

A Diffusion Process Perspective on Posterior Contraction Rates for Parameters

Wenlong Mou*, Nhat Ho†, Martin Wainwright‡, Peter Bartlett§, and Michael Jordan§

Abstract. We analyze the posterior contraction rates of parameters in Bayesian models via the Langevin diffusion process, in particular by controlling moments of the stochastic process and taking limits. Analogous to the non-asymptotic analysis of statistical M-estimators and stochastic optimization algorithms, our contraction rates depend on the structure of the population log-likelihood function, and stochastic perturbation bounds between the population and sample log-likelihood functions. Convergence rates are determined by a non-linear equation that relates the population-level structure to stochastic perturbation terms, along with a term characterizing the diffusive behavior. Based on this technique, we also prove non-asymptotic versions of a Bernstein–von Mises guarantee for the posterior. We illustrate this general theory by deriving posterior convergence rates for various concrete examples.

Key words. Bayesian inference, Diffusion processes, Posterior contraction rate, Bernstein–von Mises theorem

AMS subject classifications. 62F15, 62F12

1. Introduction. Bayesian inference is one of the central pillars of statistics. In Bayesian analysis, we first endow the parameter space with a prior distribution chosen by modelling considerations, and then apply Bayes’ rule, combining the prior with the likelihood, so as to form the posterior distribution. From a statistical perspective, this posterior is of fundamental interest, and there are various questions associated with its behavior, including its consistency as the sample size goes to infinity, and from a more refined point of view, its contraction rate in various metrics.

The earliest work on posterior consistency dates back to the seminal work of Doob [9], who exhibited conditions under which the posterior distribution is consistent for all parameters apart from a set of zero measure. Subsequent work by Freedman [13, 14] provided examples showing that this null set can be problematic for Bayesian consistency in non-parametric settings. In order to address this issue, Schwartz [40] proposed a general framework for establishing posterior consistency for both semiparametric and nonparametric models. Since then, a number of researchers have isolated conditions that are useful for studying posterior distributions [3, 54, 55].

Moving beyond posterior consistency, convergence rates for the posterior density function, along with associated parameters of models, remains an active area of research. For posterior densities, Ghosal et al. [16] gave a general testing framework for proving convergence rates for both finite and infinite dimensional models; it has been used by various researchers to analyze posterior densities for Dirichlet and nonparametric Beta mixtures [17, 18, 38, 41]. Other work [4, 58, 57] established minimax optimal rates for regression functions in nonparametric regression models. Related problems include adaptive rates for the density in nonparametric Bayesian inference [8, 15], Bayesian linear and non-linear inverse problems [33, 25], and

*Department of Statistics, University of Toronto

†Department of Statistics and Data Science, UT Austin.

‡EECS & Mathematics, Massachusetts Institute of Technology

§Department of EECS and Department of Statistics, UC Berkeley.

38 posterior contraction rates of density under misspecified models [24]. Other popular general
 39 frameworks for analyzing the density functions of posterior distributions include those of Shen
 40 and Wasserman [42], and Walker et al. [56].

41 **1.1. From frequentist to Bayesian analysis.** The focus of this paper is on posterior
 42 convergence rates for parameters—namely, how for parametric Bayesian models, the posterior
 43 distribution assigns mass to certain regions of the parameter space. Our contributions can be
 44 put into perspective by considering known results for M -estimators. In the world of frequentist
 45 statistics, estimators based on maximizing empirically-defined objective functions—known as
 46 M -estimators—play a central role. In the parametric setting, a generic M -estimator takes the
 47 form

$$48 \quad (1.1) \quad \hat{\theta}_n := \arg \max_{\theta \in \Theta} F_n(\theta) \quad \text{where } F_n(\theta) := \frac{1}{n} \sum_{i=1}^n f(\theta; X_i), \text{ with } X_i \stackrel{\text{i.i.d.}}{\sim} \mathbb{P} \text{ for } i = 1, \dots, n,$$

50 while the parameters θ range over some constraint set Θ , and the real-valued function f
 51 has domain $\Theta \times \mathcal{X}$. Maximum-likelihood is the archetypal example, obtained when f is the
 52 log-likelihood.

53 There is now a rich and well-developed theory—one which exploits ideas from both
 54 optimization theory and empirical process theory—for deriving sharp non-asymptotic bounds
 55 on the difference between the estimate $\hat{\theta}_n$ and the maximizer θ^* of the population-level
 56 objective (e.g., see the books [52, 49, 53]). This theory leverages properties of the population-
 57 level objective $F(\theta) := \mathbb{E}[f(\theta, X)]$ where the expectation is taken with respect to $X \sim \mathbb{P}$. At a
 58 high level, there are two key steps in the analysis of an M -estimator: exploiting the structure
 59 of F , and linking the behavior of the empirical objective F_n to the population objective F . In
 60 the simplest setting, the population objective is strongly concave around its unique maximum
 61 θ^* . More generally, when F is differentiable, one can consider a condition of the following type

$$62 \quad (1.2a) \quad -\langle \nabla F(\theta), \theta - \theta^* \rangle \geq \psi(\|\theta - \theta^*\|_2),$$

64 assumed to hold uniformly for all θ in a local neighborhood of θ^* . Here ψ is an increasing
 65 function on the positive real-line, with $\psi(t) = \frac{\mu}{2}t^2$ being the one obtained for a μ -strongly
 66 concave function. The second step is to relate the empirical and population objective, for
 67 instance by establishing a uniform bound on their gradients—say

$$68 \quad (1.2b) \quad \|\nabla F_n(\theta) - \nabla F(\theta)\|_2 \leq \zeta(\|\theta - \theta^*\|_2)\varepsilon_n,$$

70 where the function ζ is again defined on the positive real line, and ε_n measures the magnitude
 71 of the noise.

72 When the functions F and F_n satisfy bounds of the form (1.2a) and (1.2b), it can be shown
 73 that the estimate $\hat{\theta}_n$ satisfies a bound of the form $\|\hat{\theta}_n - \theta^*\|_2 \lesssim r_n$, where $r_n > 0$ is the largest
 74 positive solution to the inequality¹

$$75 \quad (1.3) \quad \psi(r) \leq \varepsilon_n \zeta(r).$$

¹This solution exists and is unique under mild regularity conditions on the pair (ψ, ζ) .

77 This framework is very convenient to use, since optimization theory and empirical process
 78 theory give us various tools for establishing the local growth condition (1.2a) and the stochastic
 79 perturbation bound (1.2b).

80 By using this framework with care, one can often obtain sharp results in terms of *problem*
 81 *dimension* d , in both the rate itself and sample size lower bound needed to achieve such rates.
 82 Moreover, the local growth condition (1.2a) is relatively flexible; for instance, it allows for models
 83 in which the Fisher information matrix is singular (so that the function ψ is *not* quadratic).
 84 There are many different instantiations of this general approach in past work, including various
 85 methods or establishing growth conditions and empirical process bounds [44, 34], analysis
 86 of iterative optimization algorithm [2, 12, 28, 21], as well as regularized and constrained
 87 M -estimators [27, 6].

88 **1.2. Our contributions.** Moving back to the Bayesian setup, it is natural to seek to a
 89 similarly flexible and user-friendly method for establishing finite-sample results for posterior
 90 contraction. The main contribution of this paper is to do so by using the Langevin diffusion
 91 process—a stochastic differential equation that can encode the posterior distribution—as a
 92 lens of analysis.

93 There are natural parallels between our mode of analysis, and deterministic analyses of
 94 optimization algorithms via differential equations [45, 43]. To provide such intuition, recall the
 95 M -estimator defined by the objective function (1.1). Under the given conditions, its optimum
 96 θ^* can be characterized as the limiting point of an *ordinary differential equation* known as the
 97 gradient flow, and the rate (1.3) via the gradient flow dynamics for population and empirical
 98 loss functions, respectively. Now consider the analogous approach for studying *not* the M -
 99 estimator, but rather (in the Bayesian set-up) the posterior distribution. It is well-known [37]
 100 that under mild regularity conditions, the posterior distribution can be represented as the
 101 stationary distribution of a *stochastic differential equation* known as the Langevin diffusion.
 102 Consequently, just as information about the M -estimator can be recovered by studying the
 103 gradient flow, we can recover information about the posterior distribution by studying the
 104 Langevin diffusion. In particular, we do so by leveraging stochastic calculus so as to control
 105 the moments of this diffusion process. At a high-level, our main results involving showing
 106 that, under assumptions of the form (1.2), the posterior convergence rate is governed by the
 107 inequality $\psi(r) \leq \varepsilon_n \zeta(r) + \frac{d}{n}$. By comparison to inequality (1.3), relevant for M -estimation, we
 108 see that this inequality includes an additional $\frac{d}{n}$ term: it characterizes the diffusive behavior
 109 (with dimension d and sample size n) induced from sampling from the Gibbs measure e^{-F_n} as
 110 opposed to taking its maximum.

111 With this overview in place, we now summarize the different classes of contributions that
 112 are made in this paper:

113 **Posterior contraction under one-point strong convexity.** We begin with the simplest setting,
 114 in which the population log-likelihood function is strongly concave in a global sense. Under
 115 certain regularity conditions,² we prove that the posterior contraction rate around the true
 116 parameter is $(d/n)^{1/2}$. Our technique allows us to specify precise non-asymptotic conditions
 117 on the sample size and other model properties under which a guarantee of this type holds. In

²Briefly, we require the prior distribution to be sufficiently smooth and the perturbation error between the population and empirical log-likelihood function to be well-controlled.

118 many practical examples, the results yield sharp dependence on the problem dimension.

119 *Posterior contraction under weak concavity.* We then relax our assumption from strongly
120 concave to weakly concave, and prove related guarantees. Our results allow the Fisher
121 information matrix to be degenerate, in which case the $n^{-1/2}$ convergence rate is not possible,
122 and the contraction rate is governed by the interplay of a local growth assumption and local
123 empirical process bounds. We illustrate these general results for two concrete classes of models:
124 over-specified Bayesian location Gaussian mixture models and Bayesian logistic regression
125 models.

126 *Non-asymptotic Bernstein–von Mises (BvM) results.* Our final contribution is to establish
127 two non-asymptotic BvM results for models with non-degenerate Fisher information. We
128 first derive a non-asymptotic upper bound on the Kullback–Leibler (KL) divergence between
129 the posterior distribution and the limiting Gaussian distribution. Second, we prove a non-
130 asymptotic contraction bounds for the posterior distribution that adapts to the geometry
131 of Fisher information. The bound almost matches the tail bounds satisfied by the limiting
132 Gaussian law.

133 The remainder of the paper is organized as follows. In [Section 2](#), we set up the basic frame-
134 work for Bayesian models and introduce a diffusion process that admits posterior distribution
135 as its stationary distribution. [Section 3](#) presents the main results whose proofs are in [Section 5](#).
136 [Section 4](#) is devoted to implications to concrete examples. We conclude our work with a
137 discussion in [Section 6](#) while some technical proofs are in the supplementary material [\[31\]](#).

138 *Notation.* In the paper, the expression $a_n \lesssim b_n$ will be used to denote $a_n \leq cb_n$ for some
139 positive universal constant c that does not change with n . Additionally, we write $a_n \asymp b_n$ if
140 both $a_n \lesssim b_n$ and $a_n \gtrsim b_n$ hold. For any $n \in \mathbb{N}$, we denote $[n] = \{1, 2, \dots, n\}$. The notation
141 \mathbb{S}^{d-1} stands for the unit sphere, namely, the set of vectors $u \in \mathbb{R}^d$ such that $\|u\|_2 = 1$. Given a
142 vector $\theta \in \mathbb{R}^d$ and a scalar $r > 0$, we use $\mathbb{B}(\theta, r)$ to denote the closed ball centered at θ with
143 radius r . For any subset Θ of \mathbb{R}^d , $r \geq 1$, and $\varepsilon > 0$, we denote $\mathcal{N}(\varepsilon, \Theta, \|\cdot\|_r)$ the covering
144 number of Θ under $\|\cdot\|_r$ norm, namely, the minimum number of ε -balls under $\|\cdot\|_r$ norm to
145 cover the entire set Θ . Given a positive-definite matrix $M \succ 0$, we use $\lambda_{\max}(M)$ and $\lambda_{\min}(M)$ to
146 denote its largest and smallest eigenvalue, respectively, and we use $\kappa(M) := \lambda_{\max}(M)/\lambda_{\min}(M)$
147 to denote its condition number. Finally, for any $x, y \in \mathbb{R}$, we denote $x \vee y = \max\{x, y\}$
148 and $x \wedge y = \min\{x, y\}$. Given a pair of probability distributions P and Q , such that P is
149 absolutely continuous with respect to Q . The Kullback–Leibler (KL) divergence is defined as
150 $D_{\text{KL}}(P \parallel Q) := \mathbb{E}_P[\log \frac{dP}{dQ}]$.

151 **2. Background and problem formulation.** This section is devoted to background material
152 along with formulation of the problems studied in this paper. We first set up the problem of
153 studying convergence rates for posterior distributions over parameters in [Subsection 2.1](#), and
154 provide background on its representation as the stationary distribution of a Langevin diffusion
155 process in [Subsection 2.2](#). Finally, we define the population likelihood function, and introduce
156 various smoothness conditions in [Subsection 2.3](#).

157 **2.1. Posterior contraction rates for parameters.** Consider a parametric family of distri-
158 butions $\mathcal{P}_\Theta = \{P_\theta \mid \theta \in \Theta\}$. Throughout the paper, we assume that each distribution P_θ has
159 density p_θ with respect to the Lebesgue measure. Let $X_1^n := (X_1, \dots, X_n)$ be a sequence of
160 random variables drawn i.i.d. from an underlying distribution P . In the well-specified case,

161 we have that $P = P_{\theta^*} \in \mathcal{P}_\Theta$ for some $\theta^* \in \Theta$. However, it is important to note throughout
 162 our paper, the ground truth distribution P does not have to lie in the parametric family \mathcal{P}_Θ .
 163 Instead, the posterior contraction results around the parameter θ^* hold as long as certain
 164 geometric conditions around θ^* are satisfied. These conditions are typically achieved by the
 165 parameter θ^* such that P_{θ^*} is the best approximation to P within the family. See [Section 3](#) for
 166 a concrete discussion about these conditions.

167 Given a prior π over the parameter space, we define the log-likelihood

$$168 \quad (2.1) \quad F_n(\theta) := \frac{1}{n} \sum_{i=1}^n \log p_\theta(X_i), \quad \text{along with the posterior} \quad \mathbb{Q}(\theta \mid X_1^n) := \frac{e^{nF_n(\theta)} \pi(\theta)}{\int_{\Theta} e^{nF_n(u)} \pi(u) du}.$$

170 As the sample size n increases, we expect that the posterior distribution will concentrate
 171 more of its mass over increasingly smaller neighborhoods of the true parameter θ^* . Posterior
 172 contraction rates allow us to study how quickly this concentration of mass takes place. In
 173 particular, for a given norm, we study the posterior mass of a ball of the form $\|\theta - \theta^*\| \leq \rho$ for a
 174 suitably chosen radius $\rho > 0$. For a given $\delta \in (0, 1)$, our goal is to prove statements of the form
 175 $\mathbb{Q}(\|\theta - \theta^*\| \geq \rho(n, d, \delta) \mid X_1^n) \leq \delta$, with probability at least $1 - \delta$ over the randomly drawn
 176 data X_1^n . Our interest is in the scaling of the radius $\rho(n, d, \delta)$ as a function of sample size n ,
 177 problem dimension d , and the error tolerance δ , as well as other problem-specific parameters.

178 **2.2. From diffusion processes to the posterior distribution.** The analysis of this paper
 179 relies on a well-known connection between the posterior distribution and a particular stochastic
 180 differential equation (SDE) known as the Langevin diffusion. For a parameter $\beta > 0$, the
 181 Langevin diffusion can be written as

$$182 \quad (2.2) \quad d\theta_t = -\frac{1}{2} \nabla U(\theta_t) dt + \frac{1}{\sqrt{\beta}} dB_t,$$

184 where $(B_t, t \geq 0)$ is a standard d -dimensional Brownian motion [36], and $U : \mathbb{R}^d \rightarrow \mathbb{R}$ is known
 185 as the potential function. Suppose that we impose the following regularity conditions on the
 186 potential: (a) its gradient ∇U is locally Lipschitz, and (b) its gradient satisfies the inequality
 187 $\langle \nabla U(\theta), \theta \rangle \geq c_1 \|\theta\|_2 - c_2$ for any $\theta \in \mathbb{R}^d$, for some strictly positive constants c_1, c_2 . Under
 188 these conditions, by known results on general Langevin diffusions [1], the solution to the
 189 Langevin diffusion (2.2) exists and is unique in the strong sense. Furthermore, the density of
 190 θ_t converges in \mathbb{L}^2 to the stationary distribution with density proportional to $e^{-\beta U}$.

191 In the context of Bayesian inference, we can apply this argument to the potential function
 192 $U_n(\theta) := -F_n(\theta) - n^{-1} \log \pi(\theta)$ and $\beta = n$. Doing so will require us to verify that U_n satisfies
 193 the requisite regularity conditions. Assuming this validity, we are guaranteed that the posterior
 194 distribution $\mathbb{Q}(\theta \mid X_1^n)$ is the stationary distribution of the SDE

$$195 \quad (2.3) \quad d\theta_t = \frac{1}{2} \nabla F_n(\theta_t) dt + \frac{1}{2n} \nabla \log \pi(\theta_t) dt + \frac{1}{\sqrt{n}} dB_t,$$

197 with initial condition $\theta_0 = \theta^*$. Moreover, the density of θ_t converges in \mathbb{L}^2 to the posterior
 198 density.

199 It should be noted that this SDE-based representation of the posterior underlies various
 200 algorithms for drawing samples from the posterior distribution; we refer the reader to the

classical literature [47, 48] and the recent progress [7, 10, 11] for some results in this direction. In this paper, we exploit this SDE-based representation for statistical analysis (as opposed to efficient computation). In particular, by characterizing the behavior of the process $(\theta_t, t \geq 0)$ as a function of time, we can obtain bounds on the posterior distribution by taking limits. The following proposition guarantees the convergence of the moments based on a uniform-in-time moment upper bound and a convergence in total variation distance.

Proposition 2.1. *Consider a sequence of distributions $(\pi_t)_{t \geq 0}$ on \mathbb{R}^d such that $d_{\text{TV}}(\pi_t, \pi^*) \rightarrow 0$, and suppose that $\sup_{t \geq 0} \mathbb{E}_{\pi_t} [\|X\|_2^p] < +\infty$ and $\mathbb{E}_{\pi^*} [\|X\|_2^p] < +\infty$ for any integer $p \geq 2$. We then have $\lim_{t \rightarrow +\infty} \mathbb{E}_{\pi_t} [\|X\|_2^p] = \mathbb{E}_{\pi^*} [\|X\|_2^p]$.*

See [Appendix C.1](#) in our supplementary material [31] for the proof of this proposition.

Given this limiting behavior, we can establish posterior contraction rates for the parameters by controlling the moments of the diffusion process $\{\theta_t\}_{t \geq 0}$. The main theoretical results of this paper are obtained by following this general roadmap.

2.3. From empirical to population likelihood. Before proceeding to our main results, let us introduce some additional definitions and conditions. A useful notion for our analysis is the population log-likelihood F . It corresponds to the limit of log-likelihood function F_n , as previously defined in equation (2.1), as the sample size n goes to infinity—viz.

$$(2.4) \quad F(\theta) := \mathbb{E} [\log p_\theta(X)],$$

where the expectation is taken with respect to $X \sim P_{\theta^*}$. Throughout the paper, we impose the following smoothness conditions on the log prior density $\log \pi$:

(A) There exists a non-negative constant $B \geq 0$ such that

$$\langle \nabla \log \pi(\theta), \theta - \theta^* \rangle \leq B \|\theta - \theta^*\|_2 \quad \text{for all } \theta \in \mathbb{R}^d.$$

Although the constant B in Assumption (A) can depend on θ^* , we suppress this dependence so as to keep the notation streamlined. When the function $\log \pi$ is globally Lipschitz (so that $\|\nabla \log \pi(\theta)\|_2$ is uniformly bounded), Assumption (A) is automatically satisfied. But the one-sided nature of Assumption (A) makes it flexible and allows many practical prior distributions. For example, given scalars $\alpha, \beta > 0$, for the prior distribution $\pi(\theta) \propto \exp(-\beta^{-1} \|\theta\|_2^\alpha)$, we have

$$\begin{aligned} \langle \nabla \log \pi(\theta), \theta - \theta^* \rangle &= \frac{\alpha}{\beta} \|\theta\|_2^{\alpha-2} \left\{ \langle \theta^*, \theta - \theta^* \rangle - \|\theta - \theta^*\|_2^2 \right\} \\ &\leq \begin{cases} 2^{\alpha-2} \frac{\alpha}{\beta} \|\theta^*\|_2^{\alpha-1} \cdot \|\theta - \theta^*\|_2 & \|\theta - \theta^*\|_2 \leq \|\theta^*\|_2, \\ 0 & \text{otherwise,} \end{cases} \end{aligned}$$

so that Assumption (A) is satisfied by $B = 2^{\alpha-2} \frac{\alpha}{\beta} \|\theta^*\|_2^{\alpha-1}$.

3. Main results. We now turn to our main results. In [Subsection 3.1](#), we present a result ([Theorem 3.1](#)) that establishes the posterior convergence under strong concavity. [Subsection 3.2](#) answers the same question when the population log-likelihood is only weakly concave; see the statement of [Theorem 3.2](#). Finally, in [Subsection 3.3](#), we pursue a more fine-grained direction by establishing the non-asymptotic Bernstein–von Mises theorems (see [Proposition 3.4](#) and [Theorem 3.5](#))

240 **3.1. Posterior contraction under strong concavity.** We begin with results under strong
 241 concavity conditions. For this part, the following assumptions underlie our analysis:

242 **(S.1)** There exists a scalar $\mu > 0$ such that

$$243 \quad -\langle \nabla F(\theta), \theta^* - \theta \rangle \geq \mu \|\theta - \theta^*\|_2^2 \quad \text{for any } \theta \in \mathbb{R}^d.$$

245 **(S.2)** There exist non-negative functions ε_1 and ε_2 that map from $\mathbb{N} \times (0, 1]$ to \mathbb{R}_+ such that
 246 for any radius $r > 0$ and any $\delta \in (0, 1)$, we have

$$247 \quad \sup_{\theta \in \mathbb{B}(\theta^*, r)} \|\nabla F_n(\theta) - \nabla F(\theta)\|_2 \leq \varepsilon_1(n, \delta)r + \varepsilon_2(n, \delta) \quad \text{with prob. at least } 1 - \delta.$$

249 Assumption **(S.1)** is a standard strong concavity condition of function F around θ^* , whereas
 250 Assumption **(S.2)** provides uniform control on the gradients of the population and sample
 251 log-likelihoods. It is important to note that these assumptions, along with other assumptions to
 252 follow, *do not* require the data-generating distribution P to belong to the specified parametric
 253 class. Indeed, the results throughout this paper apply to both well-specified and mis-specified
 254 models. In the latter case, the parameter θ^* is typically the KL-projection of the true model,
 255 i.e., $\theta^* \in \arg \min_{\theta \in \Theta} D_{\text{KL}}(\mathbb{P} \parallel \mathbb{P}_\theta)$.

256 Given the above assumptions, we are ready to state our first result regarding the posterior
 257 convergence rate of parameters for a strongly concave population log-likelihood:

258 **Theorem 3.1.** *Suppose that Assumptions **(A)**, **(S.1)**, and **(S.2)** hold. Then there is a*
 259 *universal constant c such that for any $\delta \in (0, 1)$ and any sample size n for which $\varepsilon_1(n, \delta) \leq \frac{\mu}{6}$,*
 260 *we have*

$$261 \quad \mathbb{Q}\left(\|\theta - \theta^*\|_2 \geq c\sqrt{\frac{d}{n\mu}} + \frac{B}{n\mu} + \frac{\varepsilon_2(n, \delta)}{\mu} + c\sqrt{\frac{\log(1/\delta)}{n\mu}} \mid X_1^n\right) \leq \delta$$

263 *with probability $1 - \delta$, taken with respect to the random observations X_1^n .*

264 See [Subsection 5.1](#) for the proof of [Theorem 3.1](#).

265 This result guarantees posterior convergence at the rate $(d/n)^{1/2}$ when the log-likelihood
 266 is strongly concave. To be clear, such rate of posterior contraction for the parameters can be
 267 derived from the asymptotic behavior of the posterior distribution via the classical Bernstein–
 268 von Mises theorem. However, the guarantee in [Theorem 3.1](#) is non-asymptotic, and provides
 269 explicit dependence of the rate on other model parameters, including B and μ , both of which
 270 might vary as a function of θ^* . At the moment, we do not know whether the dependence of
 271 these parameters is optimal. This guarantee is valid as long as the error term $\varepsilon_1(n, \delta)$ is less
 272 than an absolute constant; such a bound typically holds as long as $n \gtrsim d$. In [Theorem 3.5](#)
 273 to follow, we also provide near-optimal non-asymptotic contraction bounds on the posterior
 274 distribution that nearly match the exact shape of the posterior distribution.

275 Although our set-up is focused on simple sampling models, it should be noted that our
 276 method is sufficiently flexible so as to accommodate certain non-i.i.d. forms of sampling, along
 277 with mis-specified models. After the first version was posted, Mazumdar et al. [29] used a
 278 variant of this result to study the posterior contraction for Thompson sampling in contextual
 279 bandits. In their problem, the data are adaptively collected instead of being i.i.d., and the
 280 empirical process bound [\(S.2\)](#) can be verified using martingale concentration inequalities.

281 While this paper focuses on the contraction of posterior distribution itself, it is worth
 282 mentioning that the proof techniques of Theorem 3.1 can be extended to study the contraction
 283 behavior of discretized Langevin diffusion. In particular, by expanding the discrete-time
 284 evolution of the iterates following Subsection 5.1, we can derive recursive relations on the
 285 moment bounds for the distance between iterates and θ^* using Assumptions (S.1) and (S.2).
 286 The solution to such recursion will lead to the rates in Theorem 3.1. This analysis tool does
 287 not depend on the ergodicity of the discretized diffusion. We defer a detailed discrete-time
 288 analysis to future work.

289 **3.2. Posterior contraction under weak concavity.** Theorem 3.1 requires global strong
 290 concavity, which is relatively strong. In this section, we relax this assumption in two ways:
 291 we relax the growth condition locally around θ^* so as to allow for weak concavity, and the
 292 global behavior need not coincide with this local behavior. Weakly concave log-likelihoods
 293 arise for singular problems, for which the Fisher information matrix at the true parameter θ^*
 294 is rank-degenerate. Examples of such singular problems include Bayesian non-linear regression
 295 models with certain choices of link functions [30], as well as over-specified mixture models [39],
 296 in which the fitted mixture model has more components than the true mixture distribution.
 297 The mismatch between local and global concavity conditions exists not only in such models,
 298 but also in non-singular problems such as Bayesian logistic regression. We discuss implications
 299 of these examples in Section 4.

300 Our analysis in the weakly concave setting is based on the following assumptions:

301 **(W.1)** There exists a convex, non-decreasing function $\psi : [0, +\infty) \rightarrow \mathbb{R}$ such that

$$302 \quad -\langle \nabla F(\theta), \theta - \theta^* \rangle \geq \psi(\|\theta - \theta^*\|_2) \quad \text{for any } \theta \in \mathbb{R}^d.$$

304 Assumption (W.1) characterizes the weak concavity of the function F around the global
 305 maxima θ^* . This condition can hold when the log-likelihood is locally strongly concave around
 306 θ^* but only weakly concave in a global sense, or it can hold when the log-likelihood is weakly
 307 concave but not strongly concave. An example of the former type is the logistic regression model
 308 analyzed in Subsection 4.1, whereas an example of the latter type is given by over-specified
 309 Gaussian mixture models Subsection 4.2.

310 Our next assumption controls the deviation between the gradients of the population and
 311 sample likelihoods, and involves a failure probability $\delta \in (0, 1)$:

312 **(W.2)** There exist a function $\varepsilon : \mathbb{N} \times (0, 1] \mapsto \mathbb{R}_+$ and a non-decreasing function $\zeta : \mathbb{R} \rightarrow \mathbb{R}$
 313 with that $\zeta(0) \geq 0$ such that for any radius $r > 0$, we

$$314 \quad \sup_{\theta \in \mathbb{B}(\theta^*, r)} \|\nabla F_n(\theta) - \nabla F(\theta)\|_2 \leq \varepsilon(n, \delta)\zeta(r) \quad \text{with prob. at least } 1 - \delta.$$

316 This type of localized empirical process bounds appeared in many existing literature in the
 317 study of M -estimators [50] and iterative algorithms [2, 12]. It is important to note that the
 318 bound depends on the radius r , making it possible to yield near-optimal rates in singular
 319 mixture models [12].

320 The previous conditions involved two functions, namely ψ and ζ . We let $\xi : \mathbb{R}_+ \rightarrow \mathbb{R}$
 321 denote the inverse function of the strictly increasing function $r \mapsto r\zeta(r)$. Our third assumption
 322 imposes certain inequalities on these functions and their derivatives:

323 **(W.3)** The function $r \mapsto \psi(\xi(r))$ is convex, and ψ and ζ satisfy the differential inequalities

324
$$r\psi'(r)\zeta(r) \stackrel{(i)}{\geq} r\psi(r)\zeta'(r) + \psi(r)\zeta(r), \quad \text{and}$$

325
$$r^2\psi''(r)\zeta(r) + r\psi'(r)\zeta(r) \stackrel{(ii)}{\geq} 3\psi(r)\zeta(r) + r^2\psi(r)\zeta''(r) \quad \text{for all } r > 0.$$

327 These differential inequalities are needed controlling the moments of the diffusion process
 328 $\{\theta_t\}_{t>0}$ in equation (2.3). In our discussion of concrete examples, we provide instances for
 329 which they are satisfied.

330 Our result involves a certain fixed point equation that depends on the parameters and
 331 functions in our assumptions. In particular, for any tolerance parameter $\delta \in (0, 1)$ and
 332 sample size n , consider the following fixed point equation in the variable $z > 0$:

333 (3.1)
$$\psi(z) = \varepsilon(n, \delta)\zeta(z)z + \frac{B}{n}z + \frac{d}{n} + \frac{\log(1/\delta)}{n}.$$

335 In order to ensure that this equation has a unique positive solution, our final assumption
 336 imposes certain condition on the growth of the functions ψ and ζ :

337
 338 **(W.4)** The limit $\liminf_{z \rightarrow +\infty} \frac{\psi(z)}{z\zeta(z)}$ is strictly positive, and the sample size n and tolerance
 339 parameter $\delta \in (0, 1)$ are such that $\varepsilon(n, \delta) < \liminf_{z \rightarrow +\infty} \frac{\psi(z)}{z\zeta(z)}$.

340 With this set-up, we are now ready to state our second main result:

341 **Theorem 3.2.** *Suppose that Assumptions (A), and (W.1) — (W.3) hold. Then for any*
 342 *given sample size n and $\delta \in (0, 1)$ such that Assumption (W.4) holds, equation (3.1) has a*
 343 *unique positive solution $z^*(n, \delta)$ such that*

344 (3.2)
$$\mathbb{Q}\left(\|\theta - \theta^*\|_2 \geq z^*(n, \delta) \mid X_1^n\right) \leq \delta \quad \text{with probability } 1 - \delta \text{ w.r.t. } X_1^n.$$

346 See Subsection 5.2 for the proof of Theorem 3.2.

347 A few comments are in order. First, the convergence guarantee (3.2) depends on the weak
 348 convexity function ψ and the perturbation function ζ through the non-linear equation (3.1).
 349 In order to understand the rate, we consider the following pair of fixed-point equations

350 (3.3a)
$$\psi(z) = 2\varepsilon(n, \delta)\zeta(z)z, \quad \text{with the solution } z_{\text{mle}}^*(n, \delta)$$

351 (3.3b)
$$\psi(z) = 2\frac{B}{n}z + 2\frac{d}{n} + 2\frac{\log(1/\delta)}{n}, \quad \text{with the solution } z_{\text{pop}}^*(n, \delta).$$

353 It is easy to see that $z^*(n, \delta) \leq \max\{z_{\text{mle}}^*(n, \delta), z_{\text{pop}}^*(n, \delta)\}$.³ This establishes that the posterior
 354 contraction rates in Theorem 3.2 are fundamentally determined by two sources of errors: on the
 355 one hand, it is known (see e.g. [50]) that the solution $z_{\text{mle}}^*(n, \delta)$ to equation (3.3a) determines
 356 (up to constant factors) the rate of convergence for the maximal likelihood estimator; on the

³Suppose the converse is true. We have $\psi(z^*(n, \delta)) > 2\varepsilon(n, \delta)\zeta(z^*(n, \delta))z^*(n, \delta)$ and $\psi(z^*(n, \delta)) > 2\frac{B}{n}z^*(n, \delta) + 2\frac{d}{n} + 2\frac{\log(1/\delta)}{n}$. Taking the average of two inequalities contradicts the fact that $z^*(n, \delta)$ is the fixed point.

357 other hand, the solution $z_{\text{pop}}^*(n, \delta)$ to equation (3.3b) captures the diffusive behavior from
 358 the posterior distribution itself. In particular, the term $\frac{B}{n}z$ is usually negligible as $z_{\text{pop}}^* \ll 1$
 359 and $B = O(\sqrt{d})$, and the solution to the equation $\psi(z) = \frac{d+\log(1/\delta)}{n}$ essentially determines
 360 the contraction rate of the ‘‘population-level posterior’’ Gibbs distribution whose density is
 361 proportional to $e^{nF(\theta)}$, when the function $-\langle \nabla F(\theta), \theta - \theta^* \rangle$ locally behaves like $\psi(\|\theta - \theta^*\|_2)$.
 362 We suspect that such an additional term is unavoidable for posterior contraction results, and
 363 we defer a rigorous justification via asymptotic shape of the re-scaled posterior to future works.

364 Second, at least in general, it is not possible to compute an explicit form for the positive
 365 solution $z^*(n, \delta)$ to the non-linear equation (3.1). However, for certain forms of the function
 366 ψ and ζ , we can derive a relatively simple upper bound. For instance, given some positive
 367 parameters (α, β) such that $\alpha > \beta$, suppose that these functions are defined locally, in a
 368 interval above zero, as follows:

$$369 \quad (3.4a) \quad \psi(r) = r^{\alpha+1}, \quad \text{and} \quad \zeta(r) = r^\beta \quad \text{for all } r \text{ in some interval } [0, \bar{r}).$$

371 Moreover, suppose that the perturbation function takes the form

$$372 \quad (3.4b) \quad \varepsilon(n, \delta) = \sqrt{(d + \log(\frac{1}{\delta})) / n}.$$

374 As shown in in Section 4, these particular forms arise in several statistical models, including
 375 Bayesian logistic regression and over specified Bayesian Gaussian mixture models. Under these
 376 conditions, we have the following simple upper bound:

377 **Corollary 3.3.** *Assume that the functions ψ , ζ have the local behavior (3.4a), and the*
 378 *perturbation term $\varepsilon(n, \delta)$ has the form (3.4b). If, in addition, the global forms of ψ and ζ*
 379 *satisfy Assumption (W.3), then for sufficiently large n , the scalar $z^*(n, \delta)$ from Theorem 3.2*
 380 *satisfies the bound $z^*(n, \delta) \leq c \left(\frac{d+\log(1/\delta)}{n} \right)^{\frac{1}{2(\alpha-\beta)}} \vee \left(\frac{d+\log(1/\delta)}{n} \right)^{\frac{1}{\alpha+1}} + \left(\frac{B}{n} \right)^{\frac{1}{\alpha}}$.*

381 Note that Corollary 3.3 ensures that the posterior has the following contraction property

$$382 \quad (3.5) \quad \mathbb{Q}\left(\|\theta - \theta^*\|_2 \geq c \left(\frac{d+\log(1/\delta)}{n} \right)^{\frac{1}{2(\alpha-\beta)} \wedge \frac{1}{\alpha+1}} + \left(\frac{B}{n} \right)^{\frac{1}{\alpha}} \mid X_1^n\right) \leq \delta \quad \text{with prob. } 1 - \delta$$

384 with respect to the training data. The posterior convergence rate scales as $(d/n)^{\frac{1}{2(\alpha-\beta)}}$ when
 385 $\alpha \geq 2\beta + 1$, in which case the posterior contraction rates match the maximal likelihood. On
 386 the other hand, this rate becomes $(d/n)^{\frac{1}{\alpha+1}}$ when $\alpha < 2\beta + 1$, and the posterior contraction is
 387 slower than maximal likelihood, owing to its diffusive behavior.

388 **Theorem 3.2** and **Corollary 3.3** rely on global conditions (W.1) and (W.2). Although
 389 these conditions can be verified for many practical examples (see Section 4), they can be
 390 restrictive in some cases, especially when multiple local maxima of the population-level function
 391 F exist. Using our techniques, it is possible to prove similar results under local assumptions.
 392 In particular, suppose that these assumptions hold only in a local ball $\mathbb{B}(\theta^*, r_0)$; then, the
 393 non-asymptotic contraction rates in Theorem 3.2 and Corollary 3.3 are available as long as we
 394 can show the posterior mass $\mathbb{Q}(\mathbb{B}(\theta^*, r_0)^c \mid X_1^n)$ is small with high probability. To obtain these
 395 rates, we could apply the arguments in the proof of Theorem 3.2 to a modified distribution,

396 which matches the shape of $\mathbb{Q}(\cdot \mid X_1^n)$ inside the ball $\mathbb{B}(\theta^*, r_0)$, while exhibiting desirable
 397 growth and smoothness conditions outside. We defer the detailed arguments based on local
 398 assumptions as well as the study of the radius r_0 to future work.

399 **3.3. Non-asymptotic Bernstein–von Mises results.** In this section, we develop non-
 400 asymptotic Bernstein–von Mises results using the diffusion process (2.3). Under mild assump-
 401 tions on the population-level and empirical-level landscapes, we establish the KL divergence
 402 between the posterior distribution and the limiting Gaussian distribution, as well as near-
 403 optimal shape-dependent posterior contraction results.

404 In order to obtain the non-asymptotic Bernstein–von Mises results, we first need the
 405 following assumptions on the second order derivatives with respect to the parameters (or
 406 equivalently Hessian matrices) of the empirical and population log-likelihoods:

407 **(BvM.1)** There exists $A > 0$ such that the population log-likelihood function F satisfies the
 408 one-point Lipschitz condition:

$$409 \quad \forall \theta \in \mathbb{R}^d, \quad \|\nabla^2 F(\theta) - \nabla^2 F(\theta^*)\|_{\text{op}} \leq A \|\theta - \theta^*\|_2.$$

410 **(BvM.2)** For any $\delta > 0$, there exist non-negative functions $\varepsilon_1^{(2)}$ and $\varepsilon_2^{(2)}$ with domain $\mathbb{N} \times (0, 1]$
 412 such that

$$413 \quad \sup_{\theta \in \mathbb{B}(\theta^*, r)} \|\nabla^2 F_n(\theta) - \nabla^2 F(\theta)\|_{\text{op}} \leq \varepsilon_1^{(2)}(n, \delta)r + \varepsilon_2^{(2)}(n, \delta),$$

414 for any radius $r > 0$ with probability at least $1 - \delta$.

415 Additionally, we also impose a smoothness assumption on the prior distribution π

$$416 \quad \text{(PS)} \quad \|\nabla \log \pi(\theta_1) - \nabla \log \pi(\theta_2)\|_2 \leq L_2 \|\theta_1 - \theta_2\|_2.$$

417 The first condition **(BvM.1)** is a standard smoothness condition needed to prove quantita-
 418 tive results about asymptotic normality (e.g., the paper [35]), and satisfied by many models such
 419 as exponential family models, location density models, as well as their mixtures and hierarchical
 420 composition. The second condition **(BvM.2)** is an empirical process condition on the Hessian
 421 matrix $\nabla^2 F_n$. This condition can usually be verified using suitable concentration bounds for
 422 each θ , as well as smoothness conditions on $\nabla^2 F_n$ used in controlling metric entropies. Both
 423 assumptions are naturally needed: the limiting Gaussian law $\mathcal{N}(\hat{\theta}^{(n)}, (nH^*)^{-1})$, which depends
 424 on the population-level Hessian at the point θ^* . The shape of posterior distribution, on the
 425 other hand, depends on the sample-level Hessian $\nabla^2 F_n$ in a local neighborhood of θ^* . These
 426 two conditions are needed to relate the shape of the sample-level posterior with the matrix H^* .
 427 The condition **(PS)** on the prior distribution is relatively mild and satisfied by many practical
 428 choices including Gaussian. As before, we note that these assumptions do not require the
 429 model to be well-specified, and our non-asymptotic Bernstein–von Mises theorems applies to
 430 the mis-specified case, where θ^* is the KL-projection of the model to this parametric class.

431 Consider the MAP estimate $\hat{\theta}^{(n)} := \arg \max_{\theta \in \mathbb{R}^d} (F_n(\theta) + \frac{1}{n} \log \pi(\theta))$. Then, we have the
 432 following upper bound on the difference between the posterior distribution of the parameters
 433 and the Gaussian distribution with mean $\hat{\theta}^{(n)}$ and covariance matrix $(nH^*)^{-1}$, where $H^* :=$
 434 $-\nabla^2 F(\theta^*)$.

437 **Proposition 3.4.** *Under Assumptions (BvM.1), (BvM.2) and PS, suppose that $H^* \succ 0$,*
 438 *and that $\|\hat{\theta}^{(n)} - \theta^*\|_2 \leq \sigma \sqrt{\frac{d}{n}}$ and $\mathbb{E}_{\mathbb{Q}}(\|\theta - \theta^*\|_2^4 | X_1^n)^{1/4} \leq \sigma \sqrt{\frac{d}{n}}$ with prob. $1 - \delta$. Then there*
 439 *exists a constant c such that the KL divergence $D_{KL}(\mathbb{Q}(\cdot | X_1^n) \| \mathcal{N}(\hat{\theta}^{(n)}, (nH^*)^{-1}))$ is at most*

$$440 \quad c \cdot \frac{1}{\lambda_{\min}(H^*)} \left(\frac{A^2 d^2 \sigma^4}{n} + \frac{\varepsilon_1^{(2)}(n, \delta)^2 d^2 \sigma^4}{n} + \sigma^2 \left(\varepsilon_2^{(2)}(n, \delta)^2 + \frac{L_2^2}{n^2} \right) d \right) \quad \text{with prob. at least } 1 - 2\delta.$$

442 See [Appendix A.2](#) for the proof of this claim.

443 A few remarks are in order. First, assuming that the problem-dependent constants
 444 (A, σ, L_2) are of constant order, and that the deviation bound scales as $\varepsilon_2^{(2)}(n, \delta) = O(1/\sqrt{n})$,
 445 [Proposition 3.4](#) shows that the KL divergence between the posterior distribution and the
 446 Gaussian limit is of order $O(1/n)$; second, the non-asymptotic behavior of posterior distribution
 447 depends on the Hessian matrix $H^* = -\nabla^2 F(\theta^*)$. In the well-specified case where the data points
 448 X_1^n are i.i.d. samples from the distribution \mathbb{P}_{θ^*} , the standard Fisher-information identity $H^* =$
 449 $\mathbb{E}_{\theta^*} [\nabla \log p_{\theta^*}(X) \nabla \log p_{\theta^*}(X)^\top]$ holds true, and the Bayesian credible set is asymptotically
 450 the same as the confidence set in the frequentist sense. On the other hand, in the mis-specified
 451 models where $\theta^* = \arg \min_{\theta \in \Theta} D_{KL}(\mathbb{P} \| \mathbb{P}_\theta)$, the limiting Gaussian law is $\mathcal{N}(\hat{\theta}^{(n)}, (nH^*)^{-1})$,
 452 depending on the Hessian matrix but not the covariance of the log-likelihood. This result
 453 coincides with the asymptotic Bernstein–von Mises theorem for mis-specified parametric
 454 models [\[23\]](#), providing a non-asymptotic characterization. Using Pinsker’s inequality and
 455 Talagrand’s T_2 -inequality [\[46\]](#), the KL divergence bound can also be transformed into bounds
 456 in term of total variation and Wasserstein-2 distances, yielding a non-asymptotic $O(1/\sqrt{n})$
 457 rate of convergence.

458 We can also use the diffusion process approach to derive more fine-grained concentration
 459 bounds for the posterior distribution, with behavior matching the limiting Gaussian law. Doing
 460 so requires the following stronger version of the posterior contraction condition:

$$461 \quad (3.6) \quad \left(\mathbb{E}_{\mathbb{Q}} \left[\|\theta - \theta^*\|_2^{2p} | X_1^n \right] \right)^{1/p} \leq \frac{\sigma^2 p d}{n}, \quad \text{for all } p > 0 \text{ with probability at least } 1 - \delta.$$

463 In addition, we define the function

$$464 \quad \mathcal{H}_n(t, \delta) := (A + \varepsilon_1^{(2)}(n, \delta))^2 \cdot \frac{\sigma^4 d^2 t^2}{n^2} + \frac{\sigma d}{n} \left(\varepsilon_2^{(2)}(n, \delta)^2 + \frac{L_2^2}{n^2} + (A + \varepsilon_1^{(2)}(n, \delta))^2 \frac{\sigma d}{n} \right),$$

466 which plays the role of a higher-order term. Equipped with this notation, we have:

467 **Theorem 3.5.** *Suppose that conditions (BvM.1), (BvM.2), and (PS) are in force, the*
 468 *Hessian H^* is strictly positive definite, and the high-probability posterior contraction con-*
 469 *dition [\(3.6\)](#) holds. Then for any $\delta \in (0, 1)$, uniformly over all $\omega \in (0, 1)$ and $t > 0$, we*
 470 *have*

$$471 \quad (3.7) \quad \mathbb{Q} \left(\left\| \theta - \hat{\theta}^{(n)} \right\|_{H^*}^2 \geq (1 + \omega) \frac{d}{n} + c \frac{1 + \log \kappa(H^*)}{\omega} \left(\frac{t}{n} + \mathcal{H}_n(t, \delta) \right) \mid X_1^n \right) \leq e^{-t},$$

473 *with probability at least $1 - \delta$.*

474 See [Appendix A.1](#) for the proof of the theorem.

475 A few remarks are in order. Note that the limiting Gaussian density $\gamma_n = \mathcal{N}(0, (nH^*)^{-1})$
 476 satisfies a tail bound of the form $\gamma_n\left(\|\theta - \hat{\theta}^{(n)}\|_{H^*}^2 \geq \frac{d}{n} + \frac{t}{n}\right) \leq e^{-t/2}$ for any $t > 0$. Unless the
 477 posterior is actually Gaussian in finite samples, it cannot satisfy this bound exactly. However,
 478 [Theorem 3.5](#) provides a bound with near-matching behavior: note that the leading-order term
 479 scales $\frac{d}{n}$, matching the asymptotics with a pre-factor $1 + \omega$ that can be made arbitrarily close
 480 to 1 (at the expense of the other term). The $\frac{t}{n}$ dependency on the tail probability comes
 481 with a mild $\log \kappa(H^*)$ factor due to technical reasons. The bound also contains a high-order
 482 term $\mathcal{H}_n(t, \delta)$, which scales as $O(n^{-2})$. It is also worth noticing that the terms in [Theorem 3.5](#)
 483 depend on the tail probability $\nu = e^{-t}$ only logarithmically, allowing for very small value of
 484 ν . We can therefore use equation (3.7) to construct non-asymptotic credible sets of ellipsoid
 485 shape, adapted to the geometry of local Hessian matrix H^* .

486 *Proof outline:* The proofs of both [Proposition 3.4](#) and [Theorem 3.5](#) rely on a first-order
 487 approximation of the gradient ∇F_n . In particular, the diffusion process (2.3) can be written in
 488 the form $d\theta_t = -\frac{1}{2}H^*(\theta_t - \hat{\theta}^{(n)})dt + \frac{1}{2}e_n(\theta_t)dt + \frac{1}{2n}\log \pi(\theta_t)dt + \frac{1}{\sqrt{n}}dB_t$, where we have defined the
 489 linearization error $e_n(\theta) := \nabla F_n(\theta) + H^*(\theta - \theta^*)$. Under the smoothness assumption ([BvM.1](#))
 490 and the empirical process bound ([BvM.2](#)), one can show that $\|e_n(\theta)\|_2 \leq \|\theta - \theta^*\|_2 \cdot O(\sqrt{d/n})$
 491 with high probability. When this error term is ignored, the diffusion process is an Ornstein–
 492 Uhlenbeck process whose stationary distribution is $\mathcal{N}(\hat{\theta}^{(n)}, (nH^*)^{-1})$. Therefore, given the
 493 non-asymptotic bounds on the error $e_n(\theta)$ stated above, we can provide a non-asymptotic
 494 characterization of the distance between the stationary distribution and the limiting Gaussian
 495 law. In order to prove [Proposition 3.4](#), we use the Gaussian log-Sobolev inequality [19] to
 496 control the KL divergence, whereas proving [Theorem 3.5](#) is based on using Itô calculus to
 497 study the growth of a Lyapunov function defined using the metric induced by H^* . Full proofs
 498 for the two results are given in [Appendix A.2](#) and [Appendix A.1](#), respectively.

499 **4. Some illustrative examples.** Having developed some general theory, we now use it
 500 to derive some concrete results for two examples of interest in statistical analysis: Bayesian
 501 logistic regression and Gaussian mixture models.

502 **4.1. Bayesian logistic regression.** Logistic regression is a classical way of modelling the
 503 relationship between a binary response variable $Y \in \{-1, +1\}$ and a vector $X \in \mathbb{R}^d$ of
 504 explanatory variables (e.g., see the book [30]). In the logistic regression model, the pair (X, Y)
 505 are related by the conditional distribution

$$506 \quad (4.1) \quad \mathbb{P}(Y = 1 \mid X, \theta) = \frac{e^{\langle X, \theta \rangle}}{1 + e^{\langle X, \theta \rangle}}, \quad \text{where } \theta \in \mathbb{R}^d \text{ is a parameter vector.}$$

508 Suppose that we observe a collection $Z_1^n = \{Z_i\}_{i=1}^n$ of n i.i.d paired samples $Z_i = (X_i, Y_i)$,
 509 each generated in the following way. First, the covariate vector X_i is drawn from a standard
 510 Gaussian distribution $N(0, I_d)$, and then the binary response Y_i is drawn according to the
 511 conditional distribution $\mathbb{P}(\cdot \mid X_i, \theta^*)$ from equation (4.1), where $\theta^* \in \mathbb{R}^d$ is a fixed but unknown
 512 value of the parameter vector. Given these assumptions, the sample log-likelihood function
 513 of the samples Z_1^n takes the form $F_n^R(\theta) := \frac{1}{n} \sum_{i=1}^n \{\log \mathbb{P}(Y_i \mid X_i, \theta) + \log \phi(X_i)\}$, where ϕ
 514 denotes the density of a standard normal vector. Combining this log-likelihood with a given

515 prior π over θ yields the posterior distribution in the usual way. We assume that the prior
 516 function π satisfies Assumption **(A)**, and recall the constant B defined in this assumption.
 517 Throughout this section, we also assume that the norm $\|\theta^*\|_2$ is a universal constant independent
 518 of (n, d) , and we suppress the dependence on this parameter.

519 With this set-up, the following result establishes the posterior convergence rate of θ around
 520 θ^* , conditionally on the observations Z_1^n .

521 **Corollary 4.1.** *For any $\delta \in (0, 1)$, given $\frac{n}{\log n} \geq c'd \log(\frac{1}{\delta})$ i.i.d. samples from the Bayesian
 522 logistic regression model (4.1), we have $\mathbb{Q}\left(\|\theta - \theta^*\|_2 \geq c\left\{\sqrt{\frac{d}{n}} + \sqrt{\frac{\log(1/\delta)}{n}} + \frac{B}{n}\right\} \mid Z_1^n\right) \leq \delta$
 523 with probability $1 - \delta$ over the data Z_1^n .*

524 See [Appendix B.1](#) for the proof of this claim.

525 A few comments are in order. First, the result of [Corollary 4.1](#) shows that for Bayesian
 526 logistic regression model (4.1), the posterior convergence rate for the parameter is of the
 527 order $(d/n)^{1/2}$. Furthermore, this result also gives a concrete dependence of the rate on B
 528 characterizing the degree to which the prior is concentrated away from the true parameter. By
 529 taking the standard Gaussian prior $\pi = \mathcal{N}(0, I_d)$, we have $B \leq \|\theta\|_2$, which is bounded by a
 530 universal constant independent of the pair (n, d) .

531 It is important to note that [Corollary 4.1](#) is valid as long as the sample size n is mildly
 532 larger than the problem dimension d (up to logarithmic factors). To our knowledge, this is the
 533 first time that a sharp non-asymptotic posterior contraction result is established in this regime.

534 Let us sketch how [Theorem 3.2](#) can be applied so as to prove this corollary. Denote
 535 $F^R := \mathbb{E}[F_n^R]$ as the population-level log-likelihood function. The first step in our proof, as
 536 given in [Appendix B.1](#), is to show that there are universal constants c, c_1, c_2 such that

$$537 \quad (4.2a) \quad -\langle \nabla F^R(\theta), \theta - \theta^* \rangle \geq c_1 \begin{cases} \|\theta - \theta^*\|_2^2, & \text{for all } \|\theta - \theta^*\|_2 \leq 1 \\ \|\theta - \theta^*\|_2, & \text{otherwise} \end{cases}, \quad \text{and}$$

$$538 \quad (4.2b) \quad \sup_{\theta \in \mathbb{R}^d} \|\nabla F_n^R(\theta) - \nabla F^R(\theta)\|_2 \leq c_2 \left(\sqrt{\frac{d}{n}} + \sqrt{\frac{\log(1/\delta)}{n}} + \frac{\log(1/\delta)}{n} \right),$$

540 for any $r > 0$ with probability $1 - \delta$ as long as $\frac{n}{\log n} \geq cd \log(1/\delta)$. Using these results, we show
 541 that Assumptions [\(W.1\)](#) and [\(W.2\)](#) hold with

$$542 \quad (4.3) \quad \psi(r) = c_1 \begin{cases} r^2 & \text{for all } r \in (0, 1), \text{ and} \\ r & \text{otherwise} \end{cases}, \quad \text{and} \quad \zeta(r) = c_2 \quad \text{for all } r > 0.$$

544 We can check that the functions ψ and ζ satisfy the conditions in Assumptions [\(W.3\)](#)
 545 and [\(W.4\)](#). Therefore, applying [Theorem 3.2](#) to these functions yields the posterior contraction
 546 rate claimed in [Corollary 4.1](#). See [Appendix B.1](#) for the details.

547 **4.2. Over-specified Bayesian Gaussian mixture models.** Gaussian mixtures are widely
 548 used for modelling heterogeneous datasets; clusters in the data are naturally associated with
 549 different mixture components [26]. In fitting such models, the true number of components is
 550 generally unknown, and several approaches have been proposed to deal with this challenge.
 551 One of the most popular methods is to deliberately include a large number of components,

552 leading to what are known as overspecified Gaussian mixture models [39]. While the behavior
 553 of posterior densities in such mixture models is relatively well-understood [17], the behavior
 554 of the posterior in terms of its parametric components is not as well understood. When the
 555 covariance matrices are known and the parameter space is bounded, the location parameters
 556 have been shown to have posterior convergence rates of the order $n^{-1/4}$ in the Wasserstein-2
 557 metric [32]. However, neither the dependence on dimension d nor on the true number of
 558 components have been established.

559 In this section, we consider the behavior of overspecified Gaussian mixture models in a
 560 particular setting, and provide convergence rates for the parameters with precise dependence
 561 on the dimension d , and without requiring any boundedness assumption. In order to model the
 562 simplest form of over-specification, suppose that we fit a Bayesian location mixture model to a
 563 collection of i.i.d. samples $X_1^n = (X_1, \dots, X_n)$ drawn from a Gaussian distribution $\mathcal{N}(\theta^*, I_d)$.
 564 (For concreteness, we set $\theta^* = 0$.) We study the behavior of the Bayesian Gaussian mixture
 565 model

$$566 \quad (4.4) \quad \theta \sim \pi(\cdot), \quad V_i \in \{-1, 1\} \stackrel{\text{i.i.d.}}{\sim} \text{Cat}(1/2, 1/2), \quad X_i | V_i, \theta \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(V_i \theta, I_d),$$

568 where $\text{Cat}(1/2, 1/2)$ stands for the categorical distribution with parameters $(1/2, 1/2)$. We
 569 assume that the prior π satisfies the smoothness condition (cf. Assumption **(A)**); one example
 570 is a Gaussian distribution (over the location parameter θ). Our goal in this section is to
 571 characterize the posterior contraction rate of the location parameter θ around θ^* .

In order to do so, we first define the sample log-likelihood function F_n^G given data X_1^n .
 It has the form $F_n^G(\theta) := \frac{1}{n} \sum_{i=1}^n \log \left(\frac{1}{2} \phi(X_i; -\theta, I_d) + \frac{1}{2} \phi(X_i; \theta, I_d) \right)$, where $x \mapsto \phi(x; \theta, I_d) =$
 $(2\pi)^{-d/2} e^{-\|x-\theta\|_2^2/2}$ denotes the density of multivariate Gaussian distribution $\mathcal{N}(\theta, \sigma^2 I_d)$. Simi-
 larly, the population log-likelihood function is given by

$$F^G(\theta) := \mathbb{E}_X \left[\log \left(\frac{1}{2} \phi(X; -\theta, I_d) + \frac{1}{2} \phi(X; \theta, I_d) \right) \right],$$

572 where the outer expectation in the above display is taken with respect to $X \sim \mathcal{N}(\theta^*, I_d)$.

573 In **Appendix B.2**, we prove that there is a universal constant $c_1 > 0$ such that

$$574 \quad (4.5a) \quad -\langle \nabla F^G(\theta), \theta - \theta^* \rangle \geq \begin{cases} c_1 \|\theta - \theta^*\|_2^4, & \text{for all } \|\theta - \theta^*\|_2 \leq \sqrt{2} \\ 4c_1 \left(\|\theta - \theta^*\|_2^2 - 1 \right), & \text{otherwise} \end{cases},$$

575
 576 and moreover, there are universal constants (c, c_2) such that for any $\delta \in (0, 1)$, given a sample
 577 size $n \geq cd \log(1/\delta)$, we have

$$578 \quad (4.5b) \quad \sup_{\theta \in \mathbb{B}(\theta^*, r)} \|\nabla F_n^G(\theta) - \nabla F^G(\theta)\|_2 \leq c_2 \left(r + \frac{1}{\sqrt{n}} \right) \left(\sqrt{\frac{d}{n}} + \sqrt{\frac{\log(\log(n/\delta))}{n}} \right) \quad \text{with prob. } 1 - \delta.$$

580 Given the above results, the functions ψ and ζ in Assumptions **(W.1)** and **(W.2)** take
 581 the form

$$582 \quad (4.6) \quad \psi(r) = \begin{cases} c_1 r^4, & \text{for all } 0 < r \leq \sqrt{2} \\ 4c_1 (r^2 - 1), & \text{otherwise} \end{cases}, \quad \text{and} \quad \zeta(r) = r + \frac{1}{\sqrt{n}} \quad \text{for all } r > 0.$$

583

584 These functions satisfy the conditions of Assumptions **(W.3)** and **(W.4)**. Therefore, it leads to
 585 the following result regarding the posterior contraction rate of parameters under overspecified
 586 Bayesian location Gaussian mixtures **(4.4)**:

587 **Corollary 4.2.** *Given the overspecified Bayesian location Gaussian mixture model **(4.4)**, there*
 588 *are universal constants c, c' such that given any $\delta \in (0, 1)$ and a sample size $n \geq c'd \log(1/\delta)$,*
 589 *we have $\mathbb{Q}\left(\|\theta - \theta^*\|_2 \geq c\left(\frac{d}{n} + \frac{\log(\log(n/\delta))}{n}\right)^{1/4} + \left(\frac{B}{n}\right)^{1/3} \mid X_1^n\right) \leq \delta$ with probability $1 - \delta$ over*
 590 *the data X_1^n . Here, B is the non-negative constant in Assumption **(A)**.*

591 See **Appendix B.2** for the proof of **Corollary 4.2**.

592 The $O(n^{-1/4})$ rate of convergence in **Corollary 4.2** is consistent with the previous result
 593 with location parameters in overspecified Bayesian location Gaussian mixtures **[5, 22, 32]**, which
 594 is also known to be minimax optimal **[20]**. When taking the problem dimension into account,
 595 to our knowledge, the $(d/n)^{1/4}$ posterior contraction rate is a novel result, and matches existing
 596 analyses for frequentist methods **[12]**. Similar to the logistic regression case, **Corollary 4.2**
 597 only requires the sample size n to be mildly larger than the dimension d . The non-asymptotic
 598 posterior contraction results are also established for the first time in such a regime. Finally,
 599 our result does not require the boundedness of the parameter space, in contrast to past
 600 work **[5, 22, 32]**.

601 **5. Proofs.** In this section, we collect the proofs of the main theorems.

602 **5.1. Proof of Theorem 3.1.** Throughout the proof, in order to simplify notation, we
 603 omit the conditioning on the σ -field $\mathcal{F}_n := \sigma(X_1^n)$; it should be taken as given. Introduce the
 604 quantity $\alpha = \frac{1}{2}\mu - \varepsilon_1(n, \delta) > \frac{\mu}{6}$. Our proof relies on proving the following auxiliary bound

$$605 \quad (5.1) \quad \frac{1}{2}e^{\alpha t} \|\theta_t - \theta^*\|_2^2 \leq \frac{1}{\sqrt{n}}M_t + U_n \frac{(e^{\alpha t} - 1)}{2\alpha},$$

607 where $U_n := \frac{3B^2}{n^2} + \frac{3\varepsilon_2^2(n, \delta)}{\mu} + \frac{d}{n}$ and $M_t := \int_0^t e^{\alpha s} \langle \theta_s - \theta^*, dB_s \rangle$. By construction, the latter
 608 term is a martingale.

609 The proof of the bound **(5.1)** is given later in this section; we take it as given for the
 610 moment, and use it to prove the theorem. In order to bound the moments of martingale M_t ,
 611 for any $p \geq 4$, we invoke the Burkholder–Davis–Gundy inequality (e.g., §4.4 of the book **[36]**)
 612 to find that

$$613 \quad \mathbb{E} \left[\sup_{0 \leq t \leq T} |M_t|^{\frac{p}{2}} \right] \leq (pC)^{\frac{p}{4}} \mathbb{E} \left[[M]_T^{\frac{p}{4}} \right] = (pC)^{\frac{p}{4}} \mathbb{E} \left(\int_0^T e^{2\alpha s} \|\theta_s - \theta^*\|_2^2 ds \right)^{\frac{p}{4}}$$

$$614 \quad \leq (pC)^{\frac{p}{4}} \mathbb{E} \left(\sup_{0 \leq t \leq T} e^{\alpha t} \|\theta_t - \theta^*\|_2^2 \int_0^T e^{\alpha s} ds \right)^{\frac{p}{4}} \leq \left(\frac{pC e^{\alpha T}}{\alpha} \right)^{\frac{p}{4}} \mathbb{E} \left(\sup_{0 \leq t \leq T} e^{\alpha t} \|\theta_s - \theta^*\|_2^2 \right)^{\frac{p}{4}},$$

615
616

617 where C is a universal constant. Therefore, we arrive at the bound

$$\begin{aligned}
 618 \quad \mathbb{E} \left[\left(\sup_{0 \leq t \leq T} e^{\alpha t} \|\theta_t - \theta^*\|_2 \right)^p \right] &\leq \mathbb{E} \left(\frac{2}{\sqrt{n}} M_t \right)^{\frac{p}{2}} + \left(U_n \frac{(e^{\alpha T} - 1)}{\alpha} \right)^{\frac{p}{2}} \\
 619 \quad &\leq \left(U_n \frac{e^{\alpha T}}{\alpha} \right)^{\frac{p}{2}} + \left(\frac{pC e^{\alpha T}}{\alpha n} \right)^{\frac{p}{4}} \mathbb{E} \left(\sup_{0 \leq s \leq T} e^{\alpha s} \|\theta_s - \theta^*\|_2^2 \right)^{\frac{p}{4}}. \\
 620
 \end{aligned}$$

621 For the right hand side of the above inequality, we can relate it to the left hand side by using
 622 Young's inequality, which is given by

$$\begin{aligned}
 623 \quad \left(\frac{pC e^{\alpha T}}{\alpha n} \right)^{\frac{p}{4}} \mathbb{E} \left(\sup_{0 \leq s \leq T} e^{\alpha s} \|\theta_s - \theta^*\|_2^2 \right)^{\frac{p}{4}} &\leq \frac{1}{2} \left(\frac{pC e^{\alpha T}}{\alpha n} \right)^{\frac{p}{2}} + \frac{1}{2} \mathbb{E} \left(\sup_{0 \leq s \leq T} e^{\alpha s} \|\theta_s - \theta^*\|_2^2 \right)^{\frac{p}{2}}. \\
 624
 \end{aligned}$$

625 Putting the above results together, and let $\alpha = \frac{\mu}{2}$, we find that

$$\begin{aligned}
 626 \quad (\mathbb{E} [\|\theta_T - \theta^*\|_2^p])^{\frac{1}{p}} &\leq e^{-\alpha T} \left(\mathbb{E} \sup_{0 \leq t \leq T} (e^{\alpha t} \|\theta_t - \theta^*\|_2^p) \right)^{\frac{1}{p}} \leq C' \left(\sqrt{\frac{U_n}{\mu}} + \sqrt{\frac{2p}{n\mu}} \right), \\
 627
 \end{aligned}$$

628 for universal constant $C' > 0$. Therefore, the diffusion process (2.3) satisfies the bound

$$\begin{aligned}
 629 \quad \sup_{t \geq 0} (\mathbb{E} [\|\theta_t - \theta^*\|_2^p])^{\frac{1}{p}} &\leq c \left(\sqrt{\frac{d}{\mu n}} + \frac{B}{\mu n} + \frac{\varepsilon_2(n, \delta)}{\mu} + \sqrt{\frac{p}{n\mu}} \right) \quad \text{for any } p \geq 1. \\
 630
 \end{aligned}$$

631 Combining the above inequality with the inequality (5.3) yields the conclusion of the theorem.

632 **5.1.1. Proof of claim (5.1).** For the given choice $\alpha > 0$, an application of Itô's formula
 633 yields the decomposition

$$\begin{aligned}
 634 \quad \frac{1}{2} e^{\alpha t} \|\theta_t - \theta^*\|_2^2 &= -\frac{1}{2} \int_0^t \langle \theta^* - \theta_s, \nabla F_n(\theta_s) e^{\alpha s} \rangle ds + \frac{1}{2n} \int_0^t \langle \theta_s - \theta^*, \nabla \log \pi(\theta_s) e^{\alpha s} \rangle ds \\
 635 \quad &+ \frac{d}{2n} \int_0^t e^{\alpha s} ds + \frac{1}{\sqrt{n}} \int_0^t e^{\alpha s} \langle \theta_s - \theta^*, dB_s \rangle + \frac{1}{2} \int_0^t \alpha e^{\alpha s} \|\theta_s - \theta^*\|_2^2 ds \\
 636 \quad (5.2) \quad &= J_1 + J_2 + J_3 + J_4 + J_5.
 \end{aligned}$$

638 We begin by bounding the term J_1 in equation (5.2). Based on Assumption (S.2) regarding
 639 the perturbation error between F_n and F and the strong convexity of F , we have

$$\begin{aligned}
 640 \quad J_1 &= -\frac{1}{2} \int_0^t \langle \theta^* - \theta_s, \nabla F_n(\theta_s) e^{\alpha s} \rangle ds \\
 641 \quad &\leq -\frac{1}{2} \int_0^t \langle \theta^* - \theta_s, \nabla F(\theta_s) e^{\alpha s} \rangle ds + \frac{1}{2} \int_0^t \|\theta_s - \theta^*\|_2 \|\nabla F(\theta_s) - \nabla F_n(\theta_s)\|_2 e^{\alpha s} ds \\
 642 \quad &\leq -\frac{1}{2} \int_0^t \mu \|\theta_s - \theta^*\|_2^2 e^{\alpha s} ds + \frac{1}{2} \int_0^t \|\theta_s - \theta^*\|_2 (\varepsilon_1(n, \delta) \|\theta_s - \theta^*\|_2 + \varepsilon_2(n, \delta)) e^{\alpha s} ds \\
 643 \quad &\leq -\frac{1}{2} \int_0^t \mu \|\theta_s - \theta^*\|_2^2 e^{\alpha s} ds + \frac{1}{2} \int_0^t \|\theta_s - \theta^*\|_2^2 (\varepsilon_1(n, \delta) + \mu/3) e^{\alpha s} ds + \frac{3\varepsilon_2^2(n, \delta)}{2\mu} \int_0^t e^{\alpha s} ds. \\
 644
 \end{aligned}$$

645 The second term J_2 involving prior π can be controlled in the following way:

646

$$647 \quad J_2 = \frac{1}{2n} \int_0^t \langle \theta_s - \theta^*, \nabla \log \pi(\theta_s) e^{\alpha s} \rangle ds \leq \frac{1}{2n} \int_0^t B \|\theta_s - \theta^*\|_2 e^{\alpha s} ds$$

$$648 \quad \leq \int_0^t \frac{\mu}{6} \|\theta_s - \theta^*\|_2^2 e^{\alpha s} ds + \frac{3B^2}{n^2 \mu} \int_0^t e^{\alpha s} ds.$$

649

650 For the third term J_3 , a direct calculation leads to

$$651 \quad J_3 = \frac{d(e^{\alpha t} - 1)}{2\alpha n}.$$

652

653 Moving to the fourth term $J_4 = M_t/\sqrt{n}$, it is a martingale (since M_t is a martingale). Putting
654 the above results together and noting that $\alpha = \frac{1}{2}\mu - \varepsilon_1(n, \delta) > \frac{\mu}{6}$, we obtain the bound

$$655 \quad \frac{1}{2} e^{\alpha t} \|\theta_t - \theta^*\|_2^2 \leq \frac{1}{\sqrt{n}} M_t + U_n \frac{(e^{\alpha t} - 1)}{2\alpha}.$$

656

657 Putting together the pieces yields the claim (5.1).

658 **5.2. Proof of Theorem 3.2.** As in the proof of Theorem 3.1, we omit the conditioning on
659 $\mathcal{F}_n := \sigma(X_1^n)$. For any $p \geq 2$, we define the functions on the positive real line $(0, \infty)$ as

$$660 \quad \nu_{(p)}(r) := \psi\left(r^{\frac{1}{p-1}}\right) r^{\frac{p-2}{p-1}}, \quad \text{and} \quad \tau_{(p)}(r^{p-1}\zeta(r)) := r^{p-2}\psi(r).$$

661

662 Note that $\tau_{(p)}$ is defined implicitly; let us verify that this definition is meaningful. By
663 Assumption (W.2), the function $r \mapsto r^{p-1}\zeta(r)$ is a strictly increasing and surjective, mapping
664 from $[0, +\infty)$ to $[0, +\infty)$. Therefore, it is invertible, which ensures that the function $\tau_{(p)}$ is
665 well-defined.

666 Now we claim that for any $p \geq 2$, the functions $\nu_{(p)}$ and $\tau_{(p)}$ are convex and strictly
667 increasing, and that furthermore, the expectation $\mathbb{E}[\|\theta_t - \theta^*\|_2^p]$ is upper bounded by the
668 integral

$$669 \quad (5.3) \quad \frac{p}{2} \int_0^t \left(-R_p(s) + \varepsilon(n, \delta) \tau_{(p)}^{-1}(R_p(s)) + \frac{B}{n} \nu_{(p)}^{-1}(R_p(s)) + \frac{p-1+d}{n} \nu_{(p)}^{-1}(R_p(s))^{\frac{p-2}{p-1}} \right) ds,$$

670

671 where $R_p(s) := \mathbb{E} \left[\|\theta_s - \theta^*\|_2^{p-2} \psi(\|\theta_s - \theta^*\|_2) \right]$.

672 Taking the above claims as given for the moment, let us now complete the proof of the
673 theorem. Since for each finite $q \geq 1$, the process $(\theta_t : t \geq 0)$ converges in \mathbb{L}^q norm, the limit
674 $\lim_{t \rightarrow +\infty} R_p(t)$ exists. Since the functions $\tau_{(p)}$ and $\nu_{(p)}$ are convex and strictly increasing, their
675 inverse functions are concave. Moreover, simple calculation leads to

$$676 \quad (5.4) \quad \nabla_r \left(\nu_{(p)}^{-1}(r)^{\frac{p-2}{p-1}} \right) = \frac{p-2}{p-1} \cdot \frac{\nu_{(p)}^{-1}(r)^{-\frac{1}{p-1}}}{\nu'_{(p)}(\nu_{(p)}^{-1}(r))}.$$

677

678 Since $\nu_{(p)}$ is convex and increasing, the numerator is a decreasing positive function of r .
679 Additionally, the denominator is an increasing positive function of r . Therefore, the derivative

680 in equation (5.4) is a decreasing function of r , and the function $r \mapsto \nu_{(p)}^{-1}(r)^{\frac{p-2}{p-1}}$ is concave.
 681 Define the function

$$682 \quad \phi(r) := -r + \varepsilon(n, \delta)\tau_{(p)}^{-1}(r) + \frac{B}{n}\nu_{(p)}^{-1}(r) + \frac{p-1+d}{n}\nu_{(p)}^{-1}(r)^{\frac{p-2}{p-1}},$$

684 and observe that ϕ is concave and $\phi(0) = 0$. Let r_* be the smallest positive solution to the
 685 equation

$$686 \quad r = \varepsilon(n, \delta)\tau_{(p)}^{-1}(r) + \frac{B}{n}\nu_{(p)}^{-1}(r) + \frac{p-1+d}{n}\nu_{(p)}^{-1}(r)^{\frac{p-2}{p-1}}.$$

688 We then have $\phi(r) < 0$ for $r > r_*$ and $\phi(r) > 0$ for $r \in (0, r_*)$. By Lemma C.1, we have
 689 $\lim_{t \rightarrow +\infty} R_p(t) \leq r_*$.

690 Since $\nu_{(p)}$ is a convex and strictly increasing function, Jensen's inequality implies that

$$691 \quad (5.5) \quad R_p(t) = \mathbb{E} \left(\|\theta_t - \theta^*\|_2^{p-2} \psi(\|\theta_t - \theta^*\|_2) \right) \geq \nu_{(p)} \left(\mathbb{E} \|\theta_t - \theta^*\|_2^{p-1} \right).$$

693 Therefore, if we define $z_* := \lim_{t \rightarrow +\infty} \left(\mathbb{E} \|\theta_t - \theta^*\|_2^{p-1} \right)^{\frac{1}{p-1}}$, we have $z_*^{p-1} \leq \nu_{(p)}^{-1}(r_*)$. Hence,
 694 we arrive at the following inequality

$$695 \quad z_*^{p-2} \psi(z_*) \leq \varepsilon(n, \delta)\tau_{(p)}^{-1}(\nu_{(p)}(z_*^{p-1})) + \frac{B}{n}z_*^{p-1} + \frac{p-1+d}{n}z_*^{p-2}$$

$$696 \quad = \varepsilon(n, \delta)z_*^{p-1}\zeta(z_*) + \frac{B}{n}z_*^{p-1} + \frac{p-1+d}{n}z_*^{p-2}.$$

698 As a consequence, we find that

$$699 \quad \psi(z_*) \leq \varepsilon(n, \delta)\zeta(z_*)z_* + \frac{B + (p-1)d}{n}.$$

701 In Appendix C.4 of the supplementary material [31], we prove the existence and uniqueness of
 702 the positive solution to the non-linear equation (3.1). Given this claim, replacing p by $(p+1)$
 703 and putting the above results together yields

$$704 \quad \lim_{t \rightarrow +\infty} \left(\mathbb{E} (\|\theta_t - \theta^*\|_2^p) \right)^{\frac{1}{p}} \leq z_p^*,$$

706 where z_p^* is the unique positive solution to the following equation:

$$707 \quad \psi(z) = \varepsilon(n, \delta)\zeta(z)z + \frac{B}{n}z + \frac{p+d}{n}.$$

709 Combining the above inequality with the inequality (5.3) yields the conclusion of the theorem.
 710

711 We now return to prove our earlier claims about the behavior of the functions $\nu_{(p)}$, $\tau_{(p)}$, the
 712 moment bound (5.3), and the existence of unique positive solution to equation (3.1).

713 **5.2.1. Structure of the function $\nu_{(p)}$.** Since ψ is a convex and strictly increasing function,
714 by taking the second derivative, we find that

$$\begin{aligned} 715 \quad \nu_{(p)}''(r) &= \nabla_r^2 \left(\psi \left(r^{\frac{1}{p-1}} \right) r^{\frac{p-2}{p-1}} \right) \\ 716 \quad &= \frac{1}{p-1} r^{\frac{1}{p-1}-1} \psi'' \left(r^{\frac{1}{p-1}} \right) + \frac{1}{p-1} r^{-1} \left(\psi' \left(r^{\frac{1}{p-1}} \right) - r^{-\frac{1}{p-1}} \psi \left(r^{\frac{1}{p-1}} \right) \right) \geq 0 \\ 717 \end{aligned}$$

718 for all $r > 0$. As a consequence, the function $\nu_{(p)}$ is convex.

719 **5.2.2. Structure of the function $\tau_{(p)}$.** The proof is by calculating the second derivative of
720 the function $\tau_{(p)}$, and we make use of Assumption **(W.3)** on the functions ψ and ζ . Recall
721 that $\tau_{(p)}(r^{p-1}\zeta(r)) = r^{p-2}\psi(r)$ for any $r > 0$. Taking derivatives with respect to r on both
722 sides, we find that

$$723 \quad [(p-1)r^{p-2}\zeta(r) + r^{p-1}\zeta'(r)] \tau_{(p)}'(r^{p-1}\zeta(r)) = (p-2)r^{p-3}\psi(r) + r^{p-2}\psi'(r). \\ 724$$

725 Under the substitution $z = \zeta_{(p)}(r)$, we find that $\nabla_z \tau_{(p)}(z) = \frac{(p-2)\psi(r) + r\psi'(r)}{(p-1)r\zeta(r) + r^2\zeta'(r)}$.

726 Taking another derivative of the above term, we find that

$$727 \quad \nabla_z^2 \tau_{(p)}(z) = \left(\zeta_{(p)}'(r) \right)^{-1} \frac{g(r, p)}{\left((p-1)r\zeta(r) + r^2\zeta'(r) \right)^2}, \\ 728$$

729 where we denote

$$\begin{aligned} 730 \quad g(r, p) &:= [(p-1)r\zeta(r) + r^2\zeta'(r)] \cdot [(p-1)\psi'(r) + r\psi''(r)] \\ 731 \quad &\quad - [(p-1)\zeta(r) + (p+1)r\zeta'(r) + r^2\zeta''(r)] \cdot [(p-2)\psi(r) + r\psi'(r)]. \end{aligned}$$

733 According to Assumption **(W.3)**, the function $\tau_{(2)} = \psi_{(2)} \circ \zeta_{(2)}^{-1}$ is convex. Therefore, we have
734 $g(r, 2) \geq 0$ for any $r > 0$. Simple algebra with first order derivative of function g with respect
735 to parameter p leads to

$$\begin{aligned} 736 \quad \nabla_p(g(r, p)) &= \zeta(r) \cdot [(p-1)r\psi'(r) + r^2\psi''(r) - (p-2)\psi(r) - r\psi'(r)] \\ 737 \quad &\quad - r\zeta'(r) [(p-2)\psi(r) + r\psi'(r)] + r\psi'(r) \cdot [(p-1)\zeta(r) + r\zeta'(r)] \\ 738 \quad &\quad - \psi(r) \cdot [(p-1)\zeta(r) + (p+1)r\zeta'(r) + r^2\zeta''(r)] \\ 739 \quad &= 2(p-2) [r\psi'(r)\zeta(r) - \psi(r)\zeta(r) - r\zeta'(r)\psi(r)] \\ 740 \quad &\quad + [r^2\zeta(r)\psi''(r) + r\psi'(r)\zeta(r) - 3\psi(r)\zeta(r) - r^2\psi(r)\zeta''(r)] \geq 0 \end{aligned}$$

742 for all $r > 0$. Here the last inequality follows from Assumption **(W.3)**. Therefore, the function
743 g is increasing function in terms of p when $p \geq 2$, so that $g(r, p) \geq g(r, 2) \geq 0$ for all $r > 0$.
744 Given this inequality, we have $\frac{d^2}{dz^2} \tau_{(p)}(z) \geq 0$ for any $z \geq 0$, $p \geq 2$, i.e., the function $\tau_{(p)}(z)$ is a
745 convex function for $z = \zeta_{(p)}(r)$.

746 **5.2.3. Proof of claim (5.3).** For any $p \geq 2$, an application of Itô's formula yields the
 747 bound $\|\theta_t - \theta^*\|_2^p \leq \sum_{j=1}^5 T_j$, where

$$748 \quad (5.6a) \quad T_1 := -\frac{p}{2} \int_0^t \langle \theta^* - \theta_s, \nabla F(\theta_s) \rangle \|\theta_s - \theta^*\|_2^{p-2} ds,$$

$$749 \quad (5.6b) \quad T_2 := \frac{p}{2} \int_0^t \langle \theta^* - \theta_s, \nabla F(\theta_s) - \nabla F_n(\theta_s) \rangle \|\theta_s - \theta^*\|_2^{p-2} ds$$

$$750 \quad (5.6c) \quad T_3 := \frac{p}{2n} \int_0^t \langle \theta_s - \theta^*, \nabla \log \pi(\theta_s) \rangle \|\theta_s - \theta^*\|_2^{p-2} ds$$

$$751 \quad (5.6d) \quad T_4 := p \int_0^t \|\theta_s - \theta^*\|_2^{p-2} \langle \theta_s - \theta^*, dB_s \rangle$$

$$752 \quad (5.6e) \quad T_5 := \frac{p(p-1+d)}{2n} \int_0^t \|\theta_s - \theta^*\|_2^{p-2} ds.$$

754 We now upper bound the terms $\{T_j\}_{j=1}^5$ in terms of functionals of the quantity R_p . From the
 755 weak convexity of F guaranteed by Assumption W.1, we have

$$756 \quad (5.7a) \quad \mathbb{E}[T_1] = -\frac{p}{2} \mathbb{E} \left[\int_0^t \langle \theta^* - \theta_s, \nabla F(\theta_s) \rangle \|\theta_s - \theta^*\|_2^{p-2} ds \right] \leq -\frac{p}{2} \int_0^t R_p(s) ds.$$

758 Based on Assumption (W.2), we find that

$$759 \quad \mathbb{E}[T_2] = \frac{p}{2} \mathbb{E} \left[\int_0^t \langle \theta^* - \theta_s, \nabla F(\theta_s) - \nabla F_n(\theta_s) \rangle \|\theta_s - \theta^*\|_2^{p-2} ds \right]$$

$$760 \quad \leq \frac{p}{2} \varepsilon(n, \delta) \int_0^t \mathbb{E} \left[\|\theta_s - \theta^*\|_2^{p-1} \zeta(\|\theta_s - \theta^*\|_2) \right] ds.$$

762 Since the function $\tau_{(p)}$ is convex, invoking Jensen's inequality, we obtain the following inequalities:
 763

$$764 \quad \int_0^t \mathbb{E} \left[\|\theta_s - \theta^*\|_2^{p-1} \zeta(\|\theta_s - \theta^*\|_2) \right] ds \leq \int_0^t \tau_{(p)}^{-1} \mathbb{E} \left[\tau_{(p)} \left(\|\theta_s - \theta^*\|_2^{p-1} \zeta(\|\theta_s - \theta^*\|_2) \right) \right] ds$$

$$765 \quad = \int_0^t \tau_{(p)}^{-1} (R_p(s)) ds.$$

767 In light of the above inequalities, we have

$$768 \quad (5.7b) \quad \mathbb{E}[T_2] \leq \frac{p}{2} \varepsilon(n, \delta) \int_0^t \tau_{(p)}^{-1} (R_p(s)) ds.$$

770 Moving to T_3 in equation (5.6c), given Assumption (A) which controls the growth of prior
 771 distribution π , its expectation is bounded as

$$772 \quad \mathbb{E}[T_3] = \frac{p}{2n} \mathbb{E} \left[\int_0^t \langle \theta_s - \theta^*, \nabla \log \pi(\theta_s) \rangle \|\theta_s - \theta^*\|_2^{p-2} ds \right]$$

$$773 \quad (5.7c) \quad \leq \frac{pB}{2n} \int_0^t \mathbb{E} \left[\|\theta_s - \theta^*\|_2^{p-1} \right] ds.$$

774

775 By exploiting the bound (5.5) along with the fact that $\nu_{(p)}$ is strictly increasing on $[0, +\infty)$,
776 we find that

$$777 \quad (5.7d) \quad \int_0^t \mathbb{E} \left(\|\theta_s - \theta^*\|_2^{p-1} \right) ds \leq \int_0^t \nu_{(p)}^{-1} (R_p(s)) ds.$$

779 Combining the inequalities (5.7c) and (5.7d), we have

$$780 \quad (5.7e) \quad \mathbb{E} [T_3] \leq \frac{pB}{2n} \int_0^t \nu_{(p)}^{-1} (R_p(s)) ds.$$

782 Moving to the fourth term T_4 from equation (5.6d), we have

$$783 \quad (5.7f) \quad \mathbb{E} [T_4] = \mathbb{E} \left[\int_0^t \|\theta_s - \theta^*\|_2^{p-2} \langle \theta_s - \theta^*, dB_s \rangle \right] = 0,$$

785 where we have used the martingale structure.

786 For the last term T_5 , invoking Hölder's inequality and the bound (5.5), we have the moment
787 estimate:

$$788 \quad \mathbb{E} \left(\|\theta_s - \theta^*\|_2^{p-2} \right) \leq \left(\mathbb{E} \left[\|\theta_s - \theta^*\|_2^{p-1} \right] \right)^{\frac{p-2}{p-1}} \leq \nu_{(p)}^{-1} (R_p(s))^{\frac{p-2}{p-1}}.$$

790 Consequently, the term T_5 can be bounded in expectation as

$$791 \quad (5.7g) \quad \mathbb{E} [T_5] \leq \frac{p(p-1+d)}{2n} \int_0^t \nu_{(p)}^{-1} (R_p(s))^{\frac{p-2}{p-1}} ds.$$

793 Collecting the bounds on the expectations of the terms $\{T_j\}_{j=1}^5$ from equations (5.7a)-(5.7g),
794 respectively, yields the claim (5.3).

795 **6. Discussion.** In this paper, we described an approach for analyzing the posterior con-
796 traction rates of parameters based on the diffusion processes. Our theory depends on two
797 important features: the local growth of the population log-likelihood function F and stochastic
798 perturbation bounds between the gradient of F and the gradient of its sample counterpart
799 F_n . For strongly concave log-likelihood functions, we established posterior convergence rates
800 for parameter estimation of the order $(d/n)^{1/2}$, valid under appropriate conditions on the
801 perturbation error between ∇F_n and ∇F and sharp sample size requirements. On the other
802 hand, when the population log-likelihood function is weakly concave, our analysis shows that
803 convergence rates are more delicate: they depend on an interaction between the degree of
804 weak convexity, and the stochastic error bounds. In this setting, we proved that the posterior
805 convergence rate of parameter is upper bounded by the unique positive solution of a non-linear
806 equation determined by the previous interplay. Compared to the convergence rate of MLE,
807 the bound contains an additional term capturing the diffusive behavior of the posterior dis-
808 tribution. Finally, we demonstrated the utility of the diffusion process approach by deriving
809 non-asymptotic forms of Bernstein–von Mises results for models with non-degenerate Fisher
810 information.

811 Let us now discuss a few directions that arise naturally from our work. First, in the
 812 weakly convex setting, although we have established non-asymptotic posterior contraction
 813 bounds, the current results do not provide information on the shape of the asymptotic posterior
 814 distribution. For example, when F is locally strongly concave around θ^* , it is well-known from
 815 the Bernstein–von Mises theorem that the posterior distribution of parameter converges to a
 816 multivariate normal distribution centered at the maximum likelihood estimation (MLE) with
 817 the covariance matrix is given by $1/(nI(\theta^*))$ (e.g., see the book [51], Chapter 10.2), where
 818 $I(\theta^*)$ denotes the Fisher information matrix at θ^* . When the F is only weakly concave, the
 819 Fisher information matrix $I(\theta^*)$ is degenerate, so that the posterior distribution can no longer
 820 be approximated by a multivariate Gaussian distribution. It is interesting to consider how the
 821 diffusion approach might provide insight into the shape of the posterior in this setting.

822 Second, the contraction rates given in this paper can give information about the over-
 823 specification of the latent variable models, thereby having potential applications for model
 824 selection. As a concrete example, for the symmetric two-component Gaussian mixture model
 825 example discussed in Subsection 4.2, the posterior distribution concentrates around $\theta^* = 0$
 826 at a rate $O((d/n)^{1/4})$ in the over-specified case. On the other hand, for a non-degenerate
 827 mixture with symmetric modes at θ^* and $-\theta^*$ (with $\theta^* \neq 0$), it concentrates at the usual rate
 828 $O((d/n)^{1/2})$. Consequently, the degree of dispersion in the posterior serves as an indicator of
 829 over-specification. Furthermore, since our results are non-asymptotic, they also give guidance
 830 on how this procedure could be performed with finite sample size n . Finally, whereas this
 831 paper focused on posterior contraction for parametric models, we suspect that the diffusion
 832 process approach used here might also be fruitfully applied to non-parametric models.

833 **Acknowledgements.** This work was partially supported by NSF grant CCF-1955450 and
 834 ONR grant N00014-21-1-2842 to MJW_j

835

REFERENCES

- 836 [1] D. BAKRY, F. BARTHE, P. CATTIAUX, AND A. GUILLIN, *A simple proof of the Poincaré inequality for a*
 837 *large class of probability measures*, Electronic Communications in Probability, 13 (2008), pp. 60–66.
 838 [2] S. BALAKRISHNAN, M. J. WAINWRIGHT, AND B. YU, *Statistical guarantees for the EM algorithm: From*
 839 *population to sample-based analysis*, Annals of Statistics, 45 (2017), pp. 77–120.
 840 [3] A. BARRON, M. SCHERVISH, AND L. WASSERMAN, *The consistency of posterior distributions in nonpara-*
 841 *metric problems*, Ann. Statist, 27 (1999), pp. 536–561.
 842 [4] A. BHATTACHARYA, D. PATI, , AND D. B. DUNSON, *Anisotropic function estimation using multi-bandwidth*
 843 *Gaussian processes*, Annals of Statistics, 42 (2014), pp. 352–381.
 844 [5] J. CHEN, *Optimal rate of convergence for finite mixture models*, Annals of Statistics, 23 (1995), pp. 221–233.
 845 [6] S. CHRÉTIEN, M. CUCURINGU, G. LECUÉ, AND L. NEIRAC, *Learning with semi-definite programming:*
 846 *statistical bounds based on fixed point analysis and excess risk curvature*, Journal of Machine Learning
 847 Research, 22 (2021).
 848 [7] A. S. DALALYAN, *Theoretical guarantees for approximate sampling from smooth and log-concave densities*,
 849 Journal of the Royal Statistical Society: Series B (Statistical Methodology), 79 (2017), pp. 651–676.
 850 [8] R. DE JONGE AND J. H. VAN ZANTEN, *Adaptive nonparametric Bayesian inference using location-scale*
 851 *mixture priors*, Annals of Statistics, 38 (2010), pp. 3300–3320.
 852 [9] J. L. DOOB, *Application of the theory of martingales*, Actes du Colloque International Le Calcul des
 853 Probabilités et ses applications (Lyon, 28 Juin– 3 Juillet, 1948), (1949), pp. 23–27.
 854 [10] A. DURMUS AND E. MOULINES, *Nonasymptotic convergence analysis for the unadjusted Langevin algorithm*,
 855 The Annals of Applied Probability, 27 (2017), pp. 1551–1587.

- 856 [11] A. DURMUS AND E. MOULINES, *High-dimensional Bayesian inference via the unadjusted Langevin algorithm*,
857 Bernoulli, 25 (2019), pp. 2854–2882.
- 858 [12] R. DWIVEDI, N. HO, K. KHAMARU, M. J. WAINWRIGHT, M. I. JORDAN, AND B. YU, *Singularity,*
859 *misspecification, and the convergence rate of EM*, arXiv preprint arXiv:1810.00828, (2018).
- 860 [13] D. A. FREEDMAN, *On the asymptotic behavior of Bayes’ estimates in the discrete case*, Annals of Statistics,
861 34 (1963), p. 1386–1403.
- 862 [14] D. A. FREEDMAN, *On the asymptotic behavior of Bayes’ estimates in the discrete case.II*, Annals of
863 Statistics, 36 (1965), p. 454–456.
- 864 [15] C. GAO AND H. H. ZHOU, *Rate exact Bayesian adaptation with modified block priors*, Annals of Statistics,
865 44 (2016), pp. 318–345.
- 866 [16] S. GHOSAL, J. K. GHOSH, AND A. VAN DER VAART, *Convergence rates of posterior distributions*, Annals
867 of Statistics, 28 (2000), pp. 500–531.
- 868 [17] S. GHOSAL AND A. VAN DER VAART, *Entropies and rates of convergence for maximum likelihood and*
869 *Bayes estimation for mixtures of normal densities*, Annals of Statistics, 29 (2001), pp. 1233–1263.
- 870 [18] S. GHOSAL AND A. VAN DER VAART, *Posterior convergence rates of Dirichlet mixtures at smooth densities*,
871 Annals of Statistics, 35 (2007), pp. 697–723.
- 872 [19] L. GROSS, *Logarithmic Sobolev inequalities*, American Journal of Mathematics, 97 (1975), pp. 1061–1083.
- 873 [20] P. HEINRICH AND J. KAHN, *Strong identifiability and optimal minimax rates for finite mixture estimation*,
874 Annals of Statistics, 46 (2018), pp. 2844–2870.
- 875 [21] N. HO, K. KHAMARU, R. DWIVEDI, M. J. WAINWRIGHT, M. I. JORDAN, AND B. YU, *Instability,*
876 *computational efficiency and statistical accuracy*, arXiv preprint arXiv:2005.11411, (2020).
- 877 [22] H. ISHWARAN, L. F. JAMES, AND J. SUN, *Bayesian model selection in finite mixtures by marginal density*
878 *decompositions*, Journal of the American Statistical Association, 96 (2001), pp. 1316–1332.
- 879 [23] B. J. KLEIJN AND A. W. VAN DER VAART, *The Bernstein-von-Mises theorem under misspecification*,
880 Electronic Journal of Statistics, 6 (2012), pp. 354–381.
- 881 [24] B. J. K. KLEIJN AND A. W. VAN DER VAART, *Misspecification in infinite-dimensional Bayesian statistics*,
882 Annals of Statistics, 34 (2006), pp. 837–877.
- 883 [25] B. T. KNAPIK, A. W. VAN DER VAART, AND J. H. VAN ZANTEN, *Bayesian inverse problems with*
884 *Gaussian priors*, The Annals of Statistics, 39 (2011), pp. 2626–2657.
- 885 [26] B. LINDSAY, *Mixture Models: Theory, Geometry and Applications*, In NSF-CBMS Regional Conference
886 Series in Probability and Statistics. IMS, Hayward, CA., 1995.
- 887 [27] P.-L. LOH AND M. J. WAINWRIGHT, *Regularized M-estimators with nonconvexity: Statistical and algorithmic*
888 *theory for local optima*, Advances in Neural Information Processing Systems, 26 (2013).
- 889 [28] C. MA, K. WANG, Y. CHI, AND Y. CHEN, *Implicit regularization in nonconvex statistical estimation:*
890 *Gradient descent converges linearly for phase retrieval, matrix completion, and blind deconvolution*,
891 Foundations of Computational Mathematics, 20 (2020), pp. 451–632.
- 892 [29] E. MAZUMDAR, A. PACCHIANO, Y. MA, M. JORDAN, AND P. BARTLETT, *On approximate Thompson*
893 *sampling with Langevin algorithms*, in International Conference on Machine Learning, PMLR, 2020,
894 pp. 6797–6807.
- 895 [30] P. MCCULLAGH AND J. A. NELDER, *Generalized Linear Models*, Chapman and Hall/CRC, 1989.
- 896 [31] W. MOU, N. HO, M. J. WAINWRIGHT, P. L. BARTLETT, AND M. I. JORDAN, *Supplementary material*
897 *to “A Diffusion Process Perspective on Posterior Contraction Rates for Parameters”*, 2023. DOI:
898 [COMPLETED BY TYPESETTER].
- 899 [32] X. NGUYEN, *Convergence of latent mixing measures in finite and infinite mixture models*, Annals of
900 Statistics, 4 (2013), pp. 370–400.
- 901 [33] R. NICKL, *Bayesian non-linear statistical inverse problems*, Lecture Notes ETH Zurich, (2022).
- 902 [34] D. M. OSTROVSKII AND F. BACH, *Finite-sample analysis of M-estimators using self-concordance*, Electronic
903 Journal of Statistics, 15 (2021), pp. 326–391.
- 904 [35] B. T. POLYAK AND A. B. JUDITSKY, *Acceleration of stochastic approximation by averaging*, SIAM journal
905 on control and optimization, 30 (1992), pp. 838–855.
- 906 [36] D. REVUZ AND M. YOR, *Continuous Martingales and Brownian Motion*, vol. 293, Springer-Verlag,
907 third ed., 1999.
- 908 [37] H. RISKEN, *The Fokker-Planck Equation*, Springer, 1996.
- 909 [38] J. ROUSSEAU, *Rates of convergence for the posterior distributions of mixtures of Beta and adaptive*

- 910 *nonparametric estimation of the density*, Annals of Statistics, 38 (2010), pp. 146–180.
- 911 [39] J. ROUSSEAU AND K. MENGENSEN, *Asymptotic behaviour of the posterior distribution in overfitted*
912 *mixture models*, Journal of the Royal Statistical Society: Series B (Statistical Methodology), 73 (2011),
913 pp. 689–710.
- 914 [40] L. SCHWARTZ, *On Bayes procedures*, Zeitschrift für Wahrscheinlichkeitstheorie und Verwandte Gebiete, 4
915 (1965), pp. 10–26.
- 916 [41] W. SHEN, S. R. TOKDAR, , AND S. GHOSAL, *Adaptive Bayesian multivariate density estimation with*
917 *Dirichlet mixtures*, Biometrika, 100 (2013), p. 623–640.
- 918 [42] X. SHEN AND L. WASSERMAN, *Rates of convergence of posterior distributions.*, Annals of Statistics, 29
919 (2001), pp. 687–714.
- 920 [43] B. SHI, S. S. DU, M. I. JORDAN, AND W. J. SU, *Understanding the acceleration phenomenon via*
921 *high-resolution differential equations*, Mathematical Programming, (2021), pp. 1–70.
- 922 [44] V. SPOKOINY, *Parametric estimation. finite sample theory*, The Annals of Statistics, 40 (2012), pp. 2877–
923 2909.
- 924 [45] W. SU, S. BOYD, AND E. J. CANDÉS, *A differential equation for modeling Nesterov’s accelerated gradient*
925 *method: Theory and insights*, Journal of Machine Learning Research, 17 (2016), pp. 1–43.
- 926 [46] M. TALAGRAND, *Transportation cost for Gaussian and other product measures*, Geometric & Functional
927 Analysis GAFA, 6 (1996), pp. 587–600.
- 928 [47] D. TALAY, *Second-order discretization schemes of stochastic differential systems for the computation of*
929 *the invariant law*, Stochastics: An International Journal of Probability and Stochastic Processes, 29
930 (1990), pp. 13–36.
- 931 [48] D. TALAY AND L. TUBARO, *Expansion of the global error for numerical schemes solving stochastic*
932 *differential equations*, Stochastic analysis and applications, 8 (1990), pp. 483–509.
- 933 [49] S. VAN DE GEER, *Empirical Processes in M-estimation*, Cambridge University Press, 2000.
- 934 [50] S. VAN DE GEER, *Empirical Processes in M-estimation*, Cambridge University Press, 2000.
- 935 [51] A. W. VAN DER VAART, *Asymptotic Statistics*, Cambridge University Press, 1998.
- 936 [52] A. W. VAN DER VAART AND J. WELLNER, *Weak Convergence and Empirical Processes*, Springer-Verlag,
937 New York, NY, 1996.
- 938 [53] M. J. WAINWRIGHT, *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*, Cambridge University
939 Press, 2019.
- 940 [54] S. WALKER, *On sufficient conditions for Bayesian consistency*, Annals of Statistics, 90 (2003), pp. 482–488.
- 941 [55] S. WALKER, *New approaches to Bayesian consistency*, Annals of Statistics, 32 (2004), pp. 2028–2043.
- 942 [56] S. G. WALKER, A. LIJOI, AND I. PRUNSTER, *On rates of convergence for posterior distributions in*
943 *infinite-dimensional models*, Annals of Statistics, 35 (2007), pp. 738–746.
- 944 [57] Y. YANG AND D. B. DUNSON, *Bayesian manifold regression*, Annals of Statistics, 44 (2016), pp. 876–905.
- 945 [58] Y. YANG AND S. T. TOKDAR, *Minimax-optimal nonparametric regression in high dimensions*, Annals of
946 Statistics, 43 (2015), pp. 652–674.

1 **SUPPLEMENTARY MATERIALS: A Diffusion Process Perspective on Posterior**
2 **Contraction Rates for Parameters**

3 Wenlong Mou*, Nhat Ho†, Martin Wainwright‡, Peter Bartlett‡, and Michael Jordan‡
4

5 This supplementary material is devoted to the proofs deferred from the main paper. In
6 [Appendix A](#), we present the proofs of non-asymptotic Bernstein–von Mises theorems using
7 tools from diffusion process theory. The proofs of our main corollaries are given in [Appendix B](#),
8 whereas [Appendix C](#) is devoted to the proofs of auxiliary results.

9 **Appendix A. Proofs of non-asymptotic Bernstein–von Mises results.** In this section,
10 we collect the proofs of [Theorem 3.5](#) and [Proposition 3.4](#).

11 **A.1. Proof of Theorem 3.5.** For any fixed $T > 0$, we define the sequence of potential
12 functions $\Phi_t : \mathbb{R}^d \rightarrow \mathbb{R}$

13
$$\Phi_t(\theta) := \langle \theta - \hat{\theta}^{(n)}, H^* e^{H^*(t-T)}(\theta - \hat{\theta}^{(n)}) \rangle, \quad \text{for each } t \in [0, T].$$

14
15 Once again, we consider the diffusion process

16
$$d\theta_t = -\nabla F_n(\theta_t)dt + \frac{1}{n}\nabla \log \pi(\theta_t)dt + dB_t,$$

17
18 with the initial condition $\theta_0 = \hat{\theta}^{(n)}$. Using Itô’s formula, for $t \in [0, T]$, we have

19
$$\begin{aligned} \Phi_t(\theta_t) &= \int_0^t \frac{\partial \Phi_s}{\partial s}(\theta_s)ds - \int_0^t \left\langle \nabla \Phi_s(\theta_s), \nabla F_n(\theta_s) - \frac{\nabla \log \pi(\theta_s)}{n} \right\rangle ds \\ &\quad + \sqrt{\frac{2}{n}} \int_0^t \langle \nabla \Phi_s(\theta_s), dB_s \rangle + \frac{1}{n} \int_0^t \Delta \Phi_s(\theta_s) ds \\ &= \underbrace{\int_0^t \left(H^*(\theta_s - \hat{\theta}^{(n)}) - \nabla F_n(\theta_s) + \frac{\nabla \log \pi(\theta_s)}{n} \right)^\top H^* e^{H^*(s-T)}(\theta_s - \hat{\theta}^{(n)}) ds}_{:= I_1(t)} \\ &\quad + \underbrace{\sqrt{\frac{2}{n}} \int_0^t (\theta_s - \hat{\theta}^{(n)})^\top H^* e^{(s-T)H^*} dB_s}_{I_2(t)} + \underbrace{\frac{1}{n} \int_0^t \text{Tr} \left(H^* e^{H^*(s-T)} \right) ds}_{I_3(t)}. \end{aligned}$$

20
21
22 (A.1)
23
24 Note that the matrices H^* and $e^{(s-T)H^*}$ commute, so that we may write their product in an
25 arbitrary order.

26 Defining the linearization error

27
$$\Delta_s := (A + \varepsilon_1^{(2)}(n, \delta)) \left(\|\theta_s - \theta^*\|_2 + \|\hat{\theta}^{(n)} - \theta^*\|_2 \right) + \varepsilon_2^{(2)}(n, \delta) + \frac{L_2}{n},$$

*Department of EECS, UC Berkeley.

†Department of Statistics and Data Science, UT Austin.

‡Department of EECS and Department of Statistics, UC Berkeley.

29 we claim that the following bounds hold for each $t \in [0, T]$:

(A.2a)

$$30 \quad I_1(t) \leq \frac{2+\log \kappa(H^*)}{a} \sup_{0 \leq s \leq t} \Phi_s(\theta_s) + a \int_0^t \Delta_s^2 \left(\|\theta_s - \theta^*\|_2^2 + \|\widehat{\theta}^{(n)} - \theta^*\|_2^2 \right) e^{-\frac{\lambda_{\min}(H^*)}{2}(s-T)} ds,$$

$$33 \quad (A.2b) \quad \left(\mathbb{E} \sup_{0 \leq t \leq T} |I_2(t)|^p \right)^{1/p} \leq c \sqrt{\frac{p(1+\log \kappa(H^*))}{n}} \left(\mathbb{E} \sup_{0 \leq t \leq T} \Phi_t(\theta_t)^{p/2} \right)^{1/p}, \quad \text{and}$$

$$34 \quad (A.2c) \quad I_3(t) \leq \frac{d}{n}.$$

36 Here $c > 0$ is a universal constant. We prove all of these bounds in the subsections to follow.

37 Taking these bounds as given for the moment, let us complete the proof of the theorem.

38 By Jensen's inequality, for an even integer $p \geq 2$, the moments of the integral term in
39 equation (A.2a) can be bounded as

$$41 \quad (A.3) \quad \mathbb{E} \left(\int_0^T \Delta_s^2 \left(\|\theta_s - \theta^*\|_2^2 + \|\widehat{\theta}^{(n)} - \theta^*\|_2^2 \right) e^{-\frac{\lambda_{\min}(H^*)}{2}(s-T)} ds \right)^p$$

$$42 \quad \leq \left(\frac{c}{\lambda_{\min}(H^*)} \right)^{p-1} \cdot \mathbb{E} \int_0^T \Delta_s^{2p} \left(\|\theta_s - \theta^*\|_2^{2p} + \|\widehat{\theta}^{(n)} - \theta^*\|_2^{2p} \right) e^{-\frac{\lambda_{\min}(H^*)}{2}(s-T)} ds,$$

44 for a universal constant $c > 0$.

45 For any $\omega \in (0, 1)$, by taking supremum on both sides of the decomposition (A.1), combining
46 with the bounds (A.2a) and (A.2c), and taking $a = c \frac{2+\log \kappa(H^*)}{\omega}$, we arrive at the inequality

$$48 \quad \sup_{0 \leq t \leq T} \Phi_t(\theta_t) \leq (1 + \omega) \left(\frac{d}{n} + \sup_{0 \leq t \leq T} I_2(t) \right)$$

$$49 \quad + \frac{c(2 + \log \kappa(H^*))}{\omega} \int_0^T \Delta_t^2 \left(\|\theta_t - \theta^*\|_2^2 + \|\widehat{\theta}^{(n)} - \theta^*\|_2^2 \right) e^{-\frac{\lambda_{\min}(H^*)}{2}(t-T)} dt.$$

51 Taking p -th moment on both sides of the inequality, combining with the bounds (A.2b)
52 and (A.3), and applying Minkowski's inequality, we arrive at the bound

$$54 \quad \left(\mathbb{E} \sup_{0 \leq t \leq T} \Phi_t(\theta_t)^p \right)^{1/p} \leq (1 + \omega) \frac{d}{n} + \sqrt{\frac{cp(1 + \log \kappa(H^*))}{n}} \cdot \left(\mathbb{E} \sup_{0 \leq t \leq T} \Phi_t(\theta_t)^p \right)^{\frac{1}{2p}}$$

$$55 \quad + \frac{c(2 + \log \kappa(H^*))}{\omega \lambda_{\min}(H^*)} \left(\sup_{0 \leq t \leq T} \mathbb{E} \left[\Delta_t^{2p} \left(\|\theta_t - \theta^*\|_2^{2p} + \|\widehat{\theta}^{(n)} - \theta^*\|_2^{2p} \right) \right] \right)^{1/p}.$$

57 Substituting with the definition of the last term, and applying Young's inequality, we find that

$$58 \quad \left(\mathbb{E} \sup_{0 \leq t \leq T} \Phi_t(\theta_t)^p \right)^{1/p} \leq (1 + \omega) \frac{d}{n} + c \frac{1 + \log \kappa(H^*)}{\omega} \left(\frac{p}{n} + \frac{\mathcal{H}_n(p, \delta)}{\lambda_{\min}(H^*)} \right),$$

59

60 where the high-order term $\mathcal{H}_n(p, \delta)$ is defined as

$$\begin{aligned}
61 \quad \mathcal{H}_n(p, \delta) &:= (A + \varepsilon_1^{(2)}(n, \delta))^2 \left(\mathbb{E}_{\mathbb{Q}} \|\theta - \theta^*\|_2^{4p} \right)^{1/p} \\
62 \quad &+ \left\| \widehat{\theta}^{(n)} - \theta^* \right\|_2^2 \left(\varepsilon_2^{(2)}(n, \delta)^2 + \frac{L_2^2}{n^2} + (A + \varepsilon_1^{(2)}(n, \delta))^2 \left\| \widehat{\theta}^{(n)} - \theta^* \right\|_2^2 \right). \\
63
\end{aligned}$$

64 Putting together the pieces yields the conclusion of the theorem.

65 **A.1.1. Proof of claim (A.2a).** We first bound the term $I_1(t)$. Noting the defining identity
66 $\nabla F_n(\widehat{\theta}^{(n)}) + \frac{1}{n} \nabla \log \pi(\widehat{\theta}^{(n)}) = 0$, we have the following bound:

$$\begin{aligned}
67 \quad &\left\| H^*(\theta_s - \widehat{\theta}^{(n)}) - \nabla F_n(\theta_s) + \nabla \log \pi(\theta_s)/n \right\|_2 \\
68 \quad &= \left\| \int_0^1 \left(H^* - \nabla^2 F_n(\gamma\theta_s + (1-\gamma)\widehat{\theta}^{(n)}) + \nabla^2 \log \pi(\gamma\theta_s + (1-\gamma)\widehat{\theta}^{(n)})/n \right) (\theta_s - \widehat{\theta}^{(n)}) d\gamma \right\|_2 \\
69 \quad &\leq \int_0^1 \left\| H^* - \nabla^2 F_n(\gamma\theta_s + (1-\gamma)\widehat{\theta}^{(n)}) + \nabla^2 \log \pi(\gamma\theta_s + (1-\gamma)\widehat{\theta}^{(n)})/n \right\|_{\text{op}} \cdot \left\| \theta_s - \widehat{\theta}^{(n)} \right\|_2 d\gamma. \\
70
\end{aligned}$$

71 By Assumptions **(BvM.1)**, **(BvM.2)**, and **(PS)**, for any $\theta \in \mathbb{R}^d$, we have the bound

$$\begin{aligned}
72 \quad &\left\| H^* - \nabla^2 F_n(\theta) + \nabla^2 \log \pi(\theta)/n \right\|_{\text{op}} \\
73 \quad &\leq \left\| H^* - \nabla^2 F(\theta) \right\|_{\text{op}} + \left\| \nabla^2 F(\theta) - \nabla^2 F_n(\theta) \right\|_{\text{op}} + \left\| \nabla^2 \log \pi(\theta)/n \right\|_{\text{op}} \\
74 \quad &\leq A \left\| \theta - \theta^* \right\|_2 + \varepsilon_1^{(2)}(n, \delta) \left\| \theta - \theta^* \right\|_2 + \varepsilon_2^{(2)}(n, \delta) + \frac{L_2}{n}. \\
75
\end{aligned}$$

76 Substituting into the bound for $I_1(t)$, for any $a > 0$, we have that

$$\begin{aligned}
77 \quad I_1(t) &\leq \int_0^t \left\| (H^*)^{1/2} e^{H^*(s-t)/2} \right\|_{\text{op}} \\
78 \quad &\quad \times \left\| H^* - \nabla^2 F_n(\theta_s) + \nabla^2 \log \pi(\theta_s)/n \right\|_2 \left\| \theta_s - \widehat{\theta}^{(n)} \right\|_2 \sqrt{\Phi_s(\theta_s)} ds \\
79 \quad &\leq a^{-1} \sup_{0 \leq s \leq t} \Phi_s(\theta_s) \cdot \int_0^t \left\| (H^*)^{1/2} e^{H^*(s-T)/4} \right\|_{\text{op}}^2 ds \\
80 \quad &\quad + a \int_0^t \left\| H^* - \nabla^2 F_n(\theta_s) + \nabla^2 \log \pi(\theta_s)/n \right\|_{\text{op}}^2 \cdot \left\| \theta_s - \widehat{\theta}^{(n)} \right\|_2^2 e^{H^*(s-T)/4} \left\| \right\|_{\text{op}}^2 ds \\
81 \quad &\leq \frac{2 + \log \kappa(H^*)}{a} \sup_{0 \leq s \leq t} \Phi_s(\theta_s) \\
82 \quad &\quad + a \int_0^t \Delta_s^2 \left(\left\| \theta_s - \theta^* \right\|_2^2 + \left\| \widehat{\theta}^{(n)} - \theta^* \right\|_2^2 \right) e^{-\frac{\lambda_{\min}(H^*)}{2}(s-T)} ds. \\
83
\end{aligned}$$

84 Therefore, claim (A.2a) follows.

85 **A.1.2. Proof of claim (A.2b).** Note that $I_2(t)$ is a martingale with respect to the Brownian
86 filtration. Applying the Burkholder-Gundy-Davis inequality for an arbitrary $p \geq 2$ yields

$$\begin{aligned}
87 \quad \left(\mathbb{E} \sup_{0 \leq t \leq T} |I_2(t)|^p \right)^{1/p} &\leq c \sqrt{\frac{p}{n}} \left(\mathbb{E} \left(\int_0^T \left\| H^* e^{(t-T)H^*} (\theta_t - \hat{\theta}^{(n)}) \right\|_2^2 dt \right)^{\frac{p}{2}} \right)^{1/p} \\
88 \quad &\leq C \sqrt{\frac{p}{n}} \left(\mathbb{E} \left(\int_0^T \left\| (H^*)^{1/2} e^{\frac{t-T}{2}H^*} \right\|_{\text{op}}^2 \Phi_t(\theta_t) dt \right)^{\frac{p}{2}} \right)^{1/p} \\
89 \quad &\leq c \sqrt{\frac{p}{n}} \left(\mathbb{E} \sup_{0 \leq t \leq T} \Phi_t(\theta_t)^{p/2} \right)^{1/p} \cdot \sqrt{\int_0^T \left\| (H^*)^{1/2} e^{\frac{t-T}{2}H^*} \right\|_{\text{op}}^2 dt}.
\end{aligned}$$

91 We now observe that

$$92 \quad \left\| (H^*)^{1/2} e^{\frac{t-T}{2}H^*} \right\|_{\text{op}}^2 = \left\| H^* e^{(t-T)H^*} \right\|_{\text{op}} = \max_{i \in [d]} \left(\lambda_i(H^*) e^{(t-T)\lambda_i(H^*)} \right).$$

94 Taking the time integral leads to the bound

$$\begin{aligned}
95 \quad \int_0^T \left\| (H^*)^{1/2} e^{\frac{t-T}{2}H^*} \right\|_{\text{op}}^2 dt &\leq \int_0^{+\infty} \max_{i \in [d]} \left(\lambda_i(H^*) e^{-t\lambda_i(H^*)} \right) dt \\
96 \quad &\leq \underbrace{\int_0^{+\infty} \max_{\lambda_{\min}(H^*) \leq \lambda \leq \lambda_{\max}(H^*)} \left(\lambda e^{-t\lambda} \right) dt}_{=: J}.
\end{aligned}$$

98 We now split the integral J into three parts, thereby obtaining

$$\begin{aligned}
99 \quad J &\leq \int_0^{\lambda_{\max}(H^*)^{-1}} \lambda_{\max}(H^*) e^{-t\lambda_{\max}(H^*)} dt \\
100 \quad &\quad + \int_{\lambda_{\max}(H^*)^{-1}}^{\lambda_{\min}(H^*)^{-1}} \frac{dt}{et} + \int_{\lambda_{\min}(H^*)^{-1}}^{+\infty} \lambda_{\min}(H^*) e^{-t\lambda_{\min}(H^*)} dt \\
101 \quad (A.4) \quad &\leq 1 + \frac{1}{e} \log \frac{\lambda_{\max}(H^*)}{\lambda_{\min}(H^*)}.
\end{aligned}$$

103 Denote $\kappa(M) := \frac{\lambda_{\max}(M)}{\lambda_{\min}(M)}$ for a positive definite matrix M . Collecting the above inequalities,
104 we find that the term $I_2(t)$ is upper bounded as

$$105 \quad \left(\mathbb{E} \sup_{0 \leq t \leq T} |I_2(t)|^p \right)^{1/p} \leq c \sqrt{\frac{p(1 + \log \kappa(H^*))}{n}} \left(\mathbb{E} \sup_{0 \leq t \leq T} \Phi_t(\theta_t)^{p/2} \right)^{1/p}$$

107 for a universal constant $c > 0$. This completes the proof of the claim (A.2b).

108 **A.1.3. Proof of claim (A.2c).** Finally, the term $I_3(t)$ is straightforward to upper bound as

$$109 \quad I_3(t) \leq \frac{1}{n} \text{Tr} \left(H^* \int_0^T e^{H^*(s-T)} ds \right) \leq \frac{1}{n} \text{Tr} \left(H^* \int_0^{+\infty} e^{-sH^*} ds \right) = \frac{d}{n},$$

111 which establishes the claim (A.2c).

112 **A.2. Proof of Proposition 3.4.** We introduce the shorthand $\mu := \mathcal{N}(\widehat{\theta}^{(n)}, (nH^*)^{-1})$ for
 113 the target density. Since $H^* \succ 0$, the Gaussian log-Sobolev inequality implies that

$$114 \quad (A.5) \quad D_{\text{KL}}(\mathbb{Q}(\cdot | X_1^n) \parallel \mu) \leq \frac{1}{n\lambda_{\min}(H^*)} \int_{\mathbb{R}^d} \|\nabla \log \mathbb{Q}(\theta | X_1^n) - \nabla \log \mu(\theta)\|_2^2 \mathbb{Q}(d\theta | X_1^n).$$

116 Since μ is a Gaussian density, we find that

$$117 \quad \nabla \log \mu(\theta) = -nH^*(\theta - \widehat{\theta}^{(n)}).$$

119 For the posterior density $\mathbb{Q}(\cdot | X_1^n)$, we note that

$$120 \quad \begin{aligned} \nabla \log \mathbb{Q}(\theta | X_1^n) &= -n\nabla F_n(\theta) + \nabla \log \pi(\theta) \\ 121 \quad &= \int_0^1 \left(-n\nabla^2 F_n(\gamma\theta + (1-\gamma)\widehat{\theta}^{(n)}) + \nabla^2 \log \pi(\gamma\theta + (1-\gamma)\widehat{\theta}^{(n)}) \right) \\ 122 \quad &\quad \times (\theta - \widehat{\theta}^{(n)}) d\gamma. \end{aligned}$$

124 Putting the above equations together yields

$$125 \quad \begin{aligned} 126 \quad &\|\nabla \log \mathbb{Q}(\theta | X_1^n) - \nabla \log \mu(\theta)\|_2 \\ 127 \quad &\leq n \int_0^1 \|\nabla^2 F_n(\gamma\theta + (1-\gamma)\widehat{\theta}^{(n)}) - H^* + \nabla^2 \log \pi(\gamma\theta + (1-\gamma)\widehat{\theta}^{(n)})/n\|_{\text{op}} \cdot \|\theta - \widehat{\theta}^{(n)}\|_2 d\gamma. \end{aligned}$$

129 By Assumptions **(BvM.1)**, **(BvM.2)**, and **(PS)**, we have the bounds

$$130 \quad \begin{aligned} &\|\nabla^2 F_n(\gamma\theta + (1-\gamma)\widehat{\theta}^{(n)}) + \nabla^2 \log \pi(\gamma\theta + (1-\gamma)\widehat{\theta}^{(n)})/n - H^*\|_{\text{op}} \\ 131 \quad &\leq \|\nabla^2 F(\gamma\theta + (1-\gamma)\widehat{\theta}^{(n)}) - H^*\|_{\text{op}} \\ 132 \quad &\quad + \|\nabla^2 F_n(\gamma\theta + (1-\gamma)\widehat{\theta}^{(n)}) - \nabla^2 F(\gamma\theta + (1-\gamma)\widehat{\theta}^{(n)})\|_{\text{op}} + \frac{L_2}{n} \\ 133 \quad &\leq A \|\gamma\theta + (1-\gamma)\widehat{\theta}^{(n)} - \theta^*\|_2 + \varepsilon_1^{(2)}(n, \delta) \|\theta - \widehat{\theta}^{(n)}\|_2 + \varepsilon_2^{(2)}(n, \delta) + \frac{L_2}{n}. \end{aligned}$$

135 Substituting this bound into the bound (A.5) yields

$$136 \quad \begin{aligned} D_{\text{KL}}(\mathbb{Q}(\cdot | X_1^n) \parallel \mu) &\leq \frac{n}{\lambda_{\min}(H^*)} \left(A \cdot \mathbb{E}_{\mathbb{Q}} \left[\|\theta - \theta^*\|_2^4 | X_1^n \right] + A \|\widehat{\theta}^{(n)} - \theta^*\|_2^4 \right) \\ 137 \quad &\quad + \frac{n\varepsilon_1^{(2)}(n, \delta)}{\lambda_{\min}(H^*)} \mathbb{E}_{\mathbb{Q}} \left[\|\theta - \widehat{\theta}^{(n)}\|_2^3 | X_1^n \right] \\ 138 \quad &\quad + (\varepsilon_2^{(2)}(n, \delta) + L_2/n) \cdot \mathbb{E} \left[\|\theta - \widehat{\theta}^{(n)}\|_2^2 | X_1^n \right]. \end{aligned}$$

140 As a consequence, we obtain the conclusion of the proposition.

141 **Appendix B. Proofs of corollaries.** In this appendix, we collect the proofs of several
 142 corollaries stated in the main text and Section 4. The crux of the proofs of these corollaries
 143 involves a verification of assumptions to invoke the respective theorems. Note that the values
 144 of universal constants may change from line to line.

145 **B.1. Proof of Corollary 4.1.** We begin by verifying claim (4.2a) about the structure
 146 of the negative population log-likelihood function F^R and claim (4.2b) about the uniform
 147 perturbation error between ∇F^R and ∇F_n^R .

148 **B.1.1. Proof of claim (4.2a).** Following some algebra, we find that

$$149 \quad -F^R(\theta) = \mathbb{E} \left[-Y \log \left(1 + e^{-\langle X, \theta \rangle} \right) - (1 - Y) \log \left(1 + e^{\langle X, \theta \rangle} \right) \right]$$

$$150 \quad = -\mathbb{E} \left[\frac{1}{1 + e^{-\langle X, \theta^* \rangle}} \log \left(1 + e^{-\langle X, \theta \rangle} \right) + \frac{1}{1 + e^{\langle X, \theta^* \rangle}} \log \left(1 + e^{\langle X, \theta \rangle} \right) \right],$$

152 where the above expectations are taken with respect to $X \sim \mathcal{N}(0, \sigma^2 I_d)$ and $Y|X$ following
 153 probability distribution generated from logistic model (4.1). Taking the derivative of F^R with
 154 respect to θ yields

$$155 \quad \langle \nabla F^R(\theta), \theta^* - \theta \rangle$$

$$156 \quad = \mathbb{E} \left[\left(\frac{1 + e^{\langle X, \theta \rangle}}{1 + e^{\langle X, \theta^* \rangle}} - \frac{1 + e^{-\langle X, \theta \rangle}}{1 + e^{-\langle X, \theta^* \rangle}} \right) \frac{e^{-\langle X, \theta \rangle}}{(1 + e^{-\langle X, \theta \rangle})^2} \langle X, \theta - \theta^* \rangle \right].$$

158 By the mean value theorem, there exists ξ between 0 and $\langle X, \theta - \theta^* \rangle$ such that

$$159 \quad \frac{1 + e^{\langle X, \theta \rangle}}{1 + e^{\langle X, \theta^* \rangle}} - \frac{1 + e^{-\langle X, \theta \rangle}}{1 + e^{-\langle X, \theta^* \rangle}} = \langle X, \theta - \theta^* \rangle \left(\frac{e^{\langle X, \theta^* \rangle + \xi}}{1 + e^{\langle X, \theta^* \rangle}} + \frac{e^{-\langle X, \theta^* \rangle - \xi}}{1 + e^{-\langle X, \theta^* \rangle}} \right).$$

161 In light of the above equality, we arrive at the following inequalities:

$$162 \quad \langle \nabla F^R(\theta), \theta^* - \theta \rangle \geq \mathbb{E} \left[\inf_{|\xi| \in [0, |\langle X, \theta - \theta^* \rangle|]} \left(\frac{e^{\langle X, \theta^* \rangle + \xi}}{1 + e^{\langle X, \theta^* \rangle}} + \frac{e^{-\langle X, \theta^* \rangle - \xi}}{1 + e^{-\langle X, \theta^* \rangle}} \right) \right.$$

$$163 \quad \quad \quad \left. \times \frac{e^{-\langle X, \theta \rangle}}{(1 + e^{-\langle X, \theta \rangle})^2} |\langle X, \theta - \theta^* \rangle|^2 \right]$$

$$164 \quad \geq \mathbb{E} \left[\frac{1}{2} e^{-|\langle X, \theta - \theta^* \rangle|} \frac{e^{-\langle X, \theta \rangle}}{(1 + e^{-\langle X, \theta \rangle})^2} |\langle X, \theta - \theta^* \rangle|^2 \right]$$

$$165 \quad \geq \frac{1}{8} \mathbb{E} \left[e^{-|\langle X, \theta - \theta^* \rangle| - |\langle X, \theta \rangle|} |\langle X, \theta - \theta^* \rangle|^2 \right]$$

$$166 \quad \geq \frac{1}{8e^4} \mathbb{E} \left[\mathbf{1}_{\{|\langle X, \theta \rangle| \leq 2, |\langle X, \theta - \theta^* \rangle| \leq 2\}} |\langle X, \theta - \theta^* \rangle|^2 \right].$$

168 Since $X \sim \mathcal{N}(0, I_d)$, we have

$$169 \quad \begin{bmatrix} \langle X, \theta \rangle \\ \langle X, \theta - \theta^* \rangle \end{bmatrix} \sim \mathcal{N} \left(0, \begin{bmatrix} \|\theta\|_2^2 & \langle \theta, \theta - \theta^* \rangle \\ \langle \theta, \theta - \theta^* \rangle & \|\theta - \theta^*\|_2^2 \end{bmatrix} \right).$$

171 Given that result, direct calculation leads to

$$172 \quad \mathbb{E} \left(\mathbf{1}_{\{|\langle X, \theta \rangle| \leq 2, |\langle X, \theta - \theta^* \rangle| \leq 2\}} |\langle X, \theta - \theta^* \rangle|^2 \right)$$

$$173 \quad \geq \frac{c}{(1 + \|\theta\|_2)(1 + \|\theta - \theta^*\|_2)} \|\theta - \theta^*\|_2^2,$$

174

175 for a universal constant $c > 0$. Collecting the above results, for all θ such that $\|\theta - \theta^*\|_2 \leq 1$,
 176 we achieve that

$$177 \quad \langle \nabla F^R(\theta), \theta^* - \theta \rangle \geq \frac{c}{(1 + \|\theta\|_2)(1 + \|\theta - \theta^*\|_2)} \|\theta - \theta^*\|_2^2$$

$$178 \quad \geq c \frac{1}{1 + \|\theta^*\|_2} \|\theta - \theta^*\|_2^2.$$

180 For θ with $\|\theta - \theta^*\|_2 > 1$, let $\tilde{\theta} = \theta^* + \frac{\theta - \theta^*}{\|\theta - \theta^*\|_2}$. Then, we find that

$$181 \quad \langle \nabla F^R(\theta), \theta^* - \theta \rangle \geq \langle \nabla F^R(\tilde{\theta}), \theta^* - \theta \rangle \geq \frac{c}{2(1 + \|\theta^*\|_2)} \|\theta - \theta^*\|_2,$$

182 which yields the claim (4.2a).

184 **B.1.2. Proof of the bound (4.2b).** In this appendix, we prove the uniform bound (4.2b)
 185 between the empirical and population likelihood gradients. It suffices to establish the following
 186 stronger result:

$$187 \quad (\text{B.1}) \quad Z := \sup_{\theta \in \mathbb{R}^d} \|\nabla F_n^R(\theta) - \nabla F^R(\theta)\|_2 \leq c \left\{ \sqrt{\frac{d}{n}} + \sqrt{\frac{\log(1/\delta)}{n}} + \frac{\log(1/\delta)}{n} \right\},$$

188 with probability at least $1 - \delta$ for any $\frac{n}{\log n} \geq c_0 d \log(1/\delta)$ where c_0 is a universal constant.

190 In order to prove the claim (B.1), we exploit a concentration inequality due to Adamczak [SM1]; it gives tight tail bounds for supremum of unbounded empirical processes. Through-
 191 out our derivation, we use $\|X\|_{\psi_\alpha}$ to denote the Orlicz ψ_α norm for a random variable X , for
 192 any $\alpha \in (0, 2]$. Let us state a simplified version of a theorem due to Adamczak:

194 **Proposition B.1 (Theorem 4, [SM1], simplified version).** *Let $(x, \theta) \mapsto f(\theta; x)$ be a function*
 195 *with domain $\Theta \times \mathcal{X}$, and suppose that there is a function $\bar{F} : \mathcal{X} \rightarrow \mathbb{R}$ such that $|f(\theta, x)| \leq \bar{F}(x)$*
 196 *for any $\theta \in \Theta$. Let $X_1, X_2, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}_X$, and suppose that $\|\bar{F}\|_{\psi_\alpha} < +\infty$ for some $\alpha \leq 1$.*
 197 *Then the random variable $Z_n := \frac{1}{n} \sup_{\theta \in \Theta} |\sum_{i=1}^n f(\theta; X_i) - \mathbb{E}[f(\theta; X)]|$ satisfies the bound:*

$$198 \quad \mathbb{P}(Z_n > 2\mathbb{E}[Z_n] + t) \leq \exp\left(-\frac{t^2}{2\mathbb{E}[\bar{F}(X)^2]}\right) + 3 \exp\left(-\left(\frac{t}{c \|\max_{i \in [n]} \bar{F}(X_i)\|_{\psi_\alpha}}\right)^\alpha\right),$$

199 for a universal constant $c > 0$.

201 In order to prove the claim (B.1), we begin by writing Z as the supremum of a stochastic
 202 process. Let \mathbb{S}^{d-1} denote the Euclidean sphere in \mathbb{R}^d , and define the stochastic process

$$203 \quad Z_{u,\theta} := \left| \frac{1}{n} \sum_{i=1}^n f_{u,\theta}(X_i, Y_i) - \mathbb{E}[f_{u,\theta}(X, Y)] \right|,$$

204 where $f_{u,\theta}(x, y) = \frac{y \langle x, u \rangle e^{y \langle x, \theta \rangle}}{1 + e^{y \langle x, \theta \rangle}}$, indexed by vectors $u \in \mathbb{S}^{d-1}$ and $\theta \in \mathbb{B}(\theta^*; r)$. The outer
 205 expectation in the above display is taken with respect to (X, Y) drawn from the logistic
 206 model (4.1)
 207

208 Observe that $Z = \sup_{u \in \mathbb{S}^{d-1}} \sup_{\theta \in \mathbb{R}^d} Z_{u,\theta}$. Let $\{u^1, \dots, u^N\}$ be a $1/8$ -covering of \mathbb{S}^{d-1} in the
 209 Euclidean norm; there exists such a set with $N \leq 17^d$ elements. By a standard discretization
 210 argument (see Chapter 6, [SM5]), we have

$$211 \quad Z \leq 2 \max_{j=1, \dots, N} \sup_{\theta \in \mathbb{R}^d} Z_{u^j, \theta}.$$

213 Accordingly, the remainder of our argument focuses on bounding the random variable
 214 $V := \sup_{\theta \in \mathbb{R}^d} Z_{u,\theta}$, where the vector $u \in \mathbb{S}^{d-1}$ should be understood as arbitrary but fixed.
 215 For each $u \in \mathbb{S}^{d-1}$ fixed, we note that $\bar{F}(X, Y) = |\langle X, u \rangle|$ is an envelop function for the class
 216 $(f_{u,\theta}(X, Y))_{\theta \in \mathbb{R}^d}$. Additionally, by standard tail bounds for maximum of Gaussian random
 217 variables, we know that:

$$218 \quad \left\| \max_{1 \leq i \leq n} \bar{F}(X_i, Y_i) \right\|_{\psi_1} \leq \sqrt{\log n}.$$

220 Consequently, invoking Proposition B.1 yields that

$$221 \quad (B.2) \quad V \leq 2\mathbb{E}[V] + \sqrt{\frac{2 \log(1/\delta)}{n}} + \frac{c \log(1/\delta)}{n} \sqrt{\log n}$$

223 with probability at least $1 - \delta$.

224 Now define the symmetrized random variable

$$225 \quad V' := \sup_{\theta \in \mathbb{R}^d} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f_{\theta, u}(X_i, Y_i) \right|.$$

227 where $\{\varepsilon_i\}_{i=1}^n$ is an i.i.d. sequence of Rademacher variables. By standard symmetrization
 228 arguments, we have

$$229 \quad \mathbb{E}[V] \leq 2\mathbb{E}[V'].$$

231 We now bound the expectation of V' , first over the Rademacher variables. Consider the
 232 function class

$$233 \quad \mathcal{G} := \left\{ g_\theta : (x, y) \mapsto \langle x, u \rangle \varphi_\theta(x, y) \mid \theta \in \mathbb{R}^d \right\}.$$

235 It is clear that the function class \mathcal{G} has the envelope function $\bar{G}(x) := |\langle x, u \rangle|$. We claim that
 236 the L_2 -covering number of \mathcal{G} can be bounded as

$$237 \quad (B.3) \quad \bar{N}(t) := \sup_Q \left| \mathcal{N} \left(\mathcal{G}, \|\cdot\|_{L^2(Q)}, t \|\bar{G}\|_{L^2(Q)} \right) \right| \leq \left(\frac{1}{t} \right)^{c(d+1)} \quad \text{for all } t > 0,$$

239 where $c > 0$ is a universal constant.

240 Let us take the claim (B.3) as given for the moment, and use it to bound the ex-
 241 pectation of V' , first over the Rademacher variables. Define the empirical expectation

242 $\mathbb{P}_n(\bar{G}^2) := \frac{1}{n} \sum_{i=1}^n \langle X_i, u \rangle^2$. Invoking Dudley's entropy integral bound (e.g., Theorem 5.22,
243 [SM5]), we find that there are universal constants C, C' such that

$$244 \quad \mathbb{E}_\varepsilon[V'] = \mathbb{E}_\varepsilon \left[\sup_{g \in \mathcal{G}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i g(X_i, Y_i) \right| \right] \leq C \sqrt{\frac{\mathbb{P}_n(\bar{G}^2)}{n}} \int_0^1 \sqrt{1 + \log \bar{N}(t)} dt$$

$$245 \quad \leq C' \sqrt{\mathbb{P}_n(\bar{G}^2)} \sqrt{\frac{d}{n}}.$$

247 Up to this point, we have been conditioning on the observations $\{X_i\}_{i=1}^n$. Taking expectations
248 over them as well yields

$$249 \quad (\text{B.4}) \quad \mathbb{E}_{\varepsilon, X_1^n}[V'] \leq C' \sqrt{\frac{d}{n}} \cdot \mathbb{E}_{X_1^n} \left[\sqrt{\mathbb{P}_n(\bar{G}^2)} \right] \stackrel{(i)}{\leq} C' \sqrt{\frac{d}{n}} \cdot \sqrt{\mathbb{E}_{X_1^n} [\mathbb{P}_n(\bar{G}^2)]} \stackrel{(ii)}{=} C' \sqrt{\frac{d}{n}},$$

251 where step (i) follows from Jensen's inequality; and step (ii) uses the fact that $\mathbb{E}_{X_1^n} [\mathbb{P}_n(\bar{G}^2)] = 1$.
252 Putting together the bounds (B.2) and (B.4) yields the following bound with probability $1 - \delta$:

$$253 \quad V \leq c \sqrt{\frac{d + \log \delta^{-1}}{n}} + c \frac{\log \delta^{-1}}{n} \sqrt{\log n}.$$

255 This probability bound holds for each $u \in \mathbb{S}^{d-1}$. By taking the union bound over the $1/8$ -
256 covering set $\{u^1, \dots, u^N\}$ of \mathbb{S}^{d-1} where $N \leq 17^d$ and applying above bound with $\delta' = \delta/N$,
257 we obtain the claim (B.1) for sample size satisfying $\frac{n}{\log n} \geq cd \log(1/\delta)$.

258 **B.1.3. Proof of claim (B.3).** We consider a fixed sequence $(x_i, y_i, t_i)_{i=1}^m$ where $y_i \in \{-1, 1\}$,
259 $x_i \in \mathbb{R}^d$ and $t_i \in \mathbb{R}$ for $i \in [m]$. Now, we suppose that for any binary sequence $(z_i)_{i=1}^m \in \{0, 1\}^m$,
260 there exists $\theta \in \mathbb{R}^d$ such that

$$261 \quad z_i = \mathbb{I}[\langle X_i, u \rangle \varphi_\theta(X_i, Y_i) \geq t_i] \quad \text{for all } i \in [m].$$

263 Following some algebra, we find that

$$264 \quad y_i x_i^T \theta - \log \frac{Y_i t_i}{\langle X_i, u \rangle - Y_i t_i} \begin{cases} \geq 0 & z_i = 1 \\ < 0 & z_i = 0 \end{cases}.$$

266 Consequently, the set $\{[y_i x_i, \log(Y_i t_i / (\langle X_i, u \rangle - Y_i t_i))]\}_{i=1}^m$ of $(d+1)$ -dimensional points can be
267 shattered by linear separators. Therefore, we have $m \leq d+2$, which leads to the VC subgraph
268 dimension of \mathcal{G} to be at most $d+2$ (e.g., see the book [SM4]). As a consequence, we obtain
269 the conclusion of the claim (B.3).

270 **B.2. Proof of Corollary 4.2.** We prove Corollary 4.2 by verifying the claims (4.5a)
271 and (4.5b).

272 **B.2.1. Structure of F^G .** Direct algebra leads to the following equation

$$273 \quad \langle \nabla F^G(\theta), \theta^* - \theta \rangle = \left(\theta - \mathbb{E} \left[X \tanh \left(X^\top \theta \right) \right] \right)^\top (\theta - \theta^*)$$

$$274 \quad (\text{B.5}) \quad \geq \|\theta\|_2^2 - \|\theta\|_2 \left\| \mathbb{E} \left[X \tanh \left(X^\top \theta \right) \right] \right\|_2$$

276 where $\tanh(x) := \frac{\exp(x) - \exp(-x)}{\exp(x) + \exp(-x)}$ for all $x \in \mathbb{R}$. From Theorem 2 in Dwivedi et al. [SM3], we
277 have

$$278 \quad \left\| \mathbb{E} \left[X \tanh \left(X^\top \theta \right) \right] \right\|_2 \leq \left(1 - p + \frac{p}{1 + \frac{\|\theta\|_2^2}{2}} \right) \|\theta\|_2$$

280 for all $\theta \in \mathbb{R}^d$ where $p := \mathbb{P}(|Y| \leq 1) + \frac{1}{2}\mathbb{P}(|Y| > 1)$ where $Y \sim \mathcal{N}(0, 1)$. Plugging the above
281 inequality into equation (B.5) leads to

$$282 \quad \langle \nabla F^G(\theta), \theta^* - \theta \rangle \geq \frac{p \|\theta\|_2^4}{2 + \|\theta\|_2^2} \geq \begin{cases} \frac{p}{4} \|\theta\|_2^4, & \text{for } \|\theta\|_2 \leq \sqrt{2} \\ \frac{p}{2} (\|\theta\|_2^2 - 1), & \text{otherwise} \end{cases}.$$

284 As a consequence, we achieve the conclusion of claim (4.5a).

285 **B.2.2. Perturbation error between ∇F^G and ∇F_n^G .** Direct calculation indicates the
286 following equation:

$$287 \quad \nabla F_n^G(\theta) - \nabla F^G(\theta) = \frac{1}{n} \sum_{i=1}^n X_i \tanh(X_i^\top \theta) - \mathbb{E} \left[X \tanh \left(X^\top \theta \right) \right].$$

289 The outer expectation in the above display is taken with respect to $X \sim \mathcal{N}(\theta^*, \sigma^2 I_d)$ where
290 $\theta^* = 0$. Based on the proof argument of Lemma 1 from the paper [SM3], for each $r > 0$, we
291 have the following concentration inequality

$$292 \quad \mathbb{P} \left(\sup_{\theta \in \mathbb{B}(\theta^*, r)} \left\| \frac{1}{n} \sum_{i=1}^n X_i \tanh(X_i^\top \theta) - \mathbb{E} \left[X \tanh \left(X^\top \theta \right) \right] \right\|_2 \right. \\ \left. \leq cr \sqrt{\frac{d + \log(1/\delta)}{n}} \right) \geq 1 - \delta,$$

293 (B.6)
294 for any $\delta > 0$ as long as the sample size $n \geq c'd \log(1/\delta)$ where c and c' are universal constants.
296 For any $M \in \mathbb{N}_+$, by the concentration bound (B.6) and the union bound, we find that

$$297 \quad \mathbb{P} \left(\forall r \in [2^{-M}, 1], \sup_{\theta \in \mathbb{B}(\theta^*, r)} \|\nabla F_n^G(\theta) - \nabla F^G(\theta)\|_2 \right. \\ \left. \leq cr \sqrt{\frac{d + \log(M/\delta)}{n}} \right) \geq 1 - \delta.$$

298 (B.7)
299 On the other hand, based on the standard inequality $|\tanh(x)| \leq |x|$ for all $x \in \mathbb{R}$, we find
300 that

$$302 \quad \|\nabla F_n^G(\theta) - \nabla F^G(\theta)\|_2 \leq \frac{1}{n} \sum_{i=1}^n \|X_i\|_2 \left| \tanh \left(X_i^\top \theta \right) \right| + \mathbb{E} \left[\|X\|_2 \left| \tanh \left(X^\top \theta \right) \right| \right] \\ 303 \quad \leq \frac{1}{n} \sum_{i=1}^n \|X_i\|_2 \left| X_i^\top \theta \right| + \mathbb{E} \left[\|X\|_2 \left| X^\top \theta \right| \right] \\ 304 \quad \leq \left(\frac{1}{n} \sum_{i=1}^n \|X_i\|_2^2 + \mathbb{E} \left[\|X\|_2^2 \right] \right) \|\theta\|_2. \\ 305$$

306 Therefore, we have $\|\nabla F_n^G(\theta) - \nabla F^G(\theta)\|_2 \leq 2d\|\theta\|_2 \log(1/\delta)$ with probability $1 - \delta$. By
 307 choosing $M_1 := \log(2nd)$, based on the previous bound, we obtain that

$$308 \quad (B.8) \quad \mathbb{P}\left(\forall r < 2^{-M_1}, \sup_{\theta \in \mathbb{B}(\theta^*, r)} \|\nabla F_n^G(\theta) - \nabla F^G(\theta)\|_2 \leq \frac{\log(1/\delta)}{n}\right) \geq 1 - \delta.$$

310 Furthermore, for vector $\theta \in \mathbb{R}^d$ with large norm, by the concentration bound (B.6) combined
 311 with the union bound, for any $M' \in \mathbb{N}_+$, we find that

$$312 \quad \mathbb{P}\left(\forall r \in [1, 2^{M'}], \sup_{\theta \in \mathbb{B}(\theta^*, r)} \|\nabla F_n^G(\theta) - F^G(\theta)\|_2\right. \\
 313 \quad \left. \leq cr \sqrt{\frac{d + \log(M'/\delta)}{n}}\right) \geq 1 - \delta.$$

315 When r in the above bound is too large, we can simply use the fact that \tanh is a bounded
 316 function. We thus have the upper bound

$$317 \quad \|\nabla F_n^G(\theta) - \nabla F^G(\theta)\|_2 \leq \mathbb{E}[\|X\|_2] + \frac{1}{n} \sum_{i=1}^n \|X_i\|_2,$$

319 for any θ . Given the above bound, by choosing $M_2 := \log(2\sqrt{n})$, we obtain that

$$320 \quad \mathbb{P}\left(\forall r > 2^{M_2}, \sup_{\theta \in \mathbb{B}(\theta^*, r)} \|\nabla F_n^G(\theta) - \nabla F^G(\theta)\|_2 \leq r \sqrt{\frac{d + \log(1/\delta)}{n}}\right) \\
 321 \quad (B.9) \quad \geq \mathbb{P}\left(\mathbb{E}[\|X\|_2] + \frac{1}{n} \sum_{i=1}^n \|X_i\|_2 \leq 2^{M_2} \sqrt{\frac{d + \log(1/\delta)}{n}}\right) \geq 1 - \delta.$$

323 Putting the bounds (B.7), (B.8), and (B.9) together, for $n \geq cd \log(1/\delta)$, the following
 324 probability bound holds

$$325 \quad \mathbb{P}\left(\forall r > 0, \sup_{\theta \in \mathbb{B}(\theta^*, r)} \|\nabla F_n^G(\theta) - \nabla F^G(\theta)\|_2\right. \\
 326 \quad \left. \leq cr \sqrt{\frac{d + \log(\log n/\delta)}{n}} + \frac{\log(1/\delta)}{n}\right) \geq 1 - \delta,$$

328 which completes the proof of the claim (4.5b).

329 Appendix C. Proofs of the remaining auxiliary results.

330 In this appendix, we provide proofs of the remaining auxiliary results in the paper.

331 **C.1. Proof of Proposition 2.1.** For any $p \geq 2$, we define the quantity:

$$332 \quad R_p := \sup_{p \geq 0} (\mathbb{E}_{\pi_t} [\|X\|_2^p])^{1/p} \vee (\mathbb{E}_{\pi^*} [\|X\|_2^p])^{1/p}$$

333

334 For any given value $\bar{R} > 0$, we note the following decomposition:

$$\begin{aligned}
335 & \quad |\mathbb{E}_{\pi_t} [\|X\|_2^p] - \mathbb{E}_{\pi^*} [\|X\|_2^p]| \\
336 & \leq \int_{\mathbb{B}(0, \bar{R})} |\pi_t - \pi^*| \cdot \|x\|_2^p dx + \int_{\mathbb{B}(0, \bar{R})^c} \pi_t(x) \|x\|_2^p dx + \int_{\mathbb{B}(0, \bar{R})^c} \pi^*(x) \|x\|_2^p dx \\
337 & \leq \bar{R}^p \cdot d_{\text{TV}}(\pi_t, \pi^*) + \mathbb{E}_{\pi_t} [\|X\|_2^p \mathbf{1}_{\|X\|_2 > \bar{R}}] + \mathbb{E}_{\pi^*} [\|X\|_2^p \mathbf{1}_{\|X\|_2 > \bar{R}}] \\
338 & \leq \bar{R}^p \cdot d_{\text{TV}}(\pi_t, \pi^*) + \sqrt{\mathbb{E}_{\pi_t} [\|X\|_2^{2p}] \sqrt{\pi_t(\|X\|_2 > \bar{R})}} + \sqrt{\mathbb{E}_{\pi^*} [\|X\|_2^{2p}] \sqrt{\pi^*(\|X\|_2 > \bar{R})}} \\
339 & \leq \bar{R}^p \cdot d_{\text{TV}}(\pi_t, \pi^*) + 2R_{2p}^p \cdot R_2 / \bar{R}.
\end{aligned}$$

341 For any $\varepsilon > 0$, take $\bar{R} := \frac{\varepsilon}{2R_{2p}^p R_2}$, we have that:

$$\begin{aligned}
342 & \quad \lim_{t \rightarrow +\infty} |\mathbb{E}_{\pi_t} [\|X\|_2^p] - \mathbb{E}_{\pi^*} [\|X\|_2^p]| \leq \varepsilon, \\
343 &
\end{aligned}$$

344 which proves the claim.

345 **C.2. A limit result.** We begin with a lemma on the limiting behavior of a certain type of
346 function. The lemma is used in the proof of [Theorem 3.2](#) in [Subsection 5.2](#).

347 **Lemma C.1.** *Let ϕ be a concave and continuous function on $[0, +\infty)$ with $\phi(0) = \phi(c) = 0$
348 for some positive constant $c > 0$. Assume furthermore that $\phi(t) < 0$ for all $t \in (c, \infty)$. Suppose
349 that there exist two continuous functions $f, g : [0, +\infty) \rightarrow [0, +\infty)$ such that $\lim_{t \rightarrow +\infty} g(t)$
350 exists and $f(t) \leq \int_0^t \phi(g(s)) ds$ for all $t \geq 0$. Under these conditions, we have $\lim_{t \rightarrow +\infty} g(t) \leq c$.*

351 *Proof.* Define the limit $A := \lim_{t \rightarrow +\infty} g(t)$, which exists according to the assumptions. We
352 proceed via proof by contradiction. In particular, suppose that $A > c$. Based on the definition
353 of A , for the positive constant $\varepsilon = (A - c)/2 > 0$, we can find a sufficiently large positive
354 constant T such that $g(t) > A - \varepsilon$ for any $t \geq T$. Since the function ϕ is concave, with $\phi(c) = 0$
355 and $\phi(t) < 0$ for $t > c$, we have that ϕ is non-increasing on $[c, +\infty)$, and therefore

$$\begin{aligned}
356 & \quad \delta := \phi(c + \varepsilon) = - \sup_{s \geq c + \varepsilon} \phi(s) < 0. \\
357 &
\end{aligned}$$

358 Therefore, for all $t > T$, we arrive at the following inequalities

$$\begin{aligned}
359 & \quad 0 \leq f(t) \leq \int_0^T \phi(g(s)) ds + \int_T^t \phi(g(s)) ds \leq \int_0^T \phi(g(s)) ds - \delta(t - T). \\
360 &
\end{aligned}$$

361 By choosing $t = 1 + T + \delta^{-1} \int_0^T \phi(g(s)) ds$, the above inequality cannot hold. This yields the
362 desired contradiction, which completes the proof. ■

363 **C.3. A tail bound based on truncation.** We now state an upper deviation inequality
364 based on a truncation argument. Consider a sequence of random variables $\{Y_i\}_{i=1}^n$ satisfying
365 the moment bounds

$$\begin{aligned}
366 & \quad (\text{C.1}) \quad \mathbb{E} [|Y_i|^q] \leq (aq)^{bq} \quad \text{for all } q = 1, 2, \dots \\
367 &
\end{aligned}$$

368 where a, b are universal constants.

369 **Lemma C.2.** *Given an i.i.d. sequence of zero-mean random variables $\{Y_i\}_{i=1}^n$ satisfying the*
 370 *moment bounds (C.1), we have*

$$371 \quad \mathbb{P} \left(\frac{1}{n} \sum_{i=1}^n Y_i \geq (4a)^b \sqrt{\frac{\log 4/\delta}{n}} + \left(a \log \frac{n}{\delta} \right)^b \frac{\log 4/\delta}{n} \right) \leq \delta.$$

373 *Proof.* The proof of the lemma is a direct combination of truncation argument and
 374 Bernstein's inequality. In particular, for each $i \in [n]$, define the truncated random variable
 375 $\tilde{Y}_i := Y_i \mathbb{I}[|Y_i| \leq 3(a \log \frac{n}{\delta})^b]$. With this definition, we have

$$376 \quad \mathbb{P} \left((Y_i)_{i=1}^n \neq (\tilde{Y}_i)_{i=1}^n \right) = \mathbb{P} \left(\max_{1 \leq i \leq n} |Y_i| > 3 \left(a \log \frac{n}{\delta} \right)^b \right)$$

$$377 \quad \leq n \mathbb{P} \left(|Y_i| > 3 \left(a \log \frac{n}{\delta} \right)^b \right) \leq \frac{\delta}{2}.$$

379 Therefore, it is sufficient to study a concentration behavior of the quantity $\sum_{i=1}^n \tilde{Y}_i$. Invoking
 380 Bernstein's inequality [SM2], we obtain that

$$381 \quad \mathbb{P} \left(\frac{1}{n} \sum_{i=1}^n \tilde{Y}_i \geq \varepsilon \right) \leq 2 \exp \left(- \frac{n \varepsilon^2}{2(2a)^{2b} + \frac{2}{3} \varepsilon \cdot 3(a \log \frac{n}{\delta})^b} \right).$$

383 In order to make the RHS of the above inequality less than $\frac{\delta}{2}$, it suffices to set

$$384 \quad \varepsilon = (4a)^b \sqrt{\frac{\log(4/\delta)}{n}} + \left(a \log \frac{n}{\delta} \right)^b \frac{\log(4/\delta)}{n}.$$

386 Collecting all of the above inequalities yields the claim. ■

387 **C.4. Unique positive solution to equation (3.1).** We now establish that equation (3.1)
 388 has a unique positive solution under the stated assumptions. Define the function

$$389 \quad \vartheta(z) := \psi(z) - \left(\varepsilon(n, \delta) \zeta(z) z + \frac{Bz + d + \log(1/\delta)}{n} \right).$$

391 Since $\psi(0) = 0$, we have $\vartheta(0) < 0$. On the other hand, based on Assumption (W.4),
 392 $\liminf_{z \rightarrow +\infty} \vartheta(z) > 0$. Therefore, there exists a positive solution to the equation $\vartheta(z) = 0$.

393 Recall that $\xi : \mathbb{R}_+ \rightarrow \mathbb{R}$ is an inverse function of the strictly increasing function $z \mapsto z\zeta(z)$.
 394 Therefore, we can write the function ϑ as follows:

$$395 \quad \vartheta(z) = \tilde{\vartheta}(r) := \psi(\xi(r)) - \varepsilon(n, \delta) r - \frac{B\xi(r) + d + \log(1/\delta)}{n},$$

397 where $r = z \cdot \zeta(z)$. Given the convexity of function $r \mapsto \psi(\xi(r))$ guaranteed by Assump-
 398 tion (W.3), the functions $\tilde{\vartheta}$ and ϑ are convex. Putting the above results together, there exists
 399 a unique positive solution to equation (3.1).

- 401 [1] R. ADAMCZAK, *A tail inequality for suprema of unbounded empirical processes with applications to markov*
402 *chains*, *Electronic Journal of Probability*, 13 (2008), pp. 1000–1034.
- 403 [2] S. BOUCHERON, G. LUGOSI, AND P. MASSART, *Concentration Inequalities: A Nonasymptotic Theory of*
404 *Independence*, Oxford University Press, 2016.
- 405 [3] R. DWIVEDI, N. HO, K. KHAMARU, M. J. WAINWRIGHT, M. I. JORDAN, AND B. YU, *Singularity,*
406 *misspecification, and the convergence rate of EM*, arXiv preprint arXiv:1810.00828, (2018).
- 407 [4] A. W. VAN DER VAART AND J. A. WELLNER, *Weak Convergence and Empirical Processes: With Applications*
408 *to Statistics*, Springer-Verlag, New York, NY, 2000.
- 409 [5] M. J. WAINWRIGHT, *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*, Cambridge University
410 Press, 2019.