

On free energy barriers in Gaussian priors and failure of cold start MCMC for high-dimensional unimodal distributions

Afonso S. Bandeira, Antoine Maillard, Richard Nickl*, Sven Wang

ETH Zürich, ETH Zürich, University of Cambridge, MIT

November 22, 2022

Abstract

We exhibit examples of high-dimensional unimodal posterior distributions arising in non-linear regression models with Gaussian process priors for which MCMC methods can take an exponential run-time to enter the regions where the bulk of the posterior measure concentrates. Our results apply to worst-case initialised (‘cold start’) algorithms that are local in the sense that their step-sizes cannot be too large on average. The counter-examples hold for general MCMC schemes based on gradient or random walk steps, and the theory is illustrated for Metropolis-Hastings adjusted methods such as pCN and MALA.

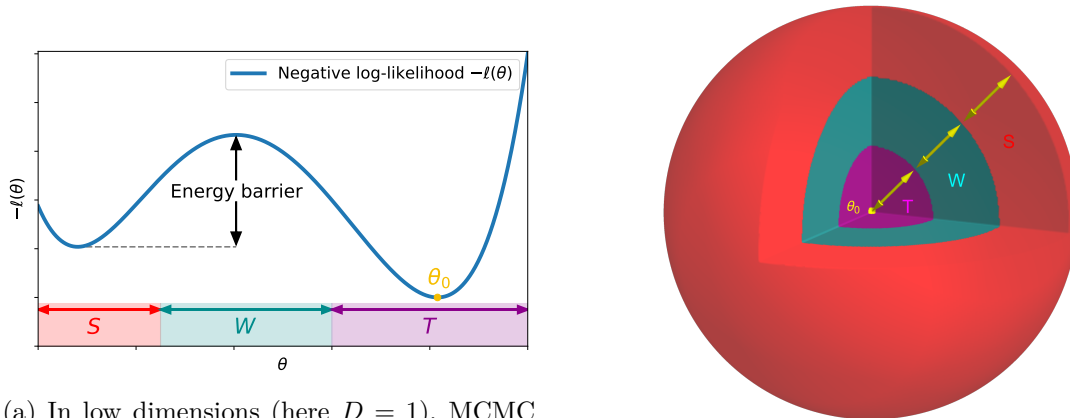
1 Introduction

Markov Chain Monte Carlo (MCMC) methods are the workhorse of Bayesian computation when closed formulas for estimators or probability distributions are not available. For this reason they have been central to the development and success of high-dimensional Bayesian statistics in the last decades, where one attempts to generate samples from some *posterior distribution* $\Pi(\cdot|\text{data})$ arising from a prior Π on D -dimensional Euclidean space and the observed data vector. MCMC methods tend to perform well in a large variety of problems, are very flexible and user-friendly, and enjoy many theoretical guarantees. Under mild assumptions, they are known to converge to their stationary ‘target’ distributions as a consequence of the ergodic theorem, albeit perhaps at a slow speed, requiring a large number of iterations to provide numerically accurate algorithms. When the target distribution is log-concave, MCMC algorithms are known to mix rapidly, even in high dimensions. But for general D -dimensional densities, we have only a restricted understanding of the scaling of the mixing time of Markov chains with D or with the ‘informativeness’ (sample size or noise level) of the data vector.

A classical source of difficulty for MCMC algorithms are multi-modal distributions. When there is a deep well in the posterior density between the starting point of an MCMC algorithm and the location where the posterior is concentrated, many MCMC algorithms are known to take an exponential time – proportional to the depth of the well – when attempting to reach the target region, even in low-dimensional settings, see Fig. 1a and also the discussion surrounding Proposition 4.2 below. However, for distributions with a single mode and when the dimension D is fixed, MCMC methods can usually be expected to perform well.

In essence this article is an attempt to explain how, in high dimensions, wells can be formed *without* multi-modality of a given posterior distribution. The difficulty in this case is volumetric, also referred to as *entropic*: while the target region contains most of the posterior mass, its (prior) volume is so small compared to the rest of the space that an MCMC algorithm may take an exponential time to find it, see Fig. 1b. This competition between ‘energy’ – here represented by the log-likelihood ℓ_N in the posterior distribution $d\Pi(\cdot|\text{data}) = \exp\{\ell_N + \log d\pi\}$ – and ‘entropy’ (related to the prior term π) has also been exploited in recent work on statistical aspects of MCMC in various high dimensional inference and statistical physics models [And89; MM09; ZK16; BGJ20a; BWZ20]. These ideas somewhat date back to the 19th century foundations of statistical mechanics [Gib73] and the

*Correspondence: nickl@maths.cam.ac.uk. RN would like to thank the Forschungsinstitut für Mathematik (FIM) at ETH Zürich for their hospitality during a sabbatical visit in spring 2022 where this research was initiated.



(a) In low dimensions (here $D = 1$), MCMC hardness usually arises because of a non-unimodal likelihood, creating an “energy barrier”, even though the maximum likelihood is attained at $\theta = \theta_0$. The MCMC algorithm is assumed to be initialised in the set \mathcal{S} containing a local maximum of the likelihood.

(b) Illustration of the arising of entropic (or volumetric) difficulties, here in dimension $D = 3$: the set of points close to θ_0 has much less volume than the set of points far away. As D increases, this phenomenon is amplified: all ratios of volumes of the three sets $\mathcal{T}, \mathcal{W}, \mathcal{S}$ scale exponentially with D .

Figure 1: Two possible sources of MCMC hardness in high dimensions: multimodal likelihoods and entropic barriers.

notion of free energy, consisting of a sum of energetic and entropic contributions which the system spontaneously attempts to minimise. The “MCMC-hardness” phenomenon described above is then akin to the meta-stable behavior of thermodynamical systems, such as glasses or supercooled liquids. As the temperature decreases, such systems can undergo a “first-order” phase transition, in which a global free energy minimum (analogous to the target region above) abruptly appears, while the system remains trapped in a suboptimal local minimum of the free energy (the starting region of the MCMC algorithm). For the system to go to thermodynamic equilibrium it must cross an extensive free energy barrier: such a crossing requires an exponentially long time, so that the system appears equilibrated on all relevant timescales, similarly to the MCMC stuck in the starting region. Classical examples include glasses and the popular experiment of rapid freezing of supercooled water (i.e. water that remained liquid at negative temperatures) after introducing a perturbation.

Inspired by recent work [BGJ20a; BWZ20; Ban+22], let us illustrate some of the volumetric phenomena which are key to our results below. We separate the parameter space into three regions (see Figures 1 and 2), which we name by common MCMC terminology. Firstly a *starting* (or initialisation) region \mathcal{S} , where an algorithm starts, secondly a *target* region \mathcal{T} where both the bulk of the posterior mass and the ground truth are situated, and thirdly an intermediate *Free-Entropy well*¹ \mathcal{W} that separates \mathcal{S} from \mathcal{T} . In our theorems, these regions will be characterised by their Euclidean distance to the ground truth parameter θ_0 generating the data. The prior volumes of the ϵ -annuli $\{\theta : r - \epsilon < \|\theta - \theta_0\|_2 \leq r\}, r > 0$, closer to the ground truth are smaller than those further out as illustrated in Fig. 1b, and in high dimensions this effect becomes quantitative in an essential way. Specifically, the trade-off between the entropic and energetic terms can happen such that the following three statements are simultaneously true.

- (i) \mathcal{T} contains “almost all” of the posterior mass.
- (ii) As one gets closer to \mathcal{T} (and thus the ground truth θ_0), the log-likelihood is strictly monotonically increasing.
- (iii) Yet \mathcal{S} still possesses exponentially more posterior mass than \mathcal{W} .

Using ‘bottleneck’ arguments from Markov chain theory (Ch.7 in [Jer03]), this means that an MCMC algorithm that starts in \mathcal{S} is expected to take an exponential time to visit \mathcal{W} . If the step size is such that it cannot “jump over” \mathcal{W} , this also implies an exponential hitting time lower bound for reaching \mathcal{T} . This is illustrated in Figure 2 for an averaged version of the model described in Section 2.

¹As classical in statistical physics, we call free entropy the negative of the free energy.

²In a physical system these regions would correspond respectively to a region including a meta-stable state, a region including the globally stable state, and a free energy barrier.

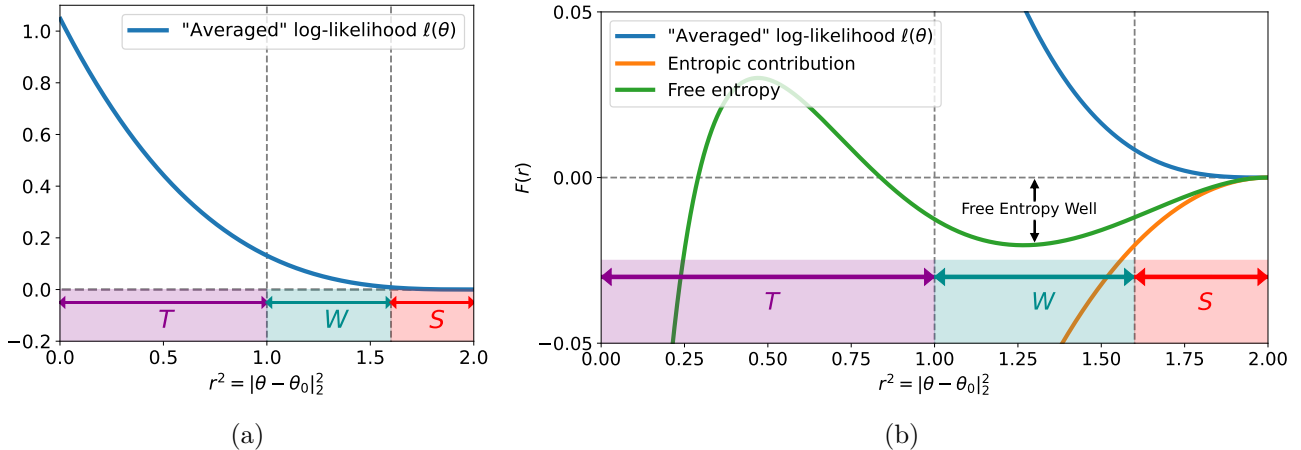


Figure 2: Illustration of a free-energy barrier (or free entropy well) arising with a unimodal posterior. The model is an “averaged” version of the spiked tensor model, with log-likelihood $\ell_n(\theta) = \lambda \langle \theta, \theta_0 \rangle^3 / 2$ and uniform prior Π on the n -dimensional unit sphere \mathbb{S}^{n-1} . θ_0 is chosen arbitrarily on \mathbb{S}^{n-1} . The posterior is $d\Pi(\theta|Y) \propto \exp\{n\ell_n(\theta)\}d\Pi(\theta)$, for $\theta \in \mathbb{S}^{n-1}$. Up to a constant, the free entropy $F(r) = (1/n) \log \int d\Pi(\theta|Y) \delta(r - \|\theta - \theta_0\|_2)$ can be decomposed as the sum of $\ell_n(\theta)$ (that only depends on $r = \|\theta - \theta_0\|_2$) and the “entropic” contribution $(1/n) \log \int d\Pi(\theta) \delta(r - \|\theta - \theta_0\|_2)$. In the figure we show $\lambda = 2.1$.

In the situation described above, the MCMC iterates never visit the region where the posterior is statistically informative, and hence yield no better inference than a random number generator. One could regard this as a ‘hardness’ result about computation of posterior distributions in high dimensions by MCMC. In this work we show that such situations can occur generically and establish hitting time lower bounds for common gradient or random walk based MCMC schemes in model problems with non-linear regression and Gaussian process priors. Before doing this, we briefly review some important results of [BGJ20a] for the problem of Tensor PCA, from which the inspiration for our work was drawn. This technique to establish lower bounds for MCMC algorithms has also recently been leveraged in [BWZ20] in the context of sparse PCA, and in [Ban+22] to establish connections between MCMC lower bounds and the Low Degree Method for algorithmic hardness predictions (see [KWB22] for an expository note on this technique).

When the target distribution is globally log-concave, pictures such as in Fig. 2 are ruled out (see also Remark 4.7) and polynomial-time mixing bounds have been shown for a variety of commonly used MCMC methods. While an exhaustive discussion would be beyond the scope of this paper, we mention here the seminal works [Dal17; DM19] which were among the first to demonstrate high-dimensional mixing of discretised Langevin methods (even upon ‘cold-start’ initialisations like the ones assumed in the present paper). In concrete non-linear regression models, polynomial-time computation guarantees were given in [NW20] under a general ‘gradient stability’ condition on the regression map which guarantees that the posterior is (with high probability) locally log-concave on a large enough region including θ_0 . While this condition can be expected to hold under natural injectivity hypotheses and was verified for an inverse problem with the Schrödinger equation in [NW20], for non-Abelian X-ray transforms in [BN21], the ‘Darcy flow’ elliptic PDE model in [Nic22] and for generalized linear models in [Alt22], all these results hinge on the existence of a suitable initialiser of the gradient MCMC scheme used. These results form part of a larger research program [Nic20; MNP19; MNP21b; MNP21a; Nic22] on algorithmic and statistical guarantees for Bayesian inversions methods [Stu10] applied to problems with partial differential equations. The present article shows that the hypothesis of existence of a suitable initialiser is – at least in principle – essential in these results if $D/N \rightarrow \kappa > 0$, and that at most ‘moderately’ high-dimensional ($D = o(N)$) MCMC implementations of Gaussian process priors may be preferable to bypass computational bottlenecks.

Our negative results apply to (worst-case initialised) Markov chains whose step-sizes cannot be too large with high probability. As we show this includes many commonly used algorithms (such as pCN and MALA) whose dynamics are of a ‘local’ nature. There are a variety of MCMC methods developed recently, such as piece-wise deterministic Markov processes, boomerang or zig-zag samplers [Fea+18; BVD18; Bie+20; WR20] which may not fall into our framework. While we are not aware of any rigorous results that would establish polynomial hitting or mixing times of these algorithms for high-dimensional posterior distributions such as those exhibited here, it is of great interest to study whether our computational hardness barriers can be overcome by ‘non-local’ methods. There is some

empirical evidence that this may be possible. For instance, in the numerical simulation of models of supercooled liquids [SGB22], methods such as swap Monte Carlo [GP01] have been observed to equilibrate to low-temperature distributions which were not reachable by local approaches. Another example is given by the planted clique problem [Jer92]: this model is conjectured to possess a large algorithmically hard phase, and local Monte Carlo methods are known to fail far from the conjectured algorithmic threshold [GZ19; AFF21; CMZ22]. On the other hand, non-local exchange Monte Carlo methods (such as parallel tempering [HN96]), have been numerically observed to perform significantly better [Ang18].

2 The spiked tensor model: an illustrative example

In this section, we present (a simplified version of) results obtained mostly in [BGJ20a]. First some notation. For any $n \geq 1$, we denote by $\mathbb{S}^{n-1} = \{\theta \in \mathbb{R}^n : \|\theta\|_2 = 1\}$ the Euclidean unit sphere in n dimensions. For $\theta, \theta' \in \mathbb{R}^n$ we denote $\theta \otimes \theta' = (\theta_i \theta'_j)_{1 \leq i, j \leq n} \in \mathbb{R}^{n^2}$ their tensor product.

Spiked tensor estimation is a synthetic model to study tensor PCA, and corresponds to a Gaussian Additive Model with a low-rank prior. More formally, it can be defined as follows [RM14].

Definition 2.1 (Spiked tensor model). *Let $p \geq 3$ denote the order of the tensor. The observations Y and the parameter θ are generated according to the following joint probability distribution:*

$$d\mathbb{Q}(Y, \theta) = \frac{1}{(2\pi)^{n^p/2}} \exp\left\{-\frac{1}{2}\|Y - \sqrt{n}\lambda\theta^{\otimes p}\|_2^2\right\} d\Pi(\theta) dY. \quad (1)$$

Here, dY denotes the Lebesgue measure on the space $(\mathbb{R}^n)^{\otimes p} = \mathbb{R}^{n^p}$ of p -tensors of size n . Π is the uniform probability measure on \mathbb{S}^{n-1} , and $\lambda \geq 0$ is the signal-to-noise ratio (SNR) parameter. In particular, the posterior distribution $\Pi(\theta|Y)$ is:

$$d\Pi(\theta|Y) = \frac{1}{\mathcal{Z}_Y} \exp(\ell_{n,Y}(\theta)) d\Pi(\theta), \quad (2)$$

in which \mathcal{Z}_Y is a normalization, and we defined the log-likelihood (up to additive constants) as:

$$\ell_{n,Y}(\theta) = \frac{1}{2} \sqrt{n} \lambda \langle \theta^{\otimes p}, Y \rangle. \quad (3)$$

In the following, we study the model from Definition 2.1 via the prism of statistical inference. In particular, we will study the posterior $\Pi(\theta|Y)$ for a fixed³ “data tensor” Y . Since such a tensor was generated according to the marginal of (1), we parameterise it as $Y = \lambda\sqrt{n}\theta_0^{\otimes p} + Z$, with Z a p -tensor with i.i.d. $\mathcal{N}(0, 1)$ coordinates, and θ_0 a “ground-truth” vector uniformly-sampled in \mathbb{S}^{n-1} . The goal of our inference task is to recover information on the low-rank perturbation $\theta_0^{\otimes p}$ (or equivalently on the vector θ_0 , possibly up to a global sign depending on the parity of p) from the posterior distribution $\Pi(\cdot|Y)$.

Crucially, we are interested in the limit of the model of Definition 2.1 as $n \rightarrow \infty$. In particular, all our statements, although sometimes non-asymptotic, are to be interpreted as n grows. We say that an event occurs “with high probability” (w.h.p.) when its probability is $1 - \mathcal{O}_n(1)$ ⁴. Moreover, by rotation invariance, all statements are uniform over $\theta_0 \in \mathbb{S}^{n-1}$, so that said probabilities only refer to the noise tensor Z . Finally, throughout our discussion we will work with latitude intervals (or bands) on the sphere, with the North Pole taken to be θ_0 . We characterise them using inner products (correlations) $\langle \theta, \theta_0 \rangle$ for odd p , and $|\langle \theta, \theta_0 \rangle|$ for even p (since in this case θ_0 and $-\theta_0$ are indistinguishable from the point of view of the observer).

Definition 2.2 (Latitude intervals). *Assume that $p \geq 3$ is even. For $0 \leq s < t \leq 1$ we define:*

- $\mathcal{S}_s = \{\theta \in \mathbb{S}^{n-1} : |\langle \theta, \theta_0 \rangle| \leq s\}$,
- $\mathcal{W}_{s,t} = \{\theta \in \mathbb{S}^{n-1} : s < |\langle \theta, \theta_0 \rangle| \leq t\}$,
- $\mathcal{T}_t = \{\theta \in \mathbb{S}^{n-1} : t < |\langle \theta, \theta_0 \rangle|\}$.

If p is odd, we define these sets similarly, replacing $|\langle \theta, \theta_0 \rangle|$ by $\langle \theta, \theta_0 \rangle$.

Note that these sets can also be characterised using the distance to the ground-truth, e.g. $\mathcal{S}_s = \{\theta \in \mathbb{S}^{n-1} : \min\{\|\theta - \theta_0\|_2^2, \|\theta + \theta_0\|_2^2\} \geq 2(1-s)\}$ when p is even.

³Note that we assume here that the statistician has access to the distribution $\Pi(\cdot|Y)$ (and in particular to λ), a setting sometimes called *Bayes-optimal* in the literature.

⁴Often the $\mathcal{O}_n(1)$ term will be exponentially small, but we will not require such a strong control.

2.1 Posterior contraction

We can use uniform concentration of the likelihood to show that as $\lambda \rightarrow \infty$ (after taking the limit $n \rightarrow \infty$) the posterior contracts in a region infinitesimally close to the ground-truth θ_0 . We first show that a region arbitrarily close to the ground truth exponentially dominates a very large starting region:

Proposition 2.3. *For any $K > 0$ there exists $\lambda_0 > 0$ and functions $\{s(\lambda), t(\lambda)\} \in [0, 1]$ such that $s(\lambda) < t(\lambda)$, $\{s(\lambda), t(\lambda)\} \rightarrow 1$ as $\lambda \rightarrow \infty$, and for all $\lambda \geq \lambda_0$:*

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log \frac{\Pi(\mathcal{S}_{s(\lambda)}|Y)}{\Pi(\mathcal{T}_{t(\lambda)}|Y)} \leq -K, \quad \text{almost surely.} \quad (4)$$

Posterior contraction is the content of the following result:

Corollary 2.4 (Posterior contraction). *There exists $\lambda_0 > 0$ and a function $s(\lambda) \in [0, 1]$ satisfying $s(\lambda) \rightarrow 1$ as $\lambda \rightarrow \infty$, such that for all $\lambda \geq \lambda_0$:*

$$\lim_{n \rightarrow \infty} \Pi[\mathcal{T}_{s(\lambda)}|Y] = 1, \quad \text{almost surely.} \quad (5)$$

The proofs of Proposition 2.3 and Corollary 2.4 are given in Appendix A.

Remark 2.5 (Suboptimality of uniform bounds). Stronger than Corollary 2.4, it is known that there exists a sharp threshold $\lambda^*(p)$ such that for any $\lambda > \lambda^*(p)$ the posterior mean, as well as the maximum likelihood estimator, sit w.h.p. in $\mathcal{T}_{s(\lambda)}$, with $s(\lambda) > 0$, while such a statement is false for $\lambda \leq \lambda^*(p)$ [PWB20; Les+17; JLM20]. The λ_0 given by Corollary 2.4 is, on the other hand, clearly not sharp, because of the crude uniform bound used in the proof. This can easily be understood in the $p = 2$ case, corresponding to rank-one matrix estimation: uniform bounds such as the ones used here would show posterior contraction for $\lambda = \omega(1)$, while it is known through the celebrated BBP transition that the maximum likelihood estimator is already correlated with the signal for any $\lambda > 1$ [BBP05]. With more refined techniques from the study of random matrices and spin glass theory of statistical physics it is often possible to obtain precise constants for such relevant thresholds.

2.2 Algorithmic bottleneck for MCMC

Simple volume arguments, associated with an ingenious use of Markov’s inequality due to [BGJ20a] and of the rotation-invariance of the noise tensor Z , allow to get a computational hardness result for MCMC algorithms, even though the posterior contracts infinitesimally close to the ground truth as we saw in Corollary 2.4. In the context of the spiked tensor model, these computational hardness results can be found in [BGJ20a] (see in particular Section 7). We will state similar results for general non-linear regression models in Section 3: in this context we will not need to use the Markov’s inequality-based technique of [BGJ20a], and will solely rely on concentration arguments.

Recall that by Section 2.1, we can find $s(\lambda)$ such that $s(\lambda) \rightarrow 1$ as $\lambda \rightarrow \infty$ and for all λ large enough $\Pi(\mathcal{T}_{s(\lambda)}|Y) = 1 - \mathcal{O}_n(1)$. Here, we show that escaping the “initialisation” region of the MCMC algorithm is hard in a large range of λ (possibly diverging with n). In what follows, the step size of the algorithm denotes the maximal change $\|x^{t+1} - x^t\|_2$ allowed in any iteration⁵. We first state this bottleneck result informally.

Proposition 2.6 (MCMC Bottleneck, informal). *Assume that $\lambda = \mathcal{O}(n^{(p-2)/4+\eta})$ for all $\eta > 0$. Then any MCMC algorithm whose invariant distribution is $\Pi(\cdot|Y)$, and with a step size bounded by $\delta = \mathcal{O}([n\lambda^2]^{-1/p})$, will take an exponential time to get out of the “initialisation” region.*

Note that the step size condition of Proposition 2.6 is always meaningful, since our hypothesis on λ implies $[n\lambda^2]^{-1/p} = \omega(n^{-1/2})$, and many MCMC algorithms (e.g. any procedure in which a number $\mathcal{O}(1)$ of coordinates of the current iterate are changed in a single iteration) will have a step size $\mathcal{O}(n^{-1/2})$.

Remark 2.7. The results of [BGJ20a] are stated when considering for the invariant distribution of the MCMC a more general “Gibbs-type” distribution $\mathbb{G}_{\beta,Y}(dx) \propto e^{\beta H(x)} d\Pi(x)$, with $H(x) = (\sqrt{n}/2)\langle x^{\otimes p}, Y \rangle$. The case we consider here is the “Bayes-optimal” $\beta = \lambda$, for which $\mathbb{G}_{\lambda,Y} = \Pi(\cdot|Y)$. For the general distribution $\mathbb{G}_{\beta,Y}$ the conditions of Proposition 2.6 become $\beta\lambda = \mathcal{O}(n^{(p-2)/2+\eta})$ and $\delta = \mathcal{O}([n\beta\lambda]^{-1/p})$. The authors of [BGJ20a] usually consider $\beta = \mathcal{O}(1)$, so that they show the bottleneck under the condition $\lambda = \mathcal{O}(n^{(p-2)/2+\eta})$.

⁵As we will detail in the following sections, see Assumption 3.1, the statements remain true if the change is allowed to be higher than the required maximum with exponentially small probability.

More generally, $\lambda \ll n^{(p-2)/4}$ is conjectured to be a regime in which *all* polynomial-time algorithms fail to recover θ_0 [RM14; WEM19; HSS15; Hop+16; KBG17]. On the other hand, “local” methods (such as gradient-based algorithms [Sar+19b; Sar+19a; BCR20; BGJ20b; BGJ21], message-passing iterations [Les+17], or natural MCMC algorithms such as the ones of previous remark) are conjectured or known to fail in the larger range $\lambda \ll n^{(p-2)/2}$. Proposition 2.6 shows that “Bayes-optimal” MCMC algorithms fail for $\lambda \ll n^{(p-2)/4}$. To the best of our knowledge, analysing this class of algorithms in the regime $n^{(p-2)/4} \ll \lambda \ll n^{(p-2)/2}$ is still open.

Let us now state formally the key ingredient behind Proposition 2.6. It is a rewriting of the “free energy wells” result of [BGJ20a].

Lemma 2.8 (Bottleneck, formal). *Assume that $\lambda = \mathcal{O}(n^{(p-2)/4+\eta})$ for all $\eta > 0$, and let $\delta = \mathcal{O}([n\lambda^2]^{-1/p})$. Let $r(\varepsilon) = n^{-1/2+\varepsilon}$. Then for any $\varepsilon > 0$ small enough, there exists $c, C > 0$ such that for large enough n , with probability at least $1 - \exp(-cn^{2\varepsilon})$ we have:*

$$\frac{\Pi(\mathcal{S}_{r(\varepsilon)}|Y)}{\Pi(\mathcal{W}_{r(\varepsilon), r(\varepsilon)+\delta}|Y)} \geq \exp\{Cn^{2\varepsilon}\}. \quad (6)$$

Note that by simple volume arguments, $\Pi(\mathcal{S}_{r(\varepsilon)}) = 1 - \mathcal{O}_n(1)$, so that $\mathcal{S}_{r(\varepsilon)}$ contains “almost all” the mass of the uniform distribution.

One can then deduce from Lemma 2.8 hitting time lower bounds for MCMCs using a folklore bottleneck argument – see Jerrum [Jer03] – that we recall here in a simplified form (see also [BWZ20], as well as Proposition 4.4, where we will detail it further along with a short proof).

Proposition 2.9. *We fix any Y and n , and let any $0 < s < t < 1$. Let $\theta^{(0)}, \theta^{(1)}, \dots$ be a Markov chain on \mathbb{S}^{n-1} with stationary distribution $\Pi(\cdot|Y)$, and initialised from $\theta^{(0)} \sim \Pi_{\mathcal{S}_s}(\cdot|Y)$, the posterior distribution conditioned on \mathcal{S}_s . Let $\tau_t = \inf\{k \in \mathbb{N} : \theta^{(k)} \in \mathcal{T}_t\}$ be the hitting time of the Markov chain onto \mathcal{T}_t . Then, for any $k \geq 1$,*

$$\Pr(\tau_t \leq k) \leq k \frac{\Pi(\mathcal{W}_{s,t}|Y)}{\Pi(\mathcal{S}_s|Y)}. \quad (7)$$

Remark 2.10 (MCMC initialisation). Note that Lemma 2.8, combined with Proposition 2.9, shows hardness of MCMC initialised in points drawn from $\Pi_{\mathcal{S}_{r(\varepsilon)}}(\cdot|Y)$. In particular, it is easy to see that this implies (via the probabilistic method) the existence of such “hard” initialising points. While one might hope to show such negative results for more general initialisation, this remains an open problem. On the other hand, [BGJ20a] shows that there exists initialisers in $\mathcal{S}_{r(\varepsilon)}$ for which vanilla Langevin dynamics achieve non-trivial recovery of the signal even for $\lambda = \Theta_n(1)$ (a phenomenon they call “equatorial passes”).

3 Main results for non-linear regression with Gaussian priors

We now turn to the main contribution of this article, which is to exhibit some of the phenomena described in Section 2 in the context of non-linear regression models. All the theorems of this section are proven in detail in Section 4.

Consider data $Z^{(N)} =^{iid} (Y_i, X_i)_{i=1}^N$ from the random design regression model

$$Y_i = \mathcal{G}(\theta)(X_i) + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(0, 1), \quad i = 1, \dots, N, \quad (8)$$

where $\mathcal{G} : \Theta \rightarrow L^2_\mu(\mathcal{X})$ is a regression map taking values in the space $L^2(\mathcal{X}) = L^2_\mu(\mathcal{X})$ on some bounded subset \mathcal{X} of \mathbb{R}^d , and where the $X_i \sim^{iid} \mu$ are drawn uniformly on \mathcal{X} . For convenience, we assume that \mathcal{X} has Lebesgue measure $\int_{\mathcal{X}} dx = 1$. The law of the data $dP_\theta^N(z_1, \dots, z_N) = \prod_{i=1}^N dP_\theta(z_i)$ is a product measure on $(\mathbb{R} \times \mathcal{X})^N$, with associated expectation operator E_θ^N . Here θ varies in some parameter space

$$\Theta \subseteq \mathbb{R}^D, \quad \frac{D}{N} \simeq \kappa \geq 0,$$

and $\theta_0 \in \Theta$ is a ‘ground truth’ (we could use ‘mis-specified’ θ_0 and project it onto Θ). We will primarily consider the case where $\kappa > 0$ and $\Theta = \mathbb{R}^D$, and consider high-dimensional asymptotics where D (and then also N) diverge to infinity, even though some aspects of our proofs do not rely

on these assumptions. We will say that events A_N hold with high probability if $P_{\theta_0}^N(A_N) \rightarrow 1$ as $N \rightarrow \infty$, and we will use the same terminology later when it involves the law of some Markov chain.

Let Π be a prior (Borel probability measure) on Θ so that given the data $Z^{(N)}$ the posterior measure is the ‘Gibbs’-type distribution

$$d\Pi(\theta|Z^{(N)}) = \frac{e^{\ell_N(\theta)} d\Pi(\theta)}{\int_{\Theta} e^{\ell_N(\theta)} d\Pi(\theta)}, \quad \theta \in \Theta, \quad (9)$$

where

$$\ell_N(\theta) = -\frac{1}{2} \sum_{i=1}^N |Y_i - \mathcal{G}(\theta)(X_i)|^2, \quad \ell(\theta) = \mathbb{E}_{\theta_0}^N \ell_N(\theta), \quad \theta \in \Theta.$$

3.1 Hardness examples for posterior computation with Gaussian priors

We are concerned here with the question of whether one can sample from the Gibbs’ measure (9) by MCMC algorithms. The priors will be Gaussian, so the ‘source’ of the difficulty will arise from the log-likelihood function ℓ_N . On the one hand, recent work ([Dal17; DM19; NW20; BN21; Nic22]) has demonstrated that if $\ell_N(\theta)$ is ‘on average’ (under E_{θ_0}) log-concave, possibly only just locally near the ground truth θ_0 , then MCMC methods that are initialised into the area of log-concavity can mix towards $\Pi(\cdot|Z^{(N)})$ in polynomial time even in high-dimensional ($D \rightarrow \infty$) and ‘informative’ ($N \rightarrow \infty$) settings. In absence of such structural assumptions, however, posterior computation may be intractable, and the purpose of this section is to give some concrete examples for this with choices of \mathcal{G} that are representative for non-linear regression models.

We will provide lower bounds on the run-time of ‘worst case’ initialised MCMC in settings where the average posterior surface is not *globally* log-concave but still unimodal. Both the log-likelihood function and posterior density exhibit linear growth towards their modes, and the average log-likelihood is locally log-concave at θ_0 . In particular the Fisher information is well defined and non-singular at the ground truth.

The computational hardness does not arise from a local optimum (‘multimodality’), but from the difficulty MCMC encounters in ‘choosing’ among many high-dimensional directions when started away from the bulk of the support of the posterior measure. That such problems occur in high dimensions is related to the probabilistic structure of the prior Π , and the manifestation of ‘free energy barriers’ in the posterior distribution.

In many applications of Bayesian statistics, such as in machine learning or in non-linear inverse problems with PDEs, *Gaussian process priors* are commonly used for inference. To connect to such situations we illustrate the key ideas that follow with two canonical examples where the prior on \mathbb{R}^D is the law

$$a) \theta \sim \mathcal{N}(0, \mathbf{I}_D/D), \quad \text{or } b) \theta \sim \mathcal{N}(0, \Sigma_\alpha), \quad (10)$$

where Σ_α is the covariance matrix arising from the law of a d -dimensional Whittle-Matérn-type Gaussian random field (see Section 4.4.1 for a detailed definition). These priors represent widely popular choices in Bayesian statistical inference [RW06; GV17] and can be expected to yield consistent statistical solutions of regression problems even when $D/N \geq \kappa > 0$, see [VZ08; GV17]. In b) we can also accommodate a further ‘rescaling’ (N -dependent shrinkage) of the prior similar to what has been used in recent theory for non-linear inverse problems ([MNP21b], [NW20], [BN21]), see Remark 4.6 for details.

We will present our main results for the case where the ground truth is $\theta_0 = 0$. This streamlines notation while also being the ‘hardest’ case for negative results, since the priors from a) and b) are then already centred at the correct parameter.

To formalise our results, let us define balls

$$B_r = \{\theta \in \mathbb{R}^D : \|\theta\|_{\mathbb{R}^D} \leq r\}, \quad r > 0, \quad (11)$$

centred at $\theta_0 = 0$. We will also require the annuli

$$\Theta_{r,\varepsilon} = \{\theta \in \mathbb{R}^D : \|\theta\|_{\mathbb{R}^D} \in (r, r + \varepsilon)\}, \quad (12)$$

for $r, \varepsilon > 0$ to be chosen. To connect this to the notation in the preceding sections, the sets $\Theta_{r,\varepsilon}$ will play the role of the initialisation (or starting) region \mathcal{S} , while B_s (for suitable s) corresponds

to the target region \mathcal{T} where the posterior mass concentrates. The ‘intermediate’ region $\mathcal{W} = \Theta_{s,\eta}$ representing the ‘free-energy barrier’ is constructed in the proofs of the theorems to follow.

Our results hold for general Markov chains whose invariant measure equals the posterior measure (9), and which admit a bound on their ‘typical’ step-sizes. As step-sizes can be random, this assumption needs to be accommodated in the probabilistic framework describing the transition probabilities of the chain. Let $\mathcal{P}_N(\theta, A)$, $N \in \mathbb{N}$, (for $\theta \in \mathbb{R}^D$ and Borel sets $A \subseteq \mathbb{R}^D$), denote a sequence of Markov kernels describing the Markov chain dynamics employed for the computation of the posterior distribution $\Pi(\cdot|Z^{(N)})$. Recall that a probability measure μ on \mathbb{R}^D is called invariant for \mathcal{P}_N if $\int_{\mathbb{R}^D} \mathcal{P}_N(\theta, A) d\mu(\theta) = \mu(A)$ for all Borel sets A .

Assumption 3.1. *Let $\mathcal{P}_N(\cdot, \cdot)$ be a sequence of Markov kernels satisfying the following:*

- i) $\mathcal{P}_N(\cdot, \cdot)$ has invariant distribution $\Pi(\cdot|Z^{(N)})$ from (9).*
- ii) For some fixed $c_0 > 0$ and for sequences $L = L_N > 0$, $\eta = \eta_N > 0$, with P_0^N -probability approaching 1 as $N \rightarrow \infty$,*

$$\sup_{\theta \in B_L} \mathcal{P}_N(\theta, \{\vartheta : \|\theta - \vartheta\|_{\mathbb{R}^D} \geq \eta/2\}) \leq e^{-c_0 N}, \quad N \geq 1.$$

This assumption states that typical steps of the Markov chain are, with high probability (both under the law of the Markov chain and the randomness of the invariant ‘target’ measure), concentrated in an area of size $\eta/2$ around the current state θ , uniformly in a ball of radius L around $\theta_0 = 0$. For standard MCMC algorithms (such as pCN, MALA) whose proposal steps are based on the discretisation of some continuous-time diffusion process, such conditions can be checked, as we will show in the next subsection.

Theorem 3.2. *Let $D/N \simeq \kappa > 0$, consider the posterior (9) arising from the model (8) and a $\mathcal{N}(0, I_D/D)$ prior of density π , and let $\theta_0 = 0$. Then there exists \mathcal{G} and a fixed constant $s \in (0, 1/3)$ for which the following statements hold true.*

- i) The expected likelihood $\ell(\theta)$ is unimodal with mode 0, locally log-concave near 0, radially symmetric, Lipschitz-continuous and monotonically decreasing in $\|\theta\|_{\mathbb{R}^D}$ on \mathbb{R}^D .*
- ii) For any fixed $r > 0$, with high probability the log-likelihood $\ell_N(\theta)$ and the posterior density $\pi(\cdot|Z^{(N)})$ are monotonically decreasing in $\|\theta\|_{\mathbb{R}^D}$ on the set $\{\theta : \|\theta\|_{\mathbb{R}^D} \geq r\}$.*
- iii) We have that $\Pi(B_s|Z^{(N)}) \xrightarrow{N \rightarrow \infty} 1$ in probability.*
- iv) There exists $\varepsilon > 0$ such that for any (sequence of) Markov kernels \mathcal{P}_N on \mathbb{R}^D and associated chains $(\vartheta_k : k \geq 1)$ that satisfy Assumption 3.1 for some $c_0 > 0$, $L = 1 + \varepsilon$, sequence $\eta_N \in (0, s)$ and all $N \geq 1$ large enough, we can find an initialisation point $\vartheta_0 \in \Theta_{2/3,\varepsilon}$ such that with high probability (under the law of $Z^{(N)}$ and the Markov chain), the hitting time τ_{B_s} for ϑ_k to reach B_s (with s as in *iii*) is lower bounded as*

$$\tau_{B_s} \geq \exp(\min\{c_0, 1\}N/2).$$

The interpretation is that despite the posterior being strictly increasing in the radial variable $\|\theta\|_{\mathbb{R}^D}$ (at least for $\|\theta\|_{\mathbb{R}^D} > r$, any $r > 0$ – note that maximisers of the posterior density may deviate from the ‘ground truth’ $\theta_0 = 0$ by some asymptotically vanishing error, cf. also Proposition 4.1), MCMC algorithms started in $\Theta_{2/3,\varepsilon}$ will still take an exponential time before visiting the region B_s where the posterior mass concentrates. This is true for small enough step-size independently of D, N . The result holds also for ϑ_0 drawn from an absolutely continuous distribution on $\Theta_{2/3,\varepsilon}$ as inspection of the proof shows. Finally, we note that at the expense of more cumbersome notation, the above high probability results (and similarly in Theorem 3.3) could be made non-asymptotic, in the sense that for all $\delta > 0$ all statements hold with probability at least $1 - \delta$ for all $N \geq N_0(\delta)$ large enough, where the dependency of N_0 on δ can be made explicit.

For ‘ellipsoidally supported’ α -regular priors b), the idea is similar but the geometry of the problem changes as the prior now ‘prefers’ low-dimensional subspaces of \mathbb{R}^D , forcing the posterior closer towards the ground truth $\theta_0 = 0$. We show that if the step size is small compared to a scaling N^{-b} for $b > 0$ determined by α , then the same hardness phenomenon persists. Note that ‘small’ is only ‘polynomially small’ in N and hence algorithmic hardness does not come from exponentially small step-sizes.

Theorem 3.3. *Let $D/N \simeq \kappa > 0$, consider the posterior (9) arising from the model (8) and a $\mathcal{N}(0, \Sigma_\alpha)$ prior of density π for some $\alpha > d/2$, and let $\theta_0 = 0$. Define $b = (\alpha/d) - (1/2) > 0$. Then there exists \mathcal{G} and some fixed constant $s_b \in (0, 1/2)$ for which the following statements hold true.*

i) The expected likelihood $\ell(\theta)$ is unimodal with mode 0, locally log-concave near 0, radially symmetric, Lipschitz continuous and monotonically decreasing in $\|\theta\|_{\mathbb{R}^D}$ on \mathbb{R}^D .

ii) For any fixed $r > 0$, with high probability $\ell_N(\theta)$ is radially symmetric and decreasing in $\|\theta\|_{\mathbb{R}^D}$ on the set $\{\theta : \|\theta\|_{\mathbb{R}^D} \geq rN^{-b}\}$.

iii) Defining $s = s_b N^{-b}$, we have $\Pi(B_s | Z^{(N)}) \xrightarrow{N \rightarrow \infty} 1$ in probability.

iv) There exist positive constants $\varepsilon, C > 0$ and $\nu = \nu(\kappa, \alpha, d) > 0$ such that for any (sequence of) Markov kernels \mathcal{P}_N on \mathbb{R}^D and associated chains $(\vartheta_k : k \geq 1)$ that satisfy Assumption 3.1 for some $c_0 > 0$, $L = L_N = C\sqrt{N}$, sequence $\eta = \eta_N \in (0, s_b N^{-b})$ and all $N \geq 1$ large enough, we can find an initialisation point $\vartheta_0 \in \Theta_{N^{-b}, \varepsilon N^{-b}}$ such that with high probability (under the law of $Z^{(N)}$ and the Markov chain), the hitting time τ_{B_s} for ϑ_k to reach B_s is lower bounded as

$$\tau_{B_s} \geq \exp(\min\{c_0, \nu\}N/2).$$

Again, *iv)* holds as well for ϑ_0 drawn from an absolutely continuous distribution on $\Theta_{N^{-b}, (1+\varepsilon)N^{-b}}$. We also note that ε depends only on α, κ, d and the choice of \mathcal{G} but not on any other parameters.

Remark 3.4. As opposed to Theorem 3.2, due to the anisotropy of the prior density π , the posterior distribution is no longer radially symmetric in the preceding theorem, whence part *ii)* differs from Theorem 3.2. But a slightly weaker form of monotonicity of the posterior density $\pi(\cdot | Z^{(N)})$ still holds: the same arguments employed to prove part *ii)* of Theorem 3.2 show that $\pi(\cdot | Z^{(N)})$ is decreasing on $\{\theta : \|\theta\|_{\mathbb{R}^D} \geq rN^{-b}\}$ (any $r > 0$) along the half-lines through 0, i.e.

$$\mathbb{P}_0^N(\pi(v e | Z^{(N)}) \leq \pi(v' e | Z^{(N)}) \text{ for all } v \geq v' \geq r, e \in \mathbb{R}^D, \|e\|_{\mathbb{R}^D} = 1) \xrightarrow{N \rightarrow \infty} 1. \quad (13)$$

We note that this notion precludes the possibility of $\pi(\cdot | Z^{(N)})$ having extremal points outside of the region of dominant posterior mass, and implies that moving toward the origin will *always* increase the posterior density. As a result, many typical Metropolis-Hastings would be encouraged to accept such ‘radially inward’ moves, if they arise as a proposal. Thus, crucially, our exponential hitting time lower bound in part *iv)* arises not through multimodality, but merely through volumetric properties of high-dimensional Gaussian measures.

Remark 3.5 (On the step-size condition). One may wonder whether larger step-sizes can help to overcome the negative result presented in the last theorem. If the step-sizes are ‘time-homogeneous’ and $\gg N^{-b}$ on average, then we may hit the region where the posterior is supported at some time. This would happen ‘by chance’ and not because the data (via ℓ_N) would suggest to move there, and future proposals will likely be outside of that bulk region, so that the chain will either exit the relevant region again or become deterministic because an accept/reject step refuses to move into such directions. In this sense, a negative result for (polynomially) small step sizes gives fundamental limitations on the ability of the chain to explore the precise characteristics of the posterior distribution. We also remark that the Lipschitz-constants of $\nabla \ell(\theta)$ are of order D or D^{1+b} in the preceding theorems, respectively. A Markov chain obtained from discretising a continuous diffusion process (such as MALA discussed in the next subsection) will generally require step-sizes that are inversely proportional to that Lipschitz constant in order to inherit the dynamics from the continuous process. For such examples, Assumption 3.1 is natural. But as discussed at the end of the introduction, there exists a variety of ‘non-local’ MCMC algorithms for which this step size assumption may not be satisfied.

3.2 Implications for common MCMC methods with ‘cold-start’

The preceding general hitting time bounds apply to commonly used MCMC methods in high-dimensional statistics. We focus in particular on algorithms that are popular with PDE models and inverse problems, see, e.g., [Cot+13; Bes+17] and also [Nic22] for many more references. We illustrate this for two natural examples with Metropolis-Hastings adjusted random walk and gradient algorithms. Other examples can be generated without difficulty.

3.2.1 Preconditioned Crank-Nicolson

We first give some hardness results for the popular preconditioned Crank-Nicolson (pCN) algorithm. A dimension-free convergence analysis for pCN was given in the important paper by Hairer, Stuart, and Vollmer [HSV14] based on ideas from [HMS11]. The results in the present section show that while the mixing bounds from [HSV14] are in principle uniform in D , the implicit dependence of the constants

on the conditions on the log-likelihood-function in [HSV14] can re-introduce exponential scaling when one wants to apply the results from [HSV14] to concrete (N -dependent) posterior distributions. This confirms a conjecture about pCN made in Section 1.2.1 of [NW20].

Let \mathcal{C} denote the covariance of some Gaussian prior on \mathbb{R}^D with density π . Then the pCN algorithm for sampling from some posterior density $\pi(\theta|Z^{(N)}) \propto e^{\ell_N(\theta)}\pi(\theta)$ is given as follows. Let $(\xi_k : k \geq 1)$ be an i.i.d. sequence of $\mathcal{N}(0, \mathcal{C})$ random vectors. For initialiser $\vartheta_0 \in \mathbb{R}^D$, step size $\beta > 0$ and $k \geq 1$, the MCMC chain is then given by

1. PROPOSAL: $p_k \sim \sqrt{1 - \beta}\vartheta_{k-1} + \sqrt{\beta}\xi_k$,
2. ACCEPT-REJECT: Set

$$\vartheta_k = \begin{cases} p_k & \text{w.p. } \min\{1, e^{\ell_N(p_k) - \ell_N(\vartheta_{k-1})}\}, \\ \vartheta_{k-1} & \text{else.} \end{cases} \quad (14)$$

By standard Markov chain arguments one verifies (see [HSV14] or Ch.1 in [Nic22]) that the (unique) invariant density of $(\vartheta_k : k \geq 1)$ equals $\pi(\cdot|Z^{(N)})$.

We now give a hitting time lower bound for the pCN algorithm which holds true in the regression setting for which the main Theorems 3.2 and 3.3 (for generic Markov chains) were derived. In particular, we emphasize that the lower bounds to follow hold for the choice of regression ‘forward’ map \mathcal{G} constructed in the proofs of Theorems 3.2 and 3.3. As for the general results, we treat the two cases of $\mathcal{C} = I_D/D$ or $\mathcal{C} = \Sigma_\alpha$ separately.

Theorem 3.6. *Let ϑ_k denote the pCN Markov chain from (14).*

i) Assume the setting of Theorem 3.2 with $\mathcal{C} = I_D/D$, and let \mathcal{G} be as in Theorem 3.2. Then there exist constants $c_1, c_2, \varepsilon > 0$ such that for any $\beta \leq c_1$, there is an initialisation point $\vartheta_0 \in \Theta_{1, \varepsilon}$ such that the hitting time $\tau_{B_s} = \inf\{k : \vartheta_k \in B_s\}$ (for B_s as in (11)) satisfies with high probability (under the law of the data and of the Markov chain) as $N \rightarrow \infty$ that $\tau_{B_s} \geq \exp(c_2 D)$.

ii) Assume the setting of Theorem 3.3 with $\mathcal{C} = \Sigma_\alpha$ for $\alpha > d/2$, and let \mathcal{G} be as in Theorem 3.2. Then there exist constants $c_1, c_2, \varepsilon > 0$ such that if $\beta \leq c_1 N^{-1-2b}$ there is an initialisation point $\vartheta_0 \in \Theta_{N^{-b}, \varepsilon N^{-b}}$ such that the hitting time $\tau_{B_s} = \inf\{k : \vartheta_k \in B_s\}$ satisfies with high probability that $\tau_{B_s} \geq \exp(c_2 D)$.

3.2.2 Gradient-based Langevin algorithms

We now turn to *gradient-based* Langevin algorithms which are based on the discretization of continuous-time diffusion processes [Cot+13; Dal17]. A polynomial time convergence analysis for the *unadjusted* Langevin algorithm in the strongly log-concave case has been given in [Dal17; DM19] and also in [Che+21] for the Metropolis-adjusted case (MALA). We show here that for unimodal but not globally log-concave distributions, the MCMC scheme can take an exponential time to reach the bulk of the posterior distribution. For simplicity we focus on the Metropolis-adjusted Langevin algorithm which is defined as follows. Let $(\xi_k : k \geq 1)$ be a sequence of i.i.d. $\mathcal{N}(0, I_D)$ variables, and let $\gamma > 0$ be a step-size.

1. PROPOSAL: $p_k = \vartheta_{k-1} + \gamma \nabla \log \pi(\vartheta_{k-1}|Z^{(N)}) + \sqrt{2\gamma}\xi_k$.
2. ACCEPT-REJECT: Set

$$\vartheta_k = \begin{cases} p_k & \text{w.p. } \min\left\{1, \frac{\pi(p_k|Z^{(N)}) \exp(-\|\vartheta_{k-1} - p_k - \gamma \nabla \log \pi(p_k|Z^{(N)})\|^2)}{\pi(\vartheta_{k-1}|Z^{(N)}) \exp(-\|\vartheta_{k-1} - \vartheta_{k-1} - \gamma \nabla \log \pi(\vartheta_{k-1}|Z^{(N)})\|^2)}\right\}, \\ \vartheta_{k-1} & \text{else.} \end{cases} \quad (15)$$

Again, standard Markov chain arguments show that $\Pi(\cdot|Z^{(N)})$ is indeed the (unique) invariant distribution of $(\vartheta_k : k \geq 1)$. We note here that for the forward \mathcal{G} featuring in our results to follow, $\nabla \log \pi$ may only be well-defined (Lebesgue-) almost everywhere on \mathbb{R}^D due to our piecewise smooth choice of w , see (21) below. However, since all proposal densities involved possess a Lebesgue density, this specification almost everywhere suffices in order to propagate the Markov chain with probability 1. Alternatively one could also straightforwardly avoid this technicality by smoothing our choice of function w in (21), which we refrain from for notational ease.

Theorem 3.7. *Let ϑ_k denote the MALA Markov chain from (15).*

i) Assume the setting of Theorem 3.2, with $\mathcal{N}(0, I_D/D)$ prior, and let \mathcal{G} also be as in Theorem 3.2. There exists some $c_1, c_2, \varepsilon > 0$ such that if the step size of $(\vartheta_k : k \geq 1)$ satisfies $\gamma \leq c_1/N$, then there is an initialisation point $\vartheta_0 \in \Theta_{1, \varepsilon}$ such that the hitting time $\tau_{B_s} = \inf\{k : \vartheta_k \in B_s\}$ (for B_s as in (11)) satisfies with high probability (under the law of the data and of the Markov chain) as $N \rightarrow \infty$ that $\tau_{B_s} \geq \exp(c_2 D)$.

ii) Assume the setting of Theorem 3.3, with a $\mathcal{N}(0, \Sigma_\alpha)$ prior, and let \mathcal{G} also be as in Theorem 3.3. Then there exist some constant $c_1, c_2, \varepsilon > 0$ such that whenever $\gamma \leq c_1 N^{-1-b-2\alpha}$, there is an initialisation point $\vartheta_0 \in \Theta_{N^{-b}, \varepsilon N^{-b}}$, such that the hitting time $\tau_{B_s} = \inf\{k : \vartheta_k \in B_s\}$ satisfies with high probability (under the law of the data and of the Markov chain) that $\tau_{B_s} \geq \exp(c_2 D)$.

As mentioned in Remark 3.5, a bound on the step-size that is inversely proportional to the Lipschitz constant of $\nabla \ell$ is natural for algorithms like MALA that arise from discretisation of a continuous time Markov process, see, e.g., [DM19; Che+21]. We emphasise again that these Lipschitz constants are D - and N -dependent, so that the required bounds on γ are not unnatural. ‘Optimal’ step-size prescriptions for MALA [RR01; BPS04; MPS12; Che+21] derived for Gaussian and log-concave targets or, more generally, mean field limits (in which the posterior distribution possesses a product or mean-field structure, unlike in the models considered here) would need to be adjusted to our model classes to be comparable.

4 Proofs of the main theorems

We begin in Section 4.1 by constructing the family of regression maps \mathcal{G} underlying our results from Section 3. Sections 4.2 and 4.3 reduce the hitting time bounds from Theorems 3.2 and 3.3 (for general Markov chains) to hitting time bounds for intermediate ‘free energy barriers’ that the Markov chain needs to travel through. Subsequently, Theorems 3.3 and 3.2 are proved in Sections 4.4 and 4.5 respectively. Finally, the proofs for pCN (Theorem 3.6) and MALA (Theorem 3.7) are contained in Section 4.6.

4.1 Radially symmetric choices of \mathcal{G}

We start with our parameterisation of the map \mathcal{G} . In our regression model and since $\mathbb{E}\varepsilon^2 = 1$,

$$\ell(\theta) = -\frac{N}{2} \mathbb{E}_{\theta_0}^1 |Y - \mathcal{G}(\theta)(X)|^2 = -\frac{N}{2} \|\mathcal{G}(\theta_0) - \mathcal{G}(\theta)\|_{L^2}^2 - \frac{N}{2}, \quad \theta \in \mathbb{R}^D. \quad (16)$$

We have $\theta_0 = 0$ and by subtracting a fixed function $\mathcal{G}(0)$ from $\mathcal{G}(\theta)$ if necessary we can also assume that $\mathcal{G}(\theta_0) = 0$. In this case, since $\text{vol}(\mathcal{X}) = 1$,

$$\ell(\theta) = -\frac{N}{2} \|\mathcal{G}(\theta)\|_{L^2}^2 - \frac{N}{2}, \quad (17)$$

Take a bounded continuous function $w : [0, \infty] \rightarrow [0, \|w\|_\infty]$ with a unique minimiser $w(0) = 0$ and take \mathcal{G} of the ‘radial’ form

$$\mathcal{G}(\theta) = \sqrt{w(\|\theta\|_{\mathbb{R}^D})} \times g(x), \quad \theta \in \mathbb{R}^D, x \in \mathcal{X},$$

where

$$g : \mathcal{X} \rightarrow [g_{\min}, g_{\max}], \quad 0 < g_{\min} < g_{\max} < \infty, \quad \|g\|_{L^2_\mu(\mathcal{X})} = 1.$$

The assumption $\mathcal{G}(\theta_0) = 0$ implies $Y_i = 0 + \varepsilon_i$ under $P_{\theta_0}^N$, so that we have

$$\begin{aligned} \ell_N(\theta) &= -\frac{1}{2} \sum_{i=1}^N |\varepsilon_i - \sqrt{w(\|\theta\|)} g(X_i)|^2 \\ &= -\frac{w(\|\theta\|_{\mathbb{R}^D})}{2} \sum_{i=1}^N g^2(X_i) - \frac{1}{2} \sum_{i=1}^N \varepsilon_i^2 + \sqrt{w(\|\theta\|)} \sum_{i=1}^N \varepsilon_i g(X_i) \end{aligned} \quad (18)$$

and the average log-likelihood is

$$\ell(\theta) = \mathbb{E}_{\theta_0}^N \ell_N(\theta) = -\frac{N}{2} w(\|\theta\|_{\mathbb{R}^D}) - \frac{N}{2}, \quad \theta \in \mathbb{R}^D. \quad (19)$$

Define ϵ -annuli of Euclidean space

$$\Theta_{r,\epsilon} = \{\theta \in \mathbb{R}^D : \|\theta\|_{\mathbb{R}^D} \in (r, r + \epsilon)\}, \quad r \geq 0. \quad (20)$$

We then also set, for any $s \geq 0$, $\epsilon > 0$,

$$w_-(r, \epsilon) = \inf_{s \in (r, r + \epsilon)} w(s), \quad w_+(r, \epsilon) = \sup_{s \in (r, r + \epsilon)} w(s).$$

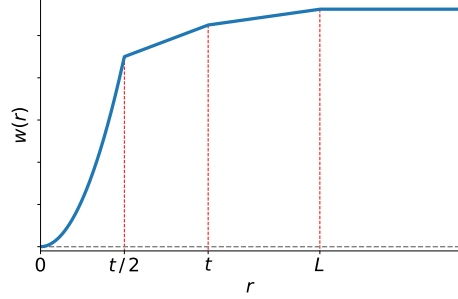
For our main theorems the map w will be monotone increasing and the preceding notation w_-, w_+ is then not necessary, but Proposition 4.2 is potentially also useful in non-monotone settings (as remarked after its proof), hence the slightly more general notation here.

The choice that \mathcal{G} is radial is convenient in the proofs, but means that the model is only identifiable up to a rotation for $\theta \neq 0$. One could easily make it identifiable by more intricate choices of \mathcal{G} , but the main point for our negative results is that the function ℓ has a unique mode at the ground truth parameter θ_0 and is identifiable there.

4.1.1 A locally log-concave, globally monotone choice of w

Define for $t < L$ and any $r > 0$ the function $w : [0, \infty) \rightarrow \mathbb{R}$ as

$$\begin{aligned} w(r) = & 4(Tr)^2 1_{[0, t/2)} \\ & + [(Tt)^2 + T(r - t/2)] 1_{[t/2, t)}(r) \\ & + [(Tt)^2 + (Tt/2) + \rho(r - t)] 1_{[t, L)}(r) \\ & + [(Tt)^2 + (Tt/2) + \rho(L - t)] 1_{[L, \infty)}(r), \end{aligned} \quad (21)$$



where $T > \rho$, are fixed constants to be chosen. Note that w is monotone increasing and

$$\|w\|_\infty = (Tt)^2 + (Tt/2) + \rho(L - t) < \infty. \quad (22)$$

The function w is quadratic near its minimum at the origin up until $t/2$, from when onwards it is piece-wise linear. In the linear regime it initially has a ‘steep’ ascent of gradient T until t , then grows more slowly with small gradient ρ from t until L , and from then on is constant. The function w is not C^∞ at the points $r = t/2, r = t, r = L$, but we can easily make it smooth by convolving with a smooth function supported in small neighbourhoods of its breakpoints r without changing the findings that follow. We abstain from this to simplify notation.

The following proposition summarises some monotonicity properties of the empirical log-likelihood function arising from the above choice of w .

Proposition 4.1. *Let w be as in (21). Then there exists $C > 0$ such that for any $r_0 > 0$ and $N \geq 1$, we have*

$$P_0^N \left(\sup_{r_0 \leq r < s \leq L} \sup_{\|\theta_s\|=s, \|\theta_r\|=r} \frac{\ell_N(\theta_s) - \ell_N(\theta_r)}{w(s) - w(r)} \leq -\frac{N}{4} \right) \geq 1 - \frac{C}{N} - \frac{C}{Nw(r_0)}.$$

In particular if $r_0 < t/2$ is such that $(Tr_0)^2 N \rightarrow \infty$ as $N \rightarrow \infty$ then the r.h.s. is $1 - o(1)$.

Recalling (19), (18) and since w is monotonically increasing, we bound

$$\begin{aligned} & P_0^N(\ell_N(\theta_s) - \ell_N(\theta_r) > (N/4)(w(r) - w(s))) \\ & = P_0^N(\ell_N(\theta_s) - \ell_N(\theta_r) - (\ell(\theta_s) - \ell(\theta_r)) > -\frac{N}{4}(w(r) - w(s))) = \\ & \Pr \left(\frac{(w(r) - w(s))}{2} \sum_{i=1}^N (g^2(X_i) - 1) + (\sqrt{w(s)} - \sqrt{w(r)}) \sum_{i=1}^N \varepsilon_i g(X_i) > \frac{N}{4}(w(s) - w(r)) \right) \\ & = \Pr \left(-\sum_{i=1}^N (g^2(X_i) - \mathbb{E}g^2(X))/2 + \frac{1}{\sqrt{w(s)} + \sqrt{w(r)}} \sum_{i=1}^N \varepsilon_i g(X_i) > \frac{N}{4} \right) \\ & \leq \Pr \left(\left| \sum_{i=1}^N (g^2(X_i) - \mathbb{E}g^2(X)) \right| > N/4 \right) + \Pr \left(\left| \sum_{i=1}^N \varepsilon_i g(X_i) \right| > \frac{2N\sqrt{w(r_0)}}{8} \right) \\ & = \mathcal{O}(1/N) + \mathcal{O}(1/(Nw(r_0))), \end{aligned}$$

using Chebyshev's inequality in the last step. Since the events in the penultimate step do not depend on $r < s \in [r_0, L]$, the result follows. \square

4.2 Bounds for posterior ratios of annuli

A key quantity in the proofs to follow will be to obtain asymptotic ($N \rightarrow \infty$) bounds of the following functional (recalling the definition of the Euclidean annuli $\Theta_{r,\varepsilon}$ from (20)),

$$\mathcal{F}_N(r, \varepsilon) = \frac{1}{N} \log \int_{\Theta_{r,\varepsilon}} e^{\ell_N(\theta)} d\Pi(\theta), \quad r \geq 0, \quad \varepsilon > 0, \quad (23)$$

in terms of the map w . As a side note, we remark that this functional has a long history in the statistical physics of glasses, in which it is often referred to as the *Franz-Parisi* potential [FP95; Ban+22].

Proposition 4.2. *Consider the regression model (8) with radially symmetric choice of \mathcal{G} from Subsection 4.1 such that $\|w\|_\infty \leq W$ for some fixed $W < \infty$ (independent of D, N), and let $\Pi = \Pi_N$ denote a sequence of prior probability measures on \mathbb{R}^D .*

i) *Suppose that for some radii $0 < s < \sigma$, constants $\varepsilon, \eta, \nu > 0$ and for all $N \geq 1$ large enough, we have*

$$\frac{1}{N} \log \frac{\Pi(\Theta_{s,\eta})}{\Pi(\Theta_{\sigma,\varepsilon})} \leq -2\nu - \frac{(w_+(\sigma, \varepsilon) - w_-(s, \eta))}{2}. \quad (24)$$

Then the posterior distribution $\Pi(\cdot|Z^{(N)})$ from (9) arising in the model (8) satisfies that with high P_0^N -probability as $N \rightarrow \infty$,

$$\frac{\Pi(\Theta_{s,\eta}|Z^{(N)})}{\Pi(\Theta_{\sigma,\varepsilon}|Z^{(N)})} \leq e^{-\nu N}. \quad (25)$$

ii) *If in addition w is monotone increasing on $[0, \infty)$ and if for some $L > 1 + \varepsilon$,*

$$\frac{1}{N} \log \frac{\Pi(B_L^c)}{\Pi(\Theta_{\sigma,\varepsilon})} \leq -2\nu, \quad (26)$$

then the posterior distribution $\Pi(\cdot|Z^{(N)})$ also satisfies (with high probability as $N \rightarrow \infty$) that

$$\frac{\Pi(B_L^c|Z^{(N)})}{\Pi(\Theta_{\sigma,\varepsilon}|Z^{(N)})} \leq e^{-\nu N}. \quad (27)$$

Remark 4.3 (The prior condition for w from (21)). If $\sigma > s > t$, for w from (21), the ‘likelihood’ term in Proposition 4.2 is

$$\frac{w_+(\sigma, \varepsilon)}{2} - \frac{w_-(s, \eta)}{2} = \frac{\rho(\sigma + \varepsilon - t) - \rho(s - t)}{2} = \frac{\rho}{2}(\sigma + \varepsilon - s) > 0, \quad (28)$$

so that if we also assume

$$Tt + \rho L = \mathfrak{c}(\sqrt{N}) \quad (29)$$

to control ω_N, ω'_N in the proof that follows, then to verify (24) it suffices to check

$$\frac{1}{N} \log \frac{\Pi(\Theta_{s,\eta})}{\Pi(\Theta_{\sigma,\varepsilon})} \leq -2\nu - \frac{\rho}{2}(\sigma + \varepsilon - s) \quad (30)$$

for all large enough N .

Proof of part i). From the definition of ℓ_N in (18) we first notice that for all $r \geq 0, \varepsilon > 0$,

$$\inf_{\theta \in \Theta_{r,\varepsilon}} \ell_N(\theta) \geq -\frac{1}{2} \sum_{i=1}^N \varepsilon_i^2 - \frac{1}{2} w_+(r, \varepsilon) \sum_{i=1}^N g^2(X_i) - \sqrt{w_+(r, \varepsilon)} \left| \sum_{i=1}^N \varepsilon_i g(X_i) \right|,$$

and

$$\sup_{\theta \in \Theta_{r,\varepsilon}} \ell_N(\theta) \leq -\frac{1}{2} \sum_{i=1}^N \varepsilon_i^2 - \frac{1}{2} w_-(r, \varepsilon) \sum_{i=1}^N g^2(X_i) + \sqrt{w_-(r, \varepsilon)} \left| \sum_{i=1}^N \varepsilon_i g(X_i) \right|,$$

We can now further bound, for our \mathcal{G} ,

$$\frac{1}{N} \log \int_{\Theta_{r,\epsilon}} e^{\ell_N(\theta)} d\Pi(\theta) \leq -\frac{1}{2N} \sum_{i=1}^N \varepsilon_i^2 - \frac{w_-(r,\epsilon)}{2N} \sum_{i=1}^N g^2(X_i) + \frac{\sqrt{w_+(r,\epsilon)}}{N} \left| \sum_{i=1}^N \varepsilon_i g(X_i) \right| + \frac{\log \Pi(\Theta_{r,\epsilon})}{N}$$

and

$$\frac{1}{N} \log \int_{\Theta_{r,\epsilon}} e^{\ell_N(\theta)} d\Pi(\theta) \geq -\frac{1}{2N} \sum_{i=1}^N \varepsilon_i^2 - \frac{w_+(r,\epsilon)}{2N} \sum_{i=1}^N g^2(X_i) - \frac{\sqrt{w_+(r,\epsilon)}}{N} \left| \sum_{i=1}^N \varepsilon_i g(X_i) \right| + \frac{\log \Pi(\Theta_{r,\epsilon})}{N}$$

We estimate $\sqrt{w_+(r,\epsilon)} \leq \bar{w}(r,\epsilon) = \max(w_+(r,\epsilon), 1)$, and noting that

$$\mathbb{E} \varepsilon_i^2 = 1 = \mathbb{E} g^2(X_i), \quad \mathbb{E} \varepsilon_i g(X_i) = 0$$

we can use Chebyshev's (or Bernstein's) inequality to construct an event of high probability such that the functional \mathcal{F}_N from (23) is bounded as

$$\mathcal{F}_N(r,\epsilon) \leq -\frac{1}{2} - \frac{w_-(r,\epsilon)}{2} + \frac{\log \Pi(\Theta_{r,\epsilon})}{N} + \omega_N(s,\eta) \quad (31)$$

and

$$\mathcal{F}_N(r,\epsilon) \geq -\frac{1}{2} - \frac{w_+(r,\epsilon)}{2} + \frac{\log \Pi(\Theta_{r,\epsilon})}{N} + \omega'_N(r,\epsilon), \quad (32)$$

where

$$\omega_N(r,\epsilon) = \mathcal{O}\left(1 + \frac{w_-(r,\epsilon) + \bar{w}(r,\epsilon)}{\sqrt{N}}\right), \quad \omega'_N(s) = \mathcal{O}\left(1 + \frac{w_+(r,\epsilon) + \bar{w}(r,\epsilon)}{\sqrt{N}}\right), \quad (33)$$

and this is uniform in all (r,ϵ) since $\|w\|_\infty \leq W$ is bounded. Using the above with (r,ϵ) chosen as (s,η) and (σ,ε) respectively, we then obtain

$$\begin{aligned} \frac{1}{N} \log \frac{\Pi(\Theta_{s,\eta} | Z^{(N)})}{\Pi(\Theta_{\sigma,\varepsilon} | Z^{(N)})} &= \mathcal{F}_N(s,\eta) - \mathcal{F}_N(\sigma,\varepsilon) \\ &\leq -\frac{w_-(s,\eta)}{2} + \frac{\log \Pi(\Theta_{s,\eta})}{N} + \frac{w_+(\sigma,\varepsilon)}{2} - \frac{\log \Pi(\Theta_{\sigma,\varepsilon})}{N} + \omega_N(s,\eta) - \omega'_N(\sigma,\varepsilon) \\ &= -\frac{w_-(s,\eta)}{2} + \frac{w_+(\sigma,\varepsilon)}{2} + \frac{1}{N} \log \frac{\Pi(\Theta_{s,\eta})}{\Pi(\Theta_{\sigma,\varepsilon})} + \omega_N(s,\eta) - \omega'_N(\sigma,\varepsilon) \end{aligned} \quad (34)$$

with high $P_{\theta_0}^N$ -probability. The result now follows from the hypothesis (24) and since the terms ω_N, ω'_N are $\mathcal{O}(1)$.

Proof of part ii). The proof of part ii) follows from an obvious modification of the previous arguments. □

In the case where $\Pi(\Theta_{s,\eta})$ and $\Pi(\Theta_{\sigma,\varepsilon})$ are comparable (so that the l.h.s. in (24) converges to zero), a local optimum at σ in the function w away from zero can verify the last inequality for 'intermediate' s such that $w(s) - w(\sigma) \leq -2\nu$. This can be used to give computational hardness results for MCMC of multi-modal distributions. But we are interested in the more challenging case of 'unimodal' examples w from (21). Before we turn to this, let us point out what can be said about the hitting times of Markov chains if the conclusion (25) of Proposition 4.2 holds.

4.3 Bounds for Markov chain hitting times

4.3.1 Hitting time bounds for intermediate sets $\Theta_{s,\eta}$

In (25) we can think of $\Theta_{\sigma,\varepsilon}$ as the 'initialisation region' (further away from θ_0) and $\Theta_{s,\eta}$ for intermediate s is the 'barrier' before we get close to $\theta_0 = 0$. The last bound permits the following classic hitting time argument, taken from [BWZ20], see also [Jer03].

Proposition 4.4. *Consider any Markov chain $(\vartheta_k : k \in \mathbb{N})$ with invariant measure $\mu = \Pi(\cdot|Z^{(N)})$ for which (25) holds. For constants $\eta < \sigma - s$, suppose ϑ_0 is started in $\Theta_{\sigma,\varepsilon}$, $\mu(\Theta_{\sigma,\varepsilon}) > 0$, drawn from the conditional distribution $\mu(\cdot|\Theta_{\sigma,\varepsilon})$, and denote by τ_s the hitting time of the Markov chain onto $\Theta_{s,\eta}$, that is, the number τ_s of iterates required until ϑ_k visits the set $\Theta_{s,\eta}$. Then*

$$\Pr(\tau_s \leq K) \leq Ke^{-\nu N}, \quad K > 0.$$

Similarly, on the event where (27) holds we have that

$$\Pr(\tau_{B_L^c} \leq K) \leq Ke^{-\nu N}, \quad K > 0.$$

Proof of Proposition 4.4 – We have

$$\begin{aligned} \Pr(\tau_s \leq K) &= \Pr(\vartheta_k \in \Theta_{s,\eta} \text{ for some } 1 \leq k \leq K | \vartheta_0 \in \Theta_{\sigma,\varepsilon}) \\ &= \frac{\Pr(\vartheta_0 \in \Theta_{\sigma,\varepsilon}, \vartheta_k \in \Theta_{s,\eta} \text{ for some } 1 \leq k \leq K)}{\mu(\Theta_{\sigma,\varepsilon})} \leq \frac{\sum_{k \leq K} \Pr(\vartheta_k \in \Theta_{s,\eta})}{\mu(\Theta_{\sigma,\varepsilon})} \\ &\leq K \frac{\mu(\Theta_{s,\eta})}{\mu(\Theta_{\sigma,\varepsilon})} \leq Ke^{-\nu N}. \end{aligned}$$

The second claim is proved analogously. \square

The last proposition holds ‘on average’ for initialisers $\vartheta_0 \sim \mu(\cdot|\Theta_{\sigma,\varepsilon})$, and since $\Pr = \mathbb{E}_{\mu(\cdot|\Theta_{\sigma,\varepsilon})} \Pr_{\vartheta_0}$ where \Pr_{ϑ_0} is the law of the Markov chain started at ϑ_0 , the hitting time inequality holds at least for one point in $\Theta_{\sigma,\varepsilon}$ since $\inf_{\vartheta_0} \Pr_{\vartheta_0} \leq \mathbb{E}_{\mu(\cdot|\Theta_{\sigma,\varepsilon})} \Pr_{\vartheta_0}$.

4.3.2 Reducing hitting times for B_s to ones for $\Theta_{s,\eta}$

We now reduce part **iv)** of Theorems 3.2 and 3.3, i.e. bounds on the hitting time of the region B_s in which the posterior contracts, to a bound for the hitting time τ_s for the annulus $\Theta_{s,\eta}$, which is controlled in Proposition 4.4. To this end, in the case of Theorem 3.2, we suppose that Propositions 4.2 and 4.4 are verified with $\nu = \sigma = 1$, some $\varepsilon > 0$ and L, s, η as in the theorem, and in the case of Theorem 3.3, we assume the same with choice $\sigma = N^{-b}$ and $\nu > 0$ given after (42) below. For c_0 from Assumption 3.1, define the events

$$A_N := \{\forall k \leq e^{(\nu \wedge c_0)N/2} : \|\vartheta_{k+1} - \vartheta_k\|_{\mathbb{R}^D} \leq \eta/2\}.$$

We can then estimate, using Assumption 3.1, that on the frequentist event on which Proposition 4.4 holds (which we apply with $K = e^{(\nu \wedge c_0)N/2} \leq e^{\nu N/2}$), under the probability law of the Markov chain we have

$$\begin{aligned} \Pr(\tau_{B_s} \leq e^{(\nu \wedge c_0)N/2}) &\leq \Pr(\tau_{B_s} \leq e^{(\nu \wedge c_0)N/2}, A_N) + \Pr(A_N^c) \\ &\leq \Pr(\tau_s \leq e^{(\nu \wedge c_0)N/2}) + \Pr(A_N^c, \tau_{B_L^c} > e^{(\nu \wedge c_0)N/2}) + \Pr(\tau_{B_L^c} \leq e^{(\nu \wedge c_0)N/2}) \\ &\leq 2e^{-\nu N/2} + e^{(\nu \wedge c_0)N/2} \sup_{\theta \in B_L} \mathcal{P}_N(\theta, \{\vartheta : \|\theta - \vartheta\|_{\mathbb{R}^D} \geq \eta/2\}) \\ &\leq 2e^{-(\nu \wedge c_0)N/2} + e^{(\nu \wedge c_0)N/2 - c_0 N} \leq 3e^{-(\nu \wedge c_0)N/2}, \end{aligned}$$

where in the second inequality we have used that on the events A_N , the Markov chain ϑ_k , when started in $\Theta_{1,\varepsilon}$, needs to pass through $\Theta_{s,\eta}$ in order to reach B_s .

4.4 Proof of Theorem 3.3

In this section, we use the results derived in the previous part of Section 4 to finish the proof of Theorem 3.3. Parts **i)** and **ii)** of the theorem follow from Proposition 4.1 and our choice of w in (21). We therefore concentrate on the proofs of part **iii)** and **iv)**. We start with proving a key lemma on small-ball estimates for truncated α -regular Gaussian priors.

4.4.1 Small ball estimates for α -regular priors

Let us first define precisely the notion of α -regular Gaussian priors. For some fixed $\alpha > d/2$, the prior Π arises as the truncated law $Law(\theta)$ of an α -regular Gaussian process with RKHS $\mathcal{H} = H^\alpha$, a Sobolev space over some bounded domain/manifold \mathcal{X} , see e.g., Sec. 6.2.1 in [Nic22] for details.

Equivalently (under the Parseval isometry) we take a Gaussian Borel measure on the usual sequence space $\ell_2 \simeq L^2$ with RKHS equal to

$$h^\alpha = \left\{ (\theta_i)_{i=1}^\infty : \sum_{i=1}^\infty i^{2(\alpha/d)} \theta_i^2 = \|\theta\|_{H^\alpha}^2 < \infty \right\}, \quad \alpha > d/2.$$

The prior Π is the truncated law of $\theta_D = (\theta_1, \dots, \theta_D)$, $D \in \mathbb{N}$.

Lemma 4.5. *Fix $z > 0$, $\alpha > d/2$ and $\kappa > 0$, and set*

$$b = \frac{\alpha}{d} - \frac{1}{2}, \quad \tau = \frac{1}{b} = \frac{2d}{2\alpha - d}.$$

Then if $D/N \simeq \kappa > 0$, there exist constants $\bar{c}_0 > c_0$ (depending on b, κ) such that for all N ($\geq N_0(z, b)$) large enough:

$$c_0(z + \kappa^{-\alpha/d} z^{-\tau/2})^{-\tau} \leq -\frac{1}{N} \log \Pi(\|\theta\|_{\mathbb{R}^D} \leq zN^{-b}) \leq \bar{c}_0 z^{-\tau}. \quad (35)$$

Proof of Lemma 4.5 – Note first that the L^2 -covering numbers of the ball $h(\alpha, B)$ of radius B in H^α satisfy the well-known two-sided estimate

$$\log \mathcal{N}(\delta, \|\cdot\|_{L^2}, h(\alpha, B)) \simeq \left(\frac{AB}{\delta} \right)^{d/\alpha}, \quad 0 < \delta < AB \quad (36)$$

for equivalence constants in \simeq depending only on d, α . The upper bound is given in Proposition 6.1.1 in [Nic22] and a lower bound can be found as well in the literature [ET96] (by injecting $H^\alpha(\mathcal{X}_0)$ into $\tilde{H}^\alpha(\mathcal{X})$ for some strict sub-domain $\mathcal{X}_0 \subset \mathcal{X}$, and using metric entropy lower bounds for the injection $H^\alpha(\mathcal{X}_0) \hookrightarrow L^2(\mathcal{X}_0)$).

Using the results about small deviation asymptotics for Gaussian measures in Banach space [LL99] – specifically Theorem 6.2.1 in [Nic22] with $a = \frac{2d}{2\alpha - d}$ – and assuming $\alpha > d/2$, this means that the concentration function of the ‘untruncated prior’ satisfies the two-sided estimate

$$-\log \Pi(\|\theta\|_{L^2} \leq \gamma) \simeq \gamma^{-\frac{2d}{2\alpha - d}} = \gamma^{-\tau}, \quad \gamma \rightarrow 0. \quad (37)$$

Here, restricting to $\gamma \in (0, 1)$, the two-sided equivalence constants depend only on α, d . Setting

$$\gamma = zN^{-b}, \quad z > 0, \quad (38)$$

and noting that $b\tau = 1$, we hence obtain that for some constants $c_l, c_u > 0$,

$$e^{-c_l z^{-\tau} N} \leq \Pi(\|\theta\|_{L^2} \leq zN^{-b}) \leq e^{-c_u z^{-\tau} N}, \quad \text{any } z > 0. \quad (39)$$

We now show that as long as $D/N \approx \kappa > 0$, one may use the above asymptotics to derive the desired small ball probabilities for the projected prior on \mathbb{R}^D .

We obviously have, by set inclusion and projection,

$$\Pi(\|\theta\|_{\mathbb{R}^D} \leq zN^{-b}) \geq \Pi(\|\theta\|_{L^2} \leq zN^{-b})$$

and hence it only remains to show the first inequality in eq. (35). The Gaussian isoperimetric theorem (Theorem 2.6.12 in [GN16]) and (39) imply that for $m \geq 4\sqrt{c_l}$ and some $c > 0$, we have that (with Φ denoting the c.d.f. for $\mathcal{N}(0, 1)$)

$$\begin{aligned} \Pi(\theta = \theta_1 + \theta_2, \|\theta_1\|_{L^2} \leq zN^{-b}, \|\theta_2\|_{h^\alpha} \leq mz^{-\tau/2}\sqrt{N}) \\ \geq \Phi(\Phi^{-1}(\Pi(\{\theta : \|\theta\| \leq zN^{-b}\})) + mz^{-\tau/2}\sqrt{N}) \\ \geq \Phi(-\sqrt{2c_l}z^{-\tau/2}\sqrt{N} + mz^{-\tau/2}\sqrt{N}) \geq 1 - e^{-cz^{-\tau}N} \end{aligned}$$

(see also the proof of Lemma 5.17 in [MNP21a] for a similar calculation). Then if the event in the last probability is denoted by I we have

$$\Pi(\|\theta_D\|_{\mathbb{R}^D} \leq zN^{-b}) \leq \Pi(\|\theta_D\|_{\mathbb{R}^D} \leq zN^{-b}, I) + e^{-cz^{-\tau}N}.$$

On I , if $D/N \rightarrow \kappa > 0$ and by the usual tail estimate for vectors in h^α , we have for some $c' > 0$ the bound

$$\|\theta - \theta_D\|_{L^2} \leq \|\theta_1\|_{L^2} + c'D^{-\alpha/d}z^{-\tau/2}\sqrt{N} \leq zN^{-b} + c'\kappa^{-\alpha/d}z^{-\tau/2}N^{-b}$$

so that for any $z > 0$,

$$\begin{aligned}\Pi(\|\theta_D\|_{\mathbb{R}^D} \leq zN^{-b}) &\leq \Pi(\|\theta\|_{L^2} \leq zN^{-b} + \|\theta - \theta_D\|_{L^2}, I) + e^{-cz^{-\tau}N} \\ &\leq \Pi(\|\theta\|_{L^2} \leq (2z + c'\kappa^{-\alpha/d}z^{-\tau/2})N^{-b}) + e^{-cz^{-\tau}N} \\ &\leq e^{-c_u(2z + c'\kappa^{-\alpha/d}z^{-\tau/2})^{-\tau}N} + e^{-cz^{-\tau}N},\end{aligned}$$

and hence the lemma follows by appropriately choosing $c_0 > 0$. \square

Remark 4.6. For statistical consistency proofs in non-linear inverse problems, often *rescaled* Gaussian priors are used to provide additional regularisation [MNP21a; NW20; BN21]. For these priors a computation analogous to the previous lemma is valid: specifically if we rescale θ by $\sqrt{N}\delta_N$, where $\delta_N = N^{-\alpha/(2\alpha+d)}$ so that $\sqrt{N}\delta_N = N^{(d/2)/(2\alpha+d)} = N^k$, then we just take $N^{-\beta+k} = N^{-b}$ in the above small ball computation, that is $-b = -\beta + k$ or $b = \beta - k$, and the same bounds (as well as the proof to follow) apply.

4.4.2 Proof of Theorem 3.3, part iv)

Lemma 4.5 and the hypotheses on η immediately imply

$$\Pi(\theta \in \Theta_{s,\eta}) = \Pi(\|\theta\|_{\mathbb{R}^D} \in (s_bN^{-b}, s_bN^{-b} + \eta)) \leq \Pi(\|\theta\|_{\mathbb{R}^D} \leq 2s_bN^{-b}) \leq e^{-c_0N(2s_b + \kappa^{-\alpha/d}(2s_b)^{-\tau/2})^{-\tau}}.$$

To lower bound $\Pi(\Theta_{N^{-b},(1+\varepsilon)N^{-b}})$, we choose ε large enough such that

$$\bar{c}_0(1 + \varepsilon)^{-\tau} < c_0(1 + \kappa^{-\alpha/d})^{-\tau},$$

which implies for all N large enough that

$$\begin{aligned}\Pi(\|\theta\|_{\mathbb{R}^D} \in (N^{-b}, (1 + \varepsilon)N^{-b})) &= \Pi(\|\theta\|_{\mathbb{R}^D} \leq (1 + \varepsilon)N^{-b}) - \Pi(\|\theta\|_{\mathbb{R}^D} \leq N^{-b}) \\ &\geq e^{-\bar{c}_0(1+\varepsilon)^{-\tau}N} - e^{-c_0(1+\kappa^{-\alpha/d})^{-\tau}N} \\ &\geq e^{-2\bar{c}_0(1+\varepsilon)^{-\tau}N}.\end{aligned}\tag{40}$$

Now, for w from (21), we choose

$$t = t_bN^{-b}, L = L_bN^{-b}, \rho \in (0, 1], \quad 0 < t_b < s_b < 1/2 < L_b < \infty, \quad T = T_bN^b,\tag{41}$$

for T_b, ρ, s_b, t_b fixed constants to be chosen, so that $\|w\|_\infty$ is bounded (uniformly in N) by a constant which depends only on T_b, L_b, ρ , whence (29) holds. Now the key inequality (30) with $s = s_bN^{-b}$ and with our choice of $\eta, \varepsilon, \sigma = N^{-b}$ will be satisfied if

$$c_0(2s_b + \kappa^{-\alpha/d}(2s_b)^{-\tau/2})^{-\tau} \geq 2\bar{c}_0(1 + \varepsilon)^{-\tau} + 2\nu + \frac{\rho}{2}N^{-b}(1 + \varepsilon - s_b).\tag{42}$$

We define ν to equal to 1/3 of the l.h.s. so that (42) will follow for the given s_b, κ, α, d by choosing ε large enough and whenever N is large enough.

Finally, let us notice that with $L = C\sqrt{N}$ for some $C \geq 2\mathbb{E}[\|\theta\|_{\ell_2}]$, where θ is the infinite Gaussian vector with RKHS h^α , we can deduce from Theorem 2.1.20 and Exercise 2.1.5 in [GN16] that

$$\Pr(\|\theta\|_{\mathbb{R}^D} \geq L) \leq 2\exp(-cC^2N/2), \quad \text{some } c > 0.$$

Thus, using also (40), choosing C large enough verifies (26). Since (40) and the a.s. boundedness of $\sup_\theta |\ell_N(\theta)|$ for ℓ_N from (18) imply that $\Pi(\Theta_{N^{-b},(1+\varepsilon)N^{-b}}|Z^{(N)}) > 0$ a.s., Proposition 4.2 and then also Proposition 4.4 apply for this prior, and the arguments from Section 4.3.2 yield the desired result.

4.4.3 Proof of Theorem 3.3, part iii)

We finish the proof of the theorem by showing point **iii**). We use the setting and choices from the previous subsection. Let us write $\mathbb{G}(A) = \int_A e^{\ell_N(\theta)} d\Pi(\theta)$ for any measurable set A . Recall the notation $B_r = \{\theta : \|\theta\|_{\mathbb{R}^D} \leq r\}, r > 0$. Repeating the argument leading to (32) with $B_{t/2}$ in place of $\Theta_{r,\varepsilon}$, and using Lemma 4.5, we have with high probability

$$\frac{1}{N} \log \mathbb{G}(B_{t/2}) \geq -\frac{1}{2} - \frac{\sup_{r \leq t_bN^{-b}/2} w(r)}{2} - \bar{c}_0\left(\frac{t_b}{2}\right)^{-\tau} + \omega'_N(t/2),$$

where $\omega'_N(t/2) = \mathcal{O}(\|w\|_\infty/\sqrt{N}) = \mathcal{O}(1)$. Likewise, we also have

$$\frac{1}{N} \log \mathbb{G}(B_s^c) \leq -\frac{1}{2} - \frac{\inf_{r \geq s_b N^{-b}} w(r)}{2} + \frac{1}{N} \log \Pi(B_s^c) + \omega''_N(s),$$

where $\omega''_N(s) = \mathcal{O}(\|w\|_\infty/\sqrt{N}) = \mathcal{O}(1)$. We can assume that $\mathbb{G}(B_s^c) > 0$. Hence, since $\Pi(B_s^c) \rightarrow 1$ in view of Lemma 4.5,

$$\begin{aligned} \frac{1}{N} \log \frac{\mathbb{G}(B_{t/2})}{\mathbb{G}(B_s^c)} &\geq -\frac{(Tt)^2}{2} - c_0 \left(\frac{t_b}{2}\right)^{-\tau} + \frac{(Tt)^2 + (Tt/2) + \rho(s-t)}{2} + \frac{1}{N} \log \Pi(B_s^c) + \mathcal{O}(1) \\ &\geq \frac{T_b t_b}{4} - c_0 \left(\frac{t_b}{2}\right)^{-\tau} + \mathcal{O}(1). \end{aligned} \quad (43)$$

Now, for $t_b < s_b$ fixed we can choose T_b large enough such that the last quantity exceeds 1 with high probability (in particular this retrospectively justifies the last $\mathcal{O}(1)$ as then $\|w\|_\infty = \mathcal{O}(1)$ for our choice of T_b). Therefore, again with high probability

$$\frac{\mathbb{G}(B_{t/2})}{\mathbb{G}(B_s^c)} \geq e^N \times (1 + \mathcal{O}(1)). \quad (44)$$

For $M_{t,s} = \{\theta : t/2 < \|\theta\|_{\mathbb{R}^D} \leq s\}$ this further implies that with high probability

$$\frac{\mathbb{G}(B_{t/2}) + \mathbb{G}(M_{t,s})}{\mathbb{G}(B_s^c)} \geq e^N \times (1 + \mathcal{O}(1)),$$

and then,

$$\begin{aligned} \Pi(B_s | Z^{(N)}) &= \frac{\mathbb{G}(B_{t/2}) + \mathbb{G}(M_{t,s})}{\mathbb{G}(B_{t/2}) + \mathbb{G}(M_{t,s}) + \mathbb{G}(B_s^c)} \\ &= \frac{\mathbb{G}(B_{t/2}) + \mathbb{G}(M_{t,s})}{(\mathbb{G}(B_{t/2}) + \mathbb{G}(M_{t,s})) \left(1 + \frac{\mathbb{G}(B_s^c)}{\mathbb{G}(B_{t/2}) + \mathbb{G}(M_{t,s})}\right)} \rightarrow 1, \end{aligned}$$

again with high probability, which is what we wanted to show.

Remark 4.7. If the map w is globally convex, say $w(s) = Ts^2/2$ for all $s > 0$, then the ‘small enough’ choice of s after (42) is still possible but then depends on the global coercivity constant T , which will prevent the previous contraction argument to work. So while the hitting time lower bound is still valid, we cannot conclude that we are never hitting regions of significant posterior probability. It is here where global log-concavity of the likelihood function helps, as it enforces a certain ‘uniform’ spread of the posterior across its support via a global coercivity constant T . In contrast the above example of w is not convex, rather it is very spiked on $(0, t/2)$ and then “flattens out”.

4.5 Proof of Theorem 3.2

The proof of Theorem 3.2 proceeds along the same lines as the one of Theorem 3.3, with scaling t, L, ρ, s, η constant in N , corresponding to $b = 0$ in N^{-b} , and replacing the volumetric Lemma 4.5 by the following basic result.

Lemma 4.8. *Let $\theta \sim \mathcal{N}(0, I_D/D)$. Let $a \in (0, 1/2)$. Then for all $D \geq D_0(a)$ large enough,*

$$-\frac{1}{D} \log \Pi(\|\theta\|_{\mathbb{R}^D} \leq z) \geq \frac{1}{2} \left(\frac{z^2}{2} - \log z - \frac{1}{2} \right), \quad \text{any } z \in (0, 1-a). \quad (45)$$

A proof of (45) is sketched in Appendix B. As a consequence of the previous lemma

$$\frac{1}{N} \log \Pi(\Theta_{s,\eta}) \leq \frac{1}{N} \log \Pi(B_{2s}) \leq \frac{\kappa}{2} (\log 2s - 2s^2 + \frac{1}{2}).$$

Moreover, to lower bound $\Pi(\Theta_{2/3,\varepsilon})$, we choose $\varepsilon > 2/3$. Then, using Theorem 2.5.7 in [GN16] as well as $\mathbb{E}\|\theta\| \leq \mathbb{E}(\|\theta\|^2)^{1/2} = 1$, and then also (45) with $z = 2/3$, we obtain that

$$\begin{aligned} \Pi(\Theta_{2/3,\varepsilon}) &\geq \Pi(\|\theta\|_{\mathbb{R}^D} - 1 \leq 1/3) \\ &\geq 1 - \Pi(\|\theta\|_{\mathbb{R}^D} \geq \mathbb{E}\|\theta\|_{\mathbb{R}^D} + 1/3) - \Pi(\|\theta\|_{\mathbb{R}^D} \leq 2/3) \\ &\geq 1 - \exp(-D/18) - \exp(-cD), \end{aligned}$$

for some fixed constant $c > 0$ given by (45), whence $\Pi(\Theta_{2/3,\varepsilon}) \rightarrow 1$ and also $N^{-1} \log \Pi(\Theta_{2/3,\varepsilon}) \rightarrow 0$. Therefore, the key inequality (30) with $\sigma = 2/3$, $\nu = 1$ holds whenever we choose $s = s_0$ small enough such that

$$-\log 2s_0 > 2\kappa^{-1} \left[2 + \frac{\rho}{2}(s_0 - 2/3 - \varepsilon) \right] + 2s_0^2 + \frac{1}{2}.$$

The rest of the detailed derivations follow the same pattern as in the proof of Theorem 3.3 and are left to the reader, including verification of (26) via an application of Theorem 2.5.7 in [GN16]. In particular, the proof of part **iii**) follows the same arguments (suppressing the N^{-b} scaling everywhere) as in Theorem 3.3.

4.6 Proofs for Section 3.2

In this section, we prove the results of Section 3.2 which detail the consequences of the general Theorems 3.2 and 3.3 for practical MCMC algorithms.

4.6.1 Proofs for pCN

Theorem 3.6 is proved by verifying the Assumption 3.1 for suitable choices of η and L , and for $c_0 = \kappa/2 > 0$.

Lemma 4.9. *Let \mathcal{P}_N denote the transition kernel of pCN from (14) with parameter $\beta > 0$.*

i) Suppose $\Pi = \mathcal{N}(0, I_D/D)$ as in Theorem 3.2, and let $L, \eta > 0$. Then for all $\beta \leq \min\{1/2, \eta/4L, \eta^2/64\}$ and all $D \geq 1$, we have (with P_0^N -probability 1)

$$\sup_{\theta \in B_L} \mathcal{P}_N(\theta, \{\vartheta : \|\theta - \vartheta\|_{\mathbb{R}^D} \geq \eta/2\}) \leq e^{-D/2}.$$

ii) Suppose $\Pi = \mathcal{N}(0, \Sigma_\alpha)$ as in Theorem 3.3, and let $L, \eta > 0$. There exists some $c > 0$ such that for all $\beta \leq \min\{1/2, \eta/4L, c\eta^2/D\}$ and all $D \geq 1$, we have (with P_0^N -probability 1)

$$\sup_{\theta \in B_L} \mathcal{P}_N(\theta, \{\vartheta : \|\theta - \vartheta\|_{\mathbb{R}^D} \geq \eta/2\}) \leq e^{-D/2}.$$

Proof of Lemma 4.9 – We begin with the proof of part **ii**). Let $\|\vartheta_k\|_{\mathbb{R}^D} \leq L$. Then using the definition of pCN and that $|\sqrt{1-\beta} - 1| \leq \beta$ for any $\beta \in [0, 1]$ (Taylor expanding $\sqrt{\cdot}$ around 1), we obtain that for any $\beta \leq \min\{1/2, \eta/4L\}$,

$$\begin{aligned} \Pr(\|\vartheta_{k+1} - \vartheta_k\|_{\mathbb{R}^D} \geq \eta/2) &\leq \Pr(\|p_{k+1} - \vartheta_k\|_{\mathbb{R}^D} \geq \eta/2) \\ &\leq \Pr(\|(\sqrt{1-\beta} - 1)\vartheta_k\|_{\mathbb{R}^D} + \sqrt{\beta}\|\xi_k\|_{\mathbb{R}^D} \geq \eta/2) \\ &\leq \Pr(\|\xi_k\|_{\mathbb{R}^D} \geq (\eta/2 - \beta L)/\sqrt{\beta}) \\ &\leq \Pr(\|\xi_k\|_{\mathbb{R}^D} \geq \frac{\eta}{4\sqrt{\beta}}) \\ &= \Pr(\|\xi_k\|_{\mathbb{R}^D} - \mathbb{E}\|\xi_k\|_{\mathbb{R}^D} \geq \frac{\eta}{4\sqrt{\beta}} - \mathbb{E}\|\xi_k\|_{\mathbb{R}^D}). \end{aligned}$$

The variables ξ_k are equal in law to a vector with components $(i^{-\alpha/d}g_i : i \leq D)$ for g_i iid $N(0, 1)$ and hence $\mathbb{E}\|\xi_k\|_{\mathbb{R}^D} \leq (\mathbb{E}\|\xi_k\|_{\mathbb{R}^D}^2)^{1/2} \leq C(\alpha, d) < \infty$ for $\alpha > d/2$. Then, for $\beta \leq c\eta^2/D$ with some sufficiently small $c > 0$ (noting that then also $\beta \leq c\eta^2$), it holds that

$$\Pr(\|\vartheta_{k+1} - \vartheta_k\|_{\mathbb{R}^D} \geq \eta/2) \leq \Pr(\|\xi_k\|_{\mathbb{R}^D} - \mathbb{E}\|\xi_k\|_{\mathbb{R}^D} \geq \frac{\eta}{8\sqrt{\beta}}) \leq \exp\left(-\frac{\eta^2}{64\beta}\right) \leq \exp(-D/2), \quad (46)$$

using, e.g. Theorem 2.5.8 in [GN16] (and representing the $\|\cdot\|_{\mathbb{R}^D}$ -norm by duality as a supremum). This completes the proof of part **ii**).

The proof of part **i**) is similar, albeit simpler, whence we leave some details to the reader. Arguing similarly as before, we obtain that for any $\beta \leq \min\{1/2, \eta/64L\}$,

$$\Pr(\|\vartheta_{k+1} - \vartheta_k\|_{\mathbb{R}^D} \geq \eta/2) \leq \Pr(\|\xi_k\|_{\mathbb{R}^D} \geq (\eta/2 - \beta L)/\sqrt{\beta}) \leq \Pr(\|g_k\|_{\mathbb{R}^D} \geq \frac{\eta\sqrt{D}}{4\sqrt{\beta}}),$$

where g_k is a $\mathcal{N}(0, I_D)$ random variable. The latter probability is bounded by a standard deviation inequality for Gaussians, see, e.g. Theorem 2.5.7 in [GN16]. Indeed, noting that $\mathbb{E}\|\xi_k\|_{\mathbb{R}^D} \leq$

$(\mathbb{E}[\|\xi_k\|_{\mathbb{R}^D}^2])^{1/2} = \sqrt{D}$, and that the one-dimensional variances satisfy $\mathbb{E}\langle g_k, v \rangle^2 = \|v\|_{\mathbb{R}^D}^2 = 1$ for any $\|v\|_{\mathbb{R}^D} = 1$, we obtain

$$\begin{aligned} \Pr(\|g_k\|_{\mathbb{R}^D} \geq \frac{\eta\sqrt{D}}{4\sqrt{\beta}}) &\leq \Pr(|\|\xi_k\|_{\mathbb{R}^D} - \mathbb{E}\|\xi_k\|_{\mathbb{R}^D}| \geq \sqrt{D}(\frac{\eta}{4\sqrt{\beta}} - 1)) \\ &\leq \exp\left(-\frac{D}{2}\left(\frac{\eta}{4\sqrt{\beta}} - 1\right)^2\right) \leq \exp\left(-\frac{D}{2}\right). \end{aligned}$$

□

Proof of Theorem 3.6 – We begin with part **ii**). Let s_b be as in Theorem 3.3 and set $\eta = \eta_N = s_b N^{-b}/2$ as well as $L = L_N C \sqrt{N}$, where C is as in Theorem 3.3. With those choices, Lemma 4.9 **ii**) implies that Assumption 3.1 is fulfilled with $c_0 = \kappa/2$, so long as β satisfies

$$\beta \leq \min\left\{\frac{1}{2}, \frac{s_b N^{-b}}{8C\sqrt{N}}, \frac{cs_b^2 N^{-2b}}{4D}\right\} \lesssim N^{-2b} D^{-1} \simeq N^{-1-2b}.$$

Hence, the desired result immediately follows from an application of Theorem 3.3 **iv**).

Part **i**) of Theorem 3.6 similarly follows from verifying Assumption 3.1 with $s \in (0, 1/3)$, L from Theorem 3.2, $\eta = s/2$ and for small enough $\beta < c_1$ (with c_1 determined by Lemma 4.9 **i**)), and subsequently applying Theorem 3.2 **iv**). □

4.6.2 Proofs for MALA

Theorem 3.7 is proved by verifying the hypotheses of Theorems 3.2 and 3.3 respectively. A key difference between pCN and MALA is that the proposal kernels for MALA, not just its acceptance probabilities, depend on the data $Z^{(N)}$ itself. Again, we begin by examining part **ii**) which regards $N(0, \Sigma_\alpha)$ priors.

Proof of Theorem 3.7, part ii) We begin by deriving a bound for the gradient $\nabla \log \pi(\cdot|Z^{(N)})$. For Lebesgue-a.e. $\theta \in \mathbb{R}^D$, recalling that $\text{vol}(\mathcal{X}) = 1$, we have that

$$\begin{aligned} \mathbb{E}_0^N[\nabla \ell_N(\theta)] &= -\frac{N}{2} w'(\|\theta\|) \frac{\theta}{\|\theta\|} \|g\|_{L^2}^2, \\ \nabla \ell_N(\theta) &= \frac{1}{2} \sum_{i=1}^N \left(\varepsilon_i - \sqrt{w(\|\theta\|)} g(X_i) \right) \frac{w'(\|\theta\|)}{2\sqrt{w(\|\theta\|)}} \frac{\theta}{\|\theta\|} g(X_i), \\ &= \frac{w'(\|\theta\|)}{4\sqrt{w(\|\theta\|)}} \frac{\theta}{\|\theta\|} \sum_{i=1}^N \varepsilon_i g(X_i) - \frac{w'(\|\theta\|)}{4} \frac{\theta}{\|\theta\|} \sum_{i=1}^N g^2(X_i). \end{aligned}$$

For any $r \in (0, t/2) \cup (t/2, t) \cup (t, L) \cup (L, \infty)$, recalling the choices for T, t, ρ in (41) we see that

$$\begin{aligned} \frac{w'(r)}{\sqrt{w(r)}} &= \frac{8Tr}{2Tr} 1_{(0, t/2)}(r) + \frac{T}{\sqrt{w(r)}} 1_{(t/2, t)}(r) + \frac{\rho}{\sqrt{w(r)}} 1_{(t, L)}(r), \\ &\lesssim 1 + N^b + 1, \end{aligned} \tag{47}$$

where, to bound the second and third term, we used that $\sqrt{w(r)} \geq Tt = t_b T_b > 0$ is bounded away from zero uniformly in N on $(t/2, \infty)$. Similarly, we have

$$\|w'\|_\infty \leq Tt/2 + T + \rho \lesssim N^b.$$

Combining the above and using Chebyshev's inequality, it follows that

$$\begin{aligned} \sup_{\theta \in \mathbb{R}^D} \|\nabla \ell_N(\theta)\|_{\mathbb{R}^D} &\lesssim N^b \left(\left| \sum_{i=1}^N \varepsilon_i g(X_i) \right| + \sum_{i=1}^N g^2(X_i) \right) \\ &\leq N^b \left(\|g\|_\infty \left| \sum_{i=1}^N \varepsilon_i \right| + \sum_{i=1}^N (g^2(X_i) - \|g\|_{L^2}^2) + N \|g\|_{L^2}^2 \right) \\ &\leq N^b (\mathcal{O}_P(\sqrt{N}) + \mathcal{O}(N)) \\ &= \mathcal{O}(N^{1+b}) + \mathcal{o}(N^{1+b}). \end{aligned}$$

Thus, the event

$$A := \left\{ \sup_{\theta \in \mathbb{R}^D} \|\nabla \ell_N(\theta)\|_{\mathbb{R}^D} \leq C' N^{1+b} \right\},$$

for some large enough $C' > 0$, has probability $P_0^N(A) \rightarrow 1$ as $N \rightarrow \infty$. We also verify that

$$\nabla \log \pi(\theta) = -\frac{1}{2} \nabla \theta^T \Sigma_\alpha^{-1} \theta = -\Sigma_\alpha^{-1} \theta, \quad (48)$$

so that with $L = L_N = C\sqrt{N}$ (for C as in Theorem 3.3) and recalling that $\Sigma_\alpha = \text{diag}(1, \dots, D^{-2\alpha})$, we obtain

$$\sup_{\|\theta\| \leq L} \|\nabla \log \pi(\theta)\|_{\mathbb{R}^D} = \sup_{\|\theta\| \leq L} \|\Sigma_\alpha^{-1} \theta\|_{\mathbb{R}^D} \lesssim D^{2\alpha} \sqrt{N} \simeq N^{2\alpha+1}.$$

Now, let s_b, C be as in Theorem 3.3 and set $\eta = \eta_N = \frac{1}{2} s_b N^{-b}$ (note that this is a permissible choice in Theorem 3.3) as well as $L = L_N = C\sqrt{N}$. Furthermore, for a small enough constant $c > 0$, let $\gamma \leq cN^{-1-2\alpha-b}$. Then since $\alpha > b$, we also have that

$$\gamma \lesssim \min\{N^{-1-2\alpha-b}, N^{-1-2b}, N^{-1/2-b}\}. \quad (49)$$

Hence, on the event A and whenever $\|\theta\|_{\mathbb{R}^D} \leq L$,

$$\gamma \|\nabla \log \pi(\vartheta_k | Z^{(N)})\|_{\mathbb{R}^D} \lesssim \gamma(N^{1+b} + N^{1+2\alpha}) \lesssim \eta.$$

Using this, (49) and choosing $c > 0$ small enough, conditional on the event A the probability $\Pr(\cdot)$ under the Markov chain satisfies

$$\begin{aligned} \Pr(\|p_{k+1} - \vartheta_k\| \geq \eta/2) &\leq \Pr(\gamma \|\nabla \log \pi(\vartheta_k | Z^{(N)})\|_{\mathbb{R}^D} \geq \eta/4) + \Pr(\sqrt{2\gamma} \|\xi_{k+1}\|_{\mathbb{R}^D} \geq \eta/4) \\ &\leq \Pr(\|\xi_{k+1}\|_{\mathbb{R}^D} \geq \frac{\eta}{4\sqrt{2\gamma}}) \\ &\leq \Pr(\|\xi_{k+1}\|_{\mathbb{R}^D} - \mathbb{E}\|\xi_{k+1}\|_{\mathbb{R}^D} \geq \sqrt{N}) \leq \exp\left(-\frac{N}{2}\right), \end{aligned}$$

where the last inequality is proved as in (46) above, using Theorem 2.5.8 in [GN16]. Thus, Assumption 3.1 is satisfied with $c_0 = 1$ and the proof is complete. \square

Proof of Theorem 3.7, part i) – The proof of part i) proceeds along the same lines, except that (47) and (48) are replaced with the bound

$$\left\| \frac{w'}{\sqrt{w}} \right\|_\infty + \|w'\|_\infty < C,$$

for some constant C independent of N , as well as the bound

$$\nabla \log \pi(\theta) = -\frac{D}{2} \nabla \|\theta\|^2 = -D\theta, \quad \sup_{\|\theta\| \leq L} \|\nabla \log \pi(\theta)\|_{\mathbb{R}^D} \simeq NL.$$

Then letting $s \in (0, 1/3)$ and $L > 0$ be as in Theorem 3.2, and fixing an arbitrary $\eta \in (0, s/2)$, the above implies that for sufficiently small constant $c > 0$ and for any $\gamma \leq c/N$, it holds that

$$\begin{aligned} \Pr(\|p_{k+1} - \vartheta_k\| \geq \eta/2) &\leq \Pr(\gamma \|\nabla \log \pi(\vartheta_k | Z^{(N)})\|_{\mathbb{R}^D} \geq \eta/4) + \Pr(\sqrt{2\gamma} \xi_{k+1} \geq \eta/4) \\ &\leq \Pr(\xi_{k+1} \geq \frac{\eta}{4\sqrt{2\gamma}}) \\ &\leq \Pr(\xi_{k+1} \geq \frac{\eta\sqrt{\kappa D}}{4\sqrt{2c}}). \end{aligned}$$

Thus, choosing $c > 0$ small enough and arguing exactly as in the last step of the proof of Theorem 3.6, part i), Assumption 3.1 is satisfied with $c_0 = 1$ and the proof is complete. \square

A Proofs of Section 2

Proof of Corollary 2.4 – We fix $K = 1$ and place ourselves under the event of Proposition 2.3, and we denote $s = s(\lambda)$ and $t = t(\lambda)$. We can decompose, since $\Pi[\mathcal{S}_s|Y] = 1 - \Pi[\mathcal{T}_s|Y]$:

$$\Pi(\mathcal{T}_s|Y) = \left[1 + \frac{\Pi(\mathcal{S}_s|Y)}{\Pi(\mathcal{T}_s|Y)}\right]^{-1}.$$

Moreover, $\Pi(\mathcal{S}_s|Y)/\Pi(\mathcal{T}_s|Y) = \Pi(\mathcal{S}_s|Y)/[\Pi(\mathcal{T}_t|Y) + \Pi(\mathcal{W}_{s,t}|Y)] \leq \Pi(\mathcal{S}_s|Y)/\Pi(\mathcal{T}_t|Y)$. Using Proposition 2.3, for $n \geq n_0(\lambda, Y)$ we have $\Pi(\mathcal{S}_s|Y)/\Pi(\mathcal{T}_s|Y) \leq \exp\{-n\}$. Therefore $\Pi[\mathcal{T}_s|Y] \geq (1 + \exp\{-n\})^{-1}$, which ends the proof. \square

The rest of this section is devoted to proving Proposition 2.3. We use a uniform bound on the injective norm of Gaussian tensors:

Lemma A.1. *For all $p \geq 3$ there exists a constant C_p , such that:*

$$\limsup_{n \rightarrow \infty} \left\{ n^{-1/2} \max_{x \in \mathbb{S}^{n-1}} |\langle x^{\otimes p}, Z \rangle| \right\} \leq C_p, \quad \text{almost surely.} \quad (50)$$

This lemma is a very crude version of much finer results: in particular the exact value of the constant μ_p such that (w.h.p.) $\max_{x \in \mathbb{S}^{n-1}} |\langle x^{\otimes p}, Z \rangle| = \sqrt{n} \mu_p (1 + o_n(1))$ has been first computed non-rigorously in [CS92], and proven in full generality in [Sub17] (see also discussions in [RM14; PWB20]). In the rest of this proof, we assume to have conditioned on eq. (50). For any $0 \leq s < t \leq 1$, we have for $n \geq n_0(Y)$:

$$\begin{aligned} \frac{\Pi(\mathcal{S}_s|Y)}{\Pi(\mathcal{T}_t|Y)} &= \frac{\int_{\mathcal{S}_s} \exp(\ell_Y(x)) d\Pi(x)}{\int_{\mathcal{T}_t} \exp(\ell_Y(x)) d\Pi(x)} \leq e^{n\lambda C_p} \frac{\int_{\mathcal{S}_s} \exp\left(\frac{n}{2} \lambda^2 \langle x, x_0 \rangle^p\right) d\Pi(x)}{\int_{\mathcal{T}_t} \exp\left(\frac{n}{2} \lambda^2 \langle x, x_0 \rangle^p\right) d\Pi(x)}, \\ &\leq \exp\left(n\lambda C_p + \frac{n\lambda^2}{2} [s^p - t^p]\right) \frac{\Pi(\mathcal{S}_s)}{\Pi(\mathcal{T}_t)}. \end{aligned} \quad (51)$$

We upper bound $\Pi(\mathcal{S}_s) \leq \Pi(\mathbb{S}^{n-1}) = 1$. To lower bound $\Pi(\mathcal{T}_t)$, we use the elementary fact (which is easy to prove using spherical coordinates):

$$\Pi(\mathcal{T}_t) = c_p I_{(1-t)/2}[(n-1)/2, (n-1)/2], \quad (52)$$

in which $I_x(a, b) = \int_0^x u^{a-1} (1-u)^{b-1} du / \int_0^1 u^{a-1} (1-u)^{b-1} du$ is the incomplete beta function, and $c_p = 1$ for odd p and $c_p = 2$ for even p . It is then elementary analysis (cf. e.g. [PWB20]) that

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \Pi(\mathcal{T}_t) = \frac{1}{2} \log(1 - t^2), \quad (53)$$

uniformly in $t \in [0, 1)$. Coming back to eq. (51), this implies that we have, for any $s < t < 1$:

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log \frac{\Pi(\mathcal{S}_s|Y)}{\Pi(\mathcal{T}_t|Y)} \leq \lambda C_p + \frac{\lambda^2}{2} [s^p - t^p] - \frac{1}{2} \log(1 - t^2). \quad (54)$$

Let $K > 0$. It is then elementary to see that it is possible to construct $0 \leq s(\lambda) < t(\lambda) < 1$ with $\lim_{\lambda \rightarrow \infty} \{s(\lambda), t(\lambda)\} = 1$, and such that the right-hand-side of eq. (54) becomes smaller than $-K$ as $\lambda \rightarrow \infty$. \square

B Small ball estimates for isotropic Gaussians

Let $\Pi = \mathcal{N}(0, I_D/D)$. In this section, we prove eq. (45), more precisely we show:

Lemma B.1. *Let $a \in (0, 1)$. Then for all $D \geq D_0(a)$ large enough, one has for all $z \in (0, 1 - a)$:*

$$-\frac{1}{D} \log \Pi(\|\theta\|_2 \leq z) \geq \frac{1}{2} \left(\frac{z^2}{2} - \log z - \frac{1}{2} \right). \quad (55)$$

Proof of Lemma B.1 – Let $f(x) = -x^2/2 + \log x + 1/2$, so that f reaches its maximum in $x = 1$, with $f(1) = 0$. By decomposition into spherical coordinates and isotropy of the Gaussian measure, one has directly:

$$\Pi(\|\theta\|_2 \leq z) = \frac{\text{vol}(\mathbb{S}^{D-1})}{(2\pi/D)^{D/2}} \int_0^z dr e^{-\frac{Dr^2}{2} + (D-1)\log r}. \quad (56)$$

Recall that $\text{vol}(\mathbb{S}^{D-1}) = 2\pi^{D/2}/\Gamma(D/2)$, so one reaches easily:

$$c_D = \frac{1}{D} \log \frac{\text{vol}(\mathbb{S}^{D-1})}{(2\pi/D)^{D/2}} - \frac{1}{2} = \frac{\log D}{2D} + \mathcal{O}(1/D). \quad (57)$$

In particular, one has for all D large enough (not depending on z):

$$\frac{1}{D} \log \Pi(\|\theta\|_2 \leq z) \leq \frac{1}{D} \log \int_0^z dr e^{-\frac{r^2}{2} + (D-1)f(r)} + c_D. \quad (58)$$

Since f is increasing on $(0, 1)$, we have for large enough D :

$$\frac{1}{D} \log \Pi(\|\theta\|_2 \leq z) \leq \left(1 - \frac{1}{D}\right) f(z) + c_D + \frac{1}{D} \log \int_0^\infty dr e^{-r^2/2}, \quad (59)$$

$$\leq \left(1 - \frac{1}{D}\right) f(z) + \frac{\log D}{D}. \quad (60)$$

Since $f(1-a) < 0$, let $D \geq D_0(a)$ large enough such that $f(1-a) \leq -2 \log D/(D-2)$. Then for all $z \leq 1-a$, one has $f(z) \leq -2 \log D/(D-2)$. Plugging it in the inequality above, we reach that for all $z \in (0, 1-a)$:

$$\frac{1}{D} \log \Pi(\|\theta\|_2 \leq z) \leq \frac{1}{2} f(z). \quad (61)$$

□

References

- [Alt22] Randolf Altmeyer. “Polynomial time guarantees for sampling based posterior inference in high-dimensional generalised linear models”. In: *arXiv preprint 2208.13296* (2022).
- [And89] Philip W Anderson. “Spin glass VI: Spin glass as cornucopia”. In: *Physics Today* 42.9 (1989), p. 9.
- [Ang18] Maria Chiara Angelini. “Parallel tempering for the planted clique problem”. In: *Journal of Statistical Mechanics: Theory and Experiment* 2018.7 (2018), p. 073404.
- [AFF21] Maria Chiara Angelini, Paolo Fachin, and Simone de Feo. “Mismatching as a tool to enhance algorithmic performances of Monte Carlo methods for the planted clique model”. In: *Journal of Statistical Mechanics: Theory and Experiment* 2021.11 (2021), p. 113406.
- [BBP05] Jinho Baik, Gérard Ben Arous, and Sandrine Péché. “Phase transition of the largest eigenvalue for nonnull complex sample covariance matrices”. In: *The Annals of Probability* 33.5 (2005), pp. 1643–1697.
- [Ban+22] Afonso S Bandeira et al. “The Franz-Parisi Criterion and Computational Trade-offs in High Dimensional Statistics”. In: *arXiv preprint arXiv:2205.09727* (2022).
- [BGJ20a] Gérard Ben Arous, Reza Gheissari, and Aukosh Jagannath. “Algorithmic thresholds for tensor PCA”. In: *The Annals of Probability* 48.4 (2020), pp. 2052–2087.

- [BGJ20b] Gérard Ben Arous, Reza Gheissari, and Aukosh Jagannath. “Bounding flows for spherical spin glass dynamics”. In: *Communications in Mathematical Physics* 373.3 (2020), pp. 1011–1048.
- [BGJ21] Gérard Ben Arous, Reza Gheissari, and Aukosh Jagannath. “Online stochastic gradient descent on non-convex losses from high-dimensional inference.” In: *J. Mach. Learn. Res.* 22 (2021), pp. 106–1.
- [BWZ20] Gérard Ben Arous, Alexander S Wein, and Ilias Zadik. “Free energy wells and overlap gap property in sparse PCA”. In: *Conference on Learning Theory*. PMLR, 2020, pp. 479–482.
- [Bes+17] Alexandros Beskos et al. “Geometric MCMC for infinite-dimensional inverse problems”. In: *J. Comput. Phys.* 335 (2017), pp. 327–351. ISSN: 0021-9991. DOI: 10.1016/j.jcp.2016.12.041.
- [Bie+20] Joris Bierkens et al. “The boomerang sampler”. In: *PMLR* (2020).
- [BCR20] Giulio Biroli, Chiara Cammarota, and Federico Ricci-Tersenghi. “How to iron out rough landscapes and get optimal performances: averaged gradient descent and its application to tensor PCA”. In: *Journal of Physics A: Mathematical and Theoretical* 53.17 (2020), p. 174003.
- [BN21] Jan Bohr and Richard Nickl. “On log-concave approximations of high-dimensional posterior measures and stability properties in non-linear inverse problems”. In: *arXiv preprint arXiv:2105.07835* (2021).
- [BVD18] Alexandre Bouchard-Côté, Sebastian J. Vollmer, and Arnaud Doucet. “The bouncy particle sampler: a nonreversible rejection-free Markov chain Monte Carlo method”. In: *J. Amer. Statist. Assoc.* 113.522 (2018), pp. 855–867. ISSN: 0162-1459. DOI: 10.1080/01621459.2017.1294075. URL: <https://doi.org/10.1080/01621459.2017.1294075>.
- [BPS04] Laird Arnault Breyer, Mauro Piccioni, and Sergio Scarlatti. “Optimal scaling of MaLa for nonlinear regression”. In: *Ann. Appl. Probab.* 14.3 (2004), pp. 1479–1505. ISSN: 1050-5164. DOI: 10.1214/105051604000000369. URL: <https://doi.org/10.1214/105051604000000369>.
- [CMZ22] Zongchen Chen, Elchanan Mossel, and Ilias Zadik. “Almost-Linear Planted Cliques Elude the Metropolis Process”. In: *arXiv preprint arXiv:2204.01911* (2022).
- [Che+21] Sinho Chewi et al. “Optimal dimension dependence of the Metropolis-Adjusted Langevin Algorithm”. In: *Conference on Learning Theory*. PMLR, 2021.
- [Cot+13] S. L. Cotter et al. “MCMC methods for functions: modifying old algorithms to make them faster”. In: *Statist. Sci.* 28.3 (2013), pp. 424–446. ISSN: 0883-4237. DOI: 10.1214/13-STS421.
- [CS92] Andrea Crisanti and H-J Sommers. “The spherical p -spin interaction spin glass model: the statics”. In: *Zeitschrift für Physik B Condensed Matter* 87.3 (1992), pp. 341–354.
- [Dal17] Arnak S. Dalalyan. “Theoretical guarantees for approximate sampling from smooth and log-concave densities”. In: *J. R. Stat. Soc. Ser. B. Stat. Methodol.* 79.3 (2017), pp. 651–676. ISSN: 1369-7412. DOI: 10.1111/rssb.12183.
- [DM19] Alain Durmus and Éric Moulines. “High-dimensional Bayesian inference via the unadjusted Langevin algorithm”. In: *Bernoulli* 25 (2019).
- [ET96] D. E. Edmunds and H. Triebel. *Function spaces, entropy numbers, differential operators*. Vol. 120. Cambridge Tracts in Mathematics. Cambridge University Press, Cambridge, 1996, pp. xii+252. ISBN: 0-521-56036-5. DOI: 10.1017/CB09780511662201.

- [Fea+18] Paul Fearnhead et al. “Piecewise deterministic Markov processes for continuous-time Monte Carlo”. In: *Statist. Sci.* 33.3 (2018), pp. 386–412. ISSN: 0883-4237. DOI: 10.1214/18-STS648. URL: <https://doi.org/10.1214/18-STS648>.
- [FP95] Silvio Franz and Giorgio Parisi. “Recipes for metastable states in spin glasses”. In: *Journal de Physique I* 5.11 (1995), pp. 1401–1415.
- [GZ19] David Gamarnik and Ilias Zadik. “The landscape of the planted clique problem: Dense subgraphs and the overlap gap property”. In: *arXiv preprint arXiv:1904.07174* (2019).
- [GV17] Subhashis Ghosal and Aad W. van der Vaart. *Fundamentals of Nonparametric Bayesian Inference*. Cambridge University Press, New York, 2017.
- [Gib73] Josiah Willard Gibbs. “A method of geometrical representation of the thermodynamic properties of substances by means of surfaces”. In: vol. 2. Transactions of the Connecticut Academy of Arts and Sciences. 1873, pp. 382–404.
- [GN16] Evarist Giné and Richard Nickl. *Mathematical foundations of infinite-dimensional statistical models*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, New York, 2016, pp. xiv+690.
- [GP01] Tomas S. Grigera and Giorgio Parisi. “Fast Monte Carlo algorithm for supercooled soft spheres”. In: *Physics Review E* 63 (2001).
- [HMS11] Martin Hairer, Jonathan Mattingly, and Martin Scheutzow. “Asymptotic coupling and a general form of Harris’ theorem with applications to stochastic delay equations”. In: *Probab. Theory Relat. Fields* 149 (2011), pp. 223–259.
- [HSV14] Martin Hairer, Andrew M. Stuart, and Sebastian J. Vollmer. “Spectral Gaps for a Metropolis-Hastings Algorithm in Infinite Dimensions”. In: *The Annals of Applied Probability* 24.6 (2014), pp. 2455–2490.
- [HSS15] Samuel B Hopkins, Jonathan Shi, and David Steurer. “Tensor principal component analysis via sum-of-square proofs”. In: *Conference on Learning Theory*. PMLR. 2015, pp. 956–1006.
- [Hop+16] Samuel B Hopkins et al. “Fast spectral algorithms from sum-of-squares proofs: tensor decomposition and planted sparse vectors”. In: *Proceedings of the forty-eighth annual ACM symposium on Theory of Computing*. 2016, pp. 178–191.
- [HN96] Koji Hukushima and Koji Nemoto. “Exchange Monte Carlo method and application to spin glass simulations”. In: *Journal of the Physical Society of Japan* 65.6 (1996), pp. 1604–1608.
- [JLM20] Aukosh Jagannath, Patrick Lopatto, and Léo Miolane. “Statistical thresholds for tensor PCA”. In: *The Annals of Applied Probability* 30.4 (2020), pp. 1910–1933.
- [Jer03] Mark Jerrum. *Counting, sampling and integrating: algorithms and complexity*. Lectures in Mathematics ETH Zürich. Birkhäuser Verlag, Basel, 2003, pp. xii+112. ISBN: 3-7643-6946-9. DOI: 10.1007/978-3-0348-8005-3.
- [Jer92] Mark Jerrum. “Large cliques elude the Metropolis process”. In: *Random Structures & Algorithms* 3.4 (1992), pp. 347–359.
- [KBG17] Chiheon Kim, Afonso S Bandeira, and Michel X Goemans. “Community detection in hypergraphs, spiked tensor models, and sum-of-squares”. In: *2017 International Conference on Sampling Theory and Applications (SampTA)*. IEEE. 2017, pp. 124–128.
- [KWB22] Dmitriy Kunisky, Alexander S Wein, and Afonso S Bandeira. “Notes on computational hardness of hypothesis testing: Predictions using the low-degree likelihood ratio”. In: *ISAAC Congress (International Society for Analysis, its Applications and Computation)*. Springer. 2022, pp. 1–50.

- [Les+17] Thibault Lesieur et al. “Statistical and computational phase transitions in spiked tensor estimation”. In: *2017 IEEE International Symposium on Information Theory (ISIT)*. IEEE. 2017, pp. 511–515.
- [LL99] Wenbo V. Li and Werner Linde. “Approximation, metric entropy and small ball estimates for Gaussian measures”. In: *Ann. Probab.* 27.3 (1999), pp. 1556–1578. ISSN: 0091-1798. DOI: 10.1214/aop/1022677459.
- [MPS12] Jonathan C. Mattingly, Natesh S. Pillai, and Andrew M. Stuart. “Diffusion limits of the random walk Metropolis algorithm in high dimensions”. In: *Ann. Appl. Probab.* 22.3 (2012), pp. 881–930. ISSN: 1050-5164. DOI: 10.1214/10-AAP754. URL: <https://doi.org/10.1214/10-AAP754>.
- [MM09] Marc Mézard and Andrea Montanari. *Information, physics, and computation*. Oxford University Press, 2009.
- [MNP21a] François Monard, Richard Nickl, and Gabriel P. Paternain. “Consistent inversion of noisy non-Abelian X-ray transforms”. In: *Comm. Pure Appl. Math.* 74.5 (2021), pp. 1045–1099.
- [MNP19] François Monard, Richard Nickl, and Gabriel P. Paternain. “Efficient nonparametric Bayesian inference for X-ray transforms”. In: *Ann. Statist.* 47.2 (2019), pp. 1113–1147.
- [MNP21b] François Monard, Richard Nickl, and Gabriel P. Paternain. “Statistical guarantees for Bayesian uncertainty quantification in nonlinear inverse problems with Gaussian process priors”. In: *Ann. Statist.* 49.6 (2021), pp. 3255–3298.
- [Nic22] Richard Nickl. *Bayesian non-linear statistical inverse problems*. ETH Zurich Lecture Notes. 2022, p. 160.
- [Nic20] Richard Nickl. “Bernstein–von Mises theorems for statistical inverse problems I: Schrödinger equation”. In: *J. Eur. Math. Soc. (JEMS)* 22.8 (2020), pp. 2697–2750.
- [NW20] Richard Nickl and Sven Wang. “On polynomial-time computation of high-dimensional posterior measures by Langevin-type algorithms”. In: *Journal of the European Mathematical Society, to appear* (2020).
- [PWB20] Amelia Perry, Alexander S Wein, and Afonso S Bandeira. “Statistical limits of spiked tensor models”. In: *Annales de l’Institut Henri Poincaré, Probabilités et Statistiques*. Vol. 56. Institut Henri Poincaré. 2020, pp. 230–264.
- [RW06] Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian processes for machine learning*. Adaptive Computation and Machine Learning. MIT Press, Cambridge, MA, 2006, pp. xviii+248. ISBN: 978-0-262-18253-9.
- [RM14] Emile Richard and Andrea Montanari. “A statistical model for tensor PCA”. In: *Advances in neural information processing systems 27* (2014).
- [RR01] Gareth O. Roberts and Jeffrey S. Rosenthal. “Optimal scaling for various Metropolis-Hastings algorithms”. In: *Statist. Sci.* 16.4 (2001), pp. 351–367. ISSN: 0883-4237. DOI: 10.1214/ss/1015346320. URL: <https://doi.org/10.1214/ss/1015346320>.
- [Sar+19a] Stefano Sarao Mannelli et al. “Passed & spurious: Descent algorithms and local minima in spiked matrix-tensor models”. In: *International Conference on Machine Learning*. PMLR. 2019, pp. 4333–4342.
- [Sar+19b] Stefano Sarao Mannelli et al. “Who is afraid of big bad minima? Analysis of gradient-flow in spiked matrix-tensor models”. In: *Advances in Neural Information Processing Systems 32* (2019).

- [SGB22] Camille Scalliet, Benjamin Guiselin, and Ludovic Berthier. “Thirty milliseconds in the life of a supercooled liquid”. In: *arXiv preprint 2207.00491* (2022).
- [Stu10] Andrew M. Stuart. “Inverse problems: a Bayesian perspective”. In: *Acta Numer.* 19 (2010), pp. 451–559. ISSN: 0962-4929. DOI: 10.1017/S0962492910000061.
- [Sub17] Eliran Subag. “The complexity of spherical p -spin models — A second moment approach”. In: *The Annals of Probability* 45.5 (2017), pp. 3385–3450.
- [VZ08] A. van der Vaart and J. H. van Zanten. “Rates of contraction of posterior distributions based on Gaussian process priors”. In: *Ann. Statist.* 36.3 (2008), pp. 1435–1463. ISSN: 0090-5364. DOI: 10.1214/009053607000000613.
- [WEM19] Alexander S Wein, Ahmed El Alaoui, and Cristopher Moore. “The Kikuchi hierarchy and tensor PCA”. In: *2019 IEEE 60th Annual Symposium on Foundations of Computer Science (FOCS)*. IEEE. 2019, pp. 1446–1468.
- [WR20] Changye Wu and Christian P. Robert. “Coordinate sampler: a non-reversible Gibbs-like MCMC sampler”. In: *Stat. Comput.* 30.3 (2020), pp. 721–730. ISSN: 0960-3174. DOI: 10.1007/s11222-019-09913-w. URL: <https://doi.org/10.1007/s11222-019-09913-w>.
- [ZK16] Lenka Zdeborová and Florent Krzakala. “Statistical physics of inference: Thresholds and algorithms”. In: *Advances in Physics* 65.5 (2016), pp. 453–552.