFT system. To take an expectation over the randomness of $\boldsymbol{\theta}^0$ and $\{\mathbf{b}^t\}_{t \geq 0}$, note that from the definition (103), we have

$$\boldsymbol{\theta}^t = \boldsymbol{\theta}^0 + \int_0^t \nabla_{\boldsymbol{\theta}} \log q(\boldsymbol{\theta}^s) \, ds + \sqrt{2}\, \mathbf{b}^t = \boldsymbol{\theta}^0 + \int_0^t \left[ \frac{1}{\sigma^2} \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}^s) + (\log g)'(\boldsymbol{\theta}^s) \right] ds + \sqrt{2}\, \mathbf{b}^t,$$

where $(\log g)'$ is applied entrywise. Then on $\mathcal{E}(C_0, C_{\mathrm{LSI}})$, by the Lipschitz continuity of $(\log g)'(\theta)$, this implies for a constant $C > 0$ that

$$d^{-1/2} \|\boldsymbol{\theta}^t\|_2 \leq \int_0^t C d^{-1/2} \|\boldsymbol{\theta}^s\|_2 \, ds + Ct + d^{-1/2} \|\boldsymbol{\theta}^0\|_2 + \sqrt{2}\, d^{-1/2} \|\mathbf{b}^t\|_2.$$

Then for any $T > 0$, Gronwall's inequality gives, for a constant $C > 0$,

$$\sup_{t \in [0,T]} d^{-1/2} \|\boldsymbol{\theta}^t\|_2 \leq C e^{CT} \left( T + d^{-1/2} \|\boldsymbol{\theta}^0\|_2 + d^{-1/2} \sup_{t \in [0,T]} \|\mathbf{b}^t\|_2 \right) \tag{167}$$

For any $p > 1$, applying

$$\left( \sup_{t \in [0,T]} d^{-1} \|\mathbf{b}^t\|_2^2 \right)^p \leq \sup_{t \in [0,T]} d^{-1} \sum_{j=1}^d |b_j^t|^{2p} \leq d^{-1} \sum_{j=1}^d \sup_{t \in [0,T]} |b_j^t|^{2p}$$

and Doob's $L^p$-maximal inequality, we have that $\langle (\sup_{t \in [0,T]} d^{-1} \|\mathbf{b}^t\|_2^2)^p \rangle$ is bounded by a $(T, p)$-dependent constant. Similarly $\langle (d^{-1} \|\boldsymbol{\theta}^0\|_2)^p \rangle$ is bounded by a $(T, p)$-dependent constant, so

$$\left\langle \left( \sup_{t \in [0,T]} d^{-1} \|\boldsymbol{\theta}^t\|_2^2 \right)^p \right\rangle \leq C_{T,p}$$

for a constant $C_{T,p} > 0$, where $\langle \cdot \rangle$ averages over $\boldsymbol{\theta}^0$ and $\{\mathbf{b}^t\}_{t \geq 0}$. Since $W_2(\frac{1}{d} \sum_{j=1}^d \delta_{(\theta_j^*, \theta_j^t)}, \mathsf{P}(\theta^*, \theta^t))^2 \leq C(d^{-1} \|\boldsymbol{\theta}^*\|_2^2 + d^{-1} \|\boldsymbol{\theta}^t\|_2^2 + \mathbb{E}(\theta^*)^2 + \mathbb{E}(\theta^t)^2)$, this implies on the event $\mathcal{E}(C_0, C_{\mathrm{LSI}})$ that

$$\left\langle \sup_{t \in [0,T]} W_2 \left( \frac{1}{d} \sum_{j=1}^d \delta_{(\theta_j^*, \theta_j^t)}, \, \mathsf{P}(\theta^*, \theta^t) \right)^{2p} \right\rangle \leq C_{T,p}' \tag{168}$$

for a different constant $C_{T,p}' > 0$. In particular, for any fixed $t \geq 0$ and $p > 1$, the squared Wasserstein-2 distance in (166) is uniformly bounded in $L^p$ and hence uniformly integrable with respect to $\langle \cdot \rangle$ for all large $n, d$, so dominated convergence implies, almost surely,

$$\lim_{n,d \to \infty} \left\langle W_2 \left( \frac{1}{d} \sum_{j=1}^d \delta_{(\theta_j^*, \theta_j^t)}, \, \mathsf{P}(\theta^*, \theta^t) \right)^2 \right\rangle = 0. \tag{169}$$

Combining (165) and (169) shows that for any fixed $t \geq 0$, almost surely,

$$\limsup_{n,d \to \infty} \left\langle W_2 \left( \frac{1}{d} \sum_{j=1}^d \delta_{(\theta_j^*, \theta_j)}, \mathsf{P}_{g_*, \omega_*; g, \omega} \right)^2 \right\rangle \leq C \left( e^{-ct} + W_2(\mathsf{P}(\theta^*, \theta^t), \, \mathsf{P}_{g_*, \omega_*; g, \omega})^2 \right).$$

By Theorem 2.5, we have

$$\lim_{t \to \infty} W_2(\mathsf{P}(\theta^*, \theta^t), \, \mathsf{P}_{g_*, \omega_*; g, \omega}) = 0$$

so taking the limit $t \to \infty$ shows (48). $\qquad \square$

To show Corollary 2.10(b) on the asymptotic free energy, we will apply an I-MMSE argument, together with the following proposition which guarantees continuity of $\mathrm{mse}, \mathrm{mse}_*$ in the noise variance $\sigma^2$. In the later proof of Theorem 2.13, we will require also continuity in the prior parameter $\alpha$; thus we establish both statements here.

**Lemma 4.12.** *Suppose Assumptions 2.1 and 2.2(b) hold. Fix any open subset $O \subset \mathbb{R}^K$, and suppose also that Assumption 2.7 holds for each $g \in \{g(\cdot, \alpha) : \alpha \in O\}$, where the constant $C_{\mathrm{LSI}} > 0$ is uniform over $\alpha \in O$. Consider any noise variance $\tilde{\sigma}^2 \geq \sigma^2$, and define $\mathrm{mse}(\tilde{\sigma}^2, \alpha), \mathrm{mse}_*(\tilde{\sigma}^2, \alpha)$ by (42) via the (approximately-TTI) DMFT limit of the Langevin dynamics (7) with a fixed prior $g(\cdot, \alpha)$ in the linear model (4) with noise variance $\tilde{\sigma}^2$.*

*Then over any compact interval $I \subset [\sigma^2, \infty)$ and compact subset $S \subset O$, $\mathrm{mse}(\tilde{\sigma}^2, \alpha), \mathrm{mse}_*(\tilde{\sigma}^2, \alpha)$ are Lipschitz functions of $(\tilde{\sigma}^2, \alpha) \in I \times S$.*

*Proof.* Consider noise/prior parameters $(s^2, \alpha)$ and $(\tilde{s}^2, \tilde{\alpha})$, where $s^2, \tilde{s}^2 \geq \sigma^2$. Let us couple the linear models with noise variances $s^2$ and $\tilde{s}^2$ by $\mathbf{y} = \mathbf{X}\boldsymbol{\theta}^* + s\mathbf{z}$ and $\tilde{\mathbf{y}} = \mathbf{X}\boldsymbol{\theta}^* + \tilde{s}\mathbf{z}$, where $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I})$. Fixing $\mathbf{X}, \boldsymbol{\theta}^*, \mathbf{z}$, let us denote

$$U(\boldsymbol{\theta}) = -\frac{1}{2s^2}\|\mathbf{X}\boldsymbol{\theta}^* + s\mathbf{z} - \mathbf{X}\boldsymbol{\theta}\|_2^2 + \sum_{j=1}^d \log g(\theta_j, \alpha)$$

so that $q(\boldsymbol{\theta}) \propto e^{U(\boldsymbol{\theta})}$ is the posterior law given $(\mathbf{X}, \mathbf{y})$ under parameters $(s^2, \alpha)$. Denote similarly $\tilde{U}(\boldsymbol{\theta})$ with $(\tilde{s}^2, \tilde{\alpha})$ in place of $(s^2, \alpha)$, and $\tilde{q}(\boldsymbol{\theta}) \propto e^{\tilde{U}(\boldsymbol{\theta})}$ as the posterior law given $(\mathbf{X}, \tilde{\mathbf{y}})$. We condition on $\mathbf{X}, \boldsymbol{\theta}^*, \mathbf{z}$ and restrict to the event

$$\mathcal{E}'(C_0, C_{\mathrm{LSI}}) = \{\|\mathbf{X}\|_{\mathrm{op}} \leq C_0, \|\boldsymbol{\theta}^*\|_2^2, \|\mathbf{z}\|_2^2 \leq C_0 d, \text{ (46) holds for both } q \text{ and } \tilde{q}\},$$

which by assumption holds a.s. for all large $n, d$. We first derive a bound on the Wasserstein-2 distance between $q$ and $\tilde{q}$, conditional on $\mathbf{X}, \boldsymbol{\theta}^*, \mathbf{z}$.

Let $\{\boldsymbol{\theta}^t\}_{t \geq 0}$ be the Langevin diffusion (103) with fixed prior $g(\cdot, \alpha)$ and stationary distribution $q(\boldsymbol{\theta})$, initialized as $\boldsymbol{\theta}^0 \sim q_0$ where $q_0$ has finite second moment and finite entropy. Let us write $\langle f(\boldsymbol{\theta}^t) \rangle$ for the expectation over $\boldsymbol{\theta}^0$ and $\{\mathbf{b}^t\}_{t \geq 0}$ defining (103), conditional on $\mathbf{X}, \boldsymbol{\theta}^*, \mathbf{z}$. We apply the following argument of [86] to bound the KL-divergence $\mathrm{D}_{\mathrm{KL}}(q_t \| \tilde{q})$ conditional on $\mathbf{X}, \boldsymbol{\theta}^*, \mathbf{z}$: Differentiating this KL-divergence in time,

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathrm{D}_{\mathrm{KL}}(q_t \| \tilde{q}) = \frac{\mathrm{d}}{\mathrm{d}t}\int q_t(\log q_t - \log \tilde{q})$$

$$= \int \left(\frac{\mathrm{d}}{\mathrm{d}t}q_t\right)(\log q_t - \log \tilde{q}) + \underbrace{\int \frac{q_t}{q_t}\left(\frac{\mathrm{d}}{\mathrm{d}t}q_t\right)}_{=0} = \int \left(\frac{\mathrm{d}}{\mathrm{d}t}q_t\right)(\log q_t - \tilde{U} + \log \tilde{Z}).$$

The law of $\boldsymbol{\theta}^t$ conditional on $\mathbf{X}, \boldsymbol{\theta}^*, \mathbf{z}$ admits a density $q_t$ which is described by the Fokker-Planck equation

$$\frac{\mathrm{d}}{\mathrm{d}t}q_t = \nabla \cdot [q_t \nabla(\log q_t - U)]$$

with initial condition $q_t|_{t=0} = q_0$. Then, applying this Fokker-Planck equation and integrating by parts, we obtain

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathrm{D}_{\mathrm{KL}}(q_t \| \tilde{q}) = -\int q_t \nabla(\log q_t - U)^\top \nabla(\log q_t - \tilde{U})$$

$$= -\int q_t \|\nabla(\log q_t - \tilde{U})\|_2^2 - \int q_t \nabla(\tilde{U} - U)^\top \nabla(\log q_t - \tilde{U})$$

$$\leq -(1/2)\int q_t \|\nabla(\log q_t - \tilde{U})\|_2^2 + (1/2)\int q_t \|\nabla(\tilde{U} - U)\|_2^2,$$

the last step applying Cauchy-Schwarz for the second term. By the LSI for $\tilde{q}$, the first term (the relative Fisher information) is lower bounded as

$$\int q_t \|\nabla(\log q_t - \tilde{U})\|_2^2 = \int q_t \left\|\nabla \log \frac{q_t}{\tilde{q}}\right\|_2^2 \geq \frac{1}{2C_{\mathrm{LSI}}}\mathrm{D}_{\mathrm{KL}}(q_t \| \tilde{q}).$$

Thus

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathrm{D}_{\mathrm{KL}}(q_t \| \tilde{q}) \leq -\frac{1}{4C_{\mathrm{LSI}}}\mathrm{D}_{\mathrm{KL}}(q_t \| \tilde{q}) + \frac{1}{2}\underbrace{\langle \|\nabla\tilde{U}(\boldsymbol{\theta}^t) - \nabla U(\boldsymbol{\theta}^t)\|_2^2 \rangle}_{:=\Delta(t)}.$$

51

Integrating this inequality shows, for some constants $C, c > 0$ depending only on $C_{\mathrm{LSI}}$ and for any $T > 0$,

$$D_{\mathrm{KL}}(q_T \| \tilde{q}) \leq C\Big( \sup_{t \in [0,T]} \Delta(t) + e^{-cT} D_{\mathrm{KL}}(q_0 \| \tilde{q}) \Big). \tag{170}$$

We now specialize (170) to the initialization $q_0 = \tilde{q}$, and bound $\Delta(t)$. We have

$$\Delta(t) \leq \left\langle \left\| \frac{1}{s^2} \mathbf{X}^\top (\mathbf{X}\boldsymbol{\theta}^* + s\mathbf{z} - \mathbf{X}\boldsymbol{\theta}^t) - \frac{1}{\tilde{s}^2} \mathbf{X}^\top (\mathbf{X}\boldsymbol{\theta}^* + \tilde{s}\mathbf{z} - \mathbf{X}\boldsymbol{\theta}^t) \right\|_2^2 + \sum_{j=1}^d \Big( \partial_\theta \log g(\theta_j^t, \alpha) - \partial_\theta \log g(\theta_j^t, \tilde{\alpha}) \Big)^2 \right\rangle.$$

Let $C, C', C'' > 0$ be constants depending on the compact sets $S, I$ of the lemma statement and changing from instance to instance. For $\alpha, \tilde{\alpha} \in S$ and $s^2, \tilde{s}^2 \in I$,

$$|s^{-2} - \tilde{s}^{-2}| \leq C|s^2 - \tilde{s}^2|, \quad |s^{-1} - \tilde{s}^{-1}| \leq C|s^2 - \tilde{s}^2|, \quad |\partial_\theta \log g(\theta; \alpha) - \partial_\theta \log g(\theta; \tilde{\alpha})| \leq C\|\alpha - \tilde{\alpha}\|_2,$$

the last inequality holding by Assumption 2.2(b). Thus

$$\Delta(t) \leq C\Big[ \|\mathbf{X}\|_{\mathrm{op}}^4 (\|\boldsymbol{\theta}^*\|_2^2 + \langle \|\boldsymbol{\theta}^t\|_2^2 \rangle) + \|\mathbf{X}\|_{\mathrm{op}}^2 \|\mathbf{z}\|_2^2 \Big] (s^2 - \tilde{s}^2)^2 + Cd\|\alpha - \tilde{\alpha}\|_2^2.$$

On the event $\mathcal{E}'(C_0, C_{\mathrm{LSI}})$, we have $\langle \|\boldsymbol{\theta}^t\|_2^2 \rangle \leq C(\langle \|\boldsymbol{\theta}^0\|_2^2 \rangle + d)$ by (121), and $\langle \|\boldsymbol{\theta}^0\|_2^2 \rangle \leq Cd$ under the initialization $q_0 = \tilde{q}$ which holds also by (121). Applying these bounds together with $\|\mathbf{X}\|_{\mathrm{op}} \leq C$, $\|\boldsymbol{\theta}^*\|_2^2 \leq Cd$, and $\|\mathbf{z}\|_2^2 \leq Cd$ by definition of $\mathcal{E}'(C_0, C_{\mathrm{LSI}})$, we have

$$\sup_{t \geq 0} \Delta(t) \leq C'd(s^2 - \tilde{s}^2)^2 + C'd\|\alpha - \tilde{\alpha}\|_2^2.$$

Applying this and $q_0 = \tilde{q}$ to (170), we have on the event $\mathcal{E}'(C_0, C_{\mathrm{LSI}})$ that

$$\sup_{t \geq 0} D_{\mathrm{KL}}(q_t \| \tilde{q}) \leq Cd(s^2 - \tilde{s}^2)^2 + Cd\|\alpha - \tilde{\alpha}\|_2^2.$$

By lower-semicontinuity of KL-divergence and the $T_2$-transportation inequality for $\tilde{q}$ implied by the LSI (c.f. [83, Theorem 9.6.1]),

$$W_2(q, \tilde{q})^2 \leq C\, D_{\mathrm{KL}}(q \| \tilde{q}) \leq C \liminf_{t \to \infty} D_{\mathrm{KL}}(q_t \| \tilde{q}) \leq C'd(s^2 - \tilde{s}^2)^2 + C'd\|\alpha - \tilde{\alpha}\|_2^2. \tag{171}$$

This gives our desired bound on the Wasserstein-2 distance between $q$ and $\tilde{q}$.

Now let $\langle f(\boldsymbol{\theta}) \rangle_q$ and $\langle f(\boldsymbol{\theta}) \rangle_{\tilde{q}}$ be the posterior expectations under $q$ (given $\mathbf{y}$) and $\tilde{q}$ (given $\tilde{\mathbf{y}}$). Then by Jensen's inequality,

$$\|\langle \boldsymbol{\theta} \rangle_q - \langle \boldsymbol{\theta} \rangle_{\tilde{q}}\|_2 \leq W_2(q, \tilde{q}).$$

Applying $|\|\mathbf{x}\|_2^2 - \|\mathbf{y}\|_2^2| \leq \|\mathbf{x} - \mathbf{y}\|_2 \cdot \|\mathbf{x} + \mathbf{y}\|_2$ and Cauchy-Schwarz, also

$$|\langle \|\boldsymbol{\theta}\|_2^2 \rangle_q - \langle \|\boldsymbol{\theta}\|_2^2 \rangle_{\tilde{q}}| \leq W_2(q, \tilde{q}) \cdot \sqrt{2\langle \|\boldsymbol{\theta}\|_2^2 \rangle_q + 2\langle \|\boldsymbol{\theta}\|_2^2 \rangle_{\tilde{q}}} \leq C\sqrt{d}\, W_2(q, \tilde{q})$$

where the last inequality applies $\langle \|\boldsymbol{\theta}\|_2 \rangle_q \leq \langle \|\boldsymbol{\theta}\|_2^2 \rangle_q^{1/2} \leq C\sqrt{d}$ on $\mathcal{E}(C_0, C_{\mathrm{LSI}})$ by (121), and similarly for $\tilde{q}$. Then, denoting by $\mathrm{MSE}(s^2, \alpha)$ and $\mathrm{MSE}(\tilde{s}^2, \tilde{\alpha})$ the values of MSE as defined in Corollary 2.10 under $q$ and $\tilde{q}$, we have

$$\begin{aligned}
|\mathrm{MSE}(s^2, \alpha) - \mathrm{MSE}(\tilde{s}^2, \tilde{\alpha})| &= \big| d^{-1} \langle \|\boldsymbol{\theta} - \langle \boldsymbol{\theta} \rangle_q\|_2^2 \rangle_q - d^{-1} \langle \|\boldsymbol{\theta} - \langle \boldsymbol{\theta} \rangle_{\tilde{q}}\|_2^2 \rangle_{\tilde{q}} \big| \\
&\leq d^{-1} \big| \langle \|\boldsymbol{\theta}\|_2^2 \rangle_q - \langle \|\boldsymbol{\theta}\|_2^2 \rangle_{\tilde{q}} \big| + d^{-1} \big| \|\langle \boldsymbol{\theta} \rangle_q\|_2^2 - \|\langle \boldsymbol{\theta} \rangle_{\tilde{q}}\|_2^2 \big| \\
&\leq \frac{C'W_2(q, \tilde{q})}{\sqrt{d}} \leq C''|s^2 - \tilde{s}^2| + C''\|\alpha - \tilde{\alpha}\|_2.
\end{aligned}$$

Similarly

$$\begin{aligned}
|\mathrm{MSE}_*(s^2, \alpha) - \mathrm{MSE}_*(\tilde{s}^2, \tilde{\alpha})| &= \big| d^{-1} \|\boldsymbol{\theta}^* - \langle \boldsymbol{\theta} \rangle_q\|_2^2 - d^{-1} \|\boldsymbol{\theta}^* - \langle \boldsymbol{\theta} \rangle_{\tilde{q}}\|_2^2 \big| \\
&\leq 2d^{-1} \|\boldsymbol{\theta}^*\|_2 \|\langle \boldsymbol{\theta} \rangle_q - \langle \boldsymbol{\theta} \rangle_{\tilde{q}}\|_2 + d^{-1} \big| \|\langle \boldsymbol{\theta} \rangle_q\|_2^2 - \|\langle \boldsymbol{\theta} \rangle_{\tilde{q}}\|_2^2 \big| \\
&\leq \frac{C'W_2(q, \tilde{q})}{\sqrt{d}} \leq C''|s^2 - \tilde{s}^2| + C''\|\alpha - \tilde{\alpha}\|_2.
\end{aligned}$$

Since $\mathcal{E}'(C_0, C_{\mathrm{LSI}})$ holds a.s. for all large $n, d$, and Corollary 2.10(a) already proven shows $\lim_{n,d\to\infty} \mathrm{MSE} = \mathrm{mse}$ and $\lim_{n,d\to\infty} \mathrm{MSE}_* = \mathrm{mse}_*$ a.s. at both $(s^2, \alpha)$ and $(\tilde{s}^2, \tilde{\alpha})$, this implies

$$|\mathrm{mse}(s^2, \alpha) - \mathrm{mse}(\tilde{s}^2, \tilde{\alpha})|, |\mathrm{mse}_*(s^2, \alpha) - \mathrm{mse}_*(\tilde{s}^2, \tilde{\alpha})| \le C|s^2 - \tilde{s}^2| + C\|\alpha - \tilde{\alpha}\|_2,$$

so $\mathrm{mse}(s^2, \alpha)$ and $\mathrm{mse}_*(s^2, \alpha)$ are locally Lipschitz as desired. $\qquad\square$

*Proof of Corollary 2.10(b).* We apply Corollary 2.10(a) and an I-MMSE relation for mismatched Gaussian channels. Write $\mathbb{E}[\cdot \mid \mathbf{X}]$ for the expectation over $(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon})$ conditional on $\mathbf{X}$ as in Proposition 4.11. Let

$$I(\mathbf{y}, \boldsymbol{\theta}^*) = \mathbb{E}\left[\log \frac{\mathsf{P}(\mathbf{y} \mid \boldsymbol{\theta}^*, \mathbf{X})}{\mathsf{P}_{g_*}(\mathbf{y} \mid \mathbf{X})} \,\middle|\, \mathbf{X}\right] = -\mathbb{E}[\log \mathsf{P}_{g_*}(\mathbf{y} \mid \mathbf{X}) \mid \mathbf{X}] - \frac{n}{2}(1 + \log 2\pi\sigma^2)$$

be the signal-observation mutual information in the linear model (4) conditional on $\mathbf{X}$, where $\mathsf{P}(\mathbf{y} \mid \boldsymbol{\theta}^*, \mathbf{X})$ is the Gaussian likelihood of $\mathbf{y}$ and $\mathsf{P}_{g_*}(\mathbf{y} \mid \mathbf{X})$ is the marginal likelihood (6) under the true prior $g_*$. Then

$$\begin{aligned}
\mathbb{E}[\log \mathsf{P}_g(\mathbf{y} \mid \mathbf{X}) \mid \mathbf{X}] &= -\mathrm{D}_{\mathrm{KL}}(\mathsf{P}_{g_*}(\mathbf{y} \mid \mathbf{X})\|\mathsf{P}_g(\mathbf{y} \mid \mathbf{X})) + \mathbb{E}[\log \mathsf{P}_{g_*}(\mathbf{y} \mid \mathbf{X}) \mid \mathbf{X}] \\
&= -\mathrm{D}_{\mathrm{KL}}(\mathsf{P}_{g_*}(\mathbf{y} \mid \mathbf{X})\|\mathsf{P}_g(\mathbf{y} \mid \mathbf{X})) - I(\mathbf{y}, \boldsymbol{\theta}^*) - \frac{n}{2}(1 + \log 2\pi\sigma^2) \qquad (172)
\end{aligned}$$

where here and throughout the proof, $\mathrm{D}_{\mathrm{KL}}(\cdot)$ denotes the KL-divergence also conditional on $\mathbf{X}$.

Let us denote the inverse noise variance by $s^{-1} = \sigma^2$ and write

$$E(s, g) = \mathbb{E}[\mathrm{YMSE}_* \mid \mathbf{X}] = n^{-1}\mathbb{E}[\|\mathbf{X}\langle\boldsymbol{\theta}\rangle - \mathbf{X}\boldsymbol{\theta}^*\|^2 \mid \mathbf{X}] \qquad (173)$$

for the expected $\mathrm{YMSE}_*$ in the linear model (4) with assumed prior $g$ and noise variance $s^{-1}$. We clarify that this means $\langle\cdot\rangle$ in (173) is the posterior average under the law

$$\mathsf{P}_g(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y}) \propto \exp\left(-\frac{s}{2}\|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 + \sum_{j=1}^d \log g(\theta_j)\right)$$

and $\mathbb{E}[\cdot \mid \mathbf{X}]$ is the expectation over $(\boldsymbol{\theta}^*, \boldsymbol{\varepsilon})$ where $\boldsymbol{\varepsilon}$ also has variance $s^{-1}$. We write also $I[s], \mathrm{D}_{\mathrm{KL}}[s]$ for the above quantities $I(\mathbf{y}, \boldsymbol{\theta}^*)$ and $\mathrm{D}_{\mathrm{KL}}(\mathsf{P}_{g_*}(\mathbf{y} \mid \mathbf{X})\|\mathsf{P}_g(\mathbf{y} \mid \mathbf{X}))$ in this model with noise variance $s^{-1}$. Then [87, Theorem 2] and [88, Eq. (24)] show the I-MMSE relations

$$\frac{\mathrm{d}}{\mathrm{d}s}I[s] = \frac{n}{2}E(s, g_*), \qquad \frac{\mathrm{d}}{\mathrm{d}s}\mathrm{D}_{\mathrm{KL}}[s] = \frac{n}{2}\big(E(s, g) - E(s, g_*)\big).$$

For any fixed $n, d$ and $\mathbf{X}$, in the limit $s \to 0$, it is direct to check that $I[s] \to 0$ and $\mathrm{D}_{\mathrm{KL}}[s] \to 0$. Thus, for $I(\mathbf{y}, \boldsymbol{\theta}^*) \equiv I[\sigma^{-2}]$ and $\mathrm{D}_{\mathrm{KL}}(\mathsf{P}_{g_*}(\mathbf{y} \mid \mathbf{X})\|\mathsf{P}_g(\mathbf{y} \mid \mathbf{X})) \equiv \mathrm{D}_{\mathrm{KL}}[\sigma^{-2}]$ in the original model with noise variance $\sigma^2$, integrating these I-MMSE relations shows

$$\mathrm{D}_{\mathrm{KL}}(\mathsf{P}_{g_*}(\mathbf{y} \mid \mathbf{X})\|\mathsf{P}_g(\mathbf{y} \mid \mathbf{X})) + I(\mathbf{y}, \boldsymbol{\theta}^*) = \frac{n}{2}\int_0^{\sigma^{-2}} E(s, g)\mathrm{d}s. \qquad (174)$$

Assumption 2.7(b) ensures that the posterior LSI (46) holds a.s. in the model with any noise variance $s^{-1} \in [\sigma^2, \infty)$. Then applying Corollary 2.10(a) already shown and the concentration of $\mathrm{YMSE}_*$ in Proposition 4.11, we have $E(s, g) \to \mathrm{ymse}_*(s, g)$ a.s. for each $s^{-1} \in (\sigma^2, \infty)$, where $\mathrm{ymse}_*(s, g)$ is defined by (42) via the DMFT limit of the Langevin dynamics (7) with fixed prior $g(\cdot)$ in the linear model with noise variance $s^{-1}$. To apply dominated convergence, we note that on the event $\|\mathbf{X}\|_{\mathrm{op}} \le C_0$, by the extension (110) of (106), we have $E(s, g) \le C$ for a constant $C > 0$ uniformly over all $s \in [0, \sigma^{-2}]$ and all $n, d$. Then, since $\|\mathbf{X}\|_{\mathrm{op}} \le C_0$ holds a.s. for all large $n, d$, taking the limit $n, d \to \infty$ and applying the bounded convergence theorem to (174) shows that almost surely,

$$\lim_{n,d\to\infty} \frac{1}{d}\Big(\mathrm{D}_{\mathrm{KL}}(\mathsf{P}_{g_*}(\mathbf{y} \mid \mathbf{X})\|\mathsf{P}_g(\mathbf{y} \mid \mathbf{X})) + I(\mathbf{y}, \boldsymbol{\theta}^*)\Big) = \frac{\delta}{2}\int_0^{\sigma^{-2}} \mathrm{ymse}_*(s, g)\mathrm{d}s. \qquad (175)$$

Let us now fix the assumed prior $g(\cdot)$, write $\mathrm{ymse}_*(s) \equiv \mathrm{ymse}_*(s, g)$, and let $(\mathrm{mse}(s), \mathrm{mse}_*(s), \omega(s), \omega_*(s))$ denote the fixed points (43) corresponding to $\mathrm{ymse}_*(s)$. Recall the marginal density $\mathsf{P}_{g,\omega}(y)$ of the scalar channel model (39), and define

$$f(\omega, \omega_*, s) = -\mathbb{E}_{g_*,\omega_*} \log \mathsf{P}_{g,\omega}(y) - \frac{1}{2}\left(2\delta + \log \frac{2\pi}{\omega} - \delta \log \frac{\delta s}{\omega} + (1-\delta)\frac{\omega}{\omega_*} + \frac{\omega}{s}\left(\frac{\omega}{\omega_*} - 2\right)\right) \tag{176}$$

$$= \underbrace{\frac{\omega}{2}\mathbb{E}\,\theta^{*2} - \mathbb{E}\log\int \exp\left(\omega\theta(\theta^* + \omega_*^{-1/2}z) - \frac{\omega}{2}\theta^2\right)g(\theta)\mathrm{d}\theta}_{:=\mathrm{I}} - \underbrace{\frac{1}{2}\left(2\delta - \delta\log\frac{\delta s}{\omega} - \frac{\delta\omega}{\omega_*} + \frac{\omega}{s}\left(\frac{\omega}{\omega_*} - 2\right)\right)}_{:=\mathrm{II}}.$$

Here, the expectations in the second line are over $\theta^* \sim g_*$ and $z \sim \mathcal{N}(0,1)$, and we have applied the explicit form of $\mathsf{P}_{g,\omega}(y)$ and evaluated $\mathbb{E}_{g_*,\omega_*}$ under the true model $y = \theta^* + \omega_*^{-1/2}z$ with some some algebraic simplification. We now claim that

$$\frac{\delta}{2}\int_0^s \mathrm{ymse}_*(t)\mathrm{d}t = f(\omega(s), \omega_*(s), s) \tag{177}$$

for all $s \in (0, \sigma^{-2})$. To show this, it suffices to check $\lim_{s\to 0} f(\omega(s), \omega_*(s), s) = 0$ and $\frac{\mathrm{d}}{\mathrm{d}s}f(\omega(s), \omega_*(s), s) = \frac{\delta}{2}\mathrm{ymse}_*(s)$, which we may do as follows:

- Let $\mathrm{MSE}(s), \mathrm{MSE}_*(s)$ denote the values of $\mathrm{MSE}, \mathrm{MSE}_*$ in a linear model with noise variance $s^{-1}$. On the event $\|\mathbf{X}\|_{\mathrm{op}} \leq C_0$, the bound (110) implies that $\mathrm{MSE}(s), \mathrm{MSE}_*(s) \leq C(1+s\|\mathbf{y}\|_2^2/d)$ for a constant $C > 0$ (independent of $s$) and for all $s^{-1} \in (\sigma^2, \infty)$. Taking the almost sure limit as $n, d \to \infty$ shows that $\mathrm{mse}(s), \mathrm{mse}_*(s) \leq C$. In particular, in the limit $s \to 0$, we have that $\mathrm{mse}(s), \mathrm{mse}_*(s)$ remain bounded, so $\omega(s), \omega_*(s) \sim \delta s$ by the fixed point relation (43). Then $\omega(s) \to 0$, $\omega_*(s) \to 0$, $\omega(s)/s \to \delta$, and $\omega(s)/\omega_*(s) \to 1$ as $s \to 0$. Applying this to (176) shows

$$\lim_{s\to 0} f(\omega(s), \omega_*(s), s) = 0.$$

- Differentiating the term I of (176) in $\omega, \omega_*$ and applying Gaussian integration-by-parts with respect to $z \sim \mathcal{N}(0,1)$, we may check that

$$\partial_\omega \mathrm{I} = \frac{1}{2}\mathbb{E}\langle(\theta^* - \theta)^2\rangle_{g,\omega} - \frac{\omega}{\omega_*}\mathbb{E}\langle(\theta - \langle\theta\rangle_{g,\omega})^2\rangle_{g,\omega},$$

$$\partial_{\omega_*}\mathrm{I} = \frac{\omega^2}{2\omega_*^2}\mathbb{E}\langle(\theta - \langle\theta\rangle_{g,\omega})^2\rangle_{g,\omega}.$$

Then at the fixed points $(\omega, \omega_*) = (\omega(s), \omega_*(s))$, we have

$$\partial_\omega \mathrm{I}|_{(\omega,\omega_*)=(\omega(s),\omega_*(s))} = \frac{1}{2}(\mathrm{mse}(s) + \mathrm{mse}_*(s)) - \frac{\omega(s)}{\omega_*(s)}\mathrm{mse}(s)$$

$$\partial_{\omega_*}\mathrm{I}|_{(\omega,\omega_*)=(\omega(s),\omega_*(s))} = \frac{\omega(s)^2}{2\omega_*(s)^2}\mathrm{mse}(s).$$

Applying $\mathrm{mse}(s) = \delta/\omega(s) - \sigma^2$ and $\mathrm{mse}_*(s) = \delta/\omega_*(s) - \sigma^2$ by (43) and comparing with the derivatives of the second term II of (176), this verifies

$$\partial_\omega f(\omega(s), \omega_*(s), s) = 0, \qquad \partial_{\omega_*}f(\omega(s), \omega_*(s), s) = 0. \tag{178}$$

Furthermore, direct calculation shows that at $(\omega, \omega_*) = (\omega(s), \omega_*(s))$,

$$\partial_s f(\omega(s), \omega_*(s), s) = \frac{\delta\sigma^2}{2} + \frac{\omega(s)\sigma^4}{2}\left(\frac{\omega(s)}{\omega_*(s)} - 2\right) = \frac{\delta}{2}\mathrm{ymse}_*(s),$$

the second equality using (44). Lemma 4.12 implies that $\mathrm{mse}(s), \mathrm{mse}_*(s), \omega(s), \omega_*(s)$ are locally Lipschitz, and hence absolutely continuous, over $s \in (0, \sigma^{-2})$. Then also $s \mapsto f(\omega(s), \omega_*(s), s)$ is absolutely

54

continuous, and we may differentiate by the chain rule to get

$$\frac{\mathrm{d}}{\mathrm{d}s} f(\omega(s), \omega_*(s), s) = \partial_\omega f(\omega(s), \omega_*(s), s) \cdot \omega'(s) + \partial_{\omega_*} f(\omega(s), \omega_*(s), s) \cdot \omega'_*(s) + \partial_s f(\omega(s), \omega_*(s), s)$$

$$= \partial_s f(\omega(s), \omega_*(s), s) = \frac{\delta}{2} \, \mathrm{ymse}_*(s).$$

Combining the above arguments verifies the claim (177).

Applying (175) and (177) to (172) and writing $(\omega, \omega_*) = (\omega(\sigma^{-2}), \omega_*(\sigma^{-2}))$ for the fixed points at the original noise variance $\sigma^2$, this shows

$$\lim_{n,d\to\infty} d^{-1} \mathbb{E}[\log \mathsf{P}_g(\mathbf{y} \mid \mathbf{X}) \mid \mathbf{X}] = -f(\omega, \omega_*, \sigma^{-2}) - \frac{\delta}{2}(1 + \log 2\pi\sigma^2).$$

Applying concentration of $d^{-1} \log \mathsf{P}_g(\mathbf{y} \mid \mathbf{X})$ with respect to $\mathbb{E}[\cdot \mid \mathbf{X}]$ which is established in Propostion 4.11, and substituting the form of $f$ in (176), this shows Corollary 2.10(b). $\qquad\square$

# 5 Analysis of empirical Bayes Langevin dynamics

In this section, we prove Theorem 2.13 on the adaptive empirical Bayes dynamics with time-varying prior parameter $\widehat{\alpha}^t$, and discuss further the examples of Section 2.4.2.

## 5.1 General analysis under uniform LSI

We introduce a few notational shorthands: Conditional on $\mathbf{X}, \boldsymbol{\theta}^*, \boldsymbol{\varepsilon}$, let

$$q_\alpha(\boldsymbol{\theta}) \equiv \mathsf{P}_{g(\cdot,\alpha)}(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y})$$

be the posterior law under the prior parameter $\alpha$. We write $\langle \cdot \rangle_\alpha$ for its posterior expectation. For $\boldsymbol{\theta} \in \mathbb{R}^d$, define

$$\bar{\mathsf{P}}_{\boldsymbol{\theta}} = \frac{1}{d} \sum_{j=1}^d \delta_{(\theta_j^*, \theta_j)}, \qquad \bar{\mathsf{P}}_\alpha = \langle \bar{\mathsf{P}}_{\boldsymbol{\theta}} \rangle_\alpha. \tag{179}$$

Thus $\bar{\mathsf{P}}_\alpha$ is a $(\mathbf{X}, \boldsymbol{\theta}^*, \boldsymbol{\varepsilon})$-dependent joint law over variables $(\theta^*, \theta)$ which satisfies

$$\mathbb{E}_{(\theta^*,\theta)\sim\bar{\mathsf{P}}_\alpha} f(\theta^*, \theta) = \frac{1}{d} \sum_{j=1}^d \langle f(\theta_j^*, \theta_j) \rangle_\alpha = \frac{1}{d} \sum_{j=1}^d \int f(\theta_j^*, \theta_j) q_\alpha(\boldsymbol{\theta}) \mathrm{d}\boldsymbol{\theta}. \tag{180}$$

We write $\theta \sim \bar{\mathsf{P}}_\alpha$ as shorthand for the $\theta$-marginal of $(\theta^*, \theta) \sim \bar{\mathsf{P}}_\alpha$.

We note that under Assumptions 2.2(b) and 2.11, all constants in (105) are uniform over $g \in \{g(\cdot, \alpha) : \alpha \in O\}$ for the bounded domain $O$ of Assumption 2.11, where a uniform bound for $|\log g(0, \alpha)|$ follows from $|\log g(0, \alpha)| \le |\log g(0, 0)| + \|\nabla_\alpha(\log g(0, 0))\|_2 \cdot \|\alpha\|_2 + C\|\alpha\|_2^2$ as implied by (19) of Assumption 2.2(b). Hence the bounds of Section 4.2 hold uniformly over $\alpha \in O$. In particular, from (106),

$$\sup_{\alpha \in O} \langle \|\boldsymbol{\theta}\|_2^2 \rangle_\alpha \le C(d + \|\mathbf{y}\|_2^2) \tag{181}$$

on an event $\{\|\mathbf{X}\|_{\mathrm{op}} \le C_0\}$ that holds a.s. for all large $n, d$.

We first prove Lemma 2.12 on the derivatives of $F, \widehat{F}$ and uniform convergence of $\widehat{F}, \nabla\widehat{F}$ over $S \subset O$.

*Proof of Lemma 2.12.* For (a), differentiating

$$\widehat{F}(\alpha) = -\frac{1}{d} \log \int \left(\frac{1}{2\pi\sigma^2}\right)^{n/2} \exp\left(-\frac{1}{2\sigma^2}\|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 + \sum_{j=1}^d \log g(\theta_j, \alpha)\right) \mathrm{d}\boldsymbol{\theta}$$

and applying the property (180), we have

$$\nabla \widehat{F}(\alpha) = -\frac{1}{d} \sum_{j=1}^{d} \langle \nabla_\alpha \log g(\theta_j, \alpha) \rangle_\alpha = -\mathbb{E}_{\theta \sim \bar{\mathsf{P}}_\alpha} \nabla_\alpha \log g(\theta, \alpha). \tag{182}$$

For the form of $\nabla F(\alpha)$, define analogously to (176)

$$f(\omega, \omega_*, \alpha) = -\mathbb{E}_{g_*, \omega_*} \log \mathsf{P}_{g(\cdot, \alpha), \omega}(y) - \frac{1}{2} \left( 2\delta + \log \frac{2\pi}{\omega} - \delta \log \frac{\delta}{\omega \sigma^2} + (1-\delta)\frac{\omega}{\omega_*} + \omega \sigma^2 \left( \frac{\omega}{\omega_*} - 2 \right) \right) \tag{183}$$

where the dependence on $\alpha$ is in $\mathsf{P}_{g(\cdot, \alpha), \omega}(y)$. For any $\alpha \in O$, let $\omega(\alpha), \omega_*(\alpha)$ be the fixed points $\omega, \omega_*$ defined by (42) via the DMFT system for the dynamics (7) with fixed prior $g \equiv g(\cdot, \alpha)$. (This DMFT system is approximately-TTI for each $\alpha \in O$ by Assumption 2.11 and Theorem 2.9, hence $\omega(\alpha), \omega_*(\alpha)$ are well-defined.) Then

$$F(\alpha) = f(\omega(\alpha), \omega_*(\alpha), \alpha) + \frac{\delta}{2}(1 + \log 2\pi\sigma^2). \tag{184}$$

By the same calculations as (178), at the fixed points $(\omega(\alpha), \omega_*(\alpha))$, we have $\partial_\omega f(\omega(\alpha), \omega_*(\alpha), \alpha) = 0$ and $\partial_{\omega_*} f(\omega(\alpha), \omega_*(\alpha), \alpha) = 0$. By Lemma 4.12, $\omega(\alpha), \omega_*(\alpha)$ are locally Lipschitz and hence absolutely continuous over $\alpha \in O$. Then $F(\alpha)$ is also absolutely continuous over $\alpha \in O$, and differentiating by the chain rule gives

$$\nabla F(\alpha) = \nabla_\alpha f(\omega, \omega_*, \alpha)\Big|_{(\omega, \omega_*)=(\omega(\alpha), \omega_*(\alpha))}$$

$$= -\nabla_\alpha \left[ \mathbb{E}_{g_*, \omega_*} \log \mathsf{P}_{g(\cdot, \alpha), \omega}(y) \right]\Big|_{(\omega, \omega_*)=(\omega(\alpha), \omega_*(\alpha))}$$

$$= -\nabla_\alpha \left[ \mathbb{E}_{g_*, \omega_*} \log \int \left( \frac{\omega}{2\pi} \right)^{1/2} \exp\left( -\frac{\omega}{2}(y-\theta)^2 + \log g(\theta, \alpha) \right) d\theta \right]\Big|_{(\omega, \omega_*)=(\omega(\alpha), \omega_*(\alpha))}$$

By definition $\mathsf{P}_\alpha$ is the joint law of $(\theta^*, \theta)$ under the generative process where $(\theta^*, y)$ are drawn from the Gaussian convolution model defining this expectation $\mathbb{E}_{g_*, \omega_*}$, and where $\theta \sim \mathsf{P}_{g(\cdot, \alpha), \omega}(\theta \mid y)$. Hence, evaluating $\nabla_\alpha$ above gives

$$\nabla F(\alpha) = -\mathbb{E}_{\theta \sim \mathsf{P}_\alpha} \nabla_\alpha \log g(\theta, \alpha).$$

For (b), let $S \subset O$ be any compact subset of the domain $O$ in Assumption 2.11, and let $Q$ be a countable dense subset of $O$. Define

$$\mathcal{E}(C_0, C_{\mathrm{LSI}}) = \{ \|\mathbf{X}\|_{\mathrm{op}} \leq C_0, \ (46) \text{ holds for } q_\alpha(\boldsymbol{\theta}) \equiv \mathsf{P}_{g(\cdot, \alpha)}(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y}) \text{ for every } \alpha \in O \}.$$

Assumptions 2.1 and 2.11 ensure for some $C_0, C_{\mathrm{LSI}} > 0$ that $\mathcal{E}(C_0, C_{\mathrm{LSI}})$ holds a.s. for all large $n, d$, where this event depends only on $\mathbf{X}$ and not on $\boldsymbol{\theta}^*, \boldsymbol{\varepsilon}$.

We restrict to the almost-sure event where the convergence statements of Corollary 2.10 and Proposition 4.11 hold for every $\alpha \in Q$, and where $\mathcal{E}(C_0, C_{\mathrm{LSI}})$ holds for all large $n, d$. Note that Corollary 2.10 shows $\widehat{F}(\alpha) \to F(\alpha)$ for each $\alpha \in Q$. To strengthen this to uniform convergence over $S$, note that Assumption 2.2(b) implies $\partial_\theta \nabla_\alpha \log g(\theta, \alpha)$ is uniformly bounded over $(\theta, \alpha) \in \mathbb{R} \times S$, so

$$\|\nabla_\alpha \log g(\theta, \alpha)\|_2 \leq \|\nabla_\alpha \log g(0, \alpha)\|_2 + C|\theta|.$$

Then, since $\nabla_\alpha \log g(0, \alpha)$ is bounded over $\alpha \in S$ by compactness of $S$, and $\sup_{\alpha \in S} \langle \|\boldsymbol{\theta}\|_2^2 \rangle_\alpha \leq Cd$ by (181), we have

$$\sup_{\alpha \in S} \|\mathbb{E}_{\theta \sim \bar{\mathsf{P}}_\alpha} \nabla_\alpha \log g(\theta, \alpha)\|_2 \leq \sup_{\alpha \in S} \frac{1}{d} \sum_{j=1}^{d} \langle \|\nabla_\alpha \log g(\theta_j, \alpha)\|_2 \rangle_\alpha$$

$$\leq \sup_{\alpha \in S} \|\nabla_\alpha \log g(0, \alpha)\|_2 + \frac{C}{d} \sum_{j=1}^{d} \langle |\theta_j| \rangle_\alpha \leq C'. \tag{185}$$

56

This shows $\nabla \widehat{F}(\alpha)$ is bounded over any compact subset $S \subset O$. Then for any compact $S \subset O$, the functions $\widehat{F}(\alpha)$ for all $n, d$ are equicontinuous in a neighborhood of each point $\alpha \in S$, and hence are uniformly equicontinuous over $S$ since a finite number of such neighborhoods cover $S$. Then by Arzela-Ascoli, the convergence $\widehat{F}(\alpha) \to F(\alpha)$ for each $\alpha \in Q$ implies uniform convergence over $\alpha \in S$.

We next show the pointwise convergence $\nabla \widehat{F}(\alpha) \to \nabla F(\alpha)$ for each $\alpha \in Q$. Recalling our definition of $\bar{\mathsf{P}}_{\boldsymbol{\theta}}$ in (179), and applying Jensen's inequality and the convexity $W_2(\lambda \mathsf{P} + (1-\lambda)\mathsf{P}', \mathsf{Q})^2 \leq \lambda W_2(\mathsf{P}, \mathsf{Q})^2 + (1-\lambda)W_2(\mathsf{P}', \mathsf{Q})^2$ of the squared Wasserstein-2 distance,

$$W_2(\bar{\mathsf{P}}_\alpha, \mathsf{P}_\alpha)^2 \leq \langle W_2(\bar{\mathsf{P}}_{\boldsymbol{\theta}}, \mathsf{P}_\alpha)^2 \rangle_\alpha.$$

For each $\alpha \in Q$, the right side converges to 0 as $n, d \to \infty$ by the statement (48) of Corollary 2.10(a). Thus $\lim_{n,d\to\infty} W_2(\bar{\mathsf{P}}_\alpha, \mathsf{P}_\alpha) = 0$. Assumption 2.2(b) ensures that $\nabla_\alpha \log g(\theta, \alpha)$ is Lipschitz in $\theta$, so this Wasserstein-2 convergence implies

$$\lim_{n,d\to\infty} \nabla \widehat{F}(\alpha) = \lim_{n,d\to\infty} \mathbb{E}_{\theta \sim \bar{\mathsf{P}}_\alpha} \nabla_\alpha \log g(\theta, \alpha) = \mathbb{E}_{\theta \sim \mathsf{P}_\alpha} \nabla_\alpha \log g(\theta, \alpha) = \nabla F(\alpha)$$

for each $\alpha \in Q$, as claimed.

To extend this to uniform convergence over any compact subset $S \subset O$, we differentiate (182) a second time. Writing $\mathrm{Var}_\alpha, \mathrm{Cov}_\alpha$ for the variance and covariance under $\langle \cdot \rangle_\alpha$,

$$\nabla^2 \widehat{F}(\alpha) = -\frac{1}{d}\left\langle \sum_{j=1}^d \nabla_\alpha^2 \log g(\theta_j, \alpha) \right\rangle_\alpha - \frac{1}{d} \mathrm{Cov}_\alpha\left[ \sum_{j=1}^d \nabla_\alpha \log g(\theta_j, \alpha) \right]. \tag{186}$$

The first term is uniformly bounded over $\alpha \in S$, by the same argument as showing boundedness of $\nabla \widehat{F}(\alpha)$ above. For the second term, on the event $\mathcal{E}(C_0, C_{\mathrm{LSI}})$, for every unit vector $v \in \mathbb{R}^K$ and $\alpha \in S$,

$$\mathrm{Var}_\alpha\left[ \sum_{j=1}^d v^\top \nabla_\alpha \log g(\theta_j, \alpha) \right] \leq (C_{\mathrm{LSI}}/2)\left\langle \sum_{j=1}^d \left( v^\top \partial_\theta \nabla_\alpha \log g(\theta_j, \alpha) \right)^2 \right\rangle_\alpha$$

by the Poincaré inequality for $q_\alpha$ implied by its LSI. Since $\partial_\theta \nabla_\alpha \log g(\theta, \alpha)$ is bounded over $\alpha \in S$, the second term of (186) is also bounded on $\mathcal{E}(C_0, C_{\mathrm{LSI}})$. Thus $\nabla^2 \widehat{F}(\alpha)$ is uniformly bounded over $\alpha \in S$ for all large $n, d$. This implies as above that for any compact $S \subset O$, the functions $\nabla \widehat{F}(\alpha)$ for all large $n, d$ are uniformly equicontinuous on $S$, so $\nabla \widehat{F}(\alpha) \to \nabla F(\alpha)$ uniformly over $\alpha \in S$. This shows part (b).

For part (c), note that if $g^* = g(\cdot, \alpha^*)$, then

$$\mathbb{E}[\widehat{F}(\alpha) \mid \mathbf{X}] - \mathbb{E}[\widehat{F}(\alpha^*) \mid \mathbf{X}] = d^{-1} \mathrm{D}_{\mathrm{KL}}(\mathsf{P}_{g(\cdot, \alpha^*)}(\mathbf{y} \mid \mathbf{X}) \| \mathsf{P}_{g(\cdot, \alpha)}(\mathbf{y} \mid \mathbf{X})) \geq 0,$$

where here $\mathrm{D}_{\mathrm{KL}}(\cdot)$ is the KL-divergence conditional on $\mathbf{X}$. Thus $\alpha^*$ is a minimizer of $\alpha \mapsto \mathbb{E}[\widehat{F}(\alpha) \mid \mathbf{X}]$ over $\mathbb{R}^K$. Applying the convergence $\widehat{F}(\alpha) - \mathbb{E}[\widehat{F}(\alpha) \mid \mathbf{X}] \to 0$ for each $\alpha \in Q$ from Proposition 4.11, we have also $\mathbb{E}[\widehat{F}(\alpha) \mid \mathbf{X}] \to F(\alpha)$ for each $\alpha \in Q$. Note that

$$\nabla_\alpha \mathbb{E}[\widehat{F}(\alpha) \mid \mathbf{X}] = -\mathbb{E}[\mathbb{E}_{\theta \sim \bar{\mathsf{P}}_\alpha} \nabla_\alpha \log g(\theta, \alpha) \mid \mathbf{X}],$$

and that $\sup_{\alpha \in S} \mathbb{E}[\langle \|\boldsymbol{\theta}\|_2^2 \rangle_\alpha \mid \mathbf{X}] \leq \mathbb{E}[\sup_{\alpha \in S} \langle \|\boldsymbol{\theta}\|_2^2 \rangle_\alpha \mid \mathbf{X}] \leq Cd$ on $\mathcal{E}(C_0, C_{\mathrm{LSI}})$, by (181). Then the argument (185) shows also that $\nabla_\alpha \mathbb{E}[\widehat{F}(\alpha) \mid \mathbf{X}]$ is uniformly bounded and equicontinuous over $\alpha \in S$, hence

$$\lim_{n,d\to\infty} \sup_{\alpha \in S} |\mathbb{E}[\widehat{F}(\alpha) \mid \mathbf{X}] - F(\alpha)| = 0.$$

Since $\alpha^*$ is a minimizer of $\mathbb{E}[\widehat{F}(\alpha) \mid \mathbf{X}]$, this implies that $F(\alpha) \geq F(\alpha^*)$ for every $\alpha \in S$. Since this holds for every compact subset $S \subset O$, this shows part (c). $\square$

We proceed to prove Theorem 2.13. Let $\{\boldsymbol{\theta}^t, \widehat{\alpha}^t\}_{t \geq 0}$ be the solution of the adaptive Langevin equations (9–10). Let $\{\alpha^t\}_{t \geq 0}$ be the (deterministic) $\alpha$-component of the DMFT limit of $\{\widehat{\alpha}^t\}_{t \geq 0}$ prescribed by Theorem 2.3(b), and consider the SDE

$$\mathrm{d}\tilde{\boldsymbol{\theta}}^t = \nabla_{\tilde{\boldsymbol{\theta}}}\left( -\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\theta}}^t\|_2^2 + \sum_{j=1}^d \log g(\tilde{\theta}_j^t, \alpha^t) \right)\mathrm{d}t + \sqrt{2}\,\mathrm{d}\mathbf{b}^t \tag{187}$$

which replaces $\widehat{\alpha}^t$ by $\alpha^t$. We couple $\{\tilde{\boldsymbol{\theta}}^t\}_{t\geq 0}$ to $\{\boldsymbol{\theta}^t, \widehat{\alpha}^t\}_{t\geq 0}$ via the same initial conditions $\tilde{\boldsymbol{\theta}}^0 = \boldsymbol{\theta}^0$ and $\alpha^0$ of Assumption 2.1, and via the same Brownian motion $\{\mathbf{b}^t\}_{t\geq 0}$.

We write $q_t$ for the density of $\tilde{\boldsymbol{\theta}}^t$ conditional on $\mathbf{X}, \boldsymbol{\theta}^*, \boldsymbol{\varepsilon}$ and averaging over $\tilde{\boldsymbol{\theta}}^0$, where $q_0 = g_0^{\otimes d}$ is the initial density of $\tilde{\boldsymbol{\theta}}^0 = \boldsymbol{\theta}^0$. In parallel to (179), we denote

$$\bar{\mathsf{P}}_{\tilde{\boldsymbol{\theta}}^t} = \frac{1}{d}\sum_{j=1}^d \delta_{(\theta_j^*, \tilde{\theta}_j^t)}, \qquad \bar{\mathsf{P}}_t = \langle \bar{\mathsf{P}}_{\tilde{\boldsymbol{\theta}}^t} \rangle \tag{188}$$

where $\langle \cdot \rangle$ is the average with respect to $\tilde{\boldsymbol{\theta}}^t \sim q_t$, i.e. the average over $\tilde{\boldsymbol{\theta}}^0$ and $\{\mathbf{b}^t\}_{t\geq 0}$. We write $\tilde{\theta}^t \sim \bar{\mathsf{P}}_t$ for the $\tilde{\theta}^t$-marginal of a sample $(\theta^*, \tilde{\theta}^t) \sim \bar{\mathsf{P}}_t$.

**Lemma 5.1.** *Under Assumptions 2.1 and 2.2(b), there exists a unique solution $\{\tilde{\boldsymbol{\theta}}^t\}_{t\geq 0}$ to (187). Letting $q_t$ be the above conditional density of $\tilde{\boldsymbol{\theta}}^t$, and letting $V(q, \alpha)$ be the Gibbs free energy (13), almost surely*

$$\limsup_{n,d\to\infty} \sup_{t\in[0,T]} \frac{\mathrm{d}}{\mathrm{d}t}\Big(V(q_t, \alpha^t) + R(\alpha^t)\Big) \leq 0. \tag{189}$$

*Proof.* Fixing $C_0 > 0$ large enough, let

$$\mathcal{E}(C_0) = \{\|\mathbf{X}\|_{\mathrm{op}} \leq C_0, \ \|\boldsymbol{\theta}^*\|_2^2, \|\boldsymbol{\varepsilon}\|_2^2 \leq C_0 d\}.$$

We restrict to the event where the almost-sure convergence statements of Theorem 2.3(b) hold and where $\mathcal{E}(C_0)$ holds for all large $n, d$.

Since $\{\alpha^t\}_{t\geq 0}$ is continuous, for each $T > 0$, there exists a compact ball $S_T$ for which $\alpha^t \in S_T$ for all $t \in [0, T]$. By Assumption 2.2(b), $(\theta, \alpha) \mapsto \partial_\theta \log g(\theta, \alpha)$ restricted to $\alpha \in S_T$ is Lipschitz. Then the drift of (187) is Lipschitz over each time horizon $[0, T]$, so (187) admits a unique solution $\{\tilde{\boldsymbol{\theta}}^t\}_{t\in[0,T]}$ (c.f. [84, Theorem II.1.2]) over $t \in [0, T]$ for every $T \geq 0$, and hence also over all $t \geq 0$. We note that

$$\frac{\mathrm{d}}{\mathrm{d}t}(\tilde{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^t) = \frac{1}{\sigma^2}\mathbf{X}^\top\mathbf{X}(\boldsymbol{\theta}^t - \tilde{\boldsymbol{\theta}}^t) + \Big[\partial_\theta \log g(\tilde{\theta}_j^t, \alpha^t) - \partial_\theta \log g(\theta_j^t, \widehat{\alpha}^t)\Big]_{j=1}^d.$$

Applying again the Lipschitz property of $(\theta, \alpha) \mapsto \partial_\theta \log g(\theta, \alpha)$ over $\alpha \in S_T$ and the bound $\|\mathbf{X}\|_{\mathrm{op}} \leq C_0$, there is a constant $C > 0$ depending on $C_0, T$ such that

$$\left\|\frac{1}{\sqrt{d}}\frac{\mathrm{d}}{\mathrm{d}t}(\tilde{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^t)\right\|_2 \leq \frac{C}{\sqrt{d}}\|\tilde{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^t\|_2 + C\|\alpha^t - \widehat{\alpha}^t\|_2.$$

Since $\sup_{t\in[0,T]}\|\alpha^t - \widehat{\alpha}^t\|_2 \to 0$ by Theorem 2.3(b), a Gronwall argument implies

$$\lim_{n,d\to\infty} \sup_{t\in[0,T]} \frac{1}{\sqrt{d}}\|\tilde{\boldsymbol{\theta}}^t - \boldsymbol{\theta}^t\|_2 = 0. \tag{190}$$

By the DMFT equation (27), the evolution of $\alpha^t$ is given by

$$\frac{\mathrm{d}}{\mathrm{d}t}\alpha^t = \mathbb{E}_{\theta^t \sim \mathsf{P}(\theta^t)}\nabla_\alpha \log g(\theta^t, \alpha^t) - \nabla R(\alpha^t) \tag{191}$$

where $\mathsf{P}(\theta^t)$ is the law of the DMFT variable $\theta^t$. The law $q_t$ of $\tilde{\boldsymbol{\theta}}^t$ satisfies the Fokker-Planck equation

$$\frac{\mathrm{d}}{\mathrm{d}t}q_t(\tilde{\boldsymbol{\theta}}) = \nabla_{\tilde{\boldsymbol{\theta}}} \cdot \left[q_t(\tilde{\boldsymbol{\theta}})\nabla_{\tilde{\boldsymbol{\theta}}}\Big(\frac{1}{2\sigma^2}\|\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\theta}}\|_2^2 - \sum_{j=1}^d \log g(\tilde{\theta}_j, \alpha^t) + \log q_t(\tilde{\boldsymbol{\theta}})\Big)\right]. \tag{192}$$

Then, using (191) and (192) to differentiate $V(q_t, \alpha^t) + R(\alpha^t)$,

$$\frac{\mathrm{d}}{\mathrm{d}t}\Big(V(q_t, \alpha^t) + R(\alpha^t)\Big) = -\frac{1}{d}\underbrace{\int \left\|\nabla_{\tilde{\boldsymbol{\theta}}}\Big(\frac{1}{2\sigma^2}\|\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\theta}}\|_2^2 - \sum_{j=1}^d \log g(\tilde{\theta}_j, \alpha^t) + \log q_t(\tilde{\boldsymbol{\theta}})\Big)\right\|_2^2 q_t(\tilde{\boldsymbol{\theta}})\mathrm{d}\tilde{\boldsymbol{\theta}}}_{:=\mathrm{FI}_t} \tag{193}$$

$$- \Big(\mathbb{E}_{\theta^t \sim \mathsf{P}(\theta^t)}\nabla_\alpha \log g(\theta^t, \alpha^t) - \nabla R(\alpha^t)\Big)^\top\Big(\int \frac{1}{d}\sum_{j=1}^d \nabla_\alpha \log g(\tilde{\theta}_j, \alpha^t)q_t(\tilde{\boldsymbol{\theta}})\mathrm{d}\tilde{\boldsymbol{\theta}} - \nabla R(\alpha^t)\Big).$$

Here, the first term $\text{FI}_t$ (the relative Fisher information) arises from differentiation in $q_t$ and integration-by-parts in $\tilde{\boldsymbol{\theta}}$, while the second term arises from differentiation in $\alpha^t$. Recalling the notation (188),

$$\int \frac{1}{d} \sum_{j=1}^d \nabla_\alpha \log g(\tilde{\theta}_j, \alpha^t) q_t(\tilde{\boldsymbol{\theta}}) \mathrm{d}\tilde{\boldsymbol{\theta}} = \mathbb{E}_{\tilde{\theta}^t \sim \bar{\mathsf{P}}_t} \nabla_\alpha \log g(\tilde{\theta}^t, \alpha^t)$$

so we may write the above as

$$
\begin{aligned}
&\frac{\mathrm{d}}{\mathrm{d}t} \Big( V(q_t, \alpha^t) + R(\alpha^t) \Big) \\
&= -\frac{1}{d} \text{FI}_t - \Big\| \mathbb{E}_{\tilde{\theta}^t \sim \bar{\mathsf{P}}_t} \nabla_\alpha \log g(\tilde{\theta}^t, \alpha^t) - \nabla R(\alpha^t) \Big\|^2 \\
&\quad + \underbrace{\Big( \mathbb{E}_{\tilde{\theta}^t \sim \bar{\mathsf{P}}_t} \nabla_\alpha \log g(\tilde{\theta}^t, \alpha^t) - \mathbb{E}_{\theta^t \sim \mathsf{P}(\theta^t)} \nabla_\alpha \log g(\theta^t, \alpha^t) \Big)^\top \Big( \mathbb{E}_{\tilde{\theta}^t \sim \bar{\mathsf{P}}_t} \nabla_\alpha \log g(\tilde{\theta}^t, \alpha^t) - \nabla R(\alpha^t) \Big)}_{:=\Delta_t}.
\end{aligned}
\tag{194}
$$

By the convexity $W_2(\lambda \mathsf{P} + (1-\lambda)\mathsf{P}', \mathsf{Q})^2 \le \lambda W_2(\mathsf{P}, \mathsf{Q})^2 + (1-\lambda)W_2(\mathsf{P}', \mathsf{Q})^2$ and Jensen's inequality,

$$\sup_{t \in [0,T]} W_2(\bar{\mathsf{P}}_t, \mathsf{P}(\theta^*, \theta^t))^2 \le \sup_{t \in [0,T]} \langle W_2(\bar{\mathsf{P}}_{\tilde{\boldsymbol{\theta}}^t}, \mathsf{P}(\theta^*, \theta^t))^2 \rangle \le \Big\langle \sup_{t \in [0,T]} W_2(\bar{\mathsf{P}}_{\tilde{\boldsymbol{\theta}}^t}, \mathsf{P}(\theta^*, \theta^t))^2 \Big\rangle, \tag{195}$$

where $\langle \cdot \rangle$ is the average over $\tilde{\boldsymbol{\theta}}^t \sim q_t$, and $\mathsf{P}(\theta^*, \theta^t)$ is the joint law of the DMFT variables $(\theta^*, \theta^t)$. By Theorem 2.3(b) and (190), for any fixed $T > 0$ we have

$$\sup_{t \in [0,T]} W_2(\bar{\mathsf{P}}_{\tilde{\boldsymbol{\theta}}^t}, \mathsf{P}(\theta^*, \theta^t))^2 \le \sup_{t \in [0,T]} 2W_2(\bar{\mathsf{P}}_{\tilde{\boldsymbol{\theta}}^t}, \bar{\mathsf{P}}_{\boldsymbol{\theta}^t})^2 + 2W_2(\bar{\mathsf{P}}_{\boldsymbol{\theta}^t}, \mathsf{P}(\theta^*, \theta^t))^2 \to 0 \tag{196}$$

almost surely as $n, d \to \infty$. The same arguments as leading to (168) show that $\sup_{t \in [0,T]} W_2(\bar{\mathsf{P}}_{\tilde{\boldsymbol{\theta}}^t}, \mathsf{P}(\theta^*, \theta^t))^2$ is uniformly integrable with respect to $\langle \cdot \rangle$ for all large $n, d$. Then applying (196) and dominated convergence to bound the right side of (195), we get

$$\lim_{n,d \to \infty} \sup_{t \in [0,T]} W_2(\bar{\mathsf{P}}_t, \mathsf{P}(\theta^*, \theta^t))^2 = 0. \tag{197}$$

Finally, applying that $(\theta, \alpha) \mapsto \nabla_\alpha \log g(\theta, \alpha)$ is uniformly Lipschitz over $\alpha \in S_T$ by Assumption 2.2(b), this Wasserstein-2 convergence implies

$$\lim_{n,d \to \infty} \sup_{t \in [0,T]} \Big| \mathbb{E}_{\theta \sim \bar{\mathsf{P}}_t} \nabla_\alpha \log g(\theta, \alpha^t) - \mathbb{E}_{\theta \sim \mathsf{P}(\theta^t)} \nabla_\alpha \log g(\theta, \alpha^t) \Big| = 0,$$

hence $\lim_{n,d \to \infty} \sup_{t \in [0,T]} |\Delta_t| = 0$ for the quantity $\Delta_t$ of (194). As the first two terms of (194) are non-positive, this shows (189). $\square$

*Proof of Theorem 2.13.* Let $S \subset O \subset \mathbb{R}^K$ be the domains of Assumption 2.11. Fixing sufficiently large constants $C_0, C_{\text{LSI}} > 0$, define

$$\mathcal{E}(C_0, C_{\text{LSI}}) = \{ \|\mathbf{X}\|_{\text{op}} \le C_0, \ \|\boldsymbol{\theta}^*\|_2^2, \|\boldsymbol{\varepsilon}\|_2^2 \le C_0 d, \text{ and (46) holds for } q_\alpha \text{ for every } \alpha \in O \}.$$

We restrict to the event where the almost-sure convergence statements of Theorem 2.3(b) and Lemma 2.12 hold, and where $\mathcal{E}(C_0, C_{\text{LSI}})$ holds for all large $n, d$. Throughout, $C, C', c > 0$ denote constants that may depend on $C_0, C_{\text{LSI}}$ and change from instance to instance.

On the event $\mathcal{E}(C_0, C_{\text{LSI}})$, we first note that by Itô's formula,

$$\mathrm{d}\|\tilde{\boldsymbol{\theta}}^t\|_2^2 = 2(\tilde{\boldsymbol{\theta}}^t)^\top \Big[ \Big( \frac{1}{\sigma^2} \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\theta}}^t) + \Big( \partial_\theta \log g(\tilde{\theta}_j^t, \alpha^t) \Big)_{j=1}^d \Big) \mathrm{d}t + \sqrt{2} \, \mathrm{d}\mathbf{b}^t \Big] + (2d)\mathrm{d}t,$$

and hence

$$\frac{\mathrm{d}}{\mathrm{d}t} \langle d^{-1} \|\tilde{\boldsymbol{\theta}}^t\|_2^2 \rangle \le C(1 + \langle d^{-1/2} \|\tilde{\boldsymbol{\theta}}^t\|_2 \rangle) + 2\Big\langle d^{-1} \sum_{j=1}^d \tilde{\theta}_j^t \cdot \partial_\theta \log g(\tilde{\theta}_j^t, \alpha^t) \Big\rangle$$

for a constant $C > 0$. Under the convexity-at-infinity condition of Assumption 2.2(b), there exist constants $C, c > 0$ for which $\theta \cdot \partial_\theta \log g(\theta, \alpha^t) \leq C|\theta| - c\theta^2$ for all $\theta \in \mathbb{R}$ and $\alpha^t \in S$. Applying this and Cauchy-Schwarz to the above, we have for some constants $C', c' > 0$ that $\frac{\mathrm{d}}{\mathrm{d}t} \langle d^{-1}\|\tilde{\boldsymbol{\theta}}^t\|_2^2 \rangle \leq C' - c' \langle d^{-1}\|\tilde{\boldsymbol{\theta}}^t\|_2^2 \rangle$. This implies on $\mathcal{E}(C_0, C_{\mathrm{LSI}})$ that

$$d^{-1}\langle \|\tilde{\boldsymbol{\theta}}^t\|_2^2 \rangle \leq C \tag{198}$$

for a constant $C > 0$ and all $t \geq 0$. The arguments leading to (168) show that for any fixed $t \geq 0$, $d^{-1}\|\tilde{\boldsymbol{\theta}}^t\|_2^2$ is uniformly integrable with respect to $\langle \cdot \rangle$ for all large $n, d$. Since $\mathcal{E}(C_0, C_{\mathrm{LSI}})$ holds a.s. for all large $n, d$, and $\lim_{n,d\to\infty} d^{-1}\|\tilde{\boldsymbol{\theta}}^t\|_2^2 = (\theta^t)^2$ a.s. by Theorem 2.3(b) and (190) where $\theta^t$ here is the $\theta$-component of the limiting DMFT system, this implies also

$$\mathbb{E}(\theta^t)^2 \leq C \tag{199}$$

for all $t \geq 0$. Furthermore, for any $s \leq t$, applying

$$\tilde{\boldsymbol{\theta}}^t - \tilde{\boldsymbol{\theta}}^s = \int_s^t \left[ \frac{1}{\sigma^2}\mathbf{X}^\top(\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\theta}}^r) + \Big(\partial_\theta \log g(\tilde{\theta}_j^r, \alpha^r)\Big)_{j=1}^d \right]\mathrm{d}r + \sqrt{2}(\mathbf{b}^t - \mathbf{b}^s)$$

and uniform Lipschitz continuity of $\theta \mapsto \partial_\theta \log g(\theta, \alpha^r)$ for $\alpha^r \in S$, we have on $\mathcal{E}(C_0, C_{\mathrm{LSI}})$ that

$$d^{-1/2}\|\tilde{\boldsymbol{\theta}}^t - \tilde{\boldsymbol{\theta}}^s\|_2 \leq \int_s^t Cd^{-1/2}\|\tilde{\boldsymbol{\theta}}^r - \tilde{\boldsymbol{\theta}}^s\|_2\,\mathrm{d}r + C(t-s)(1 + d^{-1/2}\|\tilde{\boldsymbol{\theta}}^s\|_2) + \sqrt{2}d^{-1/2}\|\mathbf{b}^t - \mathbf{b}^s\|_2.$$

Then by Gronwall's inequality,

$$d^{-1/2}\|\tilde{\boldsymbol{\theta}}^t - \tilde{\boldsymbol{\theta}}^s\|_2 \leq Ce^{C(t-s)}\Big(C(t-s)(1 + d^{-1/2}\|\tilde{\boldsymbol{\theta}}^s\|_2) + d^{-1/2}\sup_{r\in[s,t]}\|\mathbf{b}^r - \mathbf{b}^s\|_2\Big).$$

Then applying (198) and Doob's maximal inequality shows

$$d^{-1}\langle \|\tilde{\boldsymbol{\theta}}^t - \tilde{\boldsymbol{\theta}}^s\|_2^2 \rangle \leq C(t-s) \text{ for all } s \leq t \text{ with } t - s \leq 1. \tag{200}$$

We now show that for a constant $C > 0$,

$$\int_0^\infty \|\nabla F(\alpha^t) + \nabla R(\alpha^t)\|_2^2\mathrm{d}t < C. \tag{201}$$

We remind the reader that $q_t$ is the law of $\tilde{\boldsymbol{\theta}}^t$ (conditioned on $\mathbf{X}, \boldsymbol{\theta}^*, \boldsymbol{\varepsilon}$) and $q_{\alpha^t}$ is the posterior law of $\boldsymbol{\theta}$ under the prior $g \equiv g(\cdot, \alpha^t)$. On $\mathcal{E}(C_0, C_{\mathrm{LSI}})$, the LSI for $q_{\alpha^t}$ and its implied $T_2$-transportation inequality (c.f. [83, Theorem 9.6.1]) imply for the Fisher information term $\mathrm{FI}_t$ of (193) that

$$\mathrm{FI}_t \geq C_{\mathrm{LSI}}^{-1}\mathrm{D}_{\mathrm{KL}}(q_t\|q_{\alpha^t}) \geq C_{\mathrm{LSI}}^{-2}W_2(q_t, q_{\alpha^t})^2$$

for all $t \geq 0$. The average marginal distribution of coordinates of $\tilde{\boldsymbol{\theta}}^t \sim q_t$ is the $\tilde{\theta}^t$-marginal of $\bar{\mathsf{P}}_t$ defined in (188), and that of $\boldsymbol{\theta} \sim q_{\alpha^t}$ is the $\theta$-marginal of $\bar{\mathsf{P}}_{\alpha^t}$ as defined in (179). Considering the coordinatewise coupling of $\bar{\mathsf{P}}_t, \bar{\mathsf{P}}_{\alpha^t}$, we see that $W_2(\bar{\mathsf{P}}_t, \bar{\mathsf{P}}_{\alpha^t})^2 \leq d^{-1}W_2(q_t, q_{\alpha^t})^2$, so

$$d^{-1}\mathrm{FI}_t \geq C_{\mathrm{LSI}}^{-2}W_2(\bar{\mathsf{P}}_t, \bar{\mathsf{P}}_{\alpha^t})^2. \tag{202}$$

Applying this and the uniform Lipschitz continuity of $\theta \mapsto \nabla_\alpha \log g(\theta, \alpha)$ over $\alpha \in O$ guaranteed by Assumption 2.2(b),

$$\left\|\mathbb{E}_{\tilde{\theta}^t\sim\bar{\mathsf{P}}_t}\nabla_\alpha \log g(\tilde{\theta}^t, \alpha^t) - \mathbb{E}_{\theta\sim\bar{\mathsf{P}}_{\alpha^t}}\nabla_\alpha \log g(\theta, \alpha^t)\right\|_2^2 \leq C'\,W_2(\bar{\mathsf{P}}_t, \bar{\mathsf{P}}_{\alpha^t})^2 \leq Cd^{-1}\mathrm{FI}_t.$$

Then applying this as a lower bound for $d^{-1}\mathrm{FI}_t$ in (194), and applying also $C^{-1}(a-b)^2 + b^2 \geq c_0 a^2$ for a constant $c_0 > 0$ and all $a, b \in \mathbb{R}$, we get from (194) that

$$\frac{\mathrm{d}}{\mathrm{d}t}\Big(V(q_t, \alpha^t) + R(\alpha^t)\Big) \leq -c_0\left\|\mathbb{E}_{\theta\sim\bar{\mathsf{P}}_{\alpha^t}}\nabla_\alpha \log g(\theta, \alpha^t) - \nabla R(\alpha^t)\right\|^2 + \Delta_t.$$

Now note from Lemma 2.12 that

$$\mathbb{E}_{\theta \sim \bar{\mathsf{P}}_{\alpha^t}} \nabla_\alpha \log g(\theta, \alpha^t) = -\nabla \widehat{F}(\alpha^t).$$

Applying $\sup_{t \in [0,T]} \Delta_t \to 0$ and the uniform convergence $\nabla \widehat{F}(\alpha) \to \nabla F(\alpha)$ over $\alpha \in S$ from Lemma 2.12, this shows a strengthening of (189): for any $t \in [0, T]$,

$$\limsup_{n,d \to \infty} \frac{\mathrm{d}}{\mathrm{d}t} \left( V(q_t, \alpha^t) + R(\alpha^t) \right) \le -c_0 \|\nabla F(\alpha^t) + \nabla R(\alpha^t)\|_2^2.$$

Then for any $T > 0$,

$$c_0 \int_0^T \|\nabla F(\alpha^t) + \nabla R(\alpha^t)\|_2^2 \, \mathrm{d}t \le \limsup_{n,d \to \infty} V(q_0, \alpha^0) + R(\alpha^0) - V(q_T, \alpha^T) - R(\alpha^T).$$

Note that by the definition of $V(q, \alpha)$ in (13) and the conditions of finite moments and finite entropy for $g_0$ in Assumption 2.1, $V(q_0, \alpha^0) = V(g_0^{\otimes d}, \alpha^0)$ is bounded above by a constant on $\mathcal{E}(C_0, C_{\mathrm{LSI}})$ for all large $n, d$. Also by the definition (13),

$$V(q_T, \alpha^T) \ge \frac{1}{d} \mathrm{D}_{\mathrm{KL}}(q_T \| g(\cdot, \alpha^T)^{\otimes d}) + \frac{n}{2d} \log 2\pi\sigma^2 \ge \frac{n}{2d} \log 2\pi\sigma^2$$

which is bounded below by a constant for all $T$ and all large $n, d$. Then, applying also $R(\alpha^0) \le C$ and $R(\alpha^T) \ge 0$ and taking the limit $n, d \to \infty$ followed by $T \to \infty$, we obtain the claimed bound (201).

Consider the set
$$\mathrm{Crit} = \{\alpha \in S : \nabla F(\alpha) + \nabla R(\alpha) = 0\}.$$

Suppose by contradiction that $\{\alpha^t\}_{t \ge 0}$ has a limit point $\alpha^\infty \in S$ that does not belong to Crit. Lemma 2.12 implies that $\nabla F(\alpha) + \nabla R(\alpha)$ is continuous over $\alpha \in O$, so $\|\nabla F(\alpha) + \nabla R(\alpha)\|_2 > \delta$ for all $\alpha \in B_\delta(\alpha^\infty) := \{\alpha : \|\alpha - \alpha^\infty\|_2 < \delta\}$ and some $\delta > 0$. However, Assumption 2.2(b) and the DMFT equation (27) imply

$$\left\| \frac{\mathrm{d}}{\mathrm{d}t} \alpha^t \right\|_2 \le \mathbb{E}_{\theta \sim \mathsf{P}(\theta^t)} \|\nabla_\alpha \log g(\theta, \alpha^t)\|_2 + \|\nabla R(\alpha^t)\|_2 \le C(1 + \mathbb{E}|\theta^t| + \|\alpha^t\|_2) \le C' \tag{203}$$

for some constants $C, C' > 0$ and all $t \ge 0$, where the last inequality applies (199) and the assumption $\alpha^t \in S$. Then for each $t_0 \ge 0$ such that

$$\alpha^{t_0} \in B_{\delta/2}(\alpha^\infty) \tag{204}$$

we must have $\alpha^t \in B_\delta(\alpha^\infty)$ for all $t \in [t_0 - c\delta, t_0 + c\delta]$ and some constant $c > 0$. Then $\int_{t_0 - c\delta}^{t_0 + c\delta} \|\nabla F(\alpha^t) + \nabla R(\alpha^t)\|_2^2 \, \mathrm{d}t \ge 2c\delta^3$. The condition (204) must hold for infinitely many times $t_0$ because $\alpha^\infty$ is a limit point of $\{\alpha^t\}_{t \ge 0}$, but this contradicts (201). Thus we must have $\alpha^\infty \in \mathrm{Crit}$. Since this holds for every limit point $\alpha^\infty$ of $\{\alpha^t\}_{t \ge 0}$, and $S$ is compact, this implies $\lim_{t \to \infty} \mathrm{dist}(\alpha^t, \mathrm{Crit}) = 0$. If furthermore all points of Crit are isolated, then the limit point $\alpha^\infty$ of $\{\alpha^t\}_{t \ge 0}$ must be unique, and

$$\lim_{t \to \infty} \alpha^t = \alpha^\infty.$$

For the remaining statements (53), fix any $\varepsilon > 0$. Choosing $T(\varepsilon)$ such that $\|\alpha^t - \alpha^\infty\|_2 < \varepsilon/2$ for all $t > T(\varepsilon)$, we then have $\limsup_{n,d \to \infty} \|\widehat{\alpha}^t - \alpha^\infty\|_2 < \varepsilon$ by Theorem 2.3(b), showing the first statement of (53). For the second statement of (53), we note from (194) that

$$\frac{\mathrm{d}}{\mathrm{d}t} \left( V(q_t, \alpha^t) + R(\alpha^t) \right) \le -\frac{1}{d} \mathrm{FI}_t + \Delta_t.$$

Then, by the same arguments as above, for some constant $C > 0$ and every $T > 0$,

$$\limsup_{n,d \to \infty} \int_0^T d^{-1} \mathrm{FI}_t \le \limsup_{n,d \to \infty} V(q_0, \alpha^0) + R(\alpha^0) - V(q_T, \alpha^T) - R(\alpha^T) \le C.$$

Recalling (202), this implies

$$\limsup_{n,d\to\infty} \int_0^T W_2(\bar{\mathsf{P}}_t, \bar{\mathsf{P}}_{\alpha^t})^2 \mathrm{d}t \le C. \tag{205}$$

For each fixed $t \ge 0$, we have

$$\lim_{n,d\to\infty} W_2(\bar{\mathsf{P}}_t, \mathsf{P}(\theta^*, \theta^t))^2 = 0 \tag{206}$$

by (197). We have also by Jensen's inequality for the squared Wasserstein-2 distance and (48) of Corollary 2.10(a),

$$\limsup_{n,d\to\infty} W_2(\bar{\mathsf{P}}_{\alpha^t}, \mathsf{P}_{\alpha^t})^2 \le \limsup_{n,d\to\infty} \left\langle W_2(\bar{\mathsf{P}}_{\boldsymbol{\theta}}, \mathsf{P}_{\alpha^t})^2 \right\rangle_{\alpha^t} = 0 \tag{207}$$

where $\langle \cdot \rangle_{\alpha^t}$ is the average over $\boldsymbol{\theta} \sim q_{\alpha^t}$ defining $\bar{\mathsf{P}}_{\boldsymbol{\theta}}$. Then, combining (206) and (207), we have that $\lim_{n,d\to\infty} W_2(\bar{\mathsf{P}}_t, \bar{\mathsf{P}}_{\alpha^t}) = W_2(\mathsf{P}(\theta^*, \theta^t), \mathsf{P}_{\alpha^t})$. Applying this and Fatou's lemma to (205), we obtain the bound $\int_0^T W_2(\mathsf{P}(\theta^*, \theta^t), \mathsf{P}_{\alpha^t})^2 \mathrm{d}t \le C$. Since $T > 0$ is arbitrary, taking $T \to \infty$ gives

$$\int_0^\infty W_2(\mathsf{P}(\theta^*, \theta^t), \mathsf{P}_{\alpha^t})^2 \mathrm{d}t \le C. \tag{208}$$

For any $s \le t$, considering the coordinatewise coupling gives $W_2(\bar{\mathsf{P}}_s, \bar{\mathsf{P}}_t)^2 \le d^{-1}\langle \|\tilde{\boldsymbol{\theta}}^s - \tilde{\boldsymbol{\theta}}^t\|_2^2 \rangle \le C(t-s)$, where the second inequality holds for a constant $C > 0$ and all $t - s \in [0,1]$ by (200). Also

$$W_2(\bar{\mathsf{P}}_{\alpha^s}, \bar{\mathsf{P}}_{\alpha^t})^2 \le d^{-1}W_2(q_{\alpha^s}, q_{\alpha^t})^2 \le C\|\alpha^t - \alpha^s\|_2^2 \le C'(t-s)^2 \tag{209}$$

by the Wasserstein-2 Lipschitz continuity of $q_\alpha$ over $\alpha \in S$ shown in (171), and the bound (203) for $\mathrm{d}\alpha^t/\mathrm{d}t$. Then taking the limit $n, d \to \infty$ using (206) and (207), this shows

$$|W_2(\mathsf{P}(\theta^*, \theta^t), \mathsf{P}_{\alpha^t})^2 - W_2(\mathsf{P}(\theta^*, \theta^s), \mathsf{P}_{\alpha^s})^2| \le C(t-s)$$

for all $t - s \in [0,1]$. Then $t \mapsto W_2(\mathsf{P}(\theta^*, \theta^t), \mathsf{P}_{\alpha^t})^2$ is Lipschitz, so (208) implies

$$\lim_{t\to\infty} W_2(\mathsf{P}(\theta^*, \theta^t), \mathsf{P}_{\alpha^t})^2 = 0. \tag{210}$$

We have similarly to (209) that $W_2(\bar{\mathsf{P}}_{\alpha^t}, \bar{\mathsf{P}}_{\alpha^\infty})^2 \le d^{-1}W_2(q_{\alpha^t}, q_{\alpha^\infty})^2 \le C\|\alpha^t - \alpha^\infty\|_2^2$. Hence by (207), also $W_2(\mathsf{P}_{\alpha^t}, \mathsf{P}_{\alpha^\infty})^2 \le C\|\alpha^t - \alpha^\infty\|_2^2$, so

$$\lim_{t\to\infty} W_2(\mathsf{P}_{\alpha^t}, \mathsf{P}_{\alpha^\infty})^2 = 0. \tag{211}$$

Combining (210) and (211) show that for any $\varepsilon > 0$, there exists $T(\varepsilon) > 0$ such that $W_2(\mathsf{P}(\theta^*, \theta^t), \mathsf{P}_{\alpha^\infty}) < \varepsilon$ for all $t \ge T(\varepsilon)$. The second statement of (53) follows from this and the almost sure convergence $\lim_{n,d\to\infty} W_2(\frac{1}{d}\sum_j \delta_{(\theta_j^*, \theta_j^t)}, \mathsf{P}(\theta^*, \theta^t)) = 0$ ensured by Theorem 2.3(b). $\qquad\square$

## 5.2 Analysis of examples

*Analysis of Examples 2.14 and 2.15.* We prove the claims in Example 2.15 that Assumptions 2.2(b) and 2.11 hold, and that Crit consists of the unique point $\alpha = \alpha^*$. (Then these claims hold also in Example 2.14 for the Gaussian prior, which is a special case.)

Assumption 2.2(b) is immediate from the given conditions for $f(x)$. For Assumption 2.11, let us first show that there exists a compact interval $S \subset \mathbb{R}$ for which $\{\alpha^t\}_{t\ge 0}$ is confined to $S$ (for all $t \ge 0$): By Lemma 5.1 (which does not require Assumption 2.11), for each fixed $t \ge 0$, almost surely

$$\limsup_{n,d\to\infty} V(q_t, \alpha^t) - V(q_0, \alpha^0) \le 0. \tag{212}$$

By the Gibbs variational principle (12) and the lower bound $-\log g(\theta, \alpha) = f(\theta - \alpha) \ge f(0) + \frac{c_0}{2}(\theta - \alpha)^2$,

$$V(q_t, \alpha^t) \ge \widehat{F}(\alpha^t)$$

$$= -\frac{1}{d}\log \int (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2}\|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 + \sum_{j=1}^d \log g(\theta_j, \alpha^t)\right)\mathrm{d}\boldsymbol{\theta}$$

$$\ge -\frac{1}{d}\log \int (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2}\|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 - f(0)d - \sum_{j=1}^d \frac{c_0}{2}(\theta_j - \alpha^t)^2\right)\mathrm{d}\boldsymbol{\theta}$$

62

Applying $\|\mathbf{X}\|_{\mathrm{op}} \leq C$ a.s. for all large $n, d$, it is readily checked by explicit evaluation of this integral over $\boldsymbol{\theta}$ that

$$V(q_t, \alpha^t) \geq C + \frac{c_0}{2}(\alpha^t)^2 - \frac{1}{2d}\Big(\frac{\mathbf{X}^\top \mathbf{y}}{\sigma^2} + c_0 \alpha^t \mathbf{1}\Big)^\top \Big(\frac{\mathbf{X}^\top \mathbf{X}}{\sigma^2} + c_0 \mathbf{I}\Big)^{-1}\Big(\frac{\mathbf{X}^\top \mathbf{y}}{\sigma^2} + c_0 \alpha^t \mathbf{1}\Big)$$

a.s. for all large $n, d$ and a constant $C \in \mathbb{R}$ depending on $\sigma^2, \delta, f(0), c_0, \alpha^*$, where here $\mathbf{1}$ denotes the all-1's vector in $\mathbb{R}^d$. We have

$$\lim_{n,d\to\infty} \frac{1}{d}\mathbf{1}^\top \Big(\frac{\mathbf{X}^\top \mathbf{X}}{\sigma^2} + c_0 \mathbf{I}\Big)^{-1}\mathbf{1} = \sigma^2 G(-\sigma^2 c_0) < \frac{1}{c_0}$$

strictly, where $G(z) = \lim d^{-1}\operatorname{Tr}(\mathbf{X}^\top \mathbf{X} - z\mathbf{I})^{-1}$ denotes the Stieltjes transform of the Marcenko-Pastur spectral limit of $\mathbf{X}^\top \mathbf{X}$ [89, Theorem 2.5]. Applying this to lower-bound the quadratic term in $\alpha^t$ above, and applying Cauchy-Schwarz to lower-bound the linear term, we get

$$V(q_t, \alpha^t) \geq C' + c'(\alpha^t)^2$$

for some constants $C' \in \mathbb{R}$ and $c' > 0$. Now applying this and $V(q_0, \alpha^0) = V(g_0^{\otimes d}, \alpha^0) \leq C$ to (212), we deduce that $(\alpha^t)^2$ is uniformly bounded over all $t \geq 0$, i.e. there exists a compact interval $S$ for which $\alpha^t \in S$ for all $t \geq 0$, as claimed. By enlarging $S$, we may assume without loss of generality $\alpha^* \in S$. Then, taking $O$ to be any neighborhood of $S$, the remaining LSI condition of Assumption 2.11 holds by the strong convexity of $f(x)$ and Proposition 2.8.

We now show that $F(\alpha)$ is strictly convex on $O$, by showing convexity of the original negative log-likelihood $\widehat{F}(\alpha)$: Fixing sufficiently large and small constants $C_0, c > 0$, let us restrict to the event

$$\mathcal{E} = \{\|\mathbf{X}\|_{\mathrm{op}} \leq C_0,\ \|\mathbf{X}\mathbf{1}\|_2 \geq c\sqrt{d}\}$$

which holds a.s. for all large $n, d$. Recalling the form of $\nabla^2 \widehat{F}(\alpha)$ from (186) and applying this with $-\log g(\theta, \alpha) = f(\theta - \alpha)$,

$$\widehat{F}''(\alpha) = \frac{1}{d}\Big\langle \sum_{j=1}^d f''(\theta_j - \alpha) \Big\rangle_\alpha - \frac{1}{d}\operatorname{Var}_\alpha\Big[\sum_{j=1}^d f'(\theta_j - \alpha)\Big]$$

where $\langle \cdot \rangle_\alpha$ is the average under the posterior law corresponding to $g(\cdot, \alpha)$, and $\operatorname{Var}_\alpha$ is its posterior variance. Since $f(x)$ is strictly convex, the posterior density of $\theta$ is strictly log-concave for each fixed $\alpha$. Then, denoting

$$\mathbf{v}_\alpha(\boldsymbol{\theta}) = \Big(f''(\theta_j - \alpha)\Big)_{j=1}^d \in \mathbb{R}^d, \qquad \mathbf{D}_\alpha(\boldsymbol{\theta}) = \operatorname{diag}\Big(f''(\theta_j - \alpha)\Big)_{j=1}^d \in \mathbb{R}^{d\times d},$$

the Brascamp-Lieb inequality [83, Theorem 4.9.1] implies

$$\operatorname{Var}_\alpha\Big[\sum_{j=1}^d f'(\theta_j - \alpha)\Big] \leq \Big\langle \mathbf{v}_\alpha(\boldsymbol{\theta})^\top \Big(\mathbf{D}_\alpha(\boldsymbol{\theta}) + \frac{\mathbf{X}^\top \mathbf{X}}{\sigma^2}\Big)^{-1}\mathbf{v}_\alpha(\boldsymbol{\theta})\Big\rangle_\alpha$$

Observing also that $\sum_{j=1}^d f''(\theta_j - \alpha) = \mathbf{v}_\alpha(\boldsymbol{\theta})^\top \mathbf{D}_\alpha(\boldsymbol{\theta})^{-1}\mathbf{v}_\alpha(\boldsymbol{\theta})$, this shows

$$\widehat{F}''(\alpha) \geq \frac{1}{d}\Big\langle \mathbf{v}_\alpha(\boldsymbol{\theta})^\top \Big[\mathbf{D}_\alpha(\boldsymbol{\theta})^{-1} - \Big(\mathbf{D}_\alpha(\boldsymbol{\theta}) + \frac{\mathbf{X}^\top \mathbf{X}}{\sigma^2}\Big)^{-1}\Big]\mathbf{v}_\alpha(\boldsymbol{\theta})\Big\rangle_\alpha.$$

Applying the Woodbury matrix identity and $0 \preceq \sigma^2\mathbf{I} + \mathbf{X}\mathbf{D}_\alpha(\boldsymbol{\theta})^{-1}\mathbf{X}^\top \preceq C'\mathbf{I}$ on the event $\mathcal{E}$ for some constant $C' > 0$,

$$\widehat{F}''(\alpha) \geq \frac{1}{d}\Big\langle \mathbf{v}_\alpha(\boldsymbol{\theta})^\top \Big[\mathbf{D}_\alpha(\boldsymbol{\theta})^{-1}\mathbf{X}^\top \Big(\sigma^2\mathbf{I} + \mathbf{X}\mathbf{D}_\alpha(\boldsymbol{\theta})^{-1}\mathbf{X}^\top\Big)^{-1}\mathbf{X}\mathbf{D}_\alpha(\boldsymbol{\theta})^{-1}\Big]\mathbf{v}_\alpha(\boldsymbol{\theta})\Big\rangle_\alpha$$

$$\geq \frac{1}{C'd}\Big\langle \mathbf{v}_\alpha(\boldsymbol{\theta})^\top \mathbf{D}_\alpha(\boldsymbol{\theta})^{-1}\mathbf{X}^\top \mathbf{X}\mathbf{D}_\alpha(\boldsymbol{\theta})^{-1}\mathbf{v}_\alpha(\boldsymbol{\theta})\Big\rangle_\alpha = \frac{1}{C'd}\mathbf{1}^\top \mathbf{X}^\top \mathbf{X}\mathbf{1} \geq c',$$

the last inequality holding for some $c' > 0$ on $\mathcal{E}$. Thus, on $\mathcal{E}$, $\widehat{F}(\alpha) - (c'/2)\alpha^2$ is convex over $\alpha \in \mathbb{R}$. Since $\mathcal{E}$ holds a.s. for all large $n, d$ and $F(\alpha)$ is the almost-sure pointwise limit of $\widehat{F}(\alpha)$, this implies that $F(\alpha) - (c'/2)\alpha^2$ is also convex [85, Theorem 10.8], so $F(\alpha)$ is strongly convex as claimed. Lemma 2.12(c) implies that $\nabla F(\alpha^*) = 0$, i.e. $\alpha^*$ is a point of Crit, so by this convexity it is the unique point of Crit. $\square$

**Proposition 5.2.** *In the setting of Theorem 2.13, suppose $R(\alpha)$ is given by (54–55) with $\|\alpha^0\|_2 \leq D$. Then there exists a constant $C(g_*, g_0, \alpha^0) > 0$ depending only on $(g_*, g_0, \alpha^0)$ such that the DMFT process $\{\alpha^t\}_{t\geq 0}$ satisfies*

$$\|\alpha^t\|_2 \leq D + C(g_*, g_0, \alpha^0)\Big(\frac{1+\delta}{\sigma^2} + 1\Big) \text{ for all } t \geq 0.$$

*Proof.* By Lemma 5.1, for each fixed $t \geq 0$, almost surely

$$\limsup_{n,d\to\infty} \Big(V(q_t, \alpha^t) + R(\alpha^t)\Big) - \Big(V(q_0, \alpha^0) + R(\alpha^0)\Big) \leq 0. \tag{213}$$

By definition of $V(q, \alpha)$ in (13), we have

$$V(q_0, \alpha^0) = V(g_0^{\otimes d}, \alpha^0) = \frac{1}{d}\int \frac{1}{2\sigma^2}\|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|^2 \prod_{j=1}^{d} g_0(\theta_j)\mathrm{d}\theta_j + \mathrm{D_{KL}}(g_0\|g(\cdot, \alpha^0)) + \frac{n}{2d}\log 2\pi\sigma^2,$$

$$V(q_t, \alpha^t) = \frac{1}{d}\int \frac{1}{2\sigma^2}\|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|^2 q_t(\boldsymbol{\theta})\mathrm{d}\boldsymbol{\theta} + \frac{1}{d}\mathrm{D_{KL}}(q_t\|g(\cdot, \alpha^t)^{\otimes d}) + \frac{n}{2d}\log 2\pi\sigma^2.$$

Let $\Pi_{\mathbf{X}} \in \mathbb{R}^{n\times n}$ be the orthogonal projection onto the column span of $\mathbf{X}$. Then, applying the above forms with $\mathrm{D_{KL}}(q_t\|g(\cdot, \alpha^t)^{\otimes d}) \geq 0$ and noting that $\|(\mathbf{I} - \Pi_{\mathbf{X}})(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})\|^2 = \|(\mathbf{I} - \Pi_{\mathbf{X}})\mathbf{y}\|^2$ which does not depend on $\boldsymbol{\theta}$, we have

$$V(q_0, \alpha^0) - V(q_t, \alpha^t)$$

$$\leq \frac{1}{d}\int \frac{1}{2\sigma^2}\|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|^2 \prod_{j=1}^{d} g_0(\theta_j)\mathrm{d}\theta_j - \frac{1}{d}\int \frac{1}{2\sigma^2}\|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|^2 q_t(\boldsymbol{\theta})\mathrm{d}\boldsymbol{\theta} + \mathrm{D_{KL}}(g_0\|g(\cdot, \alpha^0))$$

$$= \frac{1}{d}\int \frac{1}{2\sigma^2}\|\Pi_{\mathbf{X}}(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})\|^2 \prod_{j=1}^{d} g_0(\theta_j)\mathrm{d}\theta_j - \frac{1}{d}\int \frac{1}{2\sigma^2}\|\Pi_{\mathbf{X}}(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})\|^2 q_t(\boldsymbol{\theta})\mathrm{d}\boldsymbol{\theta} + \mathrm{D_{KL}}(g_0\|g(\cdot, \alpha^0))$$

$$\leq \frac{1}{d}\int \frac{1}{2\sigma^2}\|\Pi_{\mathbf{X}}(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})\|^2 \prod_{j=1}^{d} g_0(\theta_j)\mathrm{d}\theta_j + \mathrm{D_{KL}}(g_0\|g(\cdot, \alpha^0)).$$

Let us apply

$$\|\Pi_{\mathbf{X}}(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})\|^2 \leq 2\|\Pi_{\mathbf{X}}\boldsymbol{\varepsilon}\|^2 + 2\|\mathbf{X}(\boldsymbol{\theta}^* - \boldsymbol{\theta})\|_2^2,$$

$\|\mathbf{X}\|_{\mathrm{op}}^2 \leq C(1 + \delta)$, and $\|\Pi_{\mathbf{X}}\boldsymbol{\varepsilon}\|^2 \leq C\min(n, d)\sigma^2$ for a universal constant $C > 0$ a.s. for all large $n, d$. Then, for a constant $C(g_*, g_0, \alpha^0) > 0$ depending only on $g_*, g_0, \alpha^0$,

$$\limsup_{n,d\to\infty} V(q_0, \alpha^0) - V(q_t, \alpha^t) \leq C(g_*, g_0, \alpha^0)\Big(\frac{1+\delta}{\sigma^2} + 1\Big).$$

Applying this to (213) and noting that $R(\alpha^0) = 0$ because $\|\alpha^0\| \leq D$, for every $t \geq 0$ we get

$$R(\alpha^t) \leq C(g_*, g_0, \alpha^0)\Big(\frac{1+\delta}{\sigma^2} + 1\Big).$$

The lemma follows from this bound and the condition $R(\alpha) \geq \|\alpha\| - D$ whenever $\|\alpha\| \geq D + 1$. $\qquad\square$

*Proof of Proposition 2.16.* Fix any $s^2 = \sigma^2/\delta > 0$. Throughout this proof, constants may depend on $s^2$ but not on $\delta$. Proposition 5.2 implies that there exists a constant radius $D' > 0$ (depending on $s^2$ but not on $\delta$) such that for any $\delta > 1$,

$$\alpha^t \in \mathcal{B}(D') \text{ for all } t \geq 0.$$

Set $S = \overline{\mathcal{B}(D')}$ and $O = \mathcal{B}(D' + 1)$. Then for each fixed $\alpha \in O$, Assumption 2.2(b) implies

$$C \geq -\partial_\theta^2 \log g(\theta, \alpha) \geq \begin{cases} c_0 & \text{for } |\theta| \geq r_0 \\ -C & \text{for all } \theta \in \mathbb{R} \end{cases} \tag{214}$$

for some $C, r_0, c_0 > 0$ uniformly over $\alpha \in O$. By this bound (214) and Proposition 2.8(b), for some sufficiently large $\delta_0 = \delta_0(s^2) > 0$, $\sigma^2 = \delta s^2$, and all $\delta \geq \delta_0$, the LSI (46) must hold for $g = g(\cdot, \alpha)$ and each $\alpha \in O$. This verifies Assumption 2.11.

Throughout the remainder of the proof, let $C, C', c, c' > 0$ denote constants not depending on $\delta$ that may change from instance to instance. We compare the optimization landscape of $F(\alpha)$ with that of $G_{s^2}(\alpha)$ over $O$. Let $\mathrm{mse}(\alpha), \mathrm{mse}_*(\alpha), \omega(\alpha), \omega_*(\alpha)$ be as defined by (42) and (43) for the prior $g = g(\cdot, \alpha)$. We first bound $\mathrm{mse}(\alpha), \mathrm{mse}_*(\alpha), \omega(\alpha), \omega_*(\alpha)$: Write as shorthand $\langle \cdot \rangle = \langle \cdot \rangle_{g(\cdot, \alpha), \omega(\alpha)}$ for the posterior expectation in the scalar channel model (37). We have

$$\langle (y - \theta) \rangle^2 = \frac{1}{Z} \int (y - \theta)^2 e^{-\frac{\omega(\alpha)}{2}(y-\theta)^2} g(\theta, \alpha) \mathrm{d}\theta, \qquad Z = \int e^{-\frac{\omega(\alpha)}{2}(y-\theta)^2} g(\theta, \alpha) \mathrm{d}\theta.$$

We separate the integrals over the sets $\{\theta : e^{-\frac{\omega(\alpha)}{2}(y-\theta)^2} \leq Z\}$ and $\{\theta : e^{-\frac{\omega(\alpha)}{2}(y-\theta)^2} > Z\}$, and on the latter set apply the upper bound $(y - \theta)^2 \leq -\frac{2}{\omega(\alpha)} \log Z$. This gives

$$\langle (y - \theta) \rangle^2 \leq \int (y - \theta)^2 \mathbf{1}\{e^{-\frac{\omega(\alpha)}{2}(y-\theta)^2} \leq Z\} g(\theta, \alpha) \mathrm{d}\theta - \frac{2}{\omega(\alpha)} \log Z \leq 2 \int (y - \theta)^2 g(\theta, \alpha) \mathrm{d}\theta,$$

the last inequality applying Jensen's inequality to bound $\log Z \geq \int -\frac{\omega(\alpha)}{2}(y-\theta)^2 g(\theta, \alpha) \mathrm{d}\theta$. It is clear from the lower bounds of (214) and the boundedness of $\log g(0, \alpha)$ and $\partial_\theta \log g(0, \alpha)$ over $\alpha \in \overline{O}$ that $\int \theta^2 g(\theta, \alpha) \mathrm{d}\theta < C$ for some constant $C > 0$, for all $\alpha \in O$. Thus this inequality shows

$$\langle (y - \theta) \rangle^2 \leq C(1 + y^2),$$

which implies also

$$\langle (\langle \theta \rangle - \theta) \rangle^2 \leq \langle (y - \theta) \rangle^2 \leq C(1 + y^2), \quad \langle \theta \rangle^2 \leq 2y^2 + 2(y - \langle \theta \rangle)^2 \leq 2y^2 + 2\langle (y - \theta)^2 \rangle \leq C'(1 + y^2).$$

Taking expectations over $y = \theta^* + \omega_*(\alpha)^{-1/2} z$ with $\theta^* \sim g_*$ and $z \sim \mathcal{N}(0, 1)$, we get $\mathrm{mse}(\alpha), \mathrm{mse}_*(\alpha) \leq C(1 + \omega_*(\alpha)^{-1})$. Then applying $\omega_*(\alpha)^{-1} = (\sigma^2 + \mathrm{mse}_*(\alpha))/\delta \leq s^2 + C(1 + \omega_*(\alpha)^{-1})/\delta$, for all $\delta > \delta_0$ sufficiently large, this implies $\omega_*(\alpha)^{-1} \leq C'$. This in turn shows by $\mathrm{mse}(\alpha), \mathrm{mse}_*(\alpha) \leq C(1 + \omega_*(\alpha)^{-1})$ that

$$\mathrm{mse}(\alpha), \mathrm{mse}_*(\alpha) \leq C. \tag{215}$$

Let $o_\delta(1)$ denote a quantity that converges to 0 uniformly over $\alpha \in O$ as $\delta \to \infty$ (fixing $s^2 = \sigma^2/\delta$). Then, applying (215) to the fixed point equations $\omega(\alpha) = \delta/(\sigma^2 + \mathrm{mse}(\alpha))$ and $\omega_*(\alpha) = \delta/(\sigma^2 + \mathrm{mse}_*(\alpha))$, we have

$$\omega(\alpha)^{-1} = s^2 + o_\delta(1), \qquad \omega_*(\alpha)^{-1} = s^2 + o_\delta(1). \tag{216}$$

We recall from Lemma 4.12 that $\omega(\alpha), \omega_*(\alpha)$ must be continuous functions of $\alpha \in O$. We now argue via the implicit function theorem that for all $\delta > \delta_0$ sufficiently large, these are in fact continuously-differentiable over $\alpha \in O$. For this, fix any $\alpha \in O$ and consider the map

$$f_\alpha(\omega, \omega_*) = \begin{pmatrix} \omega^{-1} - \delta^{-1}(\sigma^2 + \mathbb{E}_{g_*, \omega_*}[\langle (\theta - \langle \theta \rangle_{g(\cdot, \alpha), \omega})^2 \rangle_{g(\cdot, \alpha), \omega}]) \\ \omega_*^{-1} - \delta^{-1}(\sigma^2 + \mathbb{E}_{g_*, \omega_*}[(\theta^* - \langle \theta \rangle_{g(\cdot, \alpha), \omega})^2]) \end{pmatrix}. \tag{217}$$

Thus (42) and (43) imply that $0 = f_\alpha(\omega(\alpha), \omega_*(\alpha))$. Let us momentarily write as shorthand $\mathbb{E} = \mathbb{E}_{g_*, \omega_*}$ and $\langle \cdot \rangle = \langle \cdot \rangle_{g(\cdot, \alpha), \omega}$. Expressing $y = \theta^* + \omega_*^{-1/2} z$, $\mathbb{E}$ may be understood as the expectation over $\theta^* \sim g_*$ and $z \sim \mathcal{N}(0, 1)$. The expected posterior average $\mathbb{E}\langle \cdot \rangle$ is given explicitly by

$$\mathbb{E}\langle f(\theta) \rangle = \mathbb{E}\frac{\int f(\theta) e^{H_\alpha(\theta, \omega, \omega_*)} \mathrm{d}\theta}{\int e^{H_\alpha(\theta, \omega, \omega_*)} \mathrm{d}\theta}, \qquad H_\alpha(\theta, \omega, \omega_*) = \omega(\theta^* + \omega_*^{-1/2} z)\theta - \frac{\omega}{2}\theta^2 + \log g(\theta, \alpha),$$

and the derivatives in $(\omega, \omega_*)$ may be computed via differentiation of $H_\alpha$. Let us denote by $\kappa_j(\cdot)$ the $j^{\mathrm{th}}$ mixed cumulant associated to the posterior mean $\langle \cdot \rangle = \langle \cdot \rangle_{g(\cdot, \alpha), \omega}$, i.e.

$$\kappa_1(f(\theta)) = \langle f(\theta) \rangle, \qquad \kappa_2(f(\theta), g(\theta)) = \langle f(\theta)g(\theta) \rangle - \langle f(\theta) \rangle \langle g(\theta) \rangle,$$

etc. Then $\mathbb{E}[\langle(\theta - \langle\theta\rangle)^2\rangle] = \mathbb{E}[\kappa_2(\theta, \theta)]$ and $\mathbb{E}[(\theta^* - \langle\theta\rangle)^2] = \mathbb{E}[(\theta^* - \kappa_1(\theta))^2]$, and differentiating in $(\omega, \omega_*)$ gives

$$\partial_\omega \mathbb{E}[\langle(\theta - \langle\theta\rangle)^2\rangle] = \mathbb{E}[\kappa_3(\theta, \theta, \partial_\omega H_\alpha(\theta, \omega, \omega_*))],$$

$$\partial_{\omega_*} \mathbb{E}[\langle(\theta - \langle\theta\rangle)^2\rangle] = \mathbb{E}[\kappa_3(\theta, \theta, \partial_{\omega_*} H_\alpha(\theta, \omega, \omega_*))],$$

$$\partial_\omega \mathbb{E}[(\theta^* - \langle\theta\rangle)^2] = \mathbb{E}[-2(\theta^* - \kappa_1(\theta))\kappa_2(\theta, \partial_\omega H_\alpha(\theta, \omega, \omega_*))]$$

$$\partial_{\omega_*} \mathbb{E}[(\theta^* - \langle\theta\rangle)^2] = \mathbb{E}[-2(\theta^* - \kappa_1(\theta))\kappa_2(\theta, \partial_{\omega_*} H_\alpha(\theta, \omega, \omega_*))]$$

We note that each absolute moment $\mathbb{E}_{g_*, \omega_*(\alpha)}[\langle|\theta|^k\rangle_{g(\cdot, \alpha), \omega(\alpha)}]$ is bounded by a constant over $\alpha \in O$, by continuity of this quantity in $\alpha$ and compactness of $\overline{O}$. Then it is direct to check that each of the above four derivatives evaluated at $(\omega, \omega_*) = (\omega(\alpha), \omega_*(\alpha))$ is also bounded by a constant over $\alpha \in O$. This implies that the derivative of the map $f_\alpha(\omega, \omega_*)$ in (217) satisfies

$$\mathrm{d}_{\omega, \omega_*} f_\alpha(\omega, \omega_*)\Big|_{(\omega, \omega_*)=(\omega(\alpha), \omega_*(\alpha))} = \begin{pmatrix} -\omega^{-2} + o_\delta(1) & o_\delta(1) \\ o_\delta(1) & -\omega_*^{-2} + o_\delta(1) \end{pmatrix}\Bigg|_{(\omega, \omega_*)=(\omega(\alpha), \omega_*(\alpha))} = -s^2\,\mathbf{I} + o_\delta(1),$$
(218)

where the last equality applies (216). In particular, for $\delta > \delta_0$ sufficiently large, this derivative is invertible. Since $f_\alpha(\omega, \omega_*)$ is continuously-differentiable in $(\omega, \omega_*, \alpha)$ (where differentiability in $\alpha$ is ensured by Assumption 2.2(b) for $\log g(\theta, \alpha)$), the implicit function theorem implies that for each $\alpha_0 \in O$, there exists a unique continuously-differentiable extension of the root $(\omega(\alpha_0), \omega_*(\alpha_0))$ of $0 = f_{\alpha_0}(\omega(\alpha_0), \omega_*(\alpha_0))$ to a solution of $0 = f_\alpha(\omega, \omega_*)$ in an open neighborhood of $\alpha_0$. This extension must then coincide with $\omega(\alpha), \omega_*(\alpha)$, because Lemma 4.12 ensures that $\omega(\alpha), \omega_*(\alpha)$ are continuous in $\alpha$. Thus $\omega(\alpha), \omega_*(\alpha)$ are continuously-differentiable in $\alpha \in O$, as claimed. The implicit function theorem shows also that their first derivatives are given by

$$\begin{pmatrix} \nabla_\alpha \omega^\top \\ \nabla_\alpha \omega_*^\top \end{pmatrix} = -[\mathrm{d}_{\omega, \omega_*} f_\alpha]^{-1} \mathrm{d}_\alpha f_\alpha\Big|_{(\omega, \omega_*)=(\omega(\alpha), \omega_*(\alpha))}.$$

We may check as above that the $\alpha$-derivatives

$$\partial_{\alpha_j} \mathbb{E}[\langle(\theta - \langle\theta\rangle)^2\rangle] = \mathbb{E}[\kappa_3(\theta, \theta, \partial_{\alpha_j} H_\alpha(\theta, \omega, \omega_*))],$$

$$\partial_{\alpha_j} \mathbb{E}[(\theta^* - \langle\theta\rangle)^2] = \mathbb{E}[-2(\theta^* - \kappa_1(\theta))\kappa_2(\theta, \partial_{\alpha_j} H_\alpha(\theta, \omega, \omega_*))]$$

evaluated at $(\omega, \omega_*) = (\omega(\alpha), \omega_*(\alpha))$ are also both bounded by a constant over $\alpha \in \overline{O}$. By the definition of $f_\alpha$, this implies $\mathrm{d}_\alpha f_\alpha|_{(\omega, \omega_*)=(\omega(\alpha), \omega_*(\alpha))} = o_\delta(1)$, so together with (218), this shows also

$$\nabla_\alpha \omega(\alpha) = o_\delta(1), \quad \nabla_\alpha \omega_*(\alpha) = o_\delta(1).$$
(219)

Recall from Lemma 2.12 that

$$\nabla F(\alpha) = -\mathbb{E}_{\theta \sim \mathsf{P}_\alpha} \nabla_\alpha \log g(\theta, \alpha) = -\mathbb{E}_{g_*, \omega_*(\alpha)} \langle \nabla_\alpha \log g(\theta, \alpha)\rangle_{g(\cdot, \alpha), \omega(\alpha)}.$$
(220)

Applying continuity of $(\omega, \omega_*) \mapsto \mathbb{E}_{g_*, \omega_*}\langle\nabla_\alpha \log g(\theta, \alpha)\rangle_{g(\cdot, \alpha), \omega}$ and the approximations $\omega(\alpha)^{-1}, \omega_*(\alpha)^{-1} = s^2 + o_\delta(1)$ shown above, we have

$$\nabla F(\alpha) = -\mathbb{E}_{g_*, s^{-2}}\langle\nabla_\alpha \log g(\theta, \alpha)\rangle_{g(\cdot, \alpha), s^{-2}} + o_\delta(1) = \nabla G_{s^2}(\alpha) + o_\delta(1),$$
(221)

where $G_{s^2}(\alpha) = -\mathbb{E}_{g_*, s^{-2}}[\log \mathsf{P}_{g(\cdot, \alpha), s^{-2}}(y)]$ is the negative population log-likelihood (56) in the scalar channel model with fixed noise variance $s^2$. [Note that fixing an arbitrary point $\alpha_0 \in O$ and integrating this gradient approximation over $\alpha \in O$, this also implies

$$F(\alpha) = G(\alpha) + (F(\alpha_0) - G(\alpha_0)) + o_\delta(1),$$

i.e. $F$ approximately coincides with $G$ up to an additive shift.] Furthermore, the above continuous-differentiability of $\omega(\alpha), \omega_*(\alpha)$ and (220) imply $F(\alpha)$ is twice continuously-differentiable over $\alpha \in O$, and differentiating

$\nabla F(\alpha)$ by the chain rule gives

$$\partial_{\alpha_i}\partial_{\alpha_j}F(\alpha) = -\partial_\omega\Big(\mathbb{E}_{g_*,\omega_*(\alpha)}\langle\partial_{\alpha_i}\log g(\theta,\alpha)\rangle_{g(\cdot,\alpha),\omega(\alpha)}\Big)\cdot\partial_{\alpha_j}\omega(\alpha)$$
$$-\partial_{\omega_*}\Big(\mathbb{E}_{g_*,\omega_*(\alpha)}\langle\partial_{\alpha_i}\log g(\theta,\alpha)\rangle_{g(\cdot,\alpha),\omega(\alpha)}\Big)\cdot\partial_{\alpha_j}\omega_*(\alpha)$$
$$-\partial_{\alpha_j}\Big(\mathbb{E}_{g_*,\omega_*(\alpha)}\langle\partial_{\alpha_i}\log g(\theta,\alpha)\rangle_{g(\cdot,\alpha),\omega(\alpha)}\Big). \tag{222}$$

Writing again $\mathbb{E} = \mathbb{E}_{g_*,\omega_*}$, $\langle\cdot\rangle = \langle\cdot\rangle_{g(\cdot,\alpha),\omega}$, and $\kappa_j$ for the cumulants with respect to $\langle\cdot\rangle$, we have

$$\partial_\omega\mathbb{E}\langle\partial_{\alpha_i}\log g(\theta,\alpha)\rangle = \mathbb{E}[\kappa_2(\partial_{\alpha_i}\log g(\theta,\alpha),\partial_\omega H_\alpha(\theta,\omega,\omega_*))]$$
$$\partial_{\omega_*}\mathbb{E}\langle\partial_{\alpha_i}\log g(\theta,\alpha)\rangle = \mathbb{E}[\kappa_2(\partial_{\alpha_i}\log g(\theta,\alpha),\partial_{\omega_*}H_\alpha(\theta,\omega,\omega_*))],$$

and these are bounded at $(\omega,\omega_*) = (\omega(\alpha),\omega_*(\alpha))$ over all $\alpha \in O$. Furthermore

$$\partial_{\alpha_j}\mathbb{E}\langle\partial_{\alpha_i}\log g(\theta,\alpha)\rangle = \mathbb{E}\langle\partial_{\alpha_i}\partial_{\alpha_j}\log g(\theta,\alpha)\rangle + \mathbb{E}[\kappa_2(\partial_{\alpha_i}\log g(\theta,\alpha),\partial_{\alpha_j}\log g(\theta,\alpha))].$$

Applying these and the bounds (219) to (222),

$$\partial_{\alpha_i}\partial_{\alpha_j}F(\alpha)$$
$$= -\mathbb{E}_{g_*,\omega_*(\alpha)}\langle\partial_{\alpha_i}\partial_{\alpha_j}\log g(\theta,\alpha)\rangle_{g(\cdot,\alpha),\omega(\alpha)} - \mathbb{E}_{g_*,\omega_*(\alpha)}\operatorname{Cov}_{\langle g(\cdot,\alpha),\omega(\alpha)\rangle}(\partial_{\alpha_i}\log g(\theta,\alpha),\partial_{\alpha_j}\log g(\theta,\alpha)) + o_\delta(1)$$
$$= \underbrace{-\mathbb{E}_{g_*,s^{-2}}\langle\partial_{\alpha_i}\partial_{\alpha_j}\log g(\theta,\alpha)\rangle_{g(\cdot,\alpha),s^{-2}} - \mathbb{E}_{g_*,s^{-2}}\operatorname{Cov}_{\langle g(\cdot,\alpha),s^{-2}\rangle}(\partial_{\alpha_i}\log g(\theta,\alpha),\partial_{\alpha_j}\log g(\theta,\alpha))}_{=\partial_{\alpha_i}\partial_{\alpha_j}G_{s^2}(\alpha)} + o_\delta(1)$$

Thus we have shown
$$\nabla^2 F(\alpha) = \nabla^2 G_{s^2}(\alpha) + o_\delta(1) \tag{223}$$

where again $o_\delta(1)$ converges to 0 uniformly over $\alpha \in O$ as $\delta \to \infty$.

The approximation (221) implies that $\nabla F + \nabla R$ converges uniformly to $\nabla G_{s^2} + \nabla R$ over $\alpha \in O$, as $\delta \to \infty$. Then for all $\delta > \delta_0$ sufficiently large and for some function $\iota : [\delta_0,\infty) \to (0,\infty)$ satisfying $\iota(\delta) \to 0$ as $\delta \to \infty$, each point of $\operatorname{Crit} \cap \mathcal{B}(D) = \{\alpha \in \mathcal{B}(D) : \nabla F(\alpha) = 0\}$ must fall within a ball of radius $\iota(\delta)$ around a point of $\operatorname{Crit}_G = \{\alpha \in \mathcal{B}(D) : \nabla G_{s^2}(\alpha) = 0\}$. The approximation (223) further implies that for each such ball around a point $\alpha_0 \in \operatorname{Crit}_G$, $\nabla^2 F$ converges uniformly to $\nabla^2 G_{s^2}$ on this ball, as $\delta \to \infty$. If $\nabla^2 G_{s^2}(\alpha_0)$ is non-singular, then for all $\delta > \delta_0$ sufficiently large, an argument via the topological degree shows that there must be exactly one point of $\operatorname{Crit}$ in this ball (having the same index as $\alpha_0$ as a critical point of $G_{s^2}$) — see e.g. [90, Lemma 5]. This shows statements (1) and (2) of the proposition.

As a direct consequence of these statements, if $\alpha^*$ is the unique point of $\operatorname{Crit}_G$ and $\nabla^2 G_{s^2}(\alpha^*)$ is non-singular, then there is a unique point of $\operatorname{Crit} \cap \mathcal{B}(D)$. If furthermore $g_*(\theta) = g(\theta,\alpha^*)$, then this point of $\operatorname{Crit} \cap \mathcal{B}(D)$ must be $\alpha^*$ itself, since $\nabla F(\alpha^*) = 0$ by Lemma 2.12(c). $\qquad\square$

*Analysis of Example 2.17.* We verify Assumption 2.2(b) for Example 2.17 of the Gaussian mixture model with varying means. Let $\iota \in \{1,\ldots,K\}$ denote the mixture component of $\theta$, and let $\langle f(\iota,\theta)\rangle = \mathbb{E}[f(\iota,\theta)\,|\,\theta]$ denote the posterior average over $\iota$ given $\theta \sim \mathcal{N}(\alpha_\iota,\omega_0^{-1})$ and prior $\mathbb{P}[\iota = k] = p_k$. Let $\kappa_2(\cdot)$ denote the covariance associated to $\langle\cdot\rangle$. Then, since

$$\log g(\theta,\alpha) = \log\sum_{k=1}^K p_k\sqrt{\frac{\omega_k}{2\pi}}\exp\Big(-\frac{\omega_k}{2}(\theta - \alpha_k)^2\Big),$$

the derivatives of $\log g(\theta,\alpha)$ up to order 2 are given by

$$\partial_\theta\log g(\theta,\alpha) = \langle\omega_\iota(\alpha_\iota - \theta)\rangle, \quad \partial_\theta^2\log g(\theta,\alpha) = \kappa_2\big(\omega_\iota(\alpha_\iota - \theta),\omega_\iota(\alpha_\iota - \theta)\big) - \langle\omega_\iota\rangle$$
$$\partial_{\alpha_i}\log g(\theta,\alpha) = \omega_i(\theta - \alpha_i)\langle\mathbf{1}_{\iota=i}\rangle, \quad \partial_{\alpha_i}\partial_\theta\log g(\theta,\alpha) = \omega_i(\theta - \alpha_i)\kappa_2\big(\mathbf{1}_{\iota=i},\omega_\iota(\alpha_\iota - \theta)\big) + \omega_i\langle\mathbf{1}_{\iota=i}\rangle, \tag{224}$$
$$\partial_{\alpha_i}\partial_{\alpha_j}\log g(\theta,\alpha) = \omega_i\omega_j(\theta - \alpha_i)(\theta - \alpha_j)\kappa_2(\mathbf{1}_{\iota=i},\mathbf{1}_{\iota=j}) - \mathbf{1}_{i=j}\omega_i\langle\mathbf{1}_{\iota=i}\rangle$$

67

In particular, $|\partial_\theta \log g(\theta, \alpha)| \leq C(1 + |\theta| + |\alpha_i|)$ and $|\partial_{\alpha_i} \log g(\theta, \alpha)| \leq C(1 + |\theta| + |\alpha_i|)$, showing (19).

To bound the high-order derivatives of $\log g(\theta, \alpha)$ locally over $\alpha \in \mathbb{R}^K$, let $k_{\max} \in \{1, \ldots, K\}$ be the (unique) index corresponding to the smallest value of $\omega_k$. For any compact subset $S \subset \mathbb{R}^K$, there exist constants $B(S), c_0(S) > 0$ depending on the fixed values $\{p_1, \ldots, p_K\}$, $\{\omega_1, \ldots, \omega_K\}$ and $S$ such that for all $\alpha \in S$, we have

$$\frac{\omega_k}{2}(\theta - \alpha_k)^2 \geq \frac{\omega_{k_{\max}}}{2}(\theta - \alpha_{k_{\max}})^2 + c_0(S)\theta^2 \text{ for any } \theta > B(S) \text{ and all } k \neq k_{\max}.$$

This implies there exists a constant $C(S) > 0$ for which

$$\langle \mathbf{1}_{\iota \neq k_{\max}} \rangle \leq C(S)e^{-c_0(S)\theta^2} \text{ for all } \theta > B(S) \text{ and } \alpha \in S.$$

Let $\iota'$ denote an independent copy of $\iota$ under its posterior law given $\theta$. Then for any $\theta > B$, any $\alpha \in S$, and any $k \in \{1, \ldots, K\}$, the posterior variance of $\mathbf{1}_{\iota=k}$ is bounded as

$$\langle |\mathbf{1}_{\iota=k} - \langle \mathbf{1}_{\iota=k} \rangle|^2 \rangle \leq \langle |\mathbf{1}_{\iota=k} - \mathbf{1}_{\iota'=k}|^2 \rangle \leq 4\langle \mathbf{1}_{\iota \neq k_{\max} \text{ or } \iota' \neq k_{\max}} \rangle \leq C'(S)e^{-c_0(S)\theta^2},$$

and similarly

$$\langle |\omega_\iota(\theta - \alpha_\iota) - \langle \omega_\iota(\theta - \alpha_\iota) \rangle|^2 \rangle \leq C'(S)(1 + \theta^2)e^{-c_0(S)\theta^2}.$$

Applying these bounds and Hölder's inequality, all posterior covariances in (224) are exponentially small in $\theta^2$ for $\theta > B(S)$, implying that all derivatives of order 2 in (224) are bounded over $\alpha \in S$ and $\theta > B(S)$. Similarly they are bounded over $\alpha \in S$ and $\theta < -B(S)$, and hence also bounded uniformly over $\alpha \in S$ and $\theta \in \mathbb{R}$ since we may bound the cumulants trivially by a constant $C(S)$ for $\theta \in [-B(S), B(S)]$. The same argument bounds all mixed cumulants of $\mathbf{1}_{\iota=k}$ and $\omega_\iota(\alpha_\iota - \theta)$ of orders 3 and 4, and hence also all partial derivatives of $\log g(\theta, \alpha)$ of orders 3 and 4 over $\alpha \in S$ and $\theta \in \mathbb{R}$. These arguments show also that as $\theta \to \pm\infty$, uniformly over $\alpha \in S$, $\kappa_2(\omega_\iota(\alpha_\iota - \theta), \omega_\iota(\alpha_\iota - \theta)) \to 0$ and $\langle \omega_\iota \rangle \to \omega_{k_{\max}}$, so $\partial_\theta^2[-\log g(\theta, \alpha)] \to \omega_{k_{\max}} > 0$, verifying all statements of Assumption 2.2(b). $\square$

*Analysis of Example 2.18.* We verify Assumption 2.2(b) in Example 2.18 for the Gaussian mixture model with fixed mixture means/variances and varying weights. Again let $\iota \in \{0, \ldots, K\}$ denote the mixture component of $\theta$, and let $\langle f(\iota, \theta) \rangle = \mathbb{E}[f(\iota, \theta) \mid \theta]$ denote the posterior average over $\iota$ given $\theta \sim \mathcal{N}(\mu_\iota, \omega_\iota^{-1})$ and prior $\mathbb{P}[\iota = k] = e^{\alpha_k}/(e^{\alpha_0} + \ldots + e^{\alpha_K})$. Let $\kappa_2(\cdot)$ denote the covariance associated to $\langle \cdot \rangle$, and in addition, let $\langle \cdot \rangle_{\text{prior}}$ and $\kappa_2^{\text{prior}}$ denote the mean and covariance over $\iota$ drawn from the prior $\mathbb{P}[\iota = k] = e^{\alpha_k}/(e^{\alpha_0} + \ldots + e^{\alpha_K})$. Then, since

$$\log g(\theta, \alpha) = \log \sum_{k=0}^{K} e^{\alpha_k} \sqrt{\frac{\omega_k}{2\pi}} \exp\left(-\frac{\omega_k}{2}(\theta - \mu_k)^2\right) - \log \sum_{k=0}^{K} e^{\alpha_k},$$

the derivatives of $\log g(\theta, \alpha)$ up to order 2 are given by

$$\partial_{\alpha_i} \log g(\theta, \alpha) = \langle \mathbf{1}_{\iota=i} \rangle - \langle \mathbf{1}_{\iota=i} \rangle_{\text{prior}}, \quad \partial_{\alpha_i}\partial_{\alpha_j} \log g(\theta, \alpha) = \kappa_2(\mathbf{1}_{\iota=i}, \mathbf{1}_{\iota=j}) - \kappa_2^{\text{prior}}(\mathbf{1}_{\iota=i}, \mathbf{1}_{\iota=j}),$$
$$\partial_\theta \log g(\theta, \alpha) = \langle \omega_\iota(\mu_\iota - \theta) \rangle, \quad \partial_{\alpha_i}\partial_\theta \log g(\theta, \alpha) = \kappa_2(\mathbf{1}_{\iota=i}, \omega_\iota(\mu_\iota - \theta)), \tag{225}$$
$$\partial_\theta^2 \log g(\theta, \alpha) = \kappa_2(\omega_\iota(\mu_\iota - \theta), \omega_\iota(\mu_\iota - \theta)) - \langle \omega_\iota \rangle.$$

In particular, this shows $\sum_k \partial_{\alpha_k} \log g(\theta, \alpha) = 1 - 1 = 0$, so $\nabla_\alpha \log g(\theta, \alpha)$ always belongs to the subspace $E = \{\alpha \in \mathbb{R}^{K+1} : \alpha_0 + \ldots + \alpha_K = 0\}$. Also $\nabla R(\alpha) = r'(\|\alpha\|_2) \cdot \frac{\alpha}{\|\alpha\|_2} \in E$ if $\alpha \in E$. Furthermore, $|\partial_{\alpha_i} \log g(\theta, \alpha)| \leq C$ and $|\partial_\theta \log g(\theta, \alpha)| \leq C(1 + |\theta|)$, showing (19).

To bound the higher-order derivatives of $\log g(\theta, \alpha)$ locally over $\alpha \in E$, let $k_{\max} \in \{0, \ldots, K\}$ be the index corresponding to the smallest $\omega_k$, and among these the largest $\mu_k$ (if there are multiple $\omega_k$'s equal to the smallest value). Then for some constants $B, c_0 > 0$ depending only on the fixed values $\{\mu_0, \ldots, \mu_K\}$ and $\{\omega_0, \ldots, \omega_K\}$, we have

$$\frac{\omega_k}{2}(\theta - \mu_k)^2 \geq \frac{\omega_{k_{\max}}}{2}(\theta - \mu_{k_{\max}})^2 + c_0\theta \text{ for any } \theta > B \text{ and all } k \neq k_{\max}.$$

This implies, for any compact subset $S \subset E$, there is a constant $C(S) > 0$ for which

$$\langle \mathbf{1}_{\iota \neq k_{\max}} \rangle \leq C(S) e^{-c_0 \theta} \text{ for all } \theta > B \text{ and } \alpha \in S.$$

Then the same arguments as in the preceding example show

$$\langle |\omega_\iota - \langle \omega_\iota \rangle|^2 \rangle \leq C'(S) e^{-c_0 \theta}, \quad \langle |\omega_\iota(\mu_\iota - \theta) - \langle \omega_\iota(\mu_\iota - \theta) \rangle|^2 \rangle \leq C'(S)(1 + \theta)^2 e^{-c_0 \theta},$$

implying via Cauchy-Schwarz that each order-2 derivative in (225) is bounded over $\alpha \in S$ and $\theta > B$. Similarly it is bounded over $\alpha \in S$ and $\theta < -B$, hence also for all $\alpha \in S$ and $\theta \in \mathbb{R}$. The same argument applies to bound the mixed cumulants of $\omega_\iota$ and $\omega_\iota(\mu_\iota - \theta)$ of orders 3 and 4, and thus the partial derivatives of $\log g(\theta, \alpha)$ of orders 3 and 4. This shows also $\lim_{\theta \to \infty} \partial_\theta^2[-\log g(\theta, \alpha)] = \omega_{k_{\max}} > 0$ uniformly over $\alpha \in S$, and a similar statement holds for $\theta \to -\infty$, establishing all conditions of Assumption 2.2(b). $\qquad \square$

# A   Proof of Theorem 2.3

Theorem 2.3(a) follows immediately from [27, Theorem 2.5], upon identifying $s(\theta, \alpha)$ of [27, Theorem 2.5] as $(\log g)'(\theta)$ (with no dependence on $\alpha$) and $\mathcal{G}(\alpha, \mathsf{P}) = 0$. The required conditions of [27, Assumption 2.2] for $s(\cdot)$ hold by Assumption 2.2(a), and the conditions of [27, Assumption 2.3] for $\mathcal{G}(\cdot)$ are vacuous.

For Theorem 2.3(b), consider first the following global version of Assumption 2.2(b):

**Assumption A.1.** $\log g(\theta, \alpha)$ and $R(\alpha)$ are thrice continuously-differentiable and satisfy (19), and the conditions (20) hold for constants $C, r_0, c_0 > 0$ globally over all $\alpha \in \mathbb{R}^K$.

Under Assumption A.1, Theorem 2.3(b) again follows from [27, Theorem 2.5] upon identifying $s(\theta, \alpha) = \partial_\theta \log g(\theta, \alpha)$ and $\mathcal{G}(\alpha, \mathsf{P}) = \mathbb{E}_{\theta \sim \mathsf{P}}[\nabla_\alpha \log g(\theta, \alpha)] - \nabla R(\alpha)$, where all conditions of [27, Assumptions 2.2 and 2.3] may be checked from these conditions of Assumption A.1.

To show Theorem 2.3(b) under the weaker local conditions of Assumption 2.2(b), we may apply the following truncation argument: Note first that twice continuous-differentiability of $\log g(\theta, \alpha)$ and $R(\alpha)$ imply that $\nabla_{(\theta, \alpha)} \log g(\theta, \alpha)$ and $\nabla R(\alpha)$ are locally Lipschitz. Together with the global linear growth conditions of (19), this implies that there exists a unique (non-explosive) solution $\{(\boldsymbol{\theta}^t, \widehat{\alpha}^t)\}_{t \geq 0}$ to the joint diffusion (9–10) for all times (c.f. [91, Theorem 12.1]). Furthermore, since

$$\boldsymbol{\theta}^t = \boldsymbol{\theta}^0 + \int_0^t \left( -\frac{1}{2\sigma^2} \mathbf{X}^\top (\mathbf{X}\boldsymbol{\theta}^s - \mathbf{y}) + \left( \partial_\theta \log g(\theta_j^s, \widehat{\alpha}^s) \right)_{j=1}^d \right) \mathrm{d}s + \sqrt{2}\,\mathbf{b}^t$$

$$\widehat{\alpha}^t = \widehat{\alpha}^0 + \int_0^t \left( \frac{1}{d} \sum_{j=1}^d \nabla_\alpha \log g(\theta_j^s, \widehat{\alpha}^s) - \nabla_\alpha R(\widehat{\alpha}^s) \right) \mathrm{d}s,$$

under the growth conditions (19), this solution satisfies the bounds

$$\|\boldsymbol{\theta}^t\|_2 \leq \|\boldsymbol{\theta}^0\|_2 + C \int_0^t \left( \|\mathbf{X}\|_{\mathrm{op}}^2 \|\boldsymbol{\theta}^s\|_2 + \|\mathbf{X}\|_{\mathrm{op}} \|\mathbf{y}\|_2 + \sqrt{d} + \|\boldsymbol{\theta}^s\|_2 + \sqrt{d}\|\widehat{\alpha}^s\|_2 \right) \mathrm{d}s + \sqrt{2}\,\|\mathbf{b}^t\|_2$$

$$\|\widehat{\alpha}^t\|_2 \leq \|\widehat{\alpha}^0\|_2 + C \int_0^t \left( 1 + \|\widehat{\alpha}^s\|_2 + \|\boldsymbol{\theta}^s\|_2/\sqrt{d} \right) \mathrm{d}s$$

Fixing any $T > 0$, by the conditions of Assumption 2.2 and a standard maximal inequality for Brownian motion (see [27, Lemma 4.7]) there exists a constant $C_0 > 0$ large enough such that the event

$$\mathcal{E} = \left\{ \|\mathbf{X}\|_{\mathrm{op}} \leq C_0,\ \|\mathbf{y}\|_2 \leq C_0\sqrt{d},\ \|\boldsymbol{\theta}^0\|_2 \leq C_0\sqrt{d},\ \|\widehat{\alpha}^0\|_2 \leq C_0,\ \sup_{t \in [0,T]} \|\mathbf{b}^t\|_2 \leq C_0\sqrt{d} \right\}$$

holds a.s. for all large $n, d$. Then by a Gronwall argument, for a constant $M = M(T, C_0) > 0$,

$$\sup_{t \in [0,T]} \frac{\|\boldsymbol{\theta}^t\|_2}{\sqrt{d}} + \|\widehat{\alpha}^t\|_2 < M$$

holds on $\mathcal{E}$. Applying the conditions of Assumption 2.2(b) with $S = \{\alpha : \|\alpha\|_2 \le M\}$, there exist functions $g_M : \mathbb{R} \times \mathbb{R}^K \to \mathbb{R}$ and $R_M : \mathbb{R}^K \to \mathbb{R}$ such that $g_M(\theta, \alpha) = g(\theta, \alpha)$ and $R_M(\alpha) = R(\alpha)$ for all $\|\alpha\| \le M$, and $g_M$ and $R_M$ satisfy Assumption A.1. Let $\{(\boldsymbol{\theta}_M^t, \widehat{\alpha}_M^t)\}_{t \ge 0}$ be the solution of (9–10) defined with $g_M(\cdot)$ and $R_M(\cdot)$ in place of $g(\cdot)$ and $R(\cdot)$, and let $\boldsymbol{\eta}_M^t = \mathbf{X}\boldsymbol{\theta}_M^t$. Then as argued above, Theorem 2.3(b) holds for $\{(\boldsymbol{\theta}_M^t, \boldsymbol{\eta}_M^t, \widehat{\alpha}_M^t)\}_{t \ge 0}$, showing that a.s. as $n, d \to \infty$,

$$\frac{1}{d} \sum_{j=1}^d \delta_{\theta_j^*, \{\theta_{M,j}^t\}_{t \in [0,T]}} \overset{W_2}{\to} \mathsf{P}(\theta^*, \{\theta_M^t\}_{t \in [0,T]})$$

$$\frac{1}{n} \sum_{i=1}^n \delta_{\eta_i^*, \varepsilon_i, \{\eta_{M,i}^t\}_{t \in [0,T]}} \overset{W_2}{\to} \mathsf{P}(\eta^*, \varepsilon, \{\eta_M^t\}_{t \in [0,T]}) \tag{226}$$

$$\{\widehat{\alpha}^t\}_{t \in [0,T]} \to \{\alpha_M^t\}_{t \in [0,T]}$$

for limiting processes defined by the DMFT equations (22–28) also with $g_M(\cdot)$ and $R_M(\cdot)$ in place of $g(\cdot)$ and $R(\cdot)$. Since $\{(\boldsymbol{\theta}^t, \boldsymbol{\eta}^t, \widehat{\alpha}^t)\}_{t \in [0,T]} = \{(\boldsymbol{\theta}_M^t, \boldsymbol{\eta}_M^t, \widehat{\alpha}_M^t)\}_{t \in [0,T]}$ a.s. for all large $n, d$, this implies that (226) holds also with $\{(\boldsymbol{\theta}^t, \boldsymbol{\eta}^t, \widehat{\alpha}^t)\}_{t \in [0,T]}$ in place of $\{(\boldsymbol{\theta}_M^t, \boldsymbol{\eta}_M^t, \widehat{\alpha}_M^t)\}_{t \in [0,T]}$. Furthermore, the deterministic limit process $\{\alpha_M^t\}_{t \in [0,T]}$ must satisfy $\|\alpha_M^t\| \le M$ for all $t \in [0, T]$, so the solution up to time $T$ of the DMFT equations (22–28) with $g_M(\cdot)$ and $R_M(\cdot)$ is also a solution of these equations with $g(\cdot)$ and $R(\cdot)$. This proves Theorem 2.3(b) under Assumption 2.2(b).

# B  Correlation and response functions for a Gaussian prior

For illustration, we check Definition 2.4 explicitly for the dynamics (7) with a Gaussian prior

$$g(\theta) = \sqrt{\frac{\lambda}{2\pi}} \exp\left(-\frac{\lambda\theta^2}{2}\right).$$

Then (7) is the Ornstein-Uhlenbeck process

$$\mathrm{d}\boldsymbol{\theta}^t = \left[-\left(\frac{\mathbf{X}^\top \mathbf{X}}{\sigma^2} + \lambda\mathbf{I}\right)\boldsymbol{\theta}^t + \frac{\mathbf{X}^\top \mathbf{y}}{\sigma^2}\right]\mathrm{d}t + \sqrt{2}\,\mathrm{d}\mathbf{b}^t. \tag{227}$$

**Lemma B.1.** *Under Assumption 2.1, let*

$$\mu = \lim_{n,d \to \infty} \frac{1}{d} \sum_{i=1}^d \delta_{\lambda_i(\mathbf{X}^\top \mathbf{X}/\sigma^2)}$$

*be the almost-sure limit of the empirical eigenvalue distribution of $\mathbf{X}^\top \mathbf{X}/\sigma^2$. Then for the dynamics (227) with a fixed Gaussian prior, the corresponding DMFT system prescribed by Theorem 2.3(a) has the correlation and response functions*

$$C_\theta(t, s) = \int \left[\mathbb{E}(\theta^0)^2 \cdot e^{-(\lambda+x)(t+s)} + \frac{\mathbb{E}(\theta^*)^2 x^2 + x}{(\lambda+x)^2}(1 - e^{-(\lambda+x)t})(1 - e^{-(\lambda+x)s})\right.$$

$$\left. + \frac{1}{\lambda+x}\left(e^{-(\lambda+x)|t-s|} - e^{-(\lambda+x)(t+s)}\right)\right]\mu(\mathrm{d}x)$$

$$C_\theta(t, *) = \int \frac{\mathbb{E}(\theta^*)^2 x}{\lambda+x}(1 - e^{-(\lambda+x)t})\mu(\mathrm{d}x)$$

$$R_\theta(t, s) = \int e^{-(\lambda+x)(t-s)}\mu(\mathrm{d}x)$$

$$C_\eta(t, s) = \int \left[\mathbb{E}(\theta^0)^2 x e^{-(\lambda+x)(t+s)} + (\mathbb{E}(\theta^*)^2 x + 1)\left(\frac{x}{\lambda+x}(1 - e^{-(\lambda+x)t}) - 1\right)\left(\frac{x}{\lambda+x}(1 - e^{-(\lambda+x)s}) - 1\right)\right.$$

$$\left. + (\delta - 1) + \frac{x}{\lambda+x}\left(e^{-(\lambda+x)|t-s|} - e^{-(\lambda+x)(t+s)}\right)\right]\mu(\mathrm{d}x),$$

$$R_\eta(t, s) = \int x e^{-(\lambda+x)(t-s)}\mu(\mathrm{d}x).$$

*Proof.* Setting $\mathbf{\Lambda} = \frac{\mathbf{X}^\top \mathbf{X}}{\sigma^2} + \lambda \mathbf{I}$, the dynamics (227) have the explicit solution

$$\boldsymbol{\theta}^t = e^{-\mathbf{\Lambda}t}\boldsymbol{\theta}^0 + \mathbf{\Lambda}^{-1}\left(\mathbf{I} - e^{-\mathbf{\Lambda}t}\right)\frac{\mathbf{X}^\top \mathbf{y}}{\sigma^2} + \int_0^t e^{-\mathbf{\Lambda}(t-s)}\sqrt{2}\,\mathrm{d}\mathbf{b}^s$$

$$= e^{-\mathbf{\Lambda}t}\boldsymbol{\theta}^0 + \mathbf{\Lambda}^{-1}\left(\mathbf{I} - e^{-\mathbf{\Lambda}t}\right)\left(\frac{\mathbf{X}^\top \mathbf{X}}{\sigma^2}\boldsymbol{\theta}^* + \frac{\mathbf{X}^\top \boldsymbol{\varepsilon}}{\sigma^2}\right) + \int_0^t e^{-\mathbf{\Lambda}(t-s)}\sqrt{2}\,\mathrm{d}\mathbf{b}^s \tag{228}$$

Recall the definitions of $e_j(\boldsymbol{\theta})$ and $x_i(\boldsymbol{\theta})$ from (100) and the associated correlation and response matrices (101). Under Assumption 2.1, applying the explicit form (228) and independence of $\mathbf{X}, \boldsymbol{\theta}^0, \boldsymbol{\theta}^*, \boldsymbol{\varepsilon}$, it is direct to check that almost surely,

$$\lim_{n,d\to\infty} \frac{1}{d}\operatorname{Tr}\mathbf{C}_\theta(t,s) = \lim_{n,d\to\infty}\frac{1}{d}\langle \boldsymbol{\theta}^{t^\top}\boldsymbol{\theta}^s\rangle$$

$$= \lim_{n,d\to\infty}\frac{1}{d}\operatorname{Tr}\left(\mathbb{E}(\theta^0)^2 \cdot e^{-\mathbf{\Lambda}(t+s)} + \mathbb{E}(\theta^*)^2 \cdot \left(\frac{\mathbf{X}^\top \mathbf{X}}{\sigma^2}\right)(\mathbf{I}-e^{-\mathbf{\Lambda}t})\mathbf{\Lambda}^{-2}(\mathbf{I}-e^{-\mathbf{\Lambda}s})\left(\frac{\mathbf{X}^\top \mathbf{X}}{\sigma^2}\right)\right.$$

$$\left. + \int_0^{t\wedge s} 2e^{-\mathbf{\Lambda}(t+s-2r)}\mathrm{d}r\right) + \frac{1}{d}\operatorname{Tr}\left(\mathbb{E}\varepsilon^2 \cdot \frac{\mathbf{X}}{\sigma^2}(\mathbf{I}-e^{-\mathbf{\Lambda}t})\mathbf{\Lambda}^{-2}(\mathbf{I}-e^{-\mathbf{\Lambda}s})\frac{\mathbf{X}^\top}{\sigma^2}\right)$$

$$= \int\left[\mathbb{E}(\theta^0)^2 \cdot e^{-(\lambda+x)(t+s)} + \frac{\mathbb{E}(\theta^*)^2 x^2 + x}{(\lambda+x)^2}(1-e^{-(\lambda+x)t})(1-e^{-(\lambda+x)s})\right.$$

$$\left. + \frac{1}{\lambda+x}\left(e^{-(\lambda+x)|t-s|} - e^{-(\lambda+x)(t+s)}\right)\right]\mu(\mathrm{d}x)$$

and

$$\lim_{n,d\to\infty}\frac{1}{d}\operatorname{Tr}\mathbf{C}_\theta(t,*) = \lim_{n,d\to\infty}\frac{1}{d}\langle\boldsymbol{\theta}^{t^\top}\boldsymbol{\theta}^*\rangle = \lim_{n,d\to\infty}\frac{1}{d}\operatorname{Tr}\left(\mathbb{E}(\theta^*)^2\mathbf{\Lambda}^{-1}(\mathbf{I}-e^{-\mathbf{\Lambda}t})\frac{\mathbf{X}^\top \mathbf{X}}{\sigma^2}\right)$$

$$= \int\frac{\mathbb{E}(\theta^*)^2 x}{\lambda+x}(1-e^{-(\lambda+x)t})\mu(\mathrm{d}x).$$

Furthermore, the above form (228) for $\boldsymbol{\theta}^t$ implies

$$P_t(\boldsymbol{\theta}) = e^{-\mathbf{\Lambda}t}\boldsymbol{\theta} + +\mathbf{\Lambda}^{-1}\left(\mathbf{I}-e^{-\mathbf{\Lambda}t}\right)\left(\frac{\mathbf{X}^\top \mathbf{X}}{\sigma^2}\boldsymbol{\theta}^* + \frac{\mathbf{X}^\top \boldsymbol{\varepsilon}}{\sigma^2}\right). \tag{229}$$

Then $\nabla P_t e_j(\boldsymbol{\theta}) = \nabla[\mathbf{e}_j^\top P_t(\boldsymbol{\theta})] = e^{-\mathbf{\Lambda}t}\mathbf{e}_j$ is a constant function not depending on $\boldsymbol{\theta}$, and

$$\lim_{n,d\to\infty}\frac{1}{d}\operatorname{Tr}\mathbf{R}_\theta(t,s) = \lim_{n,d\to\infty}\frac{1}{d}\sum_{j=1}^d[\nabla e_j^\top \nabla P_{t-s}e_j](\boldsymbol{\theta}^s) = \lim_{n,d\to\infty}\frac{1}{d}\operatorname{Tr}e^{-\mathbf{\Lambda}(t-s)} = \int e^{-(\lambda+x)(t-s)}\mu(\mathrm{d}x).$$

By Theorem 4.3, this shows the forms of $C_\theta(t,s)$, $C_\theta(t,*)$, and $R_\theta(t,s)$.

From (228) and (229), we have also

$$\mathbf{X}\boldsymbol{\theta}^t - \mathbf{y} = \mathbf{X}e^{-\mathbf{\Lambda}t}\boldsymbol{\theta}^0 + \mathbf{X}\left(\mathbf{\Lambda}^{-1}\left(\mathbf{I}-e^{-\mathbf{\Lambda}t}\right)\frac{\mathbf{X}^\top \mathbf{X}}{\sigma^2} - \mathbf{I}\right)\boldsymbol{\theta}^* + \left(\mathbf{X}\mathbf{\Lambda}^{-1}\left(\mathbf{I}-e^{-\mathbf{\Lambda}t}\right)\frac{\mathbf{X}^\top}{\sigma^2} - \mathbf{I}_n\right)\boldsymbol{\varepsilon}$$

$$+ \int_0^t \mathbf{X}e^{-\mathbf{\Lambda}(t-s)}\sqrt{2}\,\mathrm{d}\mathbf{b}^s$$

and $\nabla P_t x_i(\boldsymbol{\theta}) = (\sqrt{\delta}/\sigma)\nabla[\mathbf{e}_i^\top \mathbf{X}^\top P_t(\boldsymbol{\theta})] = (\sqrt{\delta}/\sigma)e^{-\boldsymbol{\Lambda} t}\mathbf{X}^\top \mathbf{e}_i$. Then

$$\lim_{n,d\to\infty} \frac{1}{n}\operatorname{Tr}\mathbf{C}_\eta(t,s) = \lim_{n,d\to\infty} \frac{1}{n}\cdot\frac{\delta}{\sigma^2}(\mathbf{X}\boldsymbol{\theta}^t - \mathbf{y})^\top(\mathbf{X}\boldsymbol{\theta}^s - \mathbf{y})$$

$$= \lim_{n,d\to\infty} \frac{1}{d\sigma^2}\operatorname{Tr}\left(\mathbb{E}(\theta^0)^2 \cdot e^{-\boldsymbol{\Lambda} t}\mathbf{X}^\top\mathbf{X}e^{-\boldsymbol{\Lambda} s}\right.$$

$$+ \mathbb{E}(\theta^*)^2 \cdot \left(\boldsymbol{\Lambda}^{-1}\left(\mathbf{I} - e^{-\boldsymbol{\Lambda} t}\right)\frac{\mathbf{X}^\top\mathbf{X}}{\sigma^2} - \mathbf{I}\right)\mathbf{X}^\top\mathbf{X}\left(\boldsymbol{\Lambda}^{-1}\left(\mathbf{I} - e^{-\boldsymbol{\Lambda} s}\right)\frac{\mathbf{X}^\top\mathbf{X}}{\sigma^2} - \mathbf{I}\right)$$

$$+ \left.\int_0^{t\wedge s} 2\,e^{-\boldsymbol{\Lambda}(t-r)}\mathbf{X}^\top\mathbf{X}e^{-\boldsymbol{\Lambda}(s-r)}\mathrm{d}r\right)$$

$$+ \frac{1}{d\sigma^2}\operatorname{Tr}\left(\mathbb{E}\varepsilon^2 \cdot \left(\mathbf{X}\boldsymbol{\Lambda}^{-1}\left(\mathbf{I} - e^{-\boldsymbol{\Lambda} t}\right)\frac{\mathbf{X}^\top}{\sigma^2} - \mathbf{I}_n\right)^\top\left(\mathbf{X}\boldsymbol{\Lambda}^{-1}\left(\mathbf{I} - e^{-\boldsymbol{\Lambda} s}\right)\frac{\mathbf{X}^\top}{\sigma^2} - \mathbf{I}_n\right)\right)$$

$$= \int\left[\mathbb{E}(\theta^0)^2 x e^{-(\lambda+x)(t+s)} + (\mathbb{E}(\theta^*)^2 x + 1)\left(\frac{x}{\lambda+x}(1 - e^{-(\lambda+x)t}) - 1\right)\left(\frac{x}{\lambda+x}(1 - e^{-(\lambda+x)s}) - 1\right)\right.$$

$$\left.+ (\delta - 1) + \frac{x}{\lambda+x}\left(e^{-(\lambda+x)|t-s|} - e^{-(\lambda+x)(t+s)}\right)\right]\mu(\mathrm{d}x),$$

and

$$\lim_{n,d\to\infty} \frac{1}{n}\operatorname{Tr}\mathbf{R}_\eta(t,s) = \lim_{n,d\to\infty} \frac{1}{n}\sum_{i=1}^n[\nabla x_i^\top \nabla P_{t-s}x_i](\boldsymbol{\theta}^s) = \lim_{n,d\to\infty} \frac{1}{n}\cdot\frac{\delta}{\sigma^2}\sum_{i=1}^n\mathbf{e}_i^\top\mathbf{X}e^{-\boldsymbol{\Lambda}(t-s)}\mathbf{X}^\top\mathbf{e}_i$$

$$= \lim_{n,d\to\infty} \frac{1}{d\sigma^2}\operatorname{Tr}\mathbf{X}e^{-\boldsymbol{\Lambda}(t-s)}\mathbf{X}^\top = \int x e^{-(\lambda+x)(t-s)}\mu(\mathrm{d}x).$$

By Theorem 4.3, this shows the forms of $C_\eta(t,s)$ and $R_\eta(t,s)$. $\qquad\square$

From Lemma B.1 it is apparent that the approximations (30–32) and (34–35) hold with $\varepsilon(t) = Ce^{-ct}$ and

$$c_\theta^{\mathrm{init}}(s) = -\int\frac{\mathbb{E}(\theta^*)^2 x^2 + x}{(\lambda+x)^2}e^{-(\lambda+x)s}\mu(\mathrm{d}x), \qquad c_\theta^{\mathrm{tti}}(\infty) = \int\frac{\mathbb{E}(\theta^*)^2 x^2 + x}{(\lambda+x)^2}\mu(\mathrm{d}x)$$

$$c_\theta^{\mathrm{tti}}(\tau) = c_\theta^{\mathrm{tti}}(\infty) + \int\frac{1}{\lambda+x}e^{-(\lambda+x)\tau}\mu(\mathrm{d}x), \qquad r_\theta^{\mathrm{tti}}(\tau) = \int e^{-(\lambda+x)\tau}\mu(\mathrm{d}x), \qquad c_\theta(*) = \int\frac{\mathbb{E}(\theta^*)^2 x}{\lambda+x}\mu(\mathrm{d}x)$$

$$c_\eta^{\mathrm{init}}(s) = \int\frac{(\mathbb{E}(\theta^*)^2 x + 1)\lambda x}{(\lambda+x)^2}e^{-(\lambda+x)s}\mu(\mathrm{d}x), \qquad c_\eta^{\mathrm{tti}}(\infty) = \int\frac{(\mathbb{E}(\theta^*)^2 x + 1)\lambda^2}{(\lambda+x)^2}\mu(\mathrm{d}x) + \delta - 1$$

$$c_\eta^{\mathrm{tti}}(\tau) = c_\eta^{\mathrm{tti}}(\infty) + \int\frac{x}{\lambda+x}e^{-(\lambda+x)\tau}\mu(\mathrm{d}x), \qquad r_\eta^{\mathrm{tti}}(\tau) = \int x e^{-(\lambda+x)\tau}\mu(\mathrm{d}x).$$

These functions $c_\theta^{\mathrm{tti}}, c_\eta^{\mathrm{tti}}$ have the forms (33) for the positive, finite measures $\mu_\theta(\mathrm{d}a) = a^{-1}\mu(\mathrm{d}(a - \lambda))$ and $\mu_\eta(\mathrm{d}a) = [(a - \lambda)/a]\mu(\mathrm{d}(a - \lambda))$ supported on $a \in [\lambda, \infty)$. Furthermore, these functions $c_\theta^{\mathrm{tti}}, c_\eta^{\mathrm{tti}}, r_\theta^{\mathrm{tti}}, r_\eta^{\mathrm{tti}}$ satisfy the fluctuation-dissipation relations (36), verifying all conditions of Definition 2.4.

## C  Sufficient conditions for a log-Sobolev inequality

We prove Proposition 2.8 on a log-Sobolev inequality for the posterior law.

**Lemma C.1.** *Under Assumption 2.2(a), the prior density $g(\cdot)$ satisfies the LSI (17). Furthermore, consider the law*

$$\mathsf{P}(\theta) = \frac{g(\theta)e^{-\frac{a}{2}\theta^2 + b\theta}}{Z}, \qquad Z = \int g(\theta)e^{-\frac{a}{2}\theta^2 + b\theta}\mathrm{d}\theta. \tag{230}$$

*For any $a > 0$ and $b \in \mathbb{R}$, this law $\mathsf{P}(\theta)$ also satisfies the LSI (17). Both statements hold with the constant $C_{\mathrm{LSI}} = (4/c_0)\exp(8r_0^2(c_0 + C)^2/(\pi c_0))$ where $C, c_0, r_0$ are the constants of Assumption 2.2(a).*

*Proof.* Applying $x = \min(x, -c_0) + \max(x + c_0, 0)$, define

$$\ell_-(\theta) = \log g(0) + (\log g)'(0) \cdot \theta + \int_0^\theta \int_0^x \min\left((\log g)''(u), -c_0\right) du\, dx$$

$$\ell_+(\theta) = \int_0^\theta \int_0^x \max\left((\log g)''(u) + c_0, 0\right) du\, dx$$

so that $\log g(\theta) = \ell_-(\theta) + \ell_+(\theta)$. Then set

$$\tilde{\ell}_-(\theta) = -\log Z - a\theta^2 + b\theta + \ell_-(\theta)$$

so that $\log \mathsf{P}(\theta) = \tilde{\ell}_-(\theta) + \ell_+(\theta)$. By definition we have $\ell_-''(\theta) = \min((\log g)''(\theta), -c_0) \leq -c_0$ and also $\tilde{\ell}_-''(\theta) = -a + \ell_-''(\theta) \leq -c_0$. We have $(\log g)''(u) + c_0 \leq 0$ for all $|u| > r_0$ and $(\log g)''(u) + c_0 \leq c_0 + C$ for all $|u| \leq r_0$. Hence $|\ell_+'(\theta)| \leq r_0(c_0 + C)$. Thus both $\log g(\theta)$ and $\log \mathsf{P}(\theta)$ are sums of a $c_0$-strongly-log-concave potential $\ell_-(\theta)$ or $\tilde{\ell}_-(\theta)$ and a $r_0(c_0 + C)$-Lipschitz perturbation $\ell_+(\theta)$. Then [92, Lemma 2.1] shows that both laws satisfy a LSI with constant $C_{\mathrm{LSI}} = (4/c_0) \exp(8r_0^2(c_0 + C)^2/(\pi c_0))$. $\qquad\square$

*Proof of Proposition 2.8.* Under condition (a), the posterior density is strongly log-concave, satisfying

$$-\nabla^2 \log \mathsf{P}_g(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y}) = \frac{1}{\sigma^2}\mathbf{X}^\top\mathbf{X} - \mathrm{diag}\left(\left(\log g\right)''(\theta_j)\right)_{j=1}^d \succeq c_0 \mathbf{I}.$$

Hence (46) with $C_{\mathrm{LSI}} = 2/c_0$ follows from the Bakry-Emery criterion. Clearly this holds for any noise variance $\sigma^2 > 0$, verifying Assumption 2.7.

The proof under condition (b) is an adaptation of the argument of [93]; see also [10, Theorem 3.4] and [94] for similar specializations to the linear model. Under the conditions for $\mathbf{X}$ of Assumption 2.1, by the Bai-Yin law ( [95, Theorem 3.1]), for any $\varepsilon > 0$ the event

$$\mathcal{E}(\mathbf{X}) = \left\{(\sqrt{\delta} - 1)_+^2 - \varepsilon \leq \lambda_{\min}(\mathbf{X}^\top\mathbf{X}) \leq \lambda_{\max}(\mathbf{X}^\top\mathbf{X}) \leq (\sqrt{\delta} + 1)^2 + \varepsilon\right\}$$

hold a.s. for all large $n, d$ (where $\delta = \lim n/d$). Thus, choosing some sufficiently small $\varepsilon > 0$ and setting

$$\kappa = (\sqrt{\delta} - 1)_+^2 - 2\varepsilon, \quad \tau^2 = \sigma^2\left([(\sqrt{\delta} + 1)^2 - (\sqrt{\delta} - 1)_+^2 + 3\varepsilon]^{-1} - \varepsilon\right), \quad \boldsymbol{\Sigma} = \sigma^2(\mathbf{X}^\top\mathbf{X} - \kappa\mathbf{I})^{-1} - \tau^2\mathbf{I},$$

we have $\mathbf{X}^\top\mathbf{X} - \kappa\mathbf{I} \succeq \varepsilon\mathbf{I}$, $(\mathbf{X}^\top\mathbf{X} - \kappa\mathbf{I})^{-1} \succeq (\frac{\tau^2}{\sigma^2} + \varepsilon)\mathbf{I}$, and hence $\boldsymbol{\Sigma} \succeq \varepsilon\sigma^2$ on $\mathcal{E}(\mathbf{X})$. Since $\sigma^2(\mathbf{X}^\top\mathbf{X} - \kappa\mathbf{I})^{-1} = \boldsymbol{\Sigma} + \tau^2\mathbf{I}$, we have the Gaussian convolution identity

$$e^{-\frac{1}{2\sigma^2}\boldsymbol{\theta}^\top(\mathbf{X}^\top\mathbf{X} - \kappa\mathbf{I})\boldsymbol{\theta}} \propto \int e^{-\frac{1}{2\tau^2}\|\boldsymbol{\theta} - \boldsymbol{\varphi}\|_2^2} e^{-\frac{1}{2}\boldsymbol{\varphi}^\top\boldsymbol{\Sigma}^{-1}\boldsymbol{\varphi}} d\boldsymbol{\varphi}.$$

Then the posterior density $\mathsf{P}_g(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y})$ satisfies

$$\mathsf{P}_g(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y}) \propto \exp\left(-\frac{1}{2\sigma^2}\|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|^2\right)\prod_{j=1}^d g(\theta_j)$$

$$\propto \exp\left(-\frac{1}{2\sigma^2}\boldsymbol{\theta}^\top(\mathbf{X}^\top\mathbf{X} - \kappa\mathbf{I})\boldsymbol{\theta}\right)\prod_{j=1}^d g(\theta_j)\exp\left(-\frac{\kappa}{2\sigma^2}\theta_j^2 + \frac{\mathbf{x}_j^\top\mathbf{y}}{\sigma^2}\theta_j\right)$$

$$\propto \int e^{-\frac{1}{2}\boldsymbol{\varphi}^\top\boldsymbol{\Sigma}^{-1}\boldsymbol{\varphi}}\prod_{j=1}^d g(\theta_j)\exp\left(-\frac{\kappa}{2\sigma^2}\theta_j^2 - \frac{1}{2\tau^2}(\theta_j - \varphi_j)^2 + \frac{\mathbf{x}_j^\top\mathbf{y}}{\sigma^2}\theta_j\right)d\boldsymbol{\varphi}.$$

Defining

$$q_{\varphi_j}(\theta_j) = \frac{1}{Z_j(\varphi_j)}g(\theta_j)\exp\left(-\frac{\kappa}{2\sigma^2}\theta_j^2 - \frac{1}{2\tau^2}(\theta_j - \varphi_j)^2 + \frac{\mathbf{x}_j^\top\mathbf{y}}{\sigma^2}\theta_j\right),$$

$$Z_j(\varphi_j) = \int g(\theta_j)\exp\left(-\frac{\kappa}{2\sigma^2}\theta_j^2 - \frac{1}{2\tau^2}(\theta_j - \varphi_j)^2 + \frac{\mathbf{x}_j^\top\mathbf{y}}{\sigma^2}\theta_j\right)d\theta_j,$$

$$\mu(\boldsymbol{\varphi}) = \frac{e^{-\frac{1}{2}\boldsymbol{\varphi}^\top\boldsymbol{\Sigma}^{-1}\boldsymbol{\varphi}}\prod_{j=1}^d Z_j(\varphi_j)}{\int e^{-\frac{1}{2}\boldsymbol{\varphi}'^\top\boldsymbol{\Sigma}^{-1}\boldsymbol{\varphi}'}\prod_{j=1}^d Z_j(\varphi_j')d\boldsymbol{\varphi}'}$$

this gives the mixture-of-products representation

$$\mathsf{P}_g(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y}) = \int \underbrace{\prod_{j=1}^{d} q_{\varphi_j}(\theta_j)}_{:=q_{\boldsymbol{\varphi}}(\boldsymbol{\theta})} \mu(\boldsymbol{\varphi}) \mathrm{d}\boldsymbol{\varphi}.$$

Then for any $f \in C^1(\mathbb{R}^d)$,

$$\mathrm{Ent}[f(\boldsymbol{\theta})^2 \mid \mathbf{X}, \mathbf{y}] = \mathbb{E}_{\boldsymbol{\varphi} \sim \mu} \mathrm{Ent}_{\boldsymbol{\theta} \sim q_{\boldsymbol{\varphi}}} f(\boldsymbol{\theta})^2 + \mathrm{Ent}_{\boldsymbol{\varphi} \sim \mu} \mathbb{E}_{\boldsymbol{\theta} \sim q_{\boldsymbol{\varphi}}} f(\boldsymbol{\theta})^2. \tag{231}$$

For the first term of (231), note that inside the exponential defining $q_{\varphi_j}(\theta_j)$ we have $\kappa \geq -2\varepsilon$ and $\tau^2 \leq \sigma^2((1 + 3\varepsilon)^{-1} - \varepsilon)$, so the coefficient of $\theta_j^2$ is negative for sufficiently small $\varepsilon > 0$. Then by Lemma C.1, $q_{\varphi_j}(\theta_j)$ satisfies the univariate LSI (17) with constant $C_{\mathrm{LSI}} := (4/c_0) \exp(8r_0^2(c_0 + C)^2/(\pi c_0))$. So the product law $q_{\boldsymbol{\varphi}}$ satisfies the LSI with the same constant by tensorization, and

$$\mathbb{E}_{\boldsymbol{\varphi} \sim \mu} \mathrm{Ent}_{\boldsymbol{\theta} \sim q_{\boldsymbol{\varphi}}} f(\boldsymbol{\theta})^2 \leq C_{\mathrm{LSI}} \mathbb{E}_{\boldsymbol{\varphi} \sim \mu} \mathbb{E}_{\boldsymbol{\theta} \sim q_{\boldsymbol{\varphi}}} \|\nabla f(\boldsymbol{\theta})\|_2^2 = C_{\mathrm{LSI}} \mathbb{E}_{\boldsymbol{\theta} \sim q} \|\nabla f(\boldsymbol{\theta})\|_2^2. \tag{232}$$

For the second term of (231), note that

$$-\nabla_{\boldsymbol{\varphi}}^2 \log \mu(\boldsymbol{\varphi}) = \boldsymbol{\Sigma}^{-1} - \mathrm{diag}\left(\left(\log Z_j\right)''(\varphi_j)\right)_{j=1}^{d} = \boldsymbol{\Sigma}^{-1} + \mathrm{diag}\left(\frac{1}{\tau^2} - \frac{1}{\tau^4} \mathrm{Var}_{\theta_j \sim q_{\varphi_j}}[\theta_j]\right)_{j=1}^{d}.$$

The LSI for $q_{\varphi_j}$ implies $\mathrm{Var}_{\theta_j \sim q_{\varphi_j}}[\theta_j] \leq (C_{\mathrm{LSI}}/2)$ by its implied Poincaré inequality. Applying $(\sqrt{\delta} + 1)^2 - (\sqrt{\delta} - 1)_+^2 = 4\sqrt{\delta}\mathbf{1}\{\delta > 1\} + (\sqrt{\delta} + 1)^2\mathbf{1}\{\delta \leq 1\}$ and the given condition (b) for $\sigma^2$, we see that for a sufficiently small choice of $\varepsilon > 0$, we have $\tau^2 \geq C_{\mathrm{LSI}}/[2(1 - \varepsilon)]$. Then this gives $-\nabla_{\boldsymbol{\varphi}}^2 \log \mu(\boldsymbol{\varphi}) \succeq (\varepsilon/\tau^2)\mathbf{I}$. Then by the Bakry-Emery criterion, $\mu$ satisfies the LSI $\mathrm{Ent}_{\boldsymbol{\varphi} \sim \mu} f(\boldsymbol{\varphi})^2 \leq (2\tau^2/\varepsilon)\mathbb{E}_{\boldsymbol{\varphi} \sim \mu} \|\nabla f(\boldsymbol{\varphi})\|_2^2$, hence

$$\mathrm{Ent}_{\boldsymbol{\varphi} \sim \mu} \mathbb{E}_{\boldsymbol{\theta} \sim q_{\boldsymbol{\varphi}}} f(\boldsymbol{\theta})^2 \leq \frac{2\tau^2}{\varepsilon} \mathbb{E}_{\boldsymbol{\varphi} \sim \mu} \|\nabla_{\boldsymbol{\varphi}} (\mathbb{E}_{\boldsymbol{\theta} \sim q_{\boldsymbol{\varphi}}} f(\boldsymbol{\theta})^2)^{1/2}\|_2^2.$$

Denote by $q_{\boldsymbol{\varphi}^{-j}}$ the product of components of $q_{\boldsymbol{\varphi}}$ other than the $j^{\mathrm{th}}$. We compute

$$\partial_{\varphi_j}(\mathbb{E}_{\boldsymbol{\theta} \sim q_{\boldsymbol{\varphi}}} f(\boldsymbol{\theta})^2)^{1/2} = \frac{\partial_{\varphi_j} \mathbb{E}_{\boldsymbol{\theta} \sim q_{\boldsymbol{\varphi}}} f(\boldsymbol{\theta})^2}{2(\mathbb{E}_{\boldsymbol{\theta} \sim q_{\boldsymbol{\varphi}}} f(\boldsymbol{\theta})^2)^{1/2}} = \frac{\mathbb{E}_{\boldsymbol{\theta}^{-j} \sim q_{\boldsymbol{\varphi}^{-j}}} \partial_{\varphi_j} \mathbb{E}_{\theta_j \sim q_{\varphi_j}} f(\boldsymbol{\theta})^2}{2(\mathbb{E}_{\boldsymbol{\theta} \sim q_{\boldsymbol{\varphi}}} f(\boldsymbol{\theta})^2)^{1/2}} = \frac{\mathbb{E}_{\boldsymbol{\theta}^{-j} \sim q_{\boldsymbol{\varphi}^{-j}}} \mathrm{Cov}_{\theta_j \sim q_{\varphi_j}}[f(\boldsymbol{\theta})^2, \theta_j]}{2\tau^2(\mathbb{E}_{\boldsymbol{\theta} \sim q_{\boldsymbol{\varphi}}} f(\boldsymbol{\theta})^2)^{1/2}}.$$

We apply [96, Proposition 2.2]: For any law $\nu$ on $\mathbb{R}^d$ satisfying a LSI $\mathrm{Ent}_\nu f^2 \leq C_{\mathrm{LSI}} \mathbb{E}_\nu \|\nabla f\|_2^2$, and for any smooth functions $f, g : \mathbb{R}^d \to \mathbb{R}$,

$$\mathrm{Cov}_\nu[f^2, g] \leq C \sup_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\nabla g(\boldsymbol{\theta})\|_2 \cdot (\mathbb{E}_\nu f^2)^{1/2} (\mathbb{E}_\nu \|\nabla f\|_2^2)^{1/2}$$

where $C$ depends only on the LSI constant $C_{\mathrm{LSI}}$ of $\nu$. Applying this to the univariate law $\nu = q_{\varphi_j}$, followed by Cauchy-Schwarz,

$$\mathbb{E}_{\boldsymbol{\theta}^{-j} \sim q_{\boldsymbol{\varphi}^{-j}}} \mathrm{Cov}_{\theta_j \sim q_{\varphi_j}}[f(\boldsymbol{\theta})^2, \theta_j] \leq C_1 (\mathbb{E}_{\boldsymbol{\theta} \sim q_{\boldsymbol{\varphi}}} f(\boldsymbol{\theta})^2)^{1/2} (\mathbb{E}_{\boldsymbol{\theta} \sim q_{\boldsymbol{\varphi}}} (\partial_{\theta_j} f(\boldsymbol{\theta}))^2)^{1/2}$$

where $C_1$ depends only on the LSI constant $C_{\mathrm{LSI}}$ for $q_{\varphi_j}$. Summing over $j = 1, \ldots, d$ gives

$$\|\nabla_{\boldsymbol{\varphi}} (\mathbb{E}_{\boldsymbol{\theta} \sim q_{\boldsymbol{\varphi}}} f(\boldsymbol{\theta})^2)^{1/2}\|_2^2 = \sum_{j=1}^{d} \left[\partial_{\varphi_j} (\mathbb{E}_{\boldsymbol{\theta} \sim q_{\boldsymbol{\varphi}}} f(\boldsymbol{\theta})^2))^{1/2}\right]^2$$

$$\leq \sum_{j=1}^{d} \left(\frac{C_1}{2\tau^2}\right)^2 \mathbb{E}_{\boldsymbol{\theta} \sim q_{\boldsymbol{\varphi}}} (\partial_{\theta_j} f(\boldsymbol{\theta}))^2 = \left(\frac{C_1}{2\tau^2}\right)^2 \mathbb{E}_{\boldsymbol{\theta} \sim q_{\boldsymbol{\varphi}}} \|\nabla f(\boldsymbol{\theta})\|_2^2.$$

Thus

$$\mathrm{Ent}_{\boldsymbol{\varphi} \sim \mu} \mathbb{E}_{\boldsymbol{\theta} \sim q_{\boldsymbol{\varphi}}} f(\boldsymbol{\theta})^2 \leq \frac{2\tau^2}{\varepsilon} \left(\frac{C_1}{2\tau^2}\right)^2 \mathbb{E}_{\boldsymbol{\varphi} \sim \mu} \mathbb{E}_{\boldsymbol{\theta} \sim q_{\boldsymbol{\varphi}}} \|\nabla f(\boldsymbol{\theta})\|_2^2 = \frac{C_1^2}{2\varepsilon\tau^2} \mathbb{E}[\|\nabla f(\boldsymbol{\theta})\|_2^2 \mid \mathbf{X}, \mathbf{y}]. \tag{233}$$

Applying (232) and (233) to (231) completes the proof of (46), on the above event $\mathcal{E}(\mathbf{X})$. Since the given condition (b) also holds for all $\tilde{\sigma}^2 \geq \sigma^2$ when it holds for $\sigma^2$, this verifies Assumption 2.7.

Finally, under condition (c), we note that on the above event $\mathcal{E}(\mathbf{X})$ we have

$$-\nabla^2 \log \mathsf{P}_g(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y}) = \frac{1}{\sigma^2}\mathbf{X}^\top \mathbf{X} - \mathrm{diag}\left(\left(\log g\right)''(\theta_j)\right)_{j=1}^d \succeq \varepsilon \mathbf{I}$$

for all $\sigma^2 \leq [(\sqrt{\delta}-1)^2 - \varepsilon]/(C+\varepsilon)$, so (46) holds by the Bakry-Emery criterion. Choosing $\varepsilon > 0$ small enough, under condition (c) we have $[(\sqrt{\delta}-1)^2 - \varepsilon]/(C+\varepsilon) > 4C_0\sqrt{\delta}$, so that (46) holds with $C_{\mathrm{LSI}} = 2/\varepsilon$ for $\sigma^2 > [(\sqrt{\delta}-1)^2 - \varepsilon]/(C+\varepsilon)$ by the analysis of condition (b). Thus, on $\mathcal{E}(\mathbf{X})$, (46) holds with a uniform constant $C_{\mathrm{LSI}} > 0$ for all $\sigma^2 > 0$, again verifying Assumption 2.7. $\qquad\square$

# D  Auxiliary lemmas

**Lemma D.1** (Coupling of Gaussian processes). *Let $\{K(t,s)\}_{t,s\in[0,T]}$ and $\{\tilde{K}(t,s)\}_{t,s\in[0,T]}$ be two positive semidefinite covariance kernels such that for some $\varepsilon > 0$ and $C_0 > 0$*

$$\sup_{t,s\in[0,T]} |K(t,s) - \tilde{K}(t,s)| \leq \varepsilon, \tag{234}$$

*and*

$$\sup_{t,s\in[0,T]} K(t,t) + K(s,s) - 2K(t,s) \leq C_0|t-s|. \tag{235}$$

*Then there exists a coupling of the two mean-zero Gaussian processes $\{u^t\}_{t\in[0,T]}$ and $\{\tilde{u}^t\}_{t\in[0,T]}$ with covariances $\mathbb{E}[u^t u^s] = K(t,s)$ and $\mathbb{E}[\tilde{u}^t \tilde{u}^s] = \tilde{K}(t,s)$ such that*

$$\sup_{t\in[0,T]} \mathbb{E}[(u^t - \tilde{u}^t)^2] \leq (6C_0+3)\sqrt{T\varepsilon} + 15\varepsilon.$$

*Proof.* Fix $\gamma > 0$, and let $\lfloor t \rfloor = \max\{i\gamma : i \in \mathbb{Z}_+, i\gamma \leq t\}$ where $\mathbb{Z}_+ = \{0,1,2,\ldots\}$. Let $v^t = u^{\lfloor t \rfloor}$ and $\tilde{v}^t = \tilde{u}^{\lfloor t \rfloor}$ so that $\mathbb{E}[v^t v^s] = K(\lfloor t \rfloor, \lfloor s \rfloor)$ and $\mathbb{E}[\tilde{v}^t \tilde{v}^s] = \tilde{K}(\lfloor t \rfloor, \lfloor s \rfloor)$). Then by (235) and (234),

$$\sup_{t\in[0,T]} \mathbb{E}(u^t - v^t)^2 \leq C_0\gamma, \qquad \sup_{t\in[0,T]} \mathbb{E}(\tilde{u}^t - \tilde{v}^t)^2 \leq C_0\gamma + 4\varepsilon.$$

Let $X = (v^0, v^\gamma, v^{2\gamma}, \ldots, v^{\lfloor T \rfloor}) \in \mathbb{R}^N$, where here $N \leq T/\gamma+1$, and similarly let $\tilde{X} = (\tilde{v}^0, \tilde{v}^\gamma, \tilde{v}^{2\gamma}, \ldots, \tilde{v}^{\lfloor T \rfloor}) \in \mathbb{R}^N$. Let $\Sigma, \tilde{\Sigma} \in \mathbb{R}^{N\times N}$ be the covariance matrices of $X, \tilde{X}$, so $\Sigma_{ij} = K(i\gamma, j\gamma)$ and $\tilde{\Sigma}_{ij} = \tilde{K}(i\gamma, j\gamma)$. Coupling $X$ and $\tilde{X}$ by $X = \Sigma^{1/2}Z$ and $\tilde{X} = \tilde{\Sigma}^{1/2}Z$ where $Z \sim \mathcal{N}(0,I)$, for each $i = 1, \ldots, N$,

$$\mathbb{E}(X_i - \tilde{X}_i)^2 = \mathbf{e}_i^\top (\Sigma^{1/2} - \tilde{\Sigma}^{1/2})^2 \mathbf{e}_i \leq \|\Sigma^{1/2} - \tilde{\Sigma}^{1/2}\|_{\mathrm{op}}^2 \overset{(*)}{\leq} \|\Sigma - \tilde{\Sigma}\|_{\mathrm{op}} \leq N\varepsilon.$$

Here $(*)$ follows from [97, Theorem X.1.3], and the last inequality applies (234). Then we have

$$\sup_{t\in[0,T]} \mathbb{E}(u^t - \tilde{u}^t)^2 \leq 3\left[\sup_{t\in[0,T]} \mathbb{E}(u^t - v^t)^2 + \sup_{t\in[0,T]} \mathbb{E}(v^t - \tilde{v}^t)^2 + \sup_{t\in[0,T]} \mathbb{E}(\tilde{u}^t - \tilde{v}^t)^2\right]$$

$$\leq 6C_0\gamma + 12\varepsilon + 3\max_{i=1}^N \mathbb{E}(X_i - \tilde{X}_i)^2 \leq 6C_0\gamma + 12\varepsilon + 3(T/\gamma + 1)\varepsilon.$$

The conclusion follows by choosing $\gamma = \sqrt{T\varepsilon}$. $\qquad\square$

# References

[1] Greg CG Wei and Martin A Tanner. A Monte Carlo implementation of the EM algorithm and the poor man's data augmentation algorithms. *Journal of the American Statistical Association*, 85(411):699–704, 1990.

[2] Richard A Levine and George Casella. Implementations of the Monte Carlo EM algorithm. *Journal of Computational and Graphical Statistics*, 10(3):422–439, 2001.

[3] Geoffrey E Hinton. Training products of experts by minimizing contrastive divergence. *Neural computation*, 14(8):1771–1800, 2002.

[4] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265. pmlr, 2015.

[5] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32, 2019.

[6] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.

[7] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2020.

[8] Juan Kuntz, Jen Ning Lim, and Adam M Johansen. Particle algorithms for maximum likelihood training of latent variable models. In *International Conference on Artificial Intelligence and Statistics*, pages 5134–5180. PMLR, 2023.

[9] Ö Deniz Akyildiz, Francesca Romana Crucinio, Mark Girolami, Tim Johnston, and Sotirios Sabanis. Interacting particle langevin algorithm for maximum marginal likelihood estimation. *arXiv preprint arXiv:2303.13429*, 2023.

[10] Zhou Fan, Leying Guan, Yandi Shen, and Yihong Wu. Gradient flows for empirical Bayes in high-dimensional linear models. *arXiv preprint arXiv:2312.12708*, 2023.

[11] Louis Sharrock, Daniel Dodd, and Christopher Nemeth. Tuning-free maximum likelihood training of latent variable models via coin betting. In *International Conference on Artificial Intelligence and Statistics*, pages 1810–1818. PMLR, 2024.

[12] Pierre Marion, Anna Korba, Peter Bartlett, Mathieu Blondel, Valentin De Bortoli, Arnaud Doucet, Felipe Llinares-López, Courtney Paquette, and Quentin Berthet. Implicit diffusion: Efficient optimization through stochastic sampling. *arXiv preprint arXiv:2402.05468*, 2024.

[13] Charles E McCulloch. Maximum likelihood variance components estimation for binary data. *Journal of the American Statistical Association*, 89(425):330–335, 1994.

[14] Charles E McCulloch. Maximum likelihood algorithms for generalized linear mixed models. *Journal of the American statistical Association*, 92(437):162–170, 1997.

[15] Herbert Robbins. An empirical Bayes approach to statistics. In *Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Contributions to the Theory of Statistics*. The Regents of the University of California, 1956.

[16] Bradley Efron. *Large-scale inference: empirical Bayes methods for estimation, testing, and prediction.* Cambridge University Press, 2012.

[17] Theo HE Meuwissen, Ben J Hayes, and ME Goddard. Prediction of total genetic value using genome-wide dense marker maps. *Genetics*, 157(4):1819–1829, 2001.

[18] Jian Yang, S Hong Lee, Michael E Goddard, and Peter M Visscher. GCTA: a tool for genome-wide complex trait analysis. *The American Journal of Human Genetics*, 88(1):76–82, 2011.

[19] Po-Ru Loh, George Tucker, Brendan K Bulik-Sullivan, et al. Efficient Bayesian mixed-model analysis increases association power in large cohorts. *Nature genetics*, 47(3):284–290, 2015.

[20] Gerhard Moser, Sang Hong Lee, Ben J Hayes, et al. Simultaneous discovery, estimation and prediction analysis of complex traits using a Bayesian mixture model. *PLoS genetics*, 11(4):e1004969, 2015.

[21] Luke R Lloyd-Jones, Jian Zeng, Julia Sidorenko, et al. Improved polygenic prediction by Bayesian multiple regression on summary statistics. *Nature communications*, 10(1):5086, 2019.

[22] Tian Ge, Chia-Yen Chen, Yang Ni, Yen-Chen Anne Feng, and Jordan W Smoller. Polygenic prediction via Bayesian regression and continuous shrinkage priors. *Nature communications*, 10(1):1776, 2019.

[23] Jeffrey P Spence, Nasa Sinnott-Armstrong, Themistocles L Assimes, and Jonathan K Pritchard. A flexible modeling and inference framework for estimating variant effect sizes from GWAS summary statistics. *BioRxiv*, pages 2022–04, 2022.

[24] Fabio Morgante, Peter Carbonetto, Gao Wang, Yuxin Zou, Abhishek Sarkar, and Matthew Stephens. A flexible empirical Bayes approach to multivariate multiple regression, and its improved accuracy in predicting multi-tissue gene expression from genotypes. *PLoS Genetics*, 19(7):e1010539, 2023.

[25] Sumit Mukherjee, Bodhisattva Sen, and Subhabrata Sen. A mean field approach to empirical Bayes estimation in high-dimensional linear regression. *arXiv preprint arXiv:2309.16843*, 2023.

[26] Valentin De Bortoli, Alain Durmus, Marcelo Pereyra, and Ana F Vidal. Efficient stochastic optimisation by unadjusted Langevin Monte Carlo: Application to maximum marginal likelihood and empirical Bayesian estimation. *Statistics and Computing*, 31:1–18, 2021.

[27] Zhou Fan, Justin Ko, Bruno Loureiro, Yue M. Lu, and Yandi Shen. Dynamical mean-field analysis of adaptive Langevin diffusions: Propagation-of-chaos and convergence of the linear response. 2025.

[28] Michael Celentano, Chen Cheng, and Andrea Montanari. The high-dimensional asymptotics of first order methods with random data. *arXiv preprint arXiv:2112.07572*, 2021.

[29] Cedric Gerbelot, Emanuele Troiani, Francesca Mignacco, Florent Krzakala, and Lenka Zdeborova. Rigorous dynamical mean-field theory for stochastic gradient descent methods. *SIAM Journal on Mathematics of Data Science*, 6(2):400–427, 2024.

[30] Dongning Guo and Sergio Verdú. Randomly spread CDMA: Asymptotics via statistical physics. *IEEE Transactions on Information Theory*, 51(6):1983–2010, 2005.

[31] Yoshiyuki Kabashima. Inference from correlated patterns: a unified theory for perceptron learning and linear vector channels. In *Journal of Physics: Conference Series*, volume 95, page 012001. IOP Publishing, 2008.

[32] Haim Sompolinsky and Annette Zippelius. Dynamic theory of the spin-glass phase. *Physical Review Letters*, 47(5):359, 1981.

[33] Haim Sompolinsky and Annette Zippelius. Relaxational dynamics of the Edwards-Anderson model and the mean-field theory of spin-glasses. *Physical Review B*, 25(11):6860, 1982.

[34] Theodore R Kirkpatrick and Devarajan Thirumalai. Dynamics of the structural glass transition and the p-spin—interaction spin-glass model. *Physical review letters*, 58(20):2091, 1987.

[35] Andrea Crisanti, Heinz Horner, and H J Sommers. The spherical p-spin interaction spin-glass model: the dynamics. *Zeitschrift für Physik B Condensed Matter*, 92:257–271, 1993.

[36] Leticia F Cugliandolo and Jorge Kurchan. Analytical solution of the off-equilibrium dynamics of a long-range spin-glass model. *Physical Review Letters*, 71(1):173, 1993.

[37] Leticia F Cugliandolo and Jorge Kurchan. On the out-of-equilibrium relaxation of the Sherrington-Kirkpatrick model. *Journal of Physics A: Mathematical and General*, 27(17):5749, 1994.

[38] G Ben Arous and Alice Guionnet. Large deviations for Langevin spin glass dynamics. *Probability Theory and Related Fields*, 102:455–509, 1995.

[39] G Ben Arous and Alice Guionnet. Symmetric Langevin spin glass dynamics. *The Annals of Probability*, 25(3):1367–1422, 1997.

[40] Alice Guionnet. Averaged and quenched propagation of chaos for spin glass dynamics. *Probability Theory and Related Fields*, 109:183–215, 1997.

[41] Francesca Mignacco, Florent Krzakala, Pierfrancesco Urbani, and Lenka Zdeborová. Dynamical mean-field theory for stochastic gradient descent in gaussian mixture classification. *Advances in Neural Information Processing Systems*, 33:9540–9550, 2020.

[42] Stefano Sarao Mannelli, Florent Krzakala, Pierfrancesco Urbani, and Lenka Zdeborova. Passed & spurious: Descent algorithms and local minima in spiked matrix-tensor models. In *international conference on machine learning*, pages 4333–4342. PMLR, 2019.

[43] Stefano Sarao Mannelli, Giulio Biroli, Chiara Cammarota, Florent Krzakala, Pierfrancesco Urbani, and Lenka Zdeborová. Marvels and pitfalls of the Langevin algorithm in noisy high-dimensional inference. *Physical Review X*, 10(1):011057, 2020.

[44] Tengyuan Liang, Subhabrata Sen, and Pragya Sur. High-dimensional asymptotics of Langevin dynamics in spiked matrix models. *Information and Inference: A Journal of the IMA*, 12(4), 2023.

[45] Francesca Mignacco, Pierfrancesco Urbani, and Lenka Zdeborová. Stochasticity helps to navigate rough landscapes: comparing gradient-descent-based algorithms in the phase retrieval problem. *Machine Learning: Science and Technology*, 2(3):035029, 2021.

[46] Qiyang Han and Xiaocong Xu. Gradient descent inference in empirical risk minimization. *arXiv preprint arXiv:2412.09498*, 2024.

[47] Elisabeth Agoritsas, Giulio Biroli, Pierfrancesco Urbani, and Francesco Zamponi. Out-of-equilibrium dynamical mean-field equations for the perceptron model. *Journal of Physics A: Mathematical and Theoretical*, 51(8):085002, 2018.

[48] Blake Bordelon and Cengiz Pehlevan. Self-consistent dynamical field theory of kernel evolution in wide neural networks. *Advances in Neural Information Processing Systems*, 35:32240–32256, 2022.

[49] Yatin Dandi, Emanuele Troiani, Luca Arnaboldi, Luca Pesce, Lenka Zdeborová, and Florent Krzakala. The benefits of reusing batches for gradient descent in two-layer networks: Breaking the curse of information and leap exponents. *arXiv preprint arXiv:2402.03220*, 2024.

[50] Blake Bordelon, Hamza Chaudhry, and Cengiz Pehlevan. Infinite limits of multi-head transformer dynamics. *Advances in Neural Information Processing Systems*, 37:35824–35878, 2024.

[51] Blake Bordelon, Alexander Atanasov, and Cengiz Pehlevan. A dynamical model of neural scaling laws. *arXiv preprint arXiv:2402.01092*, 2024.

[52] Andrea Montanari and Pierfrancesco Urbani. Dynamical decoupling of generalization and overfitting in large two-layer networks. *arXiv preprint arXiv:2502.21269*, 2025.

[53] Ada Altieri, Giulio Biroli, and Chiara Cammarota. Dynamical mean-field theory and aging dynamics. *Journal of Physics A: Mathematical and Theoretical*, 53(37):375006, 2020.

[54] Andrea Crisanti, Silvio Franz, Jorge Kurchan, and Andrea Maiorano. Dynamical mean-field theory and the aging dynamics. In *Spin Glass Theory and Far Beyond: Replica Symmetry Breaking After 40 Years*, pages 157–186. World Scientific, 2023.

[55] G Ben Arous, Amir Dembo, and Alice Guionnet. Aging of spherical spin glasses. *Probability Theory and Related Fields*, 120:1–67, 2001.

[56] Antoine Bodin and Nicolas Macris. Rank-one matrix estimation: analytic time evolution of gradient descent dynamics. In *Conference on Learning Theory*, pages 635–678. PMLR, 2021.

[57] Toshiyuki Tanaka. A statistical-mechanics approach to large-system analysis of CDMA multiuser detectors. *IEEE Transactions on Information theory*, 48(11):2888–2910, 2002.

[58] Andrea Montanari and David Tse. Analysis of belief propagation for non-linear problems: The example of CDMA (or: How to prove Tanaka's formula). In *2006 IEEE Information Theory Workshop-ITW'06 Punta del Este*, pages 160–164. IEEE, 2006.

[59] Galen Reeves and Henry D Pfister. The replica-symmetric prediction for random linear estimation with Gaussian matrices is exact. *IEEE Transactions on Information Theory*, 65(4):2252–2283, 2019.

[60] Jean Barbier, Nicolas Macris, Mohamad Dia, and Florent Krzakala. Mutual information and optimality of approximate message-passing in random linear estimation. *IEEE Transactions on Information Theory*, 66(7):4270–4303, 2020.

[61] Jean Barbier and Nicolas Macris. The adaptive interpolation method: a simple scheme to prove replica formulas in Bayesian inference. *Probability theory and related fields*, 174:1133–1185, 2019.

[62] Jean Barbier, Florent Krzakala, Nicolas Macris, Léo Miolane, and Lenka Zdeborová. Optimal errors and phase transitions in high-dimensional generalized linear models. *Proceedings of the National Academy of Sciences*, 116(12):5451–5460, 2019.

[63] Francesco Guerra. Broken replica symmetry bounds in the mean field spin glass model. *Communications in mathematical physics*, 233:1–12, 2003.

[64] Michel Talagrand. The Parisi formula. *Annals of Mathematics*, pages 221–263, 2006.

[65] Takashi Takahashi and Yoshiyuki Kabashima. Macroscopic analysis of vector approximate message passing in a model-mismatched setting. *IEEE Transactions on Information Theory*, 68(8):5579–5600, 2022.

[66] Jean Barbier, Wei-Kuo Chen, Dmitry Panchenko, and Manuel Sáenz. Performance of Bayesian linear regression in a model with mismatch. *arXiv preprint arXiv:2107.06936*, 2021.

[67] Alice Guionnet, Justin Ko, Florent Krzakala, and Lenka Zdeborová. Estimating rank-one matrices with mismatched prior and noise: universality and large deviations. *Communications in Mathematical Physics*, 406(1):9, 2025.

[68] Yury Polyanskiy and Yihong Wu. *Information theory: From coding to learning*. Cambridge university press, 2025.

[69] 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.

[70] Ji Xu, Daniel J Hsu, and Arian Maleki. Global analysis of expectation maximization for mixtures of two gaussians. *Advances in Neural Information Processing Systems*, 29, 2016.

[71] Chi Jin, Yuchen Zhang, Sivaraman Balakrishnan, Martin J Wainwright, and Michael I Jordan. Local maxima in the likelihood of gaussian mixture models: Structural results and algorithmic consequences. *Advances in neural information processing systems*, 29, 2016.

[72] Anya Katsevich and Afonso S Bandeira. Likelihood maximization and moment matching in low SNR Gaussian mixture models. *Communications on Pure and Applied Mathematics*, 76(4):788–842, 2023.

[73] Yudong Chen, Dogyoon Song, Xumei Xi, and Yuqian Zhang. Local minima structures in gaussian mixture models. *IEEE Transactions on Information Theory*, 2024.

[74] Andreas Eberle. Reflection couplings and contraction rates for diffusions. *Probability theory and related fields*, 166:851–886, 2016.

[75] Andreas Eberle and Raphael Zimmer. Sticky couplings of multidimensional diffusions with different drifts. In *Annales de l'Institut Henri Poincaré-Probabilités et Statistiques*, volume 55, pages 2370–2394, 2019.

[76] Philip E Protter and Philip E Protter. *Stochastic differential equations*. Springer, 2005.

[77] Amir Dembo and Jean-Dominique Deuschel. Markovian perturbation, response and fluctuation dissipation theorem. *Annales de l'IHP Probabilités et statistiques*, 46(3):822–852, 2010.

[78] Xian Chen and Chen Jia. Mathematical foundation of nonequilibrium fluctuation–dissipation theorems for inhomogeneous diffusion processes with unbounded coefficients. *Stochastic Processes and their Applications*, 130(1):171–202, 2020.

[79] Jean-Michel Bismut. *Large deviations and the Malliavin calculus*, volume 45 of *Progress in Mathematics*. Birkhäuser Boston, Inc., Boston, MA, 1984.

[80] K David Elworthy and Xue-Mei Li. Formulae for the derivatives of heat semigroups. *Journal of Functional Analysis*, 125(1):252–286, 1994.

[81] Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.

[82] Sergey G Bobkov, Ivan Gentil, and Michel Ledoux. Hypercontractivity of Hamilton–Jacobi equations. *Journal de Mathématiques Pures et Appliquées*, 80(7):669–696, 2001.

[83] Dominique Bakry, Ivan Gentil, and Michel Ledoux. *Analysis and geometry of Markov diffusion operators*, volume 103. Springer, 2014.

[84] H Kunita. Stochastic differential equations and stochastic flows of diffeomorphisms. *École d'Été de Probabilités de Saint-Flour XII-1982*, pages 143–303, 1984.

[85] R.T. Rockafellar. *Convex Analysis*. Princeton Landmarks in Mathematics and Physics. Princeton University Press, 1997.

[86] Santosh Vempala and Andre Wibisono. Rapid convergence of the unadjusted Langevin algorithm: Isoperimetry suffices. *Advances in neural information processing systems*, 32, 2019.

[87] Dongning Guo, Shlomo Shamai, and Sergio Verdú. Mutual information and minimum mean-square error in Gaussian channels. *IEEE transactions on information theory*, 51(4):1261–1282, 2005.

[88] Sergio Verdú. Mismatched estimation and relative entropy. *IEEE Transactions on Information Theory*, 56(8):3712–3720, 2010.

[89] Alex Bloemendal, László Erdos, Antti Knowles, Horng-Tzer Yau, and Jun Yin. Isotropic local laws for sample covariance and generalized wigner matrices. *Electron. J. Probab*, 19(33):1–53, 2014.

[90] Song Mei, Yu Bai, and Andrea Montanari. The landscape of empirical risk for nonconvex losses. *The Annals of Statistics*, 46(6A):2747–2774, 2018.

[91] L Chris G Rogers and David Williams. *Diffusions, Markov processes, and martingales: Itô calculus*, volume 2. Cambridge university press, 2000.

[92] Jean-Baptiste Bardet, Nathaël Gozlan, Florent Malrieu, and Pierre-André Zitt. Functional inequalities for gaussian convolutions of compactly supported measures: explicit bounds and dimension dependence. *Bernoulli*, 24(1):333–353, 2018.

[93] Roland Bauerschmidt and Thierry Bodineau. A very simple proof of the LSI for high temperature spin systems. *Journal of Functional Analysis*, 276(8):2582–2588, 2019.

[94] Andrea Montanari and Yuchen Wu. Provably efficient posterior sampling for sparse linear regression via measure decomposition. *arXiv preprint arXiv:2406.19550*, 2024.

[95] Yong-Quan Yin, Zhi-Dong Bai, and Pathak R Krishnaiah. On the limit of the largest eigenvalue of the large dimensional sample covariance matrix. *Probability theory and related fields*, 78:509–521, 1988.

[96] Michel Ledoux. Logarithmic Sobolev inequalities for unbounded spin systems revisited. *Séminaire de Probabilités XXXV*, pages 167–194, 2001.

[97] Rajendra Bhatia. *Matrix analysis*, volume 169. Springer Science & Business Media, 2013.