



**HAL**  
open science

# Sampling from a log-concave distribution with compact support with proximal Langevin Monte Carlo

Nicolas Brosse, Alain Durmus, Éric Moulines, Marcelo Pereyra

## ► To cite this version:

Nicolas Brosse, Alain Durmus, Éric Moulines, Marcelo Pereyra. Sampling from a log-concave distribution with compact support with proximal Langevin Monte Carlo. Proceedings of Machine Learning Research, 2017, 65, pp.319-342. hal-01648665

**HAL Id: hal-01648665**

**<https://inria.hal.science/hal-01648665>**

Submitted on 27 Nov 2017

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Sampling from a log-concave distribution with compact support with proximal Langevin Monte Carlo

Nicolas Brosse <sup>1</sup>      Alain Durmus <sup>2</sup>      Éric Moulines <sup>3</sup>  
Marcelo Pereyra <sup>4</sup>

## Abstract

This paper presents a detailed theoretical analysis of the Langevin Monte Carlo sampling algorithm recently introduced in [DMP16] when applied to log-concave probability distributions that are restricted to a convex body  $K$ . This method relies on a regularisation procedure involving the Moreau-Yosida envelope of the indicator function associated with  $K$ . Explicit convergence bounds in total variation norm and in Wasserstein distance of order 1 are established. In particular, we show that the complexity of this algorithm given a first order oracle is polynomial in the dimension of the state space. Finally, some numerical experiments are presented to compare our method with competing MCMC approaches from the literature.

## 1 Introduction

Many statistical inference problems involve estimating parameters subject to constraints on the parameter space. In a Bayesian setting, these constraints define a posterior distribution  $\pi$  with bounded support. Some examples include truncated data problems which arise naturally in failure and survival time studies [KM05],

---

<sup>1</sup>Centre de Mathématiques Appliquées, UMR 7641, Ecole Polytechnique, France.  
nicolas.brosse@polytechnique.edu

<sup>2</sup>LTCI, Telecom ParisTech 46 rue Barrault, 75634 Paris Cedex 13, France.  
alain.durmus@telecom-paristech.fr

<sup>3</sup>Centre de Mathématiques Appliquées, UMR 7641, Ecole Polytechnique, France.  
eric.moulines@polytechnique.edu

<sup>4</sup>School of Mathematical and Computer Sciences, Heriot-Watt University, Edinburgh, EH14 4AS, U.K. m.pereyra@hw.ac.uk

ordinal data models [JA06], constrained lasso and ridge regressions [Cel+12], Latent Dirichlet Allocation [BNJ03], and non-negative matrix factorization [PBJ14]. Drawing samples from such constrained distributions is a challenging problem that has been investigated in many papers; see [GSL92], [PP14], [LS15], [BEL15]. All these works are based on efficient Markov Chain Monte Carlo methods to approximate the posterior distribution; however, with the exception of the recent work [BEL15], these methods are not theoretically well understood and do not provide any theoretical guarantees on the estimations delivered.

Recently a new MCMC method has been proposed in [DMP16] to sample from a non-smooth log-concave probability distribution on  $\mathbb{R}^d$ . This method is mainly based on a carefully designed regularised version of the target distribution  $\pi$  that enjoys a number of favourable properties that are useful for MCMC simulation. In this study, we analyse the complexity of this algorithm when applied to log-concave distributions constrained to a convex set, with a focus on complexity as the dimension of the state space increases. More precisely, we establish explicit bounds in total variation norm and in Wasserstein distance of order 1 between the iterates of the Markov kernel defined by the algorithm and the target density  $\pi$ .

The paper is organised as follows. Section 2.1 introduces the MCMC method of [DMP16]. The main complexity result is stated in Section 2.2 and compared to previous works on the subject. The proof of this result is presented in Section 3 and Section 4. The methodology is then illustrated and compared to other approaches via experiments in Section 5. Proofs are finally reported in Section 6.

## 2 The Moreau-Yosida Unadjusted Langevin Algorithm (MYULA)

### 2.1 Presentation of MYULA

Let  $\pi$  be a probability measure on  $\mathbb{R}^d$  with density w.r.t. the Lebesgue measure given for all  $x \in \mathbb{R}^d$  by  $\pi(x) = e^{-U(x)} / \int_{\mathbb{R}^d} e^{-U(y)} dy$ , where  $U : \mathbb{R}^d \rightarrow (-\infty, +\infty]$  is a measurable function. In the sequel,  $U$  will be referred to as the potential associated with  $\pi$ . Assume for the moment that  $U$  is continuously differentiable. Then, the unadjusted Langevin algorithm (ULA) introduced in [Par81] (see also [RT96]) can be used to sample from  $\pi$ . This algorithm is based on the overdamped Langevin stochastic differential equation (SDE) associated with  $U$ ,

$$dY_t = -\nabla U(Y_t)dt + \sqrt{2}dB_t, \quad (1)$$

where  $(B_t)_{t \geq 0}$  is a  $d$ -dimensional Brownian motion. Under mild assumptions on  $\nabla U$ , this SDE has a unique strong solution  $(Y_t)_{t \geq 0}$  and defines a strong Markovian

semigroup  $(P_t)_{t \geq 0}$  on  $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$  which is ergodic with respect to  $\pi$ , where  $\mathcal{B}(\mathbb{R}^d)$  is the Borel  $\sigma$ -field on  $\mathbb{R}^d$ . Since simulating exact solutions of (1) is in general computationally impossible or very hard, ULA considers the Euler-Maruyama discretization associated with (1) to approximate samples from  $\pi$ . Precisely, ULA constructs the discrete-time Markov chain  $(X_k)_{k \geq 0}$ , started at  $X_0$ , given for  $k \in \mathbb{N}$  by:

$$X_{k+1} = X_k - \gamma \nabla U(X_k) + \sqrt{2\gamma} Z_{k+1} ,$$

where  $\gamma > 0$  is the stepsize and  $(Z_k)_{k \in \mathbb{N}}$  is a sequence of i.i.d. standard Gaussian  $d$ -dimensional vectors; the process  $(X_k)_{k \geq 0}$  is used as approximate samples from  $\pi$ . However, the ULA algorithm cannot be directly applied to a distribution  $\pi$  restricted to a compact convex set. Let  $\mathsf{K} \subset \mathbb{R}^d$  be a convex body, i.e. a compact convex set with non-empty interior and  $\iota_{\mathsf{K}} : \mathbb{R}^d \rightarrow \{0, +\infty\}$  be the (convex) indicator function of  $\mathsf{K}$ , defined for  $x \in \mathbb{R}^d$  by,

$$\iota_{\mathsf{K}}(x) = \begin{cases} +\infty & \text{if } x \notin \mathsf{K}, \\ 0 & \text{if } x \in \mathsf{K}. \end{cases}$$

Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ . In this paper we consider any probability density  $\pi$  associated to a potential  $U : \mathbb{R}^d \rightarrow (-\infty, +\infty]$  of the form

$$U = f + \iota_{\mathsf{K}} , \tag{2}$$

and assume that the function  $f$  and the convex body  $\mathsf{K}$  satisfy the following assumptions. For  $x \in \mathbb{R}^d$  and  $r > 0$ , denote by  $\mathsf{B}(x, r)$  the closed ball of center  $x$  and radius  $r$ :  $\mathsf{B}(x, r) = \{y \in \mathbb{R}^d : \|y - x\| \leq r\}$ .

**H1.** (i)  $f$  is convex.

(ii)  $f$  is continuously differentiable on  $\mathbb{R}^d$  and gradient Lipschitz with Lipschitz constant  $L_f$ , i.e. for all  $x, y \in \mathbb{R}^d$

$$\|\nabla f(x) - \nabla f(y)\| \leq L_f \|x - y\| . \tag{3}$$

**H2.** There exist  $r, R > 0$ ,  $r \leq R$ , such that,

$$\mathsf{B}(0, r) \subset \mathsf{K} \subset \mathsf{B}(0, R) .$$

To apply ULA, [DMP16] suggested to carefully regularize  $U$  in such a way that 1) the convexity of  $U$  is preserved (this property is key to the theoretical analysis of the algorithm), 2) the regularisation of  $U$  is continuously differentiable and gradient Lipschitz (this regularity property is key to the algorithm's stability), and 3) the resulting approximation is close to  $\pi$  (e.g. in total variation norm). The

tool used to construct such an approximation is the Moreau-Yosida envelope of  $\iota_{\mathbf{K}}$ ,  $\iota_{\mathbf{K}}^\lambda : \mathbb{R}^d \rightarrow \mathbb{R}_+$  defined for  $x \in \mathbb{R}^d$  (see e.g. [RW98, Chapter 1 Section G]) by,

$$\iota_{\mathbf{K}}^\lambda(x) = \inf_{y \in \mathbb{R}^d} (\iota_{\mathbf{K}}(y) + (2\lambda)^{-1} \|x - y\|^2) = (2\lambda)^{-1} \|x - \text{proj}_{\mathbf{K}}(x)\|^2, \quad (4)$$

where  $\lambda > 0$  is a regularization parameter and  $\text{proj}_{\mathbf{K}}$  is the projection onto  $\mathbf{K}$ . By [RW98, Example 10.32, Theorem 9.18], the function  $\iota_{\mathbf{K}}^\lambda$  is convex and continuously differentiable with gradient given for all  $x \in \mathbb{R}^d$  by:

$$\nabla \iota_{\mathbf{K}}^\lambda(x) = \lambda^{-1}(x - \text{proj}_{\mathbf{K}}(x)). \quad (5)$$

Moreover, [RW98, Proposition 12.19] implies that  $\iota_{\mathbf{K}}^\lambda$  is  $\lambda^{-1}$ -gradient Lipschitz: for all  $x, y \in \mathbb{R}^d$ ,

$$\|\nabla \iota_{\mathbf{K}}^\lambda(x) - \nabla \iota_{\mathbf{K}}^\lambda(y)\| \leq \lambda^{-1} \|x - y\|. \quad (6)$$

Adding  $f$  to  $\iota_{\mathbf{K}}^\lambda$  under **H1** leads to the regularization  $U^\lambda : \mathbb{R}^d \rightarrow \mathbb{R}$  of the potential  $U$  defined for all  $x \in \mathbb{R}^d$  by

$$U^\lambda(x) = f(x) + \iota_{\mathbf{K}}^\lambda(x). \quad (7)$$

The following lemma shows that the probability measure  $\pi^\lambda$  on  $\mathbb{R}^d$ , with density with respect to the Lebesgue measure, also denoted by  $\pi^\lambda$  and given for all  $x \in \mathbb{R}^d$  by

$$\pi^\lambda(x) = \frac{e^{-U^\lambda(x)}}{\int_{\mathbb{R}^d} e^{-U^\lambda(s)} ds}, \quad (8)$$

is well defined. It also shows that  $U^\lambda$  has a minimizer  $x^* \in \mathbb{R}^d$ , a fact that will be used in Section 4.

**Lemma 1.** *Assume **H1-(i)** and **H2**. For all  $\lambda > 0$ ,*

- a)  $U^\lambda$  has a minimizer  $x^* \in \mathbb{R}^d$ , i.e. for all  $x \in \mathbb{R}^d$ ,  $U^\lambda(x) \geq U^\lambda(x^*)$ .
- b)  $e^{-U^\lambda}$  defines a proper density of a probability measure on  $\mathbb{R}^d$ , i.e.

$$0 < \int_{\mathbb{R}^d} e^{-U^\lambda(y)} dy < +\infty.$$

*Proof.* Note that [DMP16, Proposition 1] provides a proof in a more general case. Given the specific form of  $U^\lambda$ , a short and self-contained proof can be found in Section 6.1.  $\square$

Under **H1**, for all  $\lambda > 0$ ,  $\pi^\lambda$  is log-concave and  $U^\lambda$  is continuously differentiable by (5), with  $\nabla U^\lambda$  given for all  $x \in \mathbb{R}^d$  by

$$\nabla U^\lambda(x) = -\nabla \log \pi^\lambda(x) = \nabla f(x) + \lambda^{-1}(x - \text{proj}_K(x)). \quad (9)$$

In addition, by (6),  $\nabla U^\lambda$  is Lipschitz with constant  $L \leq L_f + \lambda^{-1}$ . Since  $U^\lambda$  is continuously differentiable, ULA is well defined. The algorithm proposed in [DMP16] then proceeds by using the Euler-Maruyama discretization of the Langevin equation associated with  $U^\lambda$ , with  $\pi^\lambda$  as proxy, to generate approximate samples from  $\pi$ . Precisely, it uses the Markov chain  $(X_k)_{k \in \mathbb{N}}$ , started at  $X_0$ , given for all  $k \in \mathbb{N}$  by

$$X_{k+1} = (1 - \frac{\gamma}{\lambda})X_k - \gamma \nabla f(X_k) + \frac{\gamma}{\lambda} \text{proj}_K(X_k) + \sqrt{2\gamma}Z_{k+1}, \quad (10)$$

where  $(Z_k)_{k \in \mathbb{N}}$  is a sequence of i.i.d. standard Gaussian  $d$ -dimensional vectors and  $\gamma > 0$  is the stepsize. Note that one iteration (10) requires a projection onto the convex body  $K$  and the evaluation of  $\nabla f$ . The kernel of the homogeneous Markov chain defined by (10) is given for  $x \in \mathbb{R}^d$  and  $A \in \mathcal{B}(\mathbb{R}^d)$  by,

$$R_\gamma(x, A) = (4\pi\gamma)^{-d/2} \int_A \exp\left(- (4\gamma)^{-1} \|y - x + \gamma \nabla U^\lambda(x)\|^2\right) dy, \quad (11)$$

where  $U^\lambda$  is defined in (7). Since the target density for the Markov chain (10) is the regularized measure  $\pi^\lambda$  and not  $\pi$ , the algorithm is named the Moreau-Yosida regularized Unadjusted Langevin Algorithm (MYULA).

## 2.2 Context and contributions

The total variation distance between two probability measures  $\mu$  and  $\nu$  is defined by  $\|\mu - \nu\|_{\text{TV}} = 2 \sup_{A \in \mathcal{B}(\mathbb{R}^d)} |\mu(A) - \nu(A)|$ . Let  $\phi, \psi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ . Denote by  $\phi = \tilde{\mathcal{O}}(\psi)$  or  $\phi = \tilde{\Omega}(\psi)$  if there exist  $C, c \geq 0$  such that for all  $t \in \mathbb{R}_+$   $\phi(t) \leq C\psi(t)(\log t)^c$  or  $\phi(t) \geq C\psi(t)(\log t)^c$  respectively. Our main result is the following:

**Theorem 2.** *Assume **H1** and **H2**. For all  $\varepsilon > 0$  and  $x \in \mathbb{R}^d$ , there exist  $\lambda > 0$  and  $\gamma \in (0, \lambda(1 + L_f^2\lambda^2)^{-1})$  such that,*

$$\|\delta_x R_\gamma^n - \pi\|_{\text{TV}} \leq \varepsilon \quad \text{for } n = \tilde{\Omega}(d^5),$$

where  $R_\gamma$  is defined in (11).

The proof of Theorem 2 follows from combining Proposition 6 and Proposition 4 below. Note that these two results imply explicit bounds between  $R_\gamma^n$  and  $\pi$  for all  $n \in \mathbb{N}$  and  $\gamma > 0$ .

The problem of sampling from a probability measure restricted to a convex compact support has been investigated in several works, mainly in the fields of

theoretical computer science and Bayesian statistics. In computer science, a line of works starting with [DF91] has studied the convergence of the ball walk and the hit-and-run algorithm towards the uniform density on a convex body  $\mathbf{K}$ , or more generally to a log-concave density. The best complexity result is achieved by [LV07, Theorem 2.1] who establishes a mixing time for these two algorithms of order  $\tilde{O}(d^4)$ . However, observe that contrary to Theorem 2, this result assumes that  $\pi$  is in near-isotropic position, i.e. there exists  $C \in \mathbb{R}_+^*$  such that for all  $u \in \mathbb{R}^d$ ,  $\|u\| = 1$ ,

$$C^{-1} \leq \int_{\mathbb{R}^d} \langle u, x \rangle^2 \pi(dx) \leq C . \quad (12)$$

Note that [LV07, Section 2.5] gives also an algorithm of complexity  $\tilde{O}(d^5)$  which provides an invertible linear map  $T$  of  $\mathbb{R}^d$  such that the measure  $\pi_T$  defined for all  $\mathbf{A} \in \mathcal{B}(\mathbb{R}^d)$  by

$$\pi_T(\mathbf{A}) = \pi(T^{-1}(\mathbf{A})) ,$$

is log-concave and near-isotropic. Also note that, unlike our method, each iteration of the ball walk or the hit-and-run algorithm requires a call to a zero-order oracle, which given  $x \in \mathbb{R}^d$ , returns the value  $U(x)$ . MYULA does not require to fulfill the condition (12) and is thus dispensed of preprocessing step. However, MYULA needs a first-order oracle which returns the value  $\nabla f(x)$  for  $x \in \mathbb{R}^d$ .

As emphasized in the introduction, probability distributions with convex compact supports or more generally with constrained parameters arise naturally in Bayesian statistics. [GSL92] includes many examples of such problems and suggests to use a Gibbs sampler, see also [RDS04]. [CSI12, Chapter 6] addresses the subject with the additional difficulty of computing normalizing constants. Recently, [PP14] adapted the Hamiltonian Monte Carlo method to sample from a truncated multivariate gaussian, and [LS15] suggested a new approach which consists in mapping the constrained domain to a sphere in an augmented space. However, these methods are not well understood from a theoretical viewpoint, and do not provide any theoretical guarantees for the estimations delivered.

Concerning the ULA algorithm, when  $U$  is continuously differentiable, the first explicit convergence bounds have been obtained by [Dal16], [DM15], [DM16]. In the constrained case  $U = f + \iota_{\mathbf{K}}$ , [BEL15] suggests a projection step in ULA i.e. to consider the Markov chain  $(\tilde{X}_k)_{k \geq 0}$ , defined for all  $k \in \mathbb{N}$  by

$$\tilde{X}_{k+1} = \text{proj}_{\mathbf{K}} \left( \tilde{X}_k - \gamma \nabla U(\tilde{X}_k) + \sqrt{2\gamma} Z_{k+1} \right) . \quad (13)$$

with  $\tilde{X}_0 = 0$ . This method is referred to as the Projected Langevin Monte Carlo (PLMC) algorithm. As in MYULA, one iteration of PLMC requires a projection onto  $\mathbf{K}$  and an evaluation of  $\nabla f$ . Let  $\tilde{R}_\gamma$  be the Markov kernel defined by (13). [BEL15] proved that for all  $\varepsilon > 0$ ,  $\|\delta_0 \tilde{R}_\gamma^n - \pi\|_{\text{TV}} \leq \varepsilon$  for  $n = \tilde{\Omega}(d^7)$  if  $\pi$  is the

uniform density on  $\mathbf{K}$  and  $n = \tilde{\Omega}(d^{12})$  if  $\pi$  is a log-concave density. Theorem 2 improves these bounds for the MYULA algorithm. Note however that the iterations of PLMC stay within the constraint set  $\mathbf{K}$  and this property can be useful in some specific problems. Nevertheless, there is a wide range of settings where this property is not particularly beneficial, for example in the case of the computation of volumes discussed in Section 5, or in Bayesian model selection where it is necessary to estimate marginal likelihoods.

### 3 Distance between $\pi$ and $\pi^\lambda$

In this section, we derive bounds between  $\pi$  and  $\pi^\lambda$  in total variation and in Wasserstein distance (recall that  $\pi$  is associated with a potential of the form (2) and  $\pi^\lambda$  is given by (8)). It is shown that the approximation error in both distances can be made arbitrarily small by adjusting the regularisation parameter  $\lambda$ .

The main quantity of interest to analyze the distance between  $\pi$  and  $\pi^\lambda$  will appear to be the integral of  $x \mapsto e^{-(2\lambda)^{-1}\|x - \text{proj}_{\mathbf{K}}(x)\|^2}$  over  $\mathbb{R}^d$ . This constant is linked to useful notions borrowed from the field of convex geometry [Kam09, Proposition 3]. Indeed, Fubini's theorem gives the following equality:

$$\begin{aligned} \int_{\mathbb{R}^d} e^{-(2\lambda)^{-1}\|x - \text{proj}_{\mathbf{K}}(x)\|^2} dx &= \int_{\mathbb{R}_+} \int_{\mathbb{R}^d} \mathbb{1}_{[\|x - \text{proj}_{\mathbf{K}}(x)\|, +\infty)}(t) \lambda^{-1} t e^{-t^2/(2\lambda)} dx dt, \\ &= \int_{\mathbb{R}_+} \text{Vol}(\mathbf{K} + \mathbf{B}(0, t)) \lambda^{-1} t e^{-t^2/(2\lambda)} dt, \end{aligned} \quad (14)$$

where  $\mathbf{A} + \mathbf{B}$  is the Minkowski sum of  $\mathbf{A}, \mathbf{B} \subset \mathbb{R}^d$ , i.e.  $\mathbf{A} + \mathbf{B} = \{x + y : x \in \mathbf{A}, y \in \mathbf{B}\}$ , and we have used in the last line that for all  $t \in \mathbb{R}_+$ ,  $\mathbf{K} + \mathbf{B}(0, t) = \{x \in \mathbb{R}^d : \|x - \text{proj}_{\mathbf{K}}(x)\| \leq t\}$ . It turns out that  $t \mapsto \text{Vol}(\mathbf{K} + \mathbf{B}(0, t))$  on  $\mathbb{R}_+$  is a polynomial. More precisely, Steiner's formula states that for all  $t \geq 0$ ,

$$\text{Vol}(\mathbf{K} + \mathbf{B}(0, t)) = \sum_{i=0}^d t^i \kappa_i \mathcal{V}_{d-i}(\mathbf{K}), \quad (15)$$

where  $\{\mathcal{V}_i(\mathbf{K})\}_{0 \leq i \leq d}$  are the intrinsic volumes of  $\mathbf{K}$ ,  $\kappa_i$  denotes the volume of the unit ball in  $\mathbb{R}^i$ , i.e.

$$\kappa_i = \pi^{i/2} / \Gamma(1 + i/2), \quad (16)$$

and  $\Gamma : \mathbb{R}_+^* \rightarrow \mathbb{R}_+^*$  is the Gamma function. We refer to [Sch13, Chapter 4.2] for this result and an introduction to this topic. Combining (14) and (15) gives:

$$\int_{\mathbb{R}^d} e^{-(2\lambda)^{-1}\|x - \text{proj}_{\mathbf{K}}(x)\|^2} dx = \sum_{i=0}^d \mathcal{V}_i(\mathbf{K}) (2\pi\lambda)^{(d-i)/2}. \quad (17)$$



This expression will provide a precise analysis of the distance in total variation and Wasserstein distance between  $\pi$  and  $\pi^\lambda$ , in particular when  $\pi$  is the uniform density on  $\mathsf{K}$ . However, in more general cases, an additional assumption on the relation between  $f$  and  $\mathsf{K}$  is necessary to bound the distance between  $\pi$  and  $\pi^\lambda$ . Under **H1-(i)** and **H2**,  $f$  has a minimum  $x_{\mathsf{K}}$  on  $\mathsf{K}$ . Define

$$\tilde{\mathsf{K}} = \{x \in \mathsf{K} \mid \mathsf{B}(x, r) \subset \mathsf{K}\} . \quad (18)$$

$\tilde{\mathsf{K}}$  has the following property.

**Lemma 3.** *Assume **H2**.  $\tilde{\mathsf{K}}$  is a non-empty convex compact set.*

*Proof.* The proof is postponed to Section 6.2. □

**H3.** (i) *There exists  $\Delta_1 > 0$  such that  $\exp(\inf_{\mathsf{K}^c}(f) - \max_{\mathsf{K}}(f)) \geq \Delta_1$ .*

(ii) *There exists  $\Delta_2 \geq 0$  such that  $0 \leq f(\text{proj}_{\tilde{\mathsf{K}}}(x_{\mathsf{K}})) - f(x_{\mathsf{K}}) \leq \Delta_2$ .*

Under **H3-(i)**, the application of Steiner's formula is possible and reveals the precise dependence of the bounds with respect to the intrinsic volumes of  $\mathsf{K}$ . A complementary view is possible under **H3-(ii)**. The obtained bounds are less precise regarding  $\mathsf{K}$  but more robust with respect to  $f$ . Note that if  $x_{\mathsf{K}} \in \tilde{\mathsf{K}}$ ,  $\Delta_2$  can be chosen equal to 0. On the other hand, if  $f$  is assumed to be  $\ell$ -Lipschitz inside  $\mathsf{K}$ ,  $\Delta_2$  is less than  $\ell R$ .

**Proposition 4.** *Assume **H1-(i)** and **H2**.*

a) *Assume **H3-(i)**. For all  $\lambda > 0$ ,*

$$\|\pi^\lambda - \pi\|_{\text{TV}} \leq 2 \left(1 + \Delta_1 \mathsf{D}(\mathsf{K}, \lambda)^{-1}\right)^{-1} , \quad (19)$$

where,

$$\mathsf{D}(\mathsf{K}, \lambda) = (\text{Vol } \mathsf{K})^{-1} \sum_{i=0}^{d-1} (2\pi\lambda)^{(d-i)/2} \mathcal{V}_i(\mathsf{K}) , \quad (20)$$

and  $\mathcal{V}_i(\mathsf{K})$  are defined in (15).

b) *In addition, assuming **H3-(i)**, for all  $\lambda \in (0, (2\pi)^{-1}(r/d)^2)$ ,*

$$\|\pi^\lambda - \pi\|_{\text{TV}} \leq 2^{3/2} \Delta_1^{-1} (\pi\lambda)^{1/2} dr^{-1} . \quad (21)$$

c) *Assume **H3-(ii)**. For all  $\lambda \in (0, 16^{-1}(r/d)^2]$ ,*

$$\|\pi^\lambda - \pi\|_{\text{TV}} \leq (4/r) \exp\left(4\lambda (\Delta_2/r)^2\right) \left\{ \sqrt{\lambda}(d + \Delta_2) + (2\lambda\Delta_2)/r \right\} . \quad (22)$$

*Proof.* The proof is postponed to Section 6.3.  $\square$

In the particular case where  $f = 0$  and  $\pi$  is the uniform density on  $\mathbf{K}$ ,  $\Delta_1$  equals 1 and the inequality (19) is in fact an equality. The dependence of the upper bound in (19) w.r.t. to  $\lambda, d, r$  is sharp. Indeed, for the cube  $\mathbf{C}$  of side  $c$ ,  $D(\mathbf{C}, \lambda)$  can be explicitly computed. [KR97, Theorem 4.2.1] gives for  $i \in \{0, \dots, d\}$ ,  $\mathcal{V}_i(\mathbf{C}) = \binom{d}{i} c^i$ , which implies:

$$D(\mathbf{C}, \lambda) = \left(1 + c^{-1} \sqrt{2\pi\lambda}\right)^d - 1 ,$$

$$\|\pi^\lambda - \pi\|_{\text{TV}} = 2 \left\{ 1 - \left(1 + c^{-1} \sqrt{2\pi\lambda}\right)^{-d} \right\} , \text{ for } U = \iota_{\mathbf{C}} .$$

For two probability measures  $\mu$  and  $\nu$  on  $\mathcal{B}(\mathbb{R}^d)$ , the Wasserstein distance of order  $p \in \mathbb{N}^*$  between  $\mu$  and  $\nu$  is defined by

$$W_p(\mu, \nu) = \left( \inf_{\zeta \in \Pi(\mu, \nu)} \int_{\mathbb{R}^d \times \mathbb{R}^d} \|x - y\|^p d\zeta(x, y) \right)^{1/p} ,$$

where  $\Pi(\mu, \nu)$  is the set of transference plans of  $\mu$  and  $\nu$ .  $\zeta$  is a transference plan of  $\mu$  and  $\nu$  if it is a probability measure on  $(\mathbb{R}^d \times \mathbb{R}^d, \mathcal{B}(\mathbb{R}^d \times \mathbb{R}^d))$  such that for all  $\mathbf{A} \in \mathcal{B}(\mathbb{R}^d)$ ,  $\zeta(\mathbf{A} \times \mathbb{R}^d) = \mu(\mathbf{A})$  and  $\zeta(\mathbb{R}^d \times \mathbf{A}) = \nu(\mathbf{A})$ .

**Proposition 5.** Assume **H1-(i)** and **H2**.

a) Assume **H3-(i)**. For all  $\lambda > 0$ ,

$$W_1(\pi, \pi^\lambda) \leq \Delta_1^{-1} \mathbf{E}(\mathbf{K}, \lambda, R) ,$$

where

$$\mathbf{E}(\mathbf{K}, \lambda, R) = (\text{Vol}(\mathbf{K}))^{-1} \sum_{i=0}^{d-1} \mathcal{V}_i(\mathbf{K}) (2\pi\lambda)^{(d-i)/2} \left\{ 2R + [\lambda(d-i+2)]^{1/2} \right\} ,$$

and  $\mathcal{V}_i(\mathbf{K})$  are defined in (15).

b) In addition, assuming **H3-(i)**, for all  $\lambda \in (0, (2\pi)^{-1} d^{-2} r^2)$ ,

$$W_1(\pi, \pi^\lambda) \leq \Delta_1^{-1} (2\pi\lambda)^{1/2} d r^{-1} \left( 2R + r (3/(2d\pi))^{1/2} \right) .$$

c) Assume **H3-(ii)**. For all  $\lambda \in (0, 16^{-1} (r/d)^2]$ ,

$$W_1(\pi, \pi^\lambda) \leq 4 \exp(4\lambda (\Delta_2/r)^2) \left\{ \sqrt{\lambda} (d + \Delta_2) (R/r) + (2\lambda \Delta_2 R)/r^2 + \sqrt{\pi\lambda} \right\} .$$

*Proof.* The proof is postponed to Section 6.4.  $\square$

Note that the bounds in Wasserstein distance between  $\pi$  and  $\pi^\lambda$  are roughly similar to those obtained in total variation norm.

## 4 Convergence analysis of MYULA

We now analyse the convergence of the Markov kernel  $R_\gamma$ , given by (11), to the target density  $\pi^\lambda$  defined in (8). For  $x \in \mathbb{R}^d$  and  $n \in \mathbb{N}$ , explicit bounds in total variation norm and in Wasserstein distance between  $\delta_x R_\gamma^n$  and  $\pi^\lambda$  are provided in Proposition 6 and Proposition 7. Because of the regularisation procedure performed in Section 2.1, the convergence analysis of MYULA (10) is an application of results of [DM15] and [DM16].

### 4.1 Convergence in total variation norm

Define  $\omega : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  for all  $r \geq 0$  by

$$\omega(r) = r^2 / \{2\Phi^{-1}(3/4)\}^2, \quad (23)$$

where  $\Phi(x) = (2\pi)^{-1/2} \int_{-\infty}^x e^{-t^2/2} dt$ .

**Proposition 6.** *Assume **H1** and **H2**. Let  $\lambda > 0$ ,  $L$  be the Lipschitz constant of  $\nabla U^\lambda$  defined in (7) and  $\bar{\gamma} \in (0, \lambda^{-1}L^{-2})$ . Then for all  $\varepsilon > 0$  and  $x \in \mathbb{R}^d$ , we get:*

$$\|\delta_x R_\gamma^n - \pi^\lambda\|_{\text{TV}} \leq \varepsilon, \quad (24)$$

provided that  $n > T\gamma^{-1}$  with

$$T = (\log\{A_2(x)\} - \log(\varepsilon/2)) / (-\log(\kappa)), \quad (25a)$$

$$\gamma \leq \frac{-d + \sqrt{d^2 + (2/3)A_1(x)\varepsilon^2(L^2T)^{-1}}}{2A_1(x)/3} \wedge \bar{\gamma}, \quad (25b)$$

where

$$A_1(x) = L^2 \left( \|x - x^*\|^2 + 2(d + 8\lambda^{-1}R^2)e^{\gamma(\lambda^{-1} - \bar{\gamma}L^2)}(\lambda^{-1} - \bar{\gamma}L^2)^{-1} \right),$$

$$\log(\kappa) = -\log(2)(4\lambda)^{-1} \left[ \log \left\{ \left( 1 + e^{(8\lambda)^{-1}\omega\{\max(1,4R)\}} \right) (1 + \max(1,4R)) \right\} + \log(2) \right]^{-1},$$

$$A_2(x) = 6 + 2^{3/2} (d\lambda + 8R^2)^{1/2} + 2(A_1(x)/L^2)^{1/2},$$

and  $x^*$  is a minimizer of  $U^\lambda$ .

*Proof.* To apply [DM15, Theorem 21], it is sufficient to check the assumption [DM15, H3], i.e. there exist  $\tilde{R} \geq 0$  and  $m > 0$  such that for all  $x, y \in \mathbb{R}^d$ ,  $\|x - y\| \geq \tilde{R}$ ,

$$\langle \nabla U^\lambda(x) - \nabla U^\lambda(y), x - y \rangle \geq m \|x - y\|^2. \quad (26)$$

|   |                               |  |                               |                                |
|---|-------------------------------|--|-------------------------------|--------------------------------|
| Upper bound on $n$ to get $\ \delta_{x^*}R_\gamma^n - \pi\ _{\text{TV}} \leq \varepsilon$ | $d \rightarrow +\infty$       | $\varepsilon \rightarrow 0$              | $R \rightarrow +\infty$       | $r \rightarrow 0$              |
| Proposition 4 and Proposition 6   | $\tilde{\mathcal{O}}(d^5)$    | $\tilde{\mathcal{O}}(\varepsilon^{-6})$  | $\tilde{\mathcal{O}}(R^4)$    | $\tilde{\mathcal{O}}(r^{-4})$  |
| [BEL15, Theorem 1] $\pi$ uniform on $\mathsf{K}$  | $\tilde{\mathcal{O}}(d^7)$    | $\tilde{\mathcal{O}}(\varepsilon^{-8})$  | $\tilde{\mathcal{O}}(R^6)$    | $\tilde{\mathcal{O}}(r^{-6})$  |
| [BEL15, Theorem 1] $\pi$ log concave  | $\tilde{\mathcal{O}}(d^{12})$ | $\tilde{\mathcal{O}}(\varepsilon^{-12})$ | $\tilde{\mathcal{O}}(R^{18})$ | $\tilde{\mathcal{O}}(r^{-18})$ |

Table 1: dependency of  $n$  on  $d, \varepsilon, R$  and  $r$  to get  $\|\delta_{x^*}R_\gamma^n - \pi\|_{\text{TV}} \leq \varepsilon$

|   |                                      |                                   |
|---|--------------------------------------|-----------------------------------|
| Upper bound on $n$ to get $\ \delta_{x^*}R_\gamma^n - \pi\ _{\text{TV}} \leq \varepsilon$ | $\Delta_1 \rightarrow 0$             | $\Delta_2 \rightarrow +\infty$    |
| Proposition 4 and Proposition 6   | $\tilde{\mathcal{O}}(\Delta_1^{-4})$ | $\tilde{\mathcal{O}}(\Delta_2^4)$ |

Table 2: dependency of  $n$  on  $\Delta_1$  and  $\Delta_2$  to get  $\|\delta_{x^*}R_\gamma^n - \pi\|_{\text{TV}} \leq \varepsilon$

By (5) and the Cauchy-Schwarz inequality, we have:

$$\langle \nabla \iota_{\mathsf{K}}^\lambda(x) - \nabla \iota_{\mathsf{K}}^\lambda(y), x - y \rangle \geq \lambda^{-1} \left( \|x - y\|^2 - 2 \left\{ \sup_{z \in \mathsf{K}} \|z\| \right\} \|x - y\| \right),$$

which implies under **H1-(i)** and **H2** that (26) holds for  $\tilde{R} = 4R$  and  $m = (2\lambda)^{-1}$ .  $\square$

Combining Proposition 4 and Proposition 6 determines the stepsize  $\gamma$  and the number of samples  $n$  to get  $\|\delta_{x^*}R_\gamma^n - \pi\|_{\text{TV}} \leq \varepsilon$ .  $\lambda$  is chosen of order  $\varepsilon^2 r^2 d^{-2} \Delta_1^2$  under **H3-(i)** and  $\varepsilon^2 r^2 \min(d^{-2}, \Delta_2^{-2})$  under **H3-(ii)**. The orders of magnitude of  $n$  in  $d, \varepsilon, R, r$  are reported in Table 1, along with the results of [BEL15]. The dependency of  $n$  towards  $\Delta_1, \Delta_2$  is presented in Table 2. A detailed table is provided in Appendix A.

## 4.2 Convergence in Wasserstein distance for strongly convex $f$

In this section,  $f$  is assumed to satisfy an additional assumption.

**H4.**  $f : \mathbb{R}^d \mapsto \mathbb{R}$  is  $m$ -strongly convex, i.e. there exists  $m > 0$  such that for all  $x, y \in \mathbb{R}^d$ ,

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + (m/2) \|x - y\|^2. \quad (27)$$

Note that under **H4**,  $U^\lambda$  defined in (7) is  $m$ -strongly convex as well. The following Proposition 7 relies on the convergence analysis in Wasserstein distance

|   |                            |   |                            |                               |
|---|----------------------------|---|----------------------------|-------------------------------|
| Upper bound on $n$ to get $W_1(\delta_{x^*}R_\gamma^n, \pi) \leq \varepsilon$ | $d \rightarrow +\infty$    | $\varepsilon \rightarrow 0$             | $R \rightarrow +\infty$    | $r \rightarrow 0$             |
| Proposition 5-c) and Proposition 7  | $\tilde{\mathcal{O}}(d^5)$ | $\tilde{\mathcal{O}}(\varepsilon^{-6})$ | $\tilde{\mathcal{O}}(R^4)$ | $\tilde{\mathcal{O}}(r^{-4})$ |

Table 3: dependency of  $n$  on  $d, \varepsilon, R$  and  $r$  to get  $W_1(\delta_{x^*}R_\gamma^n, \pi) \leq \varepsilon$

|   |                                      |                                   |
|---|--------------------------------------|-----------------------------------|
| Upper bound on $n$ to get $W_1(\delta_{x^*}R_\gamma^n, \pi) \leq \varepsilon$ | $\Delta_1 \rightarrow 0$             | $\Delta_2 \rightarrow +\infty$    |
| Proposition 5-c) and Proposition 7  | $\tilde{\mathcal{O}}(\Delta_1^{-4})$ | $\tilde{\mathcal{O}}(\Delta_2^4)$ |

Table 4: dependency of  $n$  on  $\Delta_1$  and  $\Delta_2$  to get  $W_1(\delta_{x^*}R_\gamma^n, \pi) \leq \varepsilon$

done in [DM16], which assumes that  $f$  is strongly convex. It may be possible to extend the range of validity of these results but this work goes beyond the scope of this paper.

**Proposition 7.** Assume **H1** and **H4**. Let  $\lambda > 0$ ,  $L$  be the Lipschitz constant of  $\nabla U^\lambda$  defined in (7) and  $\kappa = (2mL)(m+L)^{-1}$ . Let  $\varepsilon > 0$  and  $x \in \mathbb{R}^d$ . We have,

$$W_2(\delta_x R_\gamma^n, \pi^\lambda) \leq \varepsilon ,$$

provided that,

$$\gamma \leq \frac{m}{L^2} \left\{ -\frac{13}{12} + \left[ \left( \frac{13}{12} \right)^2 + \frac{\varepsilon^2 \kappa^2}{8md} \right]^{1/2} \right\} \wedge \frac{1}{m+L} ,$$

$$n \geq 2(\kappa\gamma)^{-1} \left\{ -\log(\varepsilon^2/4) + \log(\|x - x^*\|^2 + d/m) \right\} .$$

*Proof.* The proof is postponed to Section 6.5. □

Combining Proposition 5 and Proposition 7 determines the stepsize  $\gamma$  and the number of samples  $n$  to get  $W_1(\delta_{x^*}R_\gamma^n, \pi) \leq \varepsilon$ .  $\lambda$  is chosen of order  $\varepsilon^2 \Delta_1^2 r^2 d^{-2} R^{-2}$  under **H3-(i)** and  $\varepsilon^2 r^2 R^{-2} \min(d^{-2}, \Delta_2^{-2})$  under **H3-(ii)**. The orders of magnitude of  $n$  in  $d, \varepsilon, R, r, \Delta_1, \Delta_2$  are reported in Tables 3 and 4.

## 5 Numerical experiments

In this section we illustrate MYULA with the following three numerical experiments: computation of the volume of a high-dimensional convex set, sampling from a truncated multivariate Gaussian distribution, and Bayesian inference with

the constrained LASSO model. We benchmark our results with model-specific specialised algorithms, namely the hit-and-run algorithm [LV06] for set volume computation, the wall HMC (WHMC) [PP14] for truncated Gaussian models, and the auxiliary-variable Gibbs sampler for the Bayesian lasso model [PC08]. Where relevant we also compare with the Random Walk Metropolis Hastings (RWM) algorithm.

First we consider the computation of the volume of a high-dimensional hypercube. In a manner akin to [CV15], to apply MYULA to this problem we use an annealing strategy involving truncated Gaussian distributions whose variance is gradually increased at each step  $i \in \mathbb{N}$  of the annealing process. Precisely, for  $M \in \mathbb{N}^*$  and  $i \in \{0, \dots, M-1\}$ , the potential  $U_i$  (2) of the phase  $i$  is given for all  $x \in \mathbb{R}^d$  by,  $U_i(x) = (2\sigma_i^2)^{-1} \|x\|^2 + \iota_{\mathbf{K}}$  where  $\mathbf{K} = [-1, 1]^d$ . Observing that,

$$\frac{\int_{\mathbb{R}^d} e^{-U_{i+1}(x)} dx}{\int_{\mathbb{R}^d} e^{-U_i(x)} dx} = \pi_i(g_i) \quad , \quad g_i(x) = e^{2^{-1}(\sigma_i^{-2} - \sigma_{i+1}^{-2})\|x\|^2} \quad , \quad (28)$$

where  $\pi_i$  is the probability measure associated with  $U_i$ , the volume of  $\mathbf{K}$  is

$$\text{Vol}(\mathbf{K}) = \prod_{i=0}^{M-1} \pi_i(g_i) \int_{\mathbb{R}^d} e^{-U_0(x)} \quad ,$$

where  $U_M = \iota_{\mathbf{K}}$ . To use MYULA we consider for all  $i \in \{0, \dots, M-1\}$  the potential  $U_i^{\lambda_i}$  defined for all  $x \in \mathbb{R}^d$  by  $U_i^{\lambda_i}(x) = (2\sigma_i^2)^{-1} \|x\|^2 + \iota_{\mathbf{K}}^{\lambda_i}$  where  $\iota_{\mathbf{K}}^{\lambda_i}$  is given by (4). We choose the step-size  $\gamma_i$  proportional to  $1/\{d \max(d, \sigma_i^{-1})\}$  and the regularization parameter  $\lambda_i$  is set equal to  $2\gamma_i$ . The counterpart of (28) is then

$$\frac{\int_{\mathbb{R}^d} e^{-U_{i+1}^{\lambda_{i+1}}(x)} dx}{\int_{\mathbb{R}^d} e^{-U_i^{\lambda_i}(x)} dx} = \pi_i^{\lambda_i}(g_i^{\lambda_i}) \quad , \quad g_i^{\lambda_i}(x) = e^{2^{-1}(\sigma_i^{-2} - \sigma_{i+1}^{-2})\|x\|^2 + \iota_{\mathbf{K}}^{\lambda_i} - \iota_{\mathbf{K}}^{\lambda_{i+1}}} \quad ,$$

where  $\pi_i^{\lambda_i}$  is the probability measure associated with  $U_i^{\lambda_i}$ , and the volume of  $\mathbf{K}$  is

$$\text{Vol}(\mathbf{K}) = \prod_{i=0}^{M-1} \pi_i^{\lambda_i}(g_i^{\lambda_i}) \int_{\mathbb{R}^d} e^{-U_0^{\lambda_0}(x)} \quad ,$$

where  $U_M^{\lambda_M} = U_M = \iota_{\mathbf{K}}$ .

Figure 1 shows the volume estimates (over 10 experiments) obtained with MYULA and the hit-and-run algorithm for a unit hypercube of dimension  $d$  ranging from  $d = 10$  to  $d = 90$  (to simplify visual comparison the estimates are normalised w.r.t. the true volume). Observe that the estimates of MYULA are in agreement with the results of the hit-and-run algorithm, which serves as a benchmark for this problem. The outputs of both algorithms are at similar distances with respect to the true value 1.

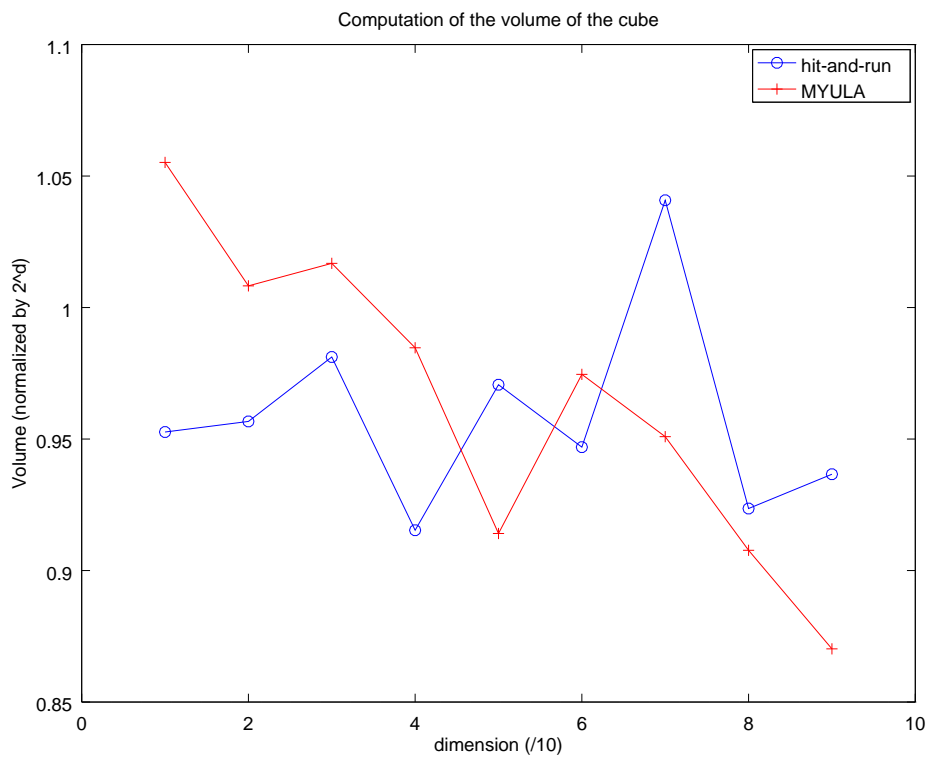


Figure 1: Computation of the volume of the cube with MYULA and hit-and-run algorithm.

Moreover, the second experiment we consider is the simulation from a  $d$ -dimensional truncated Gaussian distribution restricted on a convex set  $\mathbf{K}_d$ , with mode zero at the boundary of the set, and covariance matrix  $\Sigma$  with  $(i, j)$ th element given by  $(\Sigma)_{i,j} = 1/(1 + |i - j|)$ . Let  $\boldsymbol{\beta} \in \mathbb{R}^d$ . The potential  $U$ , given by (2) and associated with the density  $\pi(\boldsymbol{\beta})$ , is given by  $U(\boldsymbol{\beta}) = (1/2) \langle \boldsymbol{\beta}, \Sigma^{-1} \boldsymbol{\beta} \rangle + \iota_{\mathbf{K}_d}(\boldsymbol{\beta})$ . We consider three scenarios of increasing dimension:  $d = 2$  with  $\mathbf{K}_2 = [0, 5] \times [0, 1]$ ,  $d = 10$  with  $\mathbf{K}_{10} = [0, 5] \times [0, 0.5]^9$ , and  $d = 100$  with  $\mathbf{K}_{100} = [0, 5] \times [0, 0.5]^{99}$ . We generate  $10^6$  samples for MYULA,  $10^5$  samples for WHMC, and  $10^6$  samples for RWM (in all cases the initial 10% is discarded as burn-in period). Regarding algorithm parameters, we set  $\gamma = 1/1000$  and  $\lambda = 2\gamma$  for MYULA, and adjust the parameters of RWM and WHMC such that their acceptance rates are approximately 25% and 70%.

Table 5 shows the results obtained with each method for the model  $d = 2$ , and by performing 100 repetitions to obtain 95% confidence intervals. For this model we also report a solution by a cubature integration [NJ16] which provides a ground truth. Moreover, Figure 2 and Figure 3 show the results for the first three coordinates of  $\boldsymbol{\beta}$  (i.e.,  $\beta_1, \beta_2, \beta_3$ ) for  $d = 10$  and  $d = 100$  respectively. Observe the good performance of MYULA as dimensionality increases, particularly in the challenging case  $d = 100$  where it performs comparably to the specialised algorithm WHMC.

| Method | Mean   | Covariance  |
|--------|--|---|
| Truth  | $\begin{bmatrix} 0.790 \\ 0.488 \end{bmatrix}$                     | $\begin{bmatrix} 0.326 & 0.017 \\ 0.017 & 0.080 \end{bmatrix}$  |
| RWM    | $\begin{bmatrix} 0.791 \pm 0.013 \\ 0.486 \pm 0.002 \end{bmatrix}$ | $\begin{bmatrix} 0.330 \pm 0.011 & 0.017 \pm 0.002 \\ 0.017 \pm 0.002 & 0.080 \pm 0.0003 \end{bmatrix}$ |
| WHMC   | $\begin{bmatrix} 0.789 \pm 0.005 \\ 0.490 \pm 0.005 \end{bmatrix}$ | $\begin{bmatrix} 0.324 \pm 0.008 & 0.017 \pm 0.002 \\ 0.017 \pm 0.002 & 0.079 \pm 0.0007 \end{bmatrix}$ |
| MYULA  | $\begin{bmatrix} 0.758 \pm 0.052 \\ 0.484 \pm 0.016 \end{bmatrix}$ | $\begin{bmatrix} 0.309 \pm 0.038 & 0.017 \pm 0.009 \\ 0.017 \pm 0.009 & 0.088 \pm 0.002 \end{bmatrix}$  |

Table 5: Mean and covariance of  $\boldsymbol{\beta}$  in dimension 2 obtained by RWM, WHMC and MYULA.

Finally, we also report an experiment involving the analysis of a real dataset with an  $\ell_1$ -norm constrained Bayesian LASSO model (i.e. least squares regression subject to an  $\ell_1$ -ball constraint). Precisely, the observations  $Y = \{Y_1, \dots, Y_n\} \in$



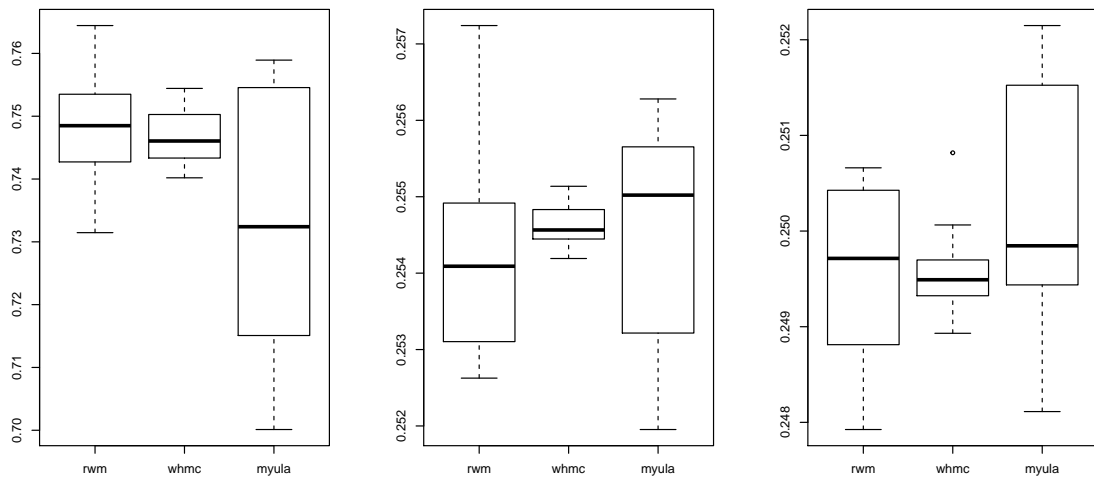


Figure 2: Boxplots of  $\beta_1, \beta_2, \beta_3$  for the truncated Gaussian variable in dimension 10.

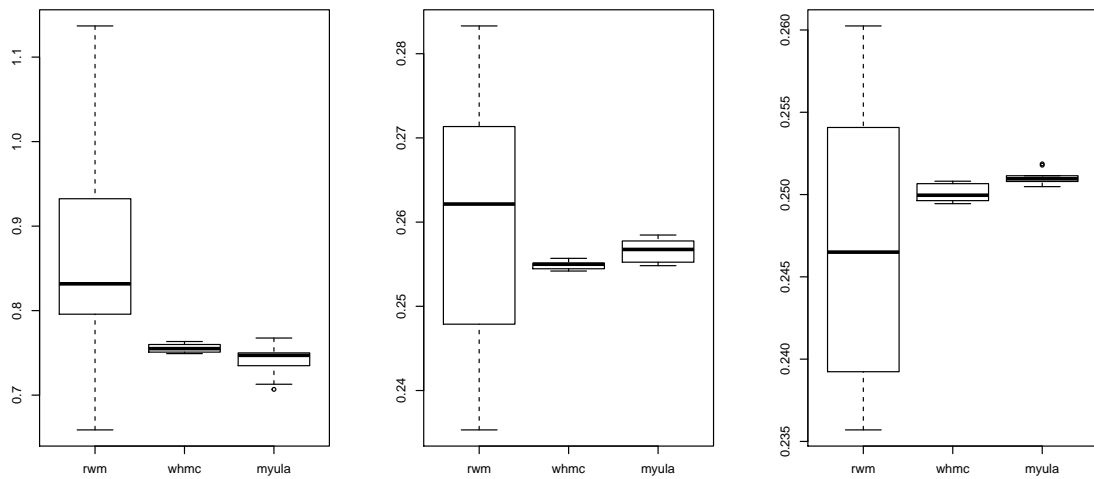


Figure 3: Boxplots of  $\beta_1, \beta_2, \beta_3$  for the truncated Gaussian variable in dimension 100.

$\mathbb{R}^n$ , for  $n \geq 1$ , are assumed to be distributed from the Gaussian distribution with mean  $X\boldsymbol{\beta}$  and covariance matrix  $\sigma^2 \mathbf{I}_n$ , where  $X \in \mathbb{R}^{n \times d}$  is the design matrix,  $\boldsymbol{\beta} \in \mathbb{R}^d$  is the regression parameter,  $\sigma^2 > 0$  and  $\mathbf{I}_n$  is the identity matrix of dimension  $n$ . The prior on  $\boldsymbol{\beta}$  is the uniform distribution over the  $\ell_1$  ball,  $B_o(0, s) = \{\boldsymbol{\beta} \in \mathbb{R}^d \mid \|\boldsymbol{\beta}\|_1 \leq s\}$ , for  $s > 0$ , where  $\|\boldsymbol{\beta}\|_1 = \sum_{i=1}^d |\beta_i|$ ,  $\beta_i$  is the  $i$ -th component of  $\boldsymbol{\beta}$ . The potential  $U^s$ , for  $s > 0$ , associated with the posterior distribution is given for all  $\boldsymbol{\beta} \in \mathbb{R}^d$  by  $U^s(\boldsymbol{\beta}) = \|Y - X\boldsymbol{\beta}\|^2 + \iota_{B_o(0,s)}(\boldsymbol{\beta})$ . We consider in our experiment the diabetes data set<sup>1</sup>, which consists in  $n = 442$  observations and  $d = 10$  explanatory variables.

Figure 4 shows the ‘‘LASSO paths’’ obtained using MYULA, the WHMC algorithm, and with the specialised Gibbs sampler of [PC08] (these paths are the posterior marginal medians associated with  $\pi^s$  for  $s = t \|\boldsymbol{\beta}^{\text{OLS}}\|_1$ ,  $t \in [0, 1]$ , and where  $\boldsymbol{\beta}^{\text{OLS}}$  is the estimate obtained by the ordinary least square regression). The dot lines represent the confidence interval at level 95%, obtained by performing 100 repetitions. MYULA estimates were obtained by using  $10^5$  samples (with the initial  $10^4$  samples discarded as burn-in period) and stepsize  $s^{3/2} \times 10^{-5}$ . WHMC estimates were obtained by using  $10^4$  samples (with the initial  $10^3$  samples discarded as burn-in period), and by adjusting parameters to achieve an acceptance rate of approximately 90%. Finally, the Gibbs sampler is targeting an unconstrained LASSO model with prior  $\boldsymbol{\beta} \mapsto (2s)^{-d} e^{-\|\boldsymbol{\beta}\|_1/s}$ , for  $s > 0$ .

## 6 Proofs

### 6.1 Proof of Lemma 1

Since  $f$  is a (proper) convex function, there exist  $a \in \mathbb{R}$ ,  $b \in \mathbb{R}^d$  such that  $f(x) \geq a + \langle b, x \rangle$  [Roc15, Theorem 23.4]. By H2 and a straightforward calculation, for  $\|x\| \geq R + 4\lambda \|b\| + 2 \{\lambda(|a| + R \|b\|)\}^{1/2}$ , we have,

$$U^\lambda(x) \geq (4\lambda)^{-1} (\|x\| - R)^2,$$

which concludes the proof.

### 6.2 Proof of Lemma 3

Under H2,  $0 \in \tilde{\mathcal{K}}$ . Let  $x_1, x_2 \in \tilde{\mathcal{K}}$  and  $t \in [0, 1]$ . We have by definition of  $\tilde{\mathcal{K}}$  (18) that  $B(tx_1 + (1-t)x_2, r) \subset tB(x_1, r) + (1-t)B(x_2, r) \subset \mathcal{K}$ , which implies that  $\tilde{\mathcal{K}}$  is convex.

<sup>1</sup><http://archive.ics.uci.edu/ml/datasets/Pima+Indians+Diabetes>

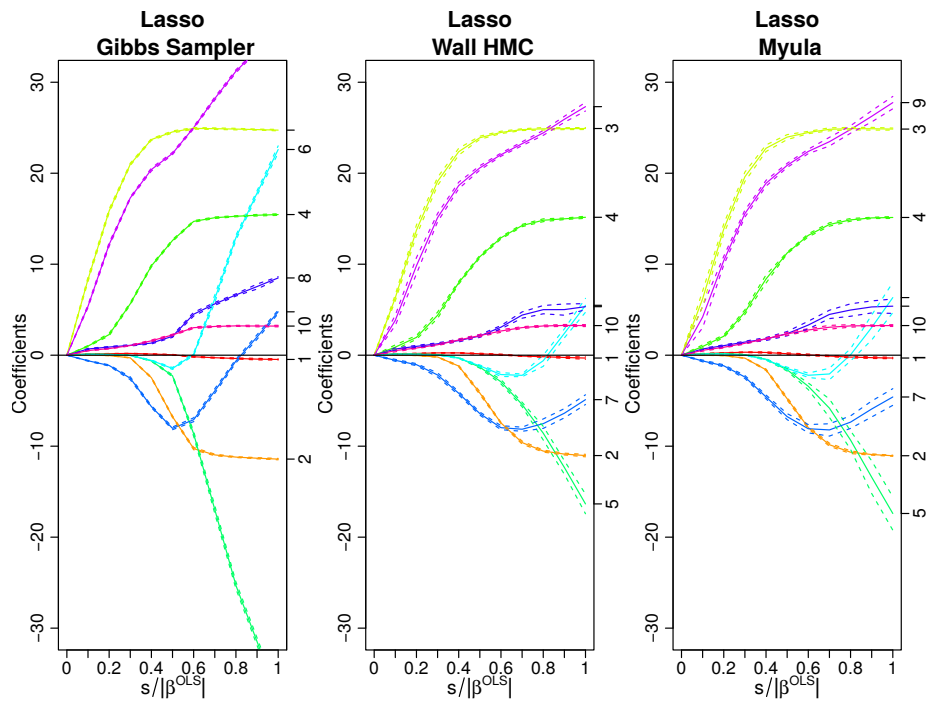


Figure 4: Lasso path for the Gibbs sampler, Wall HMC and MYULA algorithms.

To show that  $\tilde{\mathbf{K}}$  is close, it is enough to show that  $\tilde{\mathbf{K}} = \{x \in \mathbf{K} \mid \text{dist}(x, \mathbf{K}^c) \geq r\}$  where  $\text{dist}(x, \mathbf{K}^c) = \inf_{y \in \mathbf{K}^c} \|x - y\|$  since  $x \mapsto \text{dist}(x, \mathbf{K}^c)$  is Lipschitz continuous. First by definition, we have  $\tilde{\mathbf{K}} \subset \{x \in \mathbf{K} \mid \text{dist}(x, \mathbf{K}^c) \geq r\}$ . To show the converse, let  $x \in \{y \in \mathbf{K} \mid \text{dist}(y, \mathbf{K}^c) \geq r\}$ . Then,  $B_o(x, r) \subset \mathbf{K}$ , where  $B_o(x, r) = \{y \in \mathbb{R}^d \mid \|y - x\| < r\}$ , which yields  $B(x, r) \subset \mathbf{K}$  since  $\mathbf{K}$  is assumed to be close. This result then concludes the proof by definition of  $\tilde{\mathbf{K}}$ .

### 6.3 Proof of Proposition 4

a) By a direct calculation, we have:

$$\|\pi^\lambda - \pi\|_{\text{TV}} = \int_{\mathbb{R}^d} |\pi(x) - \pi^\lambda(x)| dx = 2 \left( 1 + \left\{ \int_{\mathbf{K}^c} e^{-U^\lambda(x)} dx \right\}^{-1} \int_{\mathbf{K}} e^{-f(x)} dx \right)^{-1} \quad (29)$$

$$\leq 2 \left( 1 + \exp \left( \min_{\mathbf{K}^c}(f) - \max_{\mathbf{K}}(f) \right) A \right)^{-1} . \quad (30)$$

where

$$A = \text{Vol}(\mathbf{K}) / \int_{\mathbf{K}^c} e^{-(2\lambda)^{-1} \|x - \text{proj}_{\mathbf{K}}(x)\|^2} dx . \quad (31)$$

The conclusion follows then from (17) and **H3-(i)**.

b) We give two proofs for this result, which both consist in lower bounding  $A$ . The obtained bounds are identical up to an universal constant. The first one is simpler and was suggested by a referee. The second one is more involved ; however, it has the benefit of establishing the relation between the intrinsic volumes of  $\mathbf{K}$  and the bound on the total variation norm.

Under **H2**, we have  $\mathbf{K} + B(0, t) \subset (1 + t/r)\mathbf{K}$  and using (14),

$$\begin{aligned} \int_{\mathbf{K}^c} e^{-(1/2\lambda) \|x - \text{proj}_{\mathbf{K}}(x)\|^2} dx &\leq \left\{ \int_{\mathbb{R}_+} \text{Vol}(\mathbf{K}(1 + t/r)) \lambda^{-1} t e^{-t^2/(2\lambda)} dt - \text{Vol}(\mathbf{K}) \right\} \\ &= \text{Vol}(\mathbf{K}) \left\{ \int_{\mathbb{R}_+} (1 + t/r)^d \lambda^{-1} t e^{-t^2/(2\lambda)} dt - 1 \right\} \\ &= \text{Vol}(\mathbf{K}) \sum_{i=1}^d \binom{d}{i} \left( \frac{\sqrt{2\lambda}}{r} \right)^i \Gamma(1 + i/2) \\ &\leq \text{Vol}(\mathbf{K}) \sum_{i=1}^d \left( \frac{\sqrt{2\lambda} d}{r} \right)^i , \end{aligned}$$

where the second equality follows from developping  $(1 + t/r)^d$ , making the change of variable  $t \mapsto t^2/(2\lambda)$  and using the Gamma function and the last inequality from  $\binom{d}{i}\Gamma(1 + i/2) \leq d^i$  for  $i \in \{1, \dots, d\}$ . For  $\lambda \in (0, r^2 d^{-2}/8]$ , we get

$$A^{-1} \leq \sum_{i=1}^d \left( \frac{\sqrt{2\lambda d}}{r} \right)^i \leq \frac{2\sqrt{2\lambda d}}{r}.$$

Combining it with (30) and **H3-(i)** concludes the proof.

For the second proof, it is necessary to introduce first a generalized notion of the intrinsic volumes (15), the mixed volumes. Let  $\mathcal{K}$  be the class of convex bodies of  $\mathbb{R}^d$ ,  $K_1, \dots, K_m \in \mathcal{K}$  and  $\lambda_1, \dots, \lambda_m \geq 0$ . By [Sch13, Theorem 5.1.7], there is a nonnegative symmetric function  $\mathcal{V} : (\mathcal{K})^d \rightarrow \mathbb{R}_+$ , the mixed volume, such that,

$$\text{Vol}(\lambda_1 K_1 + \dots + \lambda_m K_m) = \sum_{i_1, \dots, i_d=1}^m \lambda_{i_1} \dots \lambda_{i_d} \mathcal{V}(K_{i_1}, \dots, K_{i_d}). \quad (32)$$

Let  $m > 1$ ,  $a_1, \dots, a_m \geq 0$  and  $K_1, \dots, K_m, L$  be  $(m + 1)$  convex bodies in  $\mathbb{R}^d$  such that  $K_1 \subset L$ . By unicity of the coefficients of the polynomial in  $\lambda_1, \dots, \lambda_m$  (32) and [Sch13, p.282], we have:

$$\mathcal{V}(a_1 K_1, \dots, a_m K_m) = \left( \prod_{i=1}^m a_i \right) \mathcal{V}(K_1, \dots, K_m), \quad (33)$$

$$\mathcal{V}(K_1, K_2, \dots, K_m) \leq \mathcal{V}(L, K_2, \dots, K_m). \quad (34)$$

Denote by  $B$  the unity ball of  $\mathbb{R}^d$ ,  $B = B(0, 1)$ . Taking  $m = 2, K_1 = K, K_2 = B, \lambda_1 = 1, \lambda_2 = t$  in (32), we get:

$$\text{Vol}(K + B(0, t)) = \sum_{i=0}^d t^i \binom{d}{i} \mathcal{V}(K[d-i], B[i]), \quad (35)$$

where for a set  $A \subset \mathbb{R}^d$ , the notation  $A[i]$  means  $A$  repeated  $i$  times:  $A[i] = A, \dots, A$   $i$  times. The quermassintegrals of  $K$  are defined for  $i \in \{0, \dots, d\}$  by  $\mathcal{W}_i(K) = \mathcal{V}(K[d-i], B[i])$  [Sch13, equation 5.31]. We get then by (35) and (15),

$$\binom{d}{i} \mathcal{W}_i(K) = \kappa_i \mathcal{V}_{d-i}(K), \quad (36)$$

where  $\kappa_i$  is given by (16).

The proof consists then in identifying an upper bound on  $\mathcal{V}_i(K)(\text{Vol} K)^{-1}$  for  $i \in \{0, \dots, d\}$ . First, the sequence  $\{i! \mathcal{W}_i(K)\}_{0 \leq i \leq d}$  is shown to be log-concave, i.e. for  $i \in \{1, \dots, d-1\}$

$$(i! \mathcal{W}_i(K))^2 \geq (i+1)! \mathcal{W}_{i+1}(K) (i-1)! \mathcal{W}_{i-1}(K). \quad (37)$$

The Aleksandrov-Fenchel inequality [Sch13, equation 7.66] states, for  $i \in \{1, \dots, d-1\}$ ,

$$\mathcal{W}_i(\mathbf{K})^2 \geq \mathcal{W}_{i-1}(\mathbf{K})\mathcal{W}_{i+1}(\mathbf{K}) . \quad (38)$$

By (16),  $\kappa_i/\kappa_{i-2} = (2\pi)/i$  and the log convexity of the gamma function, we get for  $i \in \{1, \dots, d-1\}$ :

$$\frac{1}{i+1} \frac{\kappa_i}{\kappa_{i+1}} = \frac{1}{i} \frac{\kappa_{i-2}}{\kappa_{i-1}} \leq \frac{1}{i} \frac{\kappa_{i-1}}{\kappa_i} . \quad (39)$$

Combining (39), (38) and (36) shows (37).

The log-concavity of  $\{i! \mathcal{V}_i(\mathbf{K})\}_{0 \leq i \leq d}$  gives for  $i \in \{0, \dots, d-1\}$ ,

$$\frac{\mathcal{V}_i(\mathbf{K})}{\mathcal{V}_{i+1}(\mathbf{K})} \leq \frac{\mathcal{V}_{d-1}(\mathbf{K})}{\text{Vol}(\mathbf{K})} = \frac{d \mathcal{W}_1(\mathbf{K})}{2 \mathcal{W}_0(\mathbf{K})} . \quad (40)$$

Combining the definition of the quermassintegrals, (33), (34) and H2 give:

$$r \mathcal{W}_1(\mathbf{K}) = \mathcal{V}(\mathbf{K}, \dots, \mathbf{K}, \mathbf{B}(0, r)) \leq \mathcal{V}(\mathbf{K}, \dots, \mathbf{K}, \mathbf{K}) = \mathcal{W}_0(\mathbf{K}) . \quad (41)$$

By (41) and (40), we get:

$$\mathbf{D}(\mathbf{K}, \lambda) \leq \sum_{i=1}^d \left\{ dr^{-1}(\pi\lambda/2)^{1/2} \right\}^i , \quad (42)$$

where  $\mathbf{D}(\mathbf{K}, \lambda)$  is defined in (20). For all  $\lambda \in (0, 2\pi^{-1}(r/d)^2)$ , (19) gives then,

$$\|\pi^\lambda - \pi\|_{\text{TV}} \leq 2 \left\{ 1 + \exp \left( \min_{\mathbf{K}^c}(f) - \max_{\mathbf{K}}(f) \right) \left( \left\{ dr^{-1}(\pi\lambda/2)^{1/2} \right\}^{-1} - 1 \right) \right\}^{-1} .$$

Using that for all  $a, b \in \mathbb{R}_+$ ,  $b \geq 2$ ,  $(1 + a(b-1))^{-1} \leq b^{-1}/(b^{-1} + a/2)$  and H3-(i), we get for  $\lambda \in (0, 2\pi^{-1}(r/d)^2)$

$$\|\pi^\lambda - \pi\|_{\text{TV}} \leq 2^{3/2}(\pi\lambda)^{1/2} dr^{-1} \left\{ (2\pi\lambda)^{1/2} dr^{-1} + \Delta_1 \right\}^{-1} .$$

c) The proof consists in using (29) to bound  $\|\pi^\lambda - \pi\|_{\text{TV}}$ . In the first step we give an upper bound on  $\int_{\mathbb{R}^d} e^{-U^\lambda(x)} dx / \int_{\mathbf{K}} e^{-f(x)} dx$ . By Fubini's theorem, similarly to (14) we have

$$\int_{\mathbb{R}^d} e^{-U^\lambda(x)} dx \leq \int_{\mathbb{R}_+} \int_{\mathbf{K} + \mathbf{B}(0, t)} e^{-f(x)} \lambda^{-1} t e^{-t^2/(2\lambda)} dx dt . \quad (43)$$

Let  $t \geq 0$ . By definition of  $\tilde{\mathbb{K}}$ , using Lemma 3 and  $\mathbb{K} - \text{proj}_{\tilde{\mathbb{K}}}(x_{\mathbb{K}}) + \text{B}(0, t) \subset (1 + t/r)(\mathbb{K} - \text{proj}_{\tilde{\mathbb{K}}}(x_{\mathbb{K}}))$ , we have

$$\begin{aligned} \int_{\mathbb{K} + \text{B}(0, t)} e^{-f(x)} dx &= \int_{\mathbb{K} - \text{proj}_{\tilde{\mathbb{K}}}(x_{\mathbb{K}}) + \text{B}(0, t)} e^{-f(x + \text{proj}_{\tilde{\mathbb{K}}}(x_{\mathbb{K}}))} dx \\ &\leq \int_{(1+t/r)(\mathbb{K} - \text{proj}_{\tilde{\mathbb{K}}}(x_{\mathbb{K}}))} e^{-f(x + \text{proj}_{\tilde{\mathbb{K}}}(x_{\mathbb{K}}))} dx \\ &= (1 + t/r)^d \int_{\mathbb{K} - \text{proj}_{\tilde{\mathbb{K}}}(x_{\mathbb{K}})} e^{-f((1+t/r)x + \text{proj}_{\tilde{\mathbb{K}}}(x_{\mathbb{K}}))} dx. \end{aligned} \quad (44)$$

By **H1-(i)**  $f$  is convex and therefore for all  $x \in \mathbb{K} - \text{proj}_{\tilde{\mathbb{K}}}(x_{\mathbb{K}})$ ,

$$\begin{aligned} f((1 + t/r)x + \text{proj}_{\tilde{\mathbb{K}}}(x_{\mathbb{K}})) &\geq (t/r) \{f(x + \text{proj}_{\tilde{\mathbb{K}}}(x_{\mathbb{K}})) - f(\text{proj}_{\tilde{\mathbb{K}}}(x_{\mathbb{K}}))\} + f(x + \text{proj}_{\tilde{\mathbb{K}}}(x_{\mathbb{K}})) \\ &\geq -(\Delta_2 t)/r + f(x + \text{proj}_{\tilde{\mathbb{K}}}(x_{\mathbb{K}})). \end{aligned}$$

Combining it with (43) and (44), we get

$$\int_{\mathbb{R}^d} e^{-U^\lambda(x)} dx \leq \left( \int_{\mathbb{K}} e^{-f(x)} dx \right) \int_{\mathbb{R}_+} (1 + t/r)^d e^{(\Delta_2 t)/r} \lambda^{-1} t e^{-t^2/(2\lambda)} dt. \quad (45)$$

We now bound  $B = \int_{\mathbb{K}^c} e^{-U^\lambda(x)} dx / \int_{\mathbb{K}} e^{-f(x)} dx$ . Using (45) and an integration by parts, we have

$$\begin{aligned} B &\leq \int_{\mathbb{R}_+} \{(1 + t/r)^d e^{(\Delta_2 t)/r} - 1\} \lambda^{-1} t e^{-t^2/(2\lambda)} dt \\ &\leq \int_{\mathbb{R}_+} (1 + t/r)^{d-1} e^{(\Delta_2 t)/r} r^{-1} (d + \Delta_2 + (\Delta_2 t)/r) e^{-t^2/(2\lambda)} dt. \end{aligned}$$

Since for all  $t \geq 0$ ,  $(\Delta_2 t)/r - t^2/(2\lambda) \leq -t^2/(4\lambda) + 4\lambda(\Delta_2/r)^2$ , it holds

$$B \leq \frac{1}{r} \exp\left(4\lambda \left(\frac{\Delta_2}{r}\right)^2\right) \int_{\mathbb{R}_+} (1 + t/r)^{d-1} (d + \Delta_2 + (\Delta_2 t)/r) e^{-t^2/(4\lambda)} dt.$$

By developing  $(1 + t/r)^{d-1}$ , using the change of variable  $t \mapsto t^2/(4\lambda)$  and the definition of the Gamma function, we have

$$B \leq \frac{2\lambda}{r} \exp\left(4\lambda \left(\frac{\Delta_2}{r}\right)^2\right) \sum_{i=0}^{d-1} \binom{d-1}{i} \left(\frac{2\sqrt{\lambda}}{r}\right)^i \left\{ \frac{d + \Delta_2}{2\sqrt{\lambda}} \Gamma\left(\frac{1+i}{2}\right) + \frac{\Delta_2}{r} \Gamma\left(1 + \frac{i}{2}\right) \right\}.$$

Using that for all  $i \in \{0, \dots, d-1\}$ ,  $\binom{d-1}{i} \Gamma(1 + i/2) \leq d^i$ , we get for  $\lambda \in (0, 16^{-1} r^2 d^{-2}]$

$$B \leq \frac{2}{r} \exp\left(4\lambda \left(\frac{\Delta_2}{r}\right)^2\right) \left\{ \sqrt{\lambda}(d + \Delta_2) + \frac{2\lambda\Delta_2}{r} \right\},$$

which combined with (29) concludes the proof.

## 6.4 Proof of Proposition 5

a) The proof relies on a control of the Wasserstein distance by a weighted total variation. The arguments are similar to those of Proposition 4. [Vil09, Theorem 6.15] implies:

$$W_1(\pi, \pi^\lambda) \leq \int_{\mathbb{R}^d} \|x\| |\pi(x) - \pi^\lambda(x)| dx = C + D, \quad (46)$$

where

$$C = \int_{\mathbb{K}^c} \|x\| \pi^\lambda(x) dx, \quad D = \left\{ 1 - \frac{\int_{\mathbb{K}} e^{-f}}{\int_{\mathbb{R}^d} e^{-U^\lambda}} \right\} \int_{\mathbb{K}} \|x\| \pi(x) dx. \quad (47)$$

We bound these two terms separately. First using the same decomposition as in (14),  $\|x\| \leq R + \|x - \text{proj}_{\mathbb{K}}(x)\|$  and that for all  $t \in \mathbb{R}_+$ ,  $\mathbb{K} + \text{B}(0, t) = \{x \in \mathbb{R}^d : \|x - \text{proj}_{\mathbb{K}}(x)\| \leq t\}$ , we get

$$C = \left( \int_{\mathbb{R}^d} e^{-U^\lambda} \right)^{-1} \int_0^{+\infty} \int_{\mathbb{K}^c} e^{-f(x)} \|x\| t \lambda^{-1} e^{-t^2/(2\lambda)} \mathbb{1}_{[\|x - \text{proj}_{\mathbb{K}}(x)\|, +\infty)}(t) dx dt \quad (48)$$

$$\leq e^{\max_{\mathbb{K}}(f) - \min_{\mathbb{K}^c}(f)} \int_0^{+\infty} (R + t) t \lambda^{-1} e^{-t^2/(2\lambda)} \left( \frac{\text{Vol}(\mathbb{K} + \text{B}(0, t)) - \text{Vol}(\mathbb{K})}{\text{Vol}(\mathbb{K})} \right) dt. \quad (49)$$

Combining (15)-(49), **H3-(i)** and using  $\mathcal{V}_d(\mathbb{K}) = \text{Vol}(\mathbb{K})$  give

$$C \leq \Delta_1^{-1} \sum_{i=0}^{d-1} \kappa_{d-i} \frac{\mathcal{V}_i(\mathbb{K})}{\text{Vol}(\mathbb{K})} \int_0^{+\infty} (R t^{d-i+1} + t^{d-i+2}) \lambda^{-1} e^{-t^2/(2\lambda)} dt. \quad (50)$$

Using (16), for all  $k \geq 0$ ,  $\int_{\mathbb{R}_+} t^k e^{t^2/(2\lambda)} dt = (2\lambda)^{(k+1)/2} \Gamma((k+1)/2)$  and for all  $a > 1$ ,  $\Gamma(a+1/2) \leq a^{1/2} \Gamma(a)$  (by log-convexity of the Gamma function), we have

$$C \leq \Delta_1^{-1} \sum_{i=0}^{d-1} \frac{\mathcal{V}_i(\mathbb{K})}{\text{Vol}(\mathbb{K})} (2\pi\lambda)^{(d-i)/2} \left\{ R + [\lambda(d-i+2)]^{1/2} \right\}. \quad (51)$$

Regarding  $D$  defined in (47), by **H2**, **H3-(i)**, (30) and (17), we get:

$$D \leq R \Delta_1^{-1} \text{D}(\mathbb{K}, \lambda), \quad (52)$$

where  $\text{D}(\mathbb{K}, \lambda)$  is defined in (20). Combining (51) and (52) in (46) concludes the proof.



b) Using (40) and (41) in (51) gives for all  $\lambda \in (0, (2\pi)^{-1}r^2d^{-2})$

$$\begin{aligned} C &\leq \Delta_1^{-1} \sum_{i=0}^{d-1} \left( \frac{d}{r} \sqrt{\frac{\pi\lambda}{2}} \right)^{d-i} \left\{ R + [\lambda(d-i+2)]^{1/2} \right\} \\ &\leq \Delta_1^{-1} (2\pi\lambda)^{1/2} dr^{-1} \left( R + r \left( \frac{3}{2d\pi} \right)^{1/2} \right). \end{aligned}$$

Finally this bound, (52), (42) and (46) conclude the proof.

c) The proof still relies on the decomposition (46), where  $C$  and  $D$  are defined in (47). Eq. (48) gives

$$C \leq \int_0^{+\infty} (R+t)t\lambda^{-1}e^{-t^2/(2\lambda)} \left( \frac{\int_{\mathbf{K}+\mathbf{B}(0,t)} e^{-f(x)} dx}{\int_{\mathbf{K}} e^{-f(x)} dx} - 1 \right) dt.$$

Under **H 3-(ii)**, following the steps of Section 6.3-c) to upper bound the term  $\int_{\mathbf{K}+\mathbf{B}(0,t)} e^{-f(x)} dx / \int_{\mathbf{K}} e^{-f(x)} dx$ , we have

$$\begin{aligned} C &\leq \int_0^{+\infty} (R+t)t\lambda^{-1}e^{-t^2/(2\lambda)} \left( (1+t/r)^d e^{(t\Delta_2)/r} - 1 \right) dt \\ &= C_1 + C_2, \end{aligned}$$

where

$$\begin{aligned} C_1 &= R \int_0^{+\infty} t\lambda^{-1}e^{-t^2/(2\lambda)} \left( (1+t/r)^d e^{(t\Delta_2)/r} - 1 \right) dt, \\ C_2 &= \int_0^{+\infty} t^2\lambda^{-1}e^{-t^2/(2\lambda)} \left( (1+t/r)^d e^{(t\Delta_2)/r} - 1 \right) dt. \end{aligned}$$

$C_1$  is upper bounded in the same way as  $B$  in Section 6.3-c). Regarding  $C_2$ , since for all  $t \geq 0$ ,  $(\Delta_2 t)/r - t^2/(2\lambda) \leq -t^2/(4\lambda) + 4\lambda(\Delta_2/r)^2$ , developing  $(1+t/r)^d$  and using the change of variable  $t \mapsto t^2/(4\lambda)$  we get

$$\begin{aligned} C_2 &\leq e^{4\lambda(\Delta_2/r)^2} \sum_{i=0}^d \binom{d}{i} r^{-i} \int_{\mathbb{R}_+} t^{i+2} \lambda^{-1} e^{-t^2/(4\lambda)} dt \\ &\leq 4\sqrt{\lambda} e^{4\lambda(\Delta_2/r)^2} \frac{\sqrt{\pi}}{2} \sum_{i=0}^d \binom{d}{i} \left( \frac{2\sqrt{\lambda}}{r} \right)^i \Gamma\left(\frac{3}{2} + \frac{i}{2}\right). \end{aligned}$$

Using  $\binom{d}{i} \Gamma((3+i)/2) \leq (\sqrt{\pi}/2)d^i$  for  $i \in \{0, \dots, d\}$ , we have for  $\lambda \in (0, 16^{-1}r^2d^{-2}]$ ,

$$\begin{aligned} C_2 &\leq 2\sqrt{\pi\lambda}e^{4\lambda(\Delta_2/r)^2} \sum_{i=0}^d \left( \frac{2\sqrt{\lambda}d}{r} \right)^i \\ &\leq 4\sqrt{\pi\lambda}e^{4\lambda(\Delta_2/r)^2} . \end{aligned}$$

$D$  defined in (47) is upper bounded by  $RB$  where  $B$  is defined in Section 6.3-c). Combining the bounds on  $C_1, C_2, D$  gives the result.

## 6.5 Proof of Proposition 7

Assume that  $\gamma \in (0, (m+L)^{-1})$ . [DM16, Theorem 5] gives for all  $n \in \mathbb{N}^*$ :

$$W_2^2(\delta_x R_\gamma^n, \pi^\lambda) \leq 2(1 - (\kappa\gamma)/2)^n \{ \|x - x^*\|^2 + d/m \} + u(\gamma) ,$$

where,

$$u(\gamma) = 2\kappa^{-1}L^2d\gamma(\kappa^{-1} + \gamma) \left( 2 + \frac{L^2\gamma}{m} + \frac{L^2\gamma^2}{6} \right) .$$

Noting that  $\kappa\gamma \leq 1$  and  $L^2\gamma^2 \leq 1$ , it is then sufficient for  $\gamma, n$  to satisfy,

$$\begin{aligned} 4\kappa^{-2}L^2d\gamma \left( 2 + \frac{1}{6} + \frac{L^2\gamma}{m} \right) &\leq \varepsilon^2/2 , \\ 2(1 - (\kappa\gamma)/2)^n \{ \|x - x^*\|^2 + d/m \} &\leq \varepsilon^2/2 , \end{aligned}$$

which concludes the proof.

## Acknowledgments

The authors wish to express their thanks to the anonymous referees for several helpful remarks, in particular concerning a simplified proof of Proposition 4.

## References

- [BEL15] Sebastien Bubeck, Ronen Eldan, and Joseph Lehec. “Finite-time Analysis of Projected Langevin Monte Carlo”. In: *Proceedings of the 28th International Conference on Neural Information Processing Systems*. NIPS’15. Montreal, Canada: MIT Press, 2015, pp. 1243–1251. URL: <http://dl.acm.org/citation.cfm?id=2969239.2969378>.

- [BNJ03] David M Blei, Andrew Y Ng, and Michael I Jordan. “Latent dirichlet allocation”. In: *Journal of machine Learning research* 3.Jan (2003), pp. 993–1022.
- [Cel+12] Gilles Celeux et al. “Regularization in Regression: Comparing Bayesian and Frequentist Methods in a Poorly Informative Situation”. In: *Bayesian Anal.* 7.2 (June 2012), pp. 477–502. DOI: [10.1214/12-BA716](https://doi.org/10.1214/12-BA716). URL: <http://dx.doi.org/10.1214/12-BA716>.
- [CSI12] Ming-Hui Chen, Qi-Man Shao, and Joseph G Ibrahim. *Monte Carlo methods in Bayesian computation*. Springer Science & Business Media, 2012.
- [CV15] Ben Cousins and Santosh Vempala. *Computation of the volume of convex bodies*. June 2015. URL: <http://fr.mathworks.com/matlabcentral/fileexchange/43596-volume-computation-of-convex-bodies>.
- [Dal16] Arnak S Dalalyan. “Theoretical guarantees for approximate sampling from smooth and log-concave densities”. In: *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* (2016).
- [DF91] Martin Dyer and Alan Frieze. “Computing the volume of convex bodies: a case where randomness provably helps”. In: *Probabilistic combinatorics and its applications* 44 (1991), pp. 123–170.
- [DM15] A. Durmus and E. Moulines. “Non-asymptotic convergence analysis for the Unadjusted Langevin Algorithm”. In: *ArXiv e-prints* (July 2015). arXiv: [1507.05021](https://arxiv.org/abs/1507.05021) [[math.ST](#)].
- [DM16] A. Durmus and E. Moulines. “High-dimensional Bayesian inference via the Unadjusted Langevin Algorithm”. In: *ArXiv e-prints* (May 2016). arXiv: [1605.01559](https://arxiv.org/abs/1605.01559) [[math.ST](#)].
- [DMP16] A. Durmus, E. Moulines, and M. Pereyra. “Efficient Bayesian computation by proximal Markov chain Monte Carlo: when Langevin meets Moreau”. In: *ArXiv e-prints* (Dec. 2016). arXiv: [1612.07471](https://arxiv.org/abs/1612.07471) [[stat.CO](#)].
- [GSL92] A. E. Gelfand, A. F. Smith, and T.-M. Lee. “Bayesian analysis of constrained parameter and truncated data problems using Gibbs sampling”. In: *Journal of the American Statistical Association* 87.418 (1992), pp. 523–532.
- [JA06] Valen E Johnson and James H Albert. *Ordinal data modeling*. Springer Science & Business Media, 2006.
- [Kam09] Jürgen Kampf. “On weighted parallel volumes”. In: *Beiträge Algebra Geom* 50.2 (2009), pp. 495–519.

- [KM05] John P Klein and Melvin L Moeschberger. *Survival analysis: techniques for censored and truncated data*. Springer Science & Business Media, 2005.
- [KR97] Daniel A Klain and Gian-Carlo Rota. *Introduction to geometric probability*. Cambridge University Press, 1997.
- [LS15] S. Lan and B. Shahbaba. “Sampling constrained probability distributions using Spherical Augmentation”. In: *ArXiv e-prints* (June 2015). arXiv: [1506.05936](https://arxiv.org/abs/1506.05936) [stat.CO].
- [LV06] László Lovász and Santosh Vempala. “Hit-and-Run from a Corner”. In: *SIAM Journal on Computing* 35.4 (2006), pp. 985–1005. DOI: [10.1137/S009753970544727X](https://doi.org/10.1137/S009753970544727X). eprint: <http://dx.doi.org/10.1137/S009753970544727X>. URL: <http://dx.doi.org/10.1137/S009753970544727X>.
- [LV07] László Lovász and Santosh Vempala. “The Geometry of Logconcave Functions and Sampling Algorithms”. In: *Random Struct. Algorithms* 30.3 (May 2007), pp. 307–358. ISSN: 1042-9832. DOI: [10.1002/rsa.v30:3](https://doi.org/10.1002/rsa.v30:3). URL: <http://dx.doi.org/10.1002/rsa.v30:3>.
- [NJ16] Balasubramanian Narasimhan and Steven G. Johnson. *cubature: Adaptive Multivariate Integration over Hypercubes*. R package version 1.3-6. 2016. URL: <https://CRAN.R-project.org/package=cubature>.
- [Par81] G. Parisi. “Correlation functions and computer simulations”. In: *Nuclear Physics B* 180 (1981), pp. 378–384.
- [PBJ14] John Paisley, David M Blei, and Michael I Jordan. “Bayesian Nonnegative Matrix Factorization with Stochastic Variational Inference”. In: *Handbook of Mixed Membership Models and Their Applications*. Chapman and Hall/CRC, 2014, pp. 205–224.
- [PC08] T. Park and G. Casella. “The Bayesian lasso”. In: *J. Amer. Statist. Assoc.* 103.482 (2008), pp. 681–686. ISSN: 0162-1459. DOI: [10.1198/016214508000000337](https://doi.org/10.1198/016214508000000337). URL: <http://dx.doi.org/10.1198/016214508000000337>.
- [PP14] Ari Pakman and Liam Paninski. “Exact hamiltonian monte carlo for truncated multivariate gaussians”. In: *Journal of Computational and Graphical Statistics* 23.2 (2014), pp. 518–542.
- [RDS04] Gabriel Rodriguez-Yam, Richard A Davis, and Louis L Scharf. “Efficient Gibbs sampling of truncated multivariate normal with application to constrained linear regression”. In: *Unpublished manuscript* (2004).
- [Roc15] Ralph Tyrell Rockafellar. *Convex analysis*. Princeton university press, 2015.

- [RT96] G. O. Roberts and R. L. Tweedie. “Exponential convergence of Langevin distributions and their discrete approximations”. In: *Bernoulli* 2.4 (1996), pp. 341–363. ISSN: 1350-7265. DOI: [10.2307/3318418](https://doi.org/10.2307/3318418). URL: <http://dx.doi.org/10.2307/3318418>.
- [RW98] R. T. Rockafellar and R. J.-B. Wets. *Variational analysis*. Vol. 317. Grundlehren der Mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences]. Springer-Verlag, Berlin, 1998, pp. xiv+733. ISBN: 3-540-62772-3. DOI: [10.1007/978-3-642-02431-3](https://doi.org/10.1007/978-3-642-02431-3). URL: <http://dx.doi.org/10.1007/978-3-642-02431-3>.
- [Sch13] Rolf Schneider. *Convex bodies: the Brunn–Minkowski theory*. 151. Cambridge University Press, 2013.
- [Vil09] C. Villani. *Optimal transport : old and new*. Grundlehren der mathematischen Wissenschaften. Berlin: Springer, 2009. ISBN: 978-3-540-71049-3. URL: <http://opac.inria.fr/record=b1129524>.

## A Details of the orders of magnitude for Table 1 and Table 2

|                   | $d \rightarrow +\infty$ | $\varepsilon \rightarrow 0$ | $R \rightarrow +\infty$ | $r \rightarrow 0$ | $\Delta_1 \rightarrow 0$ | $\Delta_2 \rightarrow +\infty$ |
|-------------------|-------------------------|-----------------------------|-------------------------|-------------------|--------------------------|--------------------------------|
| $L, \lambda^{-1}$ | $d^2$                   | $\varepsilon^{-2}$          | 1                       | $r^{-2}$          | $\Delta_1^{-2}$          | $\Delta_2^2$                   |
| $A_1(x)$          | $d^4$                   | $\varepsilon^{-4}$          | $R^2$                   | $r^{-4}$          | $\Delta_1^{-4}$          | $\Delta_2^4$                   |
| $-\log(\kappa)$   | 1                       | 1                           | $R^{-2}$                | 1                 | 1                        | 1                              |
| $A_2(x)$          | 1                       | $\varepsilon^{-1}$          | $R$                     | $r^{-1}$          | $\Delta_1^{-1}$          | $\Delta_2$                     |
| $T$               | 1                       | $\log(\varepsilon^{-1})$    | $R^2$                   | $\log(r^{-1})$    | $\log(\Delta_1^{-1})$    | $\log(\Delta_2)$               |
| $\gamma$          | $d^{-5}$                | $\varepsilon^6$             | $R^{-2}$                | $r^{-4}$          | $\Delta_1^4$             | $\Delta_2^{-4}$                |

Table 6: dependency of  $L, A_1(x), -\log(\kappa), A_2(x), T, \gamma$  on  $d, \varepsilon, R, r, \Delta_1$  and  $\Delta_2$ .