1Efficient Bayesian computation by proximal Markov chain Monte Carlo: when2Langevin meets Moreau

 $\frac{3}{4}$

Alain Durmus*, Éric Moulines[†], and Marcelo Pereyra[‡]

5Abstract. Modern imaging methods rely strongly on Bayesian inference techniques to solve challenging imag-6 ing problems. Currently, the predominant Bayesian computation approach is convex optimisation, 7 which scales very efficiently to high dimensional image models and delivers accurate point estima-8 tion results. However, in order to perform more complex analyses, for example image uncertainty 9 quantification or model selection, it is necessary to use more computationally intensive Bayesian 10 computation techniques such as Markov chain Monte Carlo methods. This paper presents a new 11 and highly efficient Markov chain Monte Carlo methodology to perform Bayesian computation for 12high dimensional models that are log-concave and non-smooth, a class of models that is central in 13 imaging sciences. The methodology is based on a regularised unadjusted Langevin algorithm that 14 exploits tools from convex analysis, namely Moreau-Yoshida envelopes and proximal operators, to 15construct Markov chains with favourable convergence properties. In addition to scaling efficiently 16 to high dimensions, the method is straightforward to apply to models that are currently solved by 17 using proximal optimisation algorithms. We provide a detailed theoretical analysis of the proposed 18 methodology, including asymptotic and non-asymptotic convergence results with easily verifiable 19conditions, and explicit bounds on the convergence rates. The proposed methodology is demon-20strated with four experiments related to image deconvolution and tomographic reconstruction with 21total-variation and ℓ_1 priors, where we conduct a range of challenging Bayesian analyses related to 22 uncertainty quantification, hypothesis testing, and model selection in the absence of ground truth.

Key words. Mathematical imaging; inverse problems; Bayesian inference; Markov chain Monte Carlo methods;
 convex optimisation; uncertainty quantification; model selection.

25 AMS subject classifications. primary 65C40, 68U10, 62F15; secondary 65C60, 65J22.

1. Introduction. Image estimation problems are ubiquitous in science and engineering. For example, problems related to image denoising [24], deconvolution [5], compressive sensing reconstruction [12], super-resolution [30], tomographic reconstruction [27], inpainting [8], source separation [48], fusion [21], and phase retrieval [6]. The development of new theory, methodology, and algorithms for imaging problems is a focus of significant research efforts. Particularly, convex imaging problems have received a lot of attention lately, leading to major developments in this area.

Most recent works in the imaging literature adopt formal mathematical approaches to analyse problems, derive solutions, and study the underpinning algorithms. There are several mathematical frameworks available to solve imaging problems [22]. In particular, many modern methods are formulated in the Bayesian statistical framework, which relies on statistical models to represent the data observation process and the prior knowledge available, and then derives solutions by using inference techniques rooted in Bayesian decision theory [22].

[†]Centre de Mathématiques Appliquées, UMR 7641, Ecole Polytechnique (eric.moulines@polytechnique.edu).

^{*}LTCI, Telecom ParisTech (alain.durmus@telecom-paristech.fr).

[‡]Maxwell Institute for Mathematical Sciences and School of Mathematical and Computer Sciences, Heriot-Watt University (m.pereyra@hw.ac.uk).

There are currently two main approaches in Bayesian imaging methodology. The predom-39 inant approach is to use a convex formulation of the estimation problem and postulate a prior 40 distribution that is log-concave. This leads to a posterior distribution that is also log-concave, 41 and where maximum-a-posteriori (MAP) estimation can be computed efficiently by using high 42 43 dimensional convex optimisation algorithms [17]. In addition to scaling well to large settings, convex optimisation algorithms have two additional advantages that are important for prac-44 tical Bayesian computation: they are well understood theoretically and their conditions for 45 convergence are clear and simple to check; and the main algorithms are general and can be 46 applied similarly to wide range of problems. However, convex optimisation on its own cannot 47 deliver basic aspects of the Bayesian paradigm and struggles to support the complex statistical 48 analyses that are inherent to modern scientific reasoning and decision-making. 49

The second main approach in Bayesian imaging methodology is based on stochastic sim-50 ulation algorithms, namely Markov chain Monte Carlo (MCMC) algorithms. Such methods, 51which were already actively studied over two decades ago, have regained significant attention lately because of their capacity to address very challenging imaging problems that are be-53 yound the scope of optimisation-based techniques [39]. Additionally to complex models such 54as hierarchical or empirical Bayesian models, MCMC methods also enable advanced analyses 55such as hypotheses test and model selection. Unfortunately, despite great progress in high 56 dimensional MCMC methodology, solving imaging problems by stochastic simulation remains 57too expensive for applications involving moderate or large datasets. Another drawback of 5859existing MCMC methods is that the conditions for their convergence are often significantly more difficult to check than those of optimisation schemes. As a result, most practitioners 60 only assess convergence empirically. It is worth mentioning that some of these limitations can 61 be partially mitigated by resorting to variational Bayes or message passing approximations, 62 which are generally significantly more computationally efficient than stochastic simulation. 63 64 Unfortunately, such approximations are available only for specific models, and we currently 65 have little theory to analyse the approximation error involved. Similarly, it is generally difficult to provide convergence guarantees for the related algorithms, which often suffer from 66 67 local convergence issues. Observe that this is in sharp contrast with the convex optimisation 68 approach, which despite its clear limitations, is general and well understood theoretically.

In summary, convex optimisation and MCMC methods have complementary strengths and 69 weaknesses related to their computational efficiency, theoretical underpinning, and the infer-70ences they can support. As a result, it is increasingly acknowledged that the two methodolo-7172 gies should be used together. In this view, the future imaging methodological toolbox should provide a flexible framework where it is possible to perform very efficiently a first analysis 73 of a full dataset by using convex optimisation algorithms, followed by in-depth analyses by 74 MCMC simulation for specific data (e.g., particular data that will be used as evidence to support a hypothesis or a decision). Also, in this framework practitioners should be able to use MCMC algorithms to perform preliminary analyses, which then set the basis for a full 77 scale analysis with convex optimisation techniques. These could be, for example, exploratory 78 analyses with selected data aimed at calibrating the model or performing Bayesian model 7980 selection, and benchmarking analyses to assess efficient approximations (e.g., optimisationbased approximate confidence intervals [35]). Unfortunately, it is currently difficult to use 81 optimisation and MCMC methodologies in this complementary manner because optimisation 82

2

methods use predominantly non-conjugate priors that are not smooth, such as priors involving 83 the ℓ_1 or the total-variation nome, whereas MCMC methods are mainly restricted to models 84 with priors that are either conjugate to the likelihood function, or that are smooth with Lip-85 chitz gradients (the latter enables efficient high dimensional MCMC algorithms such as the 86 87 Metropolis-adjusted Langevin algorithm or Hamiltonian Monte Carlo [39]). Proximal MCMC algorithms, proposed recently in [36], are an important first step towards 88 bridging this methodological gap between convex optimisation and stochastic simulation. Un-89 like conventional high dimensional MCMC algorithms that use gradient mappings and require 90 Lipchitz differentiability, proximal MCMC algorithms draw their efficiency from convex anal-91 ysis, namely proximal mappings and Moreau-Yoshida envelopes. This allows MCMC-based 92 Bayesian computation for precisely the type of models that are solved by convex optimisation 93 (i.e., high dimensional models that are log-concave but not smooth), which in turn enables 94 advanced Bayesian analyses for these models (e.g., see [35; 2] for applications of proximal 95 MCMC to Bayesian uncertainty quantification and sparse regression). However, the proxi-96 mal MCMC algorithms presented in [36] have three shortcomings that limit their impact in 97 imaging sciences, and which this paper seeks to address. First, the conditions that guarantee 98 the convergence of the algorithms are difficult to check in practice. Second, the algorithms 99 100 assume that it is possible to compute the proximal mapping of the log-posterior distribution; in practice however this mapping is often approximated by using a forward-backward splitting 101 scheme. Third, the algorithms rely on a Metropolis-Hastings correction step to remove the 102 103 asymptotic bias introduced by the approximations and to guarantee that the Markov chains target the desired posterior distribution. Unfortunately, this correction step can degrade sig-104 105 nificantly the efficiency of the algorithms (i.e., the asymptotic bias is removed at the expense of a potentially significant increase in estimation variance and some additional bias from the 106 Markov chain's transient or burn-in regime).

107 108 This paper presents a new and significantly better proximal MCMC methodology that address all the issues of the original proximal algorithms discussed above. This new methodology 109 is highly computationally efficient and general, in that it can be applied straightforwardly to 110 111 most models currently addressed by convex optimisation (in particular, to any model that can 112 be solved by forward-backward splitting). Moreover, we provide simple theoretical conditions 113 to guarantee the convergence of the Markov chains, as well as bounds on its convergence rate. To conclude, we emphasise again that our aim is to provide a Bayesian computation method-114 115ology that complements rather than competes with modern convex optimisation, particularly 116by enabling advanced Bayesian analyses for high-dimensional models that are log-concave. The remainder of the paper is organised as follows: Section 2 defines notation, introduces 117the class of models considered, and recalls the Langevin MCMC approach that is the basis 118

of our method. In Section 3 we present the proposed MCMC method, analyse its theoretical properties in detail, provide practical implementation guidelines, and discuss connections with the original proximal MCMC algorithms described in [36]. Section 4 illustrates the methodology on four experiments related to image deconvolution and tomographic reconstruction with total-variation and ℓ_1 sparse priors, where we conduct a range of challenging Bayesian analyses

related to model comparison and uncertainty quantification. Conclusions and perspectives for

125 future work are reported in Section 5. Proofs are finally reported in Appendices A and C.

sec:bac

126 **2. Bayesian analysis and computation.**

127 **2.1. Notations and Conventions.** Denote by $\mathcal{B}(\mathbb{R}^d)$ the Borel σ -field of \mathbb{R}^d . For all 128 $A \in \mathcal{B}(\mathbb{R}^d)$, denote by Vol(A) its Lebesgue measure. Denote by $\mathbb{M}(\mathbb{R}^d)$ the set of all Borel 129 measurable functions on \mathbb{R}^d and for $f \in \mathbb{M}(\mathbb{R}^d)$, $||f||_{\infty} = \sup_{x \in \mathbb{R}^d} |f(x)|$. For μ a probability 130 measure on $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$ and $f \in \mathbb{M}(\mathbb{R}^d)$ a μ -integrable function, denote by $\mu(f)$ the integral 131 of f w.r.t. μ . For two probability measures μ and ν on $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$, the total variation norm 132 of μ and ν is defined as

133
$$\|\mu - \nu\|_{\mathrm{TV}} = \sup_{f \in \mathbb{M}(\mathbb{R}^d), \|f\|_{\infty} \le 1} \left| \int_{\mathbb{R}^d} f(x) \mathrm{d}\mu(x) - \int_{\mathbb{R}^d} f(x) \mathrm{d}\nu(x) \right|$$

134 Let $f : \mathbb{R}^d \to (-\infty, +\infty]$. If f is a Lipschitz function, namely there exists $C \ge 0$ such that for 135 all $x, y \in \mathbb{R}^d$, $|f(x) - f(y)| \le C ||x - y||$, then denote

136
$$\|f\|_{\text{Lip}} = \inf\{|f(x) - f(y)| \|x - y\|^{-1} \mid x, y \in \mathbb{R}^d, x \neq y\}.$$

f is said to be proper if there exists $x_0 \in \mathbb{R}^d$ such that $f(x_0) < +\infty$. Denote for all $M \in \mathbb{R}$, 137 $\{f \leq M\} = \{z \in \mathbb{R}^d \mid f(z) \leq M\}$. f is said to be lower semicontinuous if for all $M \in \mathbb{R}$, 138 $\{f \leq M\}$ is a closed subset of \mathbb{R}^d . For $k \geq 0$, denote by $C^k(\mathbb{R}^d)$, the set of k-times continuously 139 differentiable functions. For $f \in C^1(\mathbb{R}^d)$, denote by ∇f the gradient of f. Denote for all $q \ge 1$, the ℓ_q norm $\|\cdot\|_q$ on \mathbb{R}^d by for all $x \in \mathbb{R}^d$, $\|x\|_q = (\sum_{i=1}^d |x_i|^q)^{1/q}$. Denote by $\|\cdot\|$ the Euclidian norm on \mathbb{R}^d . For all $x \in \mathbb{R}^d$ and M > 0, denote by B(x, M), the ball centered at x of radius 140 141 142 M. For a closed convex $\mathcal{K} \subset \mathbb{R}^d$, denote by $\operatorname{proj}_{\mathcal{K}}(\cdot)$, the projection onto \mathcal{K} , and $\iota_{\mathcal{K}}$ the convex 143 indicator of \mathcal{K} defined by $\iota_{\mathcal{K}}(x) = 0$ if $x \in \mathcal{K}$, and $\iota_{\mathcal{K}}(x) = +\infty$ otherwise. In the sequel, we take the convention that $\inf \emptyset = \infty$, $1/\infty = 0$ and for $n, p \in \mathbb{N}$, n < p then $\sum_{p=0}^{n} = 0$ and 144 145146 $\prod_{n=1}^{n} = 1.$

2.2. Imaging inverse problems. We consider inverse problems where we seek to estimate an unknown quantity $x \in \mathbb{R}^d$ from an observation y, related to x by a forward statistical model with likelihood function p(y|x). Following a Bayesian approach, we use prior knowledge about x to reduce the uncertainty and deliver accurate estimation results [22]. Precisely, we specify a prior distribution p(x) promoting expected properties (e.g., sparsity, piecewise regularity, or smoothness), and combine observed and prior information by using Bayes' theorem, leading to the posterior distribution [40]

$$\pi(x) \triangleq p(x|y) = \frac{p(y|x)p(x)}{\int_{\mathbb{R}^d} p(y|x)p(x)\mathrm{d}x}$$

147 that we henceforth denote as π , and which models our knowledge about x after observing y. 148 In this paper we focus on inverse problems that are convex. We assume that π is log-concave,

149 i.e.

posterior (1)
$$\pi(x) = \frac{\mathrm{e}^{-U(x)}}{\int_{\mathbb{R}^d} \mathrm{e}^{-U(s)} \mathrm{d}s},$$

for some measurable function $U: \mathbb{R}^d \to (-\infty, +\infty]$ satisfying the following condition.

 $\frac{\text{H1. } U = f + g, \text{ where } f : \mathbb{R}^d \to \mathbb{R} \text{ and } g : \mathbb{R}^d \to (-\infty, +\infty] \text{ are two lower bounded}}{functions satisfying:}$

15(i) f is convex, continuously differentiable, and gradient Lipschitz with Lipschitz constant L_f , 155 i.e. for all $x, y \in \mathbb{R}^d$

eq:gradient-Liþő6 (2

Lipse (2) $\|\nabla f(x) - \nabla f(y)\| \le L_f \|x - y\|$.

item: assum-priorit) g is proper, convex and lower semi continuous (l.s.c).

Notice that the class (1) comprises many important models that are used extensively in modern imaging sciences. Particularly, models of the form $U(x) = ||y - Ax||^2/2\sigma^2 + \phi(Bx)$ for some linear operators A, B, and convex regulariser ϕ that is typically non-smooth, and which may also encode convex constraints on the parameter space. In such cases f(x) = $||y - Ax||^2/2\sigma^2$ and $g(x) = \phi(Bx)$ for instance.

163 When x is high-dimensional, drawing inferences from π directly is generally not possible. 164 Instead we use summaries, particularly point estimators, that capture some of the information 165 about π that is relevant for the application considered [40]. In particular, modern statistical 166 imaging methodology relies strongly on the maximum-a-posteriori (MAP) estimator defined 167 by:

$$\hat{x}_{MAP} = \underset{x \in \mathbb{R}^d}{\arg \max} \pi(x) = \underset{x \in \mathbb{R}^d}{\arg \min} U(x) ,$$

169 which can often be computed efficiently, even in very large problems, by using proximal 170 convex optimisation algorithms [10; 32]. From the practitioner's viewpoint, this is a main 171 advantage w.r.t. most other summaries that require high-dimensional integration w.r.t. π , 172 which is generally significantly more computationally expensive [39].

173However, in its raw form, mathematical imaging based on optimisation struggles to support complex statistical analyses. For example, such methods are typically unable to as-174sess the uncertainty in the solutions delivered and to support uncertainty quantification and 175176decision-making procedures (e.g. hypothesis tests). Similarly, they have difficulty checking and comparing alternative mathematical models intrinsically (i.e., without ground truth avail-177able). To perform such advanced (often Bayesian) analyses and deliver the full richness of the 178statistical paradigm it is necessary to use Monte Carlo stochastic simulation algorithms [17]. 179As mentioned previously, the high-dimensionality and the lack of smoothness of π pose 180 181 important challenges from a Bayesian computation viewpoint. This paper presents a new MCMC methodology to tackle this problem. The proposed methodology is general, robust, 182theoretically sound, and computationally efficient, and can be applied straightforwardly to any 183 model satisfying (1) that can be addressed by using proximal convex optimisation (particularly 184by using the gradient of f and the proximal operator of g, similarly to forward-backward 185186splitting algorithms).

Finally, we mention at this point some recent works that also consider new MCMC methods to sample from non-smooth posterior distributions with ℓ_1 priors, which is a specific subclass of (1). Most of these works consider Gibbs sampling strategies based either on auxiliary variables [34] or on direct simulation from the univariate conditional densities involved [25; 26]. An alternative strategy is to use a non-linear transformations to change the ℓ_1 prior into a Gaussian distribution, which then enables using the randomize-then-optimise (RTO) method of [4] to generate samples (see [47] for details). Similarly to our methodology, RTO combines optimisation and sampling steps, albeit in a completely different way (precisely, RTO simulates

195 high-dimensional Gaussian vectors by minimising a loss function with random parameters).

196 **2.3. Bayesian computation: unadjusted and Metropolis-adjusted Langevin algorithms.** 197 The MCMC method proposed in this paper is derived from the discretization of overdamped 198 Langevin diffusions. Let $\overline{U} : \mathbb{R}^d \to \mathbb{R}$ be a continuously differentiable function and consider 199 the Langevin stochastic differential equations (SDE) given by

q:langevin-200 (4)
$$\mathrm{d}\mathbf{X}_t = -\nabla \bar{U}(\mathbf{X}_t)\mathrm{d}t + \sqrt{2}\mathrm{d}B_t^d ,$$

where $(B_t^d)_{t\geq 0}$ is a *d*-dimensional Brownian motion. Under additional mild assumptions, this equation has a unique strong solution. In addition if $\int_{\mathbb{R}} e^{-\bar{U}(x)} dx < \infty$, then $\bar{\pi}(x) \propto e^{-\bar{U}(x)}$ is the unique invariant distribution of the semi-group associated with the Langevin SDE, see [23]. Consequently, if we could solve (4) and let $t \to \infty$, this would provide samples from $\bar{\pi}$ useful for Bayesian computation. Since it is possible to analytically solve (4) only in very specific cases, we consider a discrete-time Euler-Maruyama approximation and obtain the following Markov chain $(X_k)_{k\geq 0}$: for all $k \geq 0$

efinition-Eule
$$08$$
 (5)

-e

de

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

where $\gamma > 0$ is a given step size and $(Z_k)_{k\geq 1}$ is a sequence of i.i.d. *d*-dimensional standard Gaussian random variables. This scheme has been first introduced in molecular dynamics by [15] and [33], and then popularized in artificial intelligence by [18], [19] and in computational statistics by [31] and [42]. Following [42], this algorithm is referred to as the Unadjusted Langevin Algorithm (ULA).

214In Bayesian computation, the samples $(X_k)_{k>0}$ generated by ULA (5) are used to estimate probabilities and expectations w.r.t. $\bar{\pi}$. This scheme has attracted significant attention 215216recently, in particular for high-dimensional problems were most Monte Carlo methods strug-217gle. Theory for ULA advanced significantly recently with the development of non-asymptotic 218bounds in total variation distance between $\bar{\pi}$ and the marginal laws of the Markov chain $(X_k)_{k\geq 0}$ defined by ULA [11; 13], with explicit dependence on the stepsize γ and the dimen-219 sion d (see Subsection 3.2). These new theoretical results are important because they provide 220 221 estimation accuracy guarantees for ULA, as well as valuable new insights into the convergence properties of the algorithm. In particular, they establish that if \overline{U} is convex and gradient Lip-222chitz, then ULA's convergence properties deteriorate at most polynomially as d increases. 223 Remarkably, if in addition \overline{U} is strongly convex, then it deteriorates at most linearly with d, 224confirming the empirical evidence that ULA is a highly computationally efficient method to 225226 sample in high-dimensional settings.

It is worth emphasising at this point that this deep understanding of ULA is very recent. Indeed, without a proper theoretical underpinning, ULA has been traditionally regarded as unreliable and rarely applied directly in statistics or statistical image processing. Instead, most applications reported in the literature adopt a safe approach and complement ULA with a Metropolis-Hastings correction step targeting $\bar{\pi}$, as recommended by [44] and [42]. This

correction guarantees that the resulting Metropolis Adjusted Langevin Algorithm (MALA) generates a reversible Markov chain with respect to $\bar{\pi}$, and therefore eliminates the asymptotic bias. And perhaps more importantly, it places ULA within the sound theoretical framework of Metropolis-Hasting algorithms. For sufficiently smooth densities MALA inherits the good convergence properties of ULA and scales efficiently to high-dimensional settings [42].

Unfortunately, neither ULA nor MALA are well defined for non-smooth target densities, 237 which strongly limits their application to modern mathematical imaging problems. In fact, 238 both theory and experimental evidence show that ULA and MALA often run into difficulties 239if π is not sufficiently regular. For example, when $\nabla \log \pi$ is not Lipchitz continuous ULA is 240 generally explosive and MALA is not geometrically ergodic (see [42; 36, Figure 2]). Similarly, 241when $\nabla \log \pi$ is subdifferentiable and therefore, at least from a purely algorithmic viewpoint, 242the algorithms could still be applied, the theory underpinning the ULA and MALA collapses 243 and even the convergence of the time-continuous Langevin diffusion driving the algorithms 244becomes unclear. Moreover, many applications involve constraints on the parameter space 245and then π is supported only on a bounded convex set \mathcal{K} . In such case, $\nabla \log \pi$ is bounded 246on \mathcal{K} and infinite or not defined outside \mathcal{K} . Then it is not possible to use ULA, and MALA 247 typically behaves very poorly (the algorithm gets "stuck" whenever the proposal drives the 248249 Markov chain outside \mathcal{K}). Following a proximal MCMC approach [36], in the following section we present a new ULA that exploits tools from convex calculus and proximal optimisation to 250address these issues, and sample efficiently from high-dimensional log-concave densities of the 251252form H1 that are beyond the scope of conventional ULAs and MALAs. nore-yosida-regul

3. Proximal MCMC: Moreau-Yosida regularised Unadjusted Langevin Algorithm.

3.1. Proposed method. A central idea in this work is to replace the non-smooth potential U with a carefully designed smooth approximation U^{λ} which, by construction, has the following two key properties: 1) its Euler-Maruyama discrete-time approximations are always stable and have favourable convergence properties, and 2) we can make $\pi^{\lambda} \propto e^{-U^{\lambda}}$ arbitrarily close to π by adjusting an approximation parameter $\lambda > 0$.

In a manner akin to [36], we define such approximations by using Moreau-Yosida envelopes [9] which we recall below. Let $g : \mathbb{R}^d \to (-\infty, +\infty)$ be a l.s.c convex function and $\lambda > 0$. The λ -Moreau-Yosida envelope of g is a carefully regularised approximation of g given by

-eq:id-MY-env62 (6)
$$g^{\lambda}(x) = \min_{y \in \mathbb{R}^d} \left\{ g(y) + (2\lambda)^{-1} \|x - y\|^2 \right\} ,$$

where λ is a regularisation parameter that controls a trade-off between the regularity properties of g^{λ} and the approximation error involved. Remarkably, by [43, Example 10.32, Theorem 9.18], the approximation g^{λ} inherits the convexity of g and is always continuously differentiable, even if g is not. In fact, g^{λ} is gradient Lipshitz [43, Proposition 12.19]: for all $x, y \in \mathbb{R}^d$,

noreau'yosidà67 (7)
$$\left\| \nabla g^{\lambda}(x) - \nabla g^{\lambda}(y) \right\| \le \lambda^{-1} \left\| x - y \right\| .$$

268 The gradient is given by for all $x \in \mathbb{R}^d$

ssec:MYULAI

$$\overline{\text{tion-grad-prox}^{\lambda}(x)} = \lambda^{-1} \left(x - \text{prox}_{g}^{\lambda}(x) \right) ,$$

270 where 271 (9)

(10)

 $\operatorname{prox}_{g}^{\lambda}(x) = \operatorname*{arg\,min}_{y \in \mathbb{R}^{d}} \left\{ g(y) + (2\lambda)^{-1} \left\| x - y \right\|^{2} \right\} ,$

is the proximal operator of g [9]. This operator is used extensively in imaging methods based on convex optimisation, where it is generally computed efficiently by using a specialised algorithm [10; 32]. Indeed, similarly to gradient mappings, $\operatorname{prox}_{g}^{\lambda}$ also moves points in the direction of the minimum of g (by an amount related to the value of λ), and has many properties that are useful for devising fixed-point methods [9].

In addition, g^{λ} envelopes g from below: for all $x \in \mathbb{R}^d$, $g^{\lambda}(x) \leq g(x)$, and since for 278 $0 < \lambda < \lambda'$ and $x, y \in \mathbb{R}^d$, $g(y) + (2\lambda')^{-1} ||x - y||^2 \leq g(y) + (2\lambda)^{-1} ||x - y||^2$, we get that for 279 all $x \in \mathbb{R}^d g^{\lambda'}(x) \leq g^{\lambda}(x)$. By [43, Theorem 1.25], g^{λ} converges pointwise to g as λ goes to 0, 280 i.e. for all $x \in \mathbb{R}^d$,

 $\lim_{\lambda \to 0} \mathbf{g}^{\lambda}(x) = \mathbf{g}(x) \; .$

:limit-d-lambda81

-eq:prox`g71

Hence, g^{λ} provides a convex and smooth approximation to g that we can make arbitrarily close to g by adjusting the value of λ .

So under H1, if g is not continuously differentiable, but the proximity operator associated with g is available, we can consider sampling algorithms that use the λ -Moreau-Yosida envelope g^{λ} instead of g. Here we propose to replace the potential U with the approximation $U^{\lambda} : \mathbb{R}^{d} \to \mathbb{R}$ defined for all $x \in \mathbb{R}^{d}$ by

$$U^{\lambda}(x) = g^{\lambda}(x) + f(x) ,$$

which we will use to define a surrogate target density $\pi^{\lambda} \propto e^{-U^{\lambda}}$. We will see that such approximation is endowed with very useful regularity and approximation accuracy properties.

Proposition 1 below implies that the probability measure π^{λ} on \mathbb{R}^d , with density with respect to the Lebesgue measure, also denoted by π^{λ} and given for all $x \in \mathbb{R}^d$ by

$$\pi^{\lambda}(x) = rac{\mathrm{e}^{-U^{\lambda}(x)}}{\int_{\mathbb{R}^d} \mathrm{e}^{-U^{\lambda}(s)} \mathrm{d}s}$$

is well defined, log-concave, Lipschitz continuously differentiable, and as close to π as required.

H2. Assume that one of these two conditions holds:

sum:integrable $g(i) e^{-g}$ is integrable with respect to the Lebesgue measure.

assum:lipschitz(ii) g is Lipschitz.

Proposition 1. Assume H1 and H2.

29a) For all $\lambda > 0$, π^{λ} defines a proper density of a probability measure on \mathbb{R}^d , i.e.

297

sum:integrabil@@2

ite-measure-MW

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

29b) For all $\lambda > 0$, π^{λ} is log-concave and continuously differentiable with

on-grad-prox¹099 (11) $\nabla U^{\lambda}(x) = -\nabla \log \pi^{\lambda}(x) = \nabla f(x) + \lambda^{-1}(x - \operatorname{prox}_{g}^{\lambda}(x)) .$

300 In addition, ∇U^{λ} is Lipschitz with constant $L \leq L_f + \lambda^{-1}$.

This manuscript is for review purposes only.

302

oo:dist`TV`MY30¢)

$$\lim_{\lambda \to 0} \|\pi^{\lambda} - \pi\|_{\mathrm{TV}} = 0$$

o:dist'TV'MY30d) If **H**2-(*ii*) then for all $\lambda > 0$,

304

306

$$\|\pi^{\lambda} - \pi\|_{\mathrm{TV}} \le \lambda \, \|g\|_{\mathrm{Lip}}^2$$

Proof. The proof is postponed to Appendix A. 305

Figure 1 shows the approximations of two non-smooth densities that satisfy H_1 : 1. the Laplace density $\pi(x) = (1/2) \exp(|x|)$, for which

The approximation π^{λ} converges to π as $\lambda \downarrow 0$ in total variation norm, i.e.

308

item:caseUbilf(

tem:caseLaplace7

$$\pi^{\lambda}(x) = \frac{\exp\left\{(\lambda/2 - |x|)\mathbb{1}_{\{|x| \ge \lambda\}} - (x^2/(2\lambda))\mathbb{1}_{\{|x| < \lambda\}}\right\}}{2\left\{e^{-\lambda/2} + (2\pi/\lambda)^{1/2}(\mathbf{\Phi}(\lambda^{1/2}) - 1/2)\right\}}$$

309

where Φ is the cumulative function of the standard normal distribution. 2. the uniform density $\pi(x) = (1/2) \exp(-\iota_{[-1,1]}(x))$, for which

311
$$\pi^{\lambda}(x) = \left\{2 + \sqrt{2\pi\lambda}\right\}^{-1} \exp\left[\left\{-\max(|x| - 1, 0)\right\}^2 / (2\lambda)\right]$$

We observe that the approximations are smooth and converge to π as λ decreases, as described 312

by Proposition 1. Also for these two examples, analytic expressions for $\|\pi - \pi^{\lambda}\|_{TV}$ can be 313 found, and Figure 2 shows $\|\pi - \pi^{\lambda}\|_{TV}$ as a function of $\lambda > 0$. Notice that in the case of the 314

Laplace density $\|\pi - \pi^{\lambda}\|_{TV}$ goes to 0 quadratically in λ as λ goes to 0, which is faster than

315the linear bound given in Proposition 1-d). Also note that this bound does not apply to the

316 uniform density, and in this case $\|\pi - \pi^{\lambda}\|_{\text{TV}}$ vanishes at rate $\sqrt{\lambda}$.



Figure 1. Density plots for the Laplace (a) and uniform (b) distributions (solid black), and their smooth approximations π^{λ} for $\lambda = 1, 0.1, 0.01$ (dashed blue and green, and solid red).

FigMoreauApprox

317

We now make two key observations. First, Proposition 1 shows that ∇U^{λ} is gradient 318Lipschitz and therefore it guarantees that the Langevin SDE constructed with U^{λ} converges 319



Figure 2. Total variation norm between π and its smooth approximation π^{λ} as function of λ .

FigMoreauApprox-TV

to π^{λ} as $t \to \infty$ (formally, it guarantees that the Langevin SDE associated with π^{λ} admits a unique strong solution $(\mathbf{X}_{t}^{\lambda})_{t\geq 0}$ and π^{λ} is the unique stationary distribution of the semigroup). More importantly, as it will be seen below, it implies that the ULA chain derived from a Euler-

323 Maruyama discretisation of this Langevin diffusion will be, by construction, well behaved and

324 useful for Monte Carlo integration with respect to π^{λ} .

Second, Proposition 1 also establishes that λ controls the estimation bias involved in performing estimations with π^{λ} as a substitute of π . This approximation error can be made arbitrarily small, and is bounded explicitly by $\lambda ||g||^2_{\text{Lip}}$ when g is Lipschitz.

We are now in a position to present the new MCMC methodology proposed in this work, which is essentially an application of ULA to π^{λ} . Precisely, given $\lambda > 0$ and a stepsize $\gamma > 0$, we use an Euler-Maruyama approximation of $(\mathbf{X}_t^{\lambda})_{t\geq 0}$, and obtain the following Markov chain $(X_k^{\mathrm{M}})_{k\geq 0}$: for all $k \geq 0$

(12) MYULA:
$$X_{k+1}^{\mathrm{M}} = (1 - \frac{\gamma}{\lambda})X_k^{\mathrm{M}} - \gamma \nabla f(X_k^{\mathrm{M}}) + \frac{\gamma}{\lambda} \operatorname{pros}_g^{\lambda}(X_k^{\mathrm{M}}) + \sqrt{2\gamma}Z_{k+1}$$
,

where $\{Z_k, k \in \mathbb{N}^*\}$ is a sequence of i.i.d. *d* dimensional standard Gaussian random variables. This algorithm will be referred to as the *Moreau-Yosida Unadjusted Langevin Algorithm* (MYULA), and is summarised in Algorithm 1 below (see Subsection 3.3 for guidelines for setting the values of γ and λ). Note that the stationary distribution of the MYULA sequence $\{X_k^{\mathrm{M}}, k \in \mathbb{N}\}$ is different from the target distribution π^{λ} , and depends on the stepsize $\gamma > 0$. Nevertheless, we show in Subsection 3.2 that, choosing λ and γ appropriately, the samples are very close to π .

Besides, to compute the expectation of a function $h : \mathbb{R}^d \to \mathbb{R}$ under π from $\{X_k^{\mathrm{M}}; 0 \le k \le n\}$, an optional importance sampling step might be used to correct the regularization. This step amounts to approximate $\int_{\mathbb{R}^d} h(x)\pi(x) \mathrm{d}x$ by the weighted sum

$$\underbrace{\operatorname{ce`sampling}}_{\text{ce`sampling}} (13) \qquad \qquad \operatorname{S}_n(h) = \sum_{k=0}^n \omega_{k,n} h(X_k) , \text{ with } \omega_{k,n} = \left\{ \sum_{\ell=0}^n \mathrm{e}^{\bar{g}^{\lambda}(X_{\ell}^{\mathrm{M}})} \right\}^{-1} \mathrm{e}^{\bar{g}^{\lambda}(X_{k}^{\mathrm{M}})}$$

q:def-MYRULA32

tan

344 where for all $x \in \mathbb{R}^d$

$$\bar{g}^{\lambda}(x) = g^{\lambda}(x) - g(x) = g(\operatorname{prox}_{g}^{\lambda}(x)) - g(x) + (2\lambda)^{-1} \left\| x - \operatorname{prox}_{g}^{\lambda}(x) \right\|^{2}$$

To remove this asymptotic bias, we can add an Hastings-Metropolis step, which will produce a Markov chain $\{\tilde{X}_k^{\lambda}, k \in \mathbb{N}\}$ which is reversible this time with respect to π^{λ} and use similarly an importance sampling step to correct for the bias introduced by smoothing. This algorithm will be called the *Moreau-Yosida Regularized Metropolis-adjusted Langevin Algorithm* (MYMALA).

The focus of this work is on MYULA without importance sampling or Metropolis-Hastings correction. A study of MYMALA is currently in progress and will be reported separately.

Algorithm 1 Moreau-Yoshida unadjusted Langevin algorithm (MYULA) set $X_0^{\mathrm{M}} \in \mathbb{R}^d$, $\lambda > 0$, $\gamma \in (0, \lambda/(\lambda L_f + 1)]$, $n \in \mathbb{N}$ for k = 0 : n do $Z_{k+1} \sim \mathcal{N}(0, \mathbb{I}_d)$ $X_{k+1}^{\mathrm{M}} = (1 - \frac{\gamma}{\lambda})X_k^{\mathrm{M}} - \gamma \nabla f(X_k^{\mathrm{M}}) + \frac{\gamma}{\lambda} \operatorname{prox}_g^{\lambda}(X_k^{\mathrm{M}}) + \sqrt{2\gamma}Z_{k+1}$ end for

For illustration, Figure 3 shows the sample approximations of the univariate Laplace and 353 uniform distributions of Figure 1 (the true densities are depicted in solid blue for comparison). 354The histograms were generated using 10⁴ iterations of MYULA with parameters $\lambda = 10^{-3}$ 355 and $\delta = 2\lambda$. Observe that the samples provide a good approximation of the desired target 356 densities, particularly of the uniform distribution which is beyond the scope of the conventional 357 ULA. Also observe that the approximation of the uniform distribution has Gaussian tails, as 358 per the Lipchitz differentiability of π^{λ} approximation (see Proposition 1 and Figure 1). From 359 Proposition 1, this error can be made arbitrarily low by adjusting the value of λ . 360



Figure 3. MYULA sample approximations of the Laplace (a) and uniform (b) distributions (true density in solid blue). Histograms computed with 10^4 samples generated using $\lambda = 10^{-3}$ and $\delta = 2\lambda$.

FigMYULAApprox

361 Finally, similarly to ULA, it is possible to adapt MYULA to use a stochastic gradient

strategy based on an unbiased estimator of ∇f (see [46] for details). This can be useful in 362 applications that involve very large datasets for which computing the exact gradient would be 363 too computationally expensive, for example machine learning applications. The specialisation 364 of MYULA to such problems is currently under investigation and will be reported separately. 365 vergence analysis

3.2. Theoretical convergence analysis of MYULA. In this section we present a detailed theoretical analysis of MYULA implemented with fixed regularization parameter $\lambda > 0$ and 367 step-size $\gamma > 0$. We first establish that the chains generated by MYULA converge geomet-368 rically fast to an approximation of π that is controlled by λ and γ , and which can be made 369 arbitrarily close to π . More importantly, we also establish non-asymptotic bounds for the 370 estimation error of MYULA with a finite number of iterations. This enables an analysis of 371 372 the behaviour of MYULA as the dimensionality of the model increases, as well as deriving 373 practical guidelines for setting λ and γ for specific models.

First, under H 1, it has been observed that q^{λ} is λ^{-1} -gradient Lipschitz, which im-374 plies that U^{λ} is gradient Lipschitz as well: there exists $L \geq 0$ such that for all $x, y \in \mathbb{R}^d$, 375 $\left\|\nabla U^{\lambda}(x) - \nabla U^{\lambda}(y)\right\| \leq L \left\|x - y\right\|$ and 376

$$L \le L_f + \lambda^{-1}$$
 .

Of course, this bound strongly depends on the decomposition of U in a smooth and a non-378 smooth part, which is arbitrary and therefore can be pessimistic (for instance, if U is contin-379 uously differentiable, q can be chosen to be 0 which implies $U^{\lambda} = U$ and $L = L_f$). 380 381

We assume first the following assumption on the potential U^{λ} .

H3. There exist a minimizer x^* of U^{λ} , $\eta_c > 0$ and $R_c \ge 0$ such that for all $x \in \mathbb{R}^d$, $||x - x^{\star}|| \geq \mathbf{R}_{c},$

rexpo potential
$$84$$
 (15)

383

stilipiuilambda77

ssum:potentia

inition'R'kernel92

(14)

$$U^{\lambda}(x) - U^{\lambda}(x^{\star}) \ge \eta_{c} \|x - x^{\star}\| .$$

Note that in fact H3 always holds under H1 and H2, since by Lemma 4 and Proposition 1 there 385exist $C_1, C_2 > 0$ such that $U^{\lambda}(x) \ge C_1 ||x|| - C_2$. Therefore, since U^{λ} is continuous on \mathbb{R}^d , there 386 exists a minimizer x^* of U^{λ} and (15) holds with $\eta_c \leftarrow C_1/2$ and $R_c \leftarrow 2(C_2 + ||x^*|| + U^{\lambda}(x^*))/C_1$. 387 However, these constants are non quantitative, and that is why we introduce H_3 to derive 388 quantitative bounds. 389

Consider the Markov kernel R_{γ} associated to the Euler-Maruyama discretization (12) 390 given, for all $A \in \mathcal{B}(\mathbb{R}^d)$ and $x \in \mathbb{R}^d$ by 391

(16)
$$R_{\gamma}(x,\mathsf{A}) = (4\pi\gamma)^{-d/2} \int_{\mathsf{A}} \exp\left(-(4\gamma)^{-1} \left\|y - x + \gamma\nabla U^{\lambda}(x)\right\|^{2}\right) \mathrm{d}y \,.$$

The sequence $(X_n^{\mathrm{M}})_{n\geq 0}$ defined by (12) is a homogeneous Markov chain associated with the 393 Markov kernel R_{γ} . Therefore for all $n \in \mathbb{N}$, $n \geq 1$, and $x \in \mathbb{R}^d$, the distribution of X_n^{M} started 394 at x is $R^n_{\gamma}(x, \cdot)$ defined by induction for all $\mathsf{A} \in \mathcal{B}(\mathbb{R}^d)$ by 395

396
$$R_{\gamma}^{n}(x,\mathsf{A}) = \int_{\mathsf{A}} R_{\gamma}^{n-1}(x,\mathrm{d}y) R_{\gamma}(y,\mathsf{A})$$

397

It is easily seen that under \mathbf{H}_1 , since U^{λ} is continuously differentiable, R_{γ} is irreducible with respect to the Lebesgue measure, all compact sets are 1-small and the kernel is strongly aperiodic. In addition under \mathbf{H}_3 , since U is also convex then [13, Proposition 13] shows that R_{γ} satisfies a Foster-Lyapunov drift condition, i.e. for all $\bar{\gamma} \in (0, L], \gamma \in (0, \bar{\gamma}]$ and for all $x \in \mathbb{R}^d$,

403
$$R_{\gamma}V_{\rm c}(x) \le \varrho_{\rm c}^{\gamma}V_{\rm c}(x) + b_{\rm c}\gamma ,$$

404 where

405 (17a)
$$V_{\rm c}(x) = \exp\left\{ (\eta_{\rm c}/4) \left(\|x - x^{\star}\|^2 + 1 \right)^{1/2} \right\}$$

406 (17b)
$$\varrho_{\rm c} = {\rm e}^{-2^{-4}\eta_{\rm c}^2(2^{1/2}-1)}$$
, ${\rm a}_{\rm c} = \max(1, 2d/\eta_{\rm c}, {\rm R}_{\rm c})$

408 (17c)
$$b_{\rm c} = \{(\eta_{\rm c}/4)(d + (\eta_{\rm c}\bar{\gamma}/4)) - \log(\varrho_{\rm c})\} e^{\eta_{\rm c}({\rm a}_{\rm c}^2 + 1)^{1/2}/4 + (\eta_{\rm c}\bar{\gamma}/4)(d + (\eta_{\rm c}\bar{\gamma}/4))}$$

 $\frac{1}{\text{eq:convex'drift9}}{410} \quad \text{By [29, Theorem 16.0.1], } R_{\gamma} \text{ has a unique invariant distribution } \pi_{\gamma}^{\lambda} \text{ and is } V_{\text{c}}\text{-uniformly} \\ \text{geometrically ergodic: there exists } \kappa_{\text{c}} \in (0, 1) \text{ and } C_{\text{c}} \geq 0 \text{ such that all } n \geq 0 \text{ and } x \in \mathbb{R}^{d},$

411
$$\|R_{\gamma}^{n}(x,\cdot) - \pi_{\gamma}^{\lambda}\|_{\mathrm{TV}} \leq C_{\mathrm{c}} V_{\mathrm{c}}(x) \kappa_{\mathrm{c}}^{n}$$

412 Note π_{γ}^{λ} is different from π^{λ} , nevertheless the following result shows that choosing γ small 413 enough, the ULA generates samples very close to the distribution π^{λ} .

414 We are now ready to present our main theoretical result: a non-asymptotic bound of 415 the total-variation distance between π and the marginal laws of the samples generated by 416 MYULA. Denote in the following by $\omega : \mathbb{R}_+ \to \mathbb{R}_+$ the function given for all $r \ge 0$ by

-eq:Fsmall[7] (18)
$$\omega(r) = r^2 / \left\{ 2 \Phi^{-1}(3/4) \right\}^2$$
.

418

c-stepsize-convV

Theorem 2 ([13, Corollary 19]). Assume H1 and H3. Let $\bar{\gamma} \in (0, L^{-1}]$. For all $\varepsilon > 0$ and $x \in \mathbb{R}^d$, we have

$$||R_{\gamma}^{n}(x,\cdot) - \pi||_{\mathrm{TV}} \leq \varepsilon ,$$

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

420
$$T = \max\left\{32\eta_{\rm c}^{-2}\log\left(8\varepsilon^{-1}A_1(x)\right), \log(16\varepsilon^{-1})/(-\log(\kappa))\right\}$$

421
422
$$\gamma \le \frac{-d + \sqrt{d^2 + (2/3)A_2(x)\varepsilon^2(L^2T)^{-1}}}{2A_2(x)/3} \wedge \bar{\gamma} ,$$

14

423 where
$$\alpha_{c} = \max(1, 4d/\eta_{c}, R_{c})$$

424 $\beta_{c} = (\eta_{c}/4) \left[\eta_{c}\alpha_{c}/4 + d\right] \max\left\{1, (\alpha_{c}^{2} + 1)^{-1/2} \exp(\eta_{c}(\alpha_{c}^{2} + 1)^{1/2}/4)\right\}$
425 $A_{1}(x) = (1/2)(V_{c}(x) + b_{c}(-\varrho_{c}^{\gamma}\log(\varrho_{c}))^{-1} + 8\eta_{c}^{-2}\beta_{c}) + 16\eta_{c}^{-2}\beta_{c}e^{32^{-1}\eta_{c}^{2}\omega\left\{(8/\eta_{c})\log(32\eta_{c}^{-2}\beta_{c})\right\}}$
426 $A_{2}(x) = L^{2}\left(4\eta_{c}^{-1}\left[1 + \log\left\{V_{c}(x) + b_{c}(-\varrho_{c}^{\gamma}\log(\varrho_{c}))^{-1}\right\}\right]\right)^{2}$
427 $\log(\kappa) = -\log(2)(\eta_{c}^{2}/32)\left[\log\left\{8\eta_{c}^{-2}\beta_{c}\left(3 + 4\eta_{c}^{-2}e^{32^{-1}\eta_{c}^{2}\omega\left\{(8/\eta_{c})\log(32\eta_{c}^{-2}\beta_{c})\right\}\right)\right\} + \log(2)\right]^{-1},$

431

437

440

utsideBallDriftV

dec-stepsize-StV

423

 a_c, ρ_c, b_c, V_c are defined in (17) and ω in (18). 430

Proof. The proof follows from combining [13, Lemma 4, Theorem 14, Theorem 16].

432 This result implies that the number of iteration to reach a precision target ε is, at worse, of order $d^5 \log^2(\varepsilon^{-1})\varepsilon^{-2}$ for this class of models. Significantly more precise bounds can be 433 obtained under more stringent assumption on U^{λ} . In particular, we consider the case where 434 U^{λ} is strongly convex outside some ball; see [14]. 435

H4. There exist $R_s \ge 1$ and m > 0, such that for all $x, y \in \mathbb{R}^d$, $||x - y|| \ge R_s$,

$$\left\langle \nabla U^{\lambda}(x) - \nabla U^{\lambda}(y), x - y \right\rangle \ge m \left\| x - y \right\|^{2}$$

Of course, in the case where f is strongly convex then this assumption holds. 438

Theorem 3 ([13, Lemma 4, Theorem 21]). Assume H1 and H4. Let $\bar{\gamma} \in (0, L^{-1}]$. Then for all $\varepsilon > 0$, we get $||R_{\gamma}^n(x, \cdot) - \pi||_{\text{TV}} \leq \varepsilon$ provided that $n > T\gamma^{-1}$ with

441

$$T = \left(\log\{A_1(x)\} - \log(\varepsilon/2)\right) / (-\log(\kappa))$$
442
443

$$\gamma \le \frac{-d + \sqrt{d^2 + (2/3)A_2(x)\varepsilon^2(L^2T)^{-1}}}{2A_2(x)/3} \wedge \bar{\gamma} ,$$

where 444

445
$$A_{1}(x) = 5 + \left(d/m + R_{s}^{2}\right)^{1/2} + (A_{1}(x)/L^{2})^{1/2}$$
446
$$A_{2}(x) = L^{2} \left(\|x - x^{\star}\|^{2} + 2(d + mR_{s}^{2})(e^{-\gamma(2m + \bar{\gamma}L^{2})}/(2m + \bar{\gamma}L^{2}))^{-1}\right)$$
447
$$\log(\kappa) = -(\log(2)m/2) \left[\log\left\{\left(1 + e^{m\omega\{\max(1,R_{s})\}/4}\right)(1 + \max(1,R_{s}))\right\} + \log(2)\right]^{-1},$$

447
448
$$\log(\kappa) = -(\log(2)m/2) \left[\log\left\{ \left(1 + e^{m\omega\{\max(1, R_s)\}/4} \right) \left(1 + \max(1, R_s) \right) \right\} + \log(2) \left(1 + \log(2) + \log(2)$$

and ω is given in (18). 449

This result implies that the worst minimal number of iterations to achieve a precision level 450 $\varepsilon > 0$ is this time of order $d \log(d) \log^2(\varepsilon^{-1}) \varepsilon^{-2}$. 451guidelines

3.3. Selection of λ and γ . We now discuss practical guidelines for setting the values 452for λ and for γ . As mentioned previously, our aim is to provide an efficient computation 453methodology that can be applied straightforwardly to any model satisfying H1. Hence, 454

 $455 \\ 456$

rather than seeking optimal values for specific models, we focus on general rules that are simple, robust, and which only involve tractable quantities such as Lipschitz constants.

First, by Theorem 2, γ should take its value in the range $\gamma \in (0, \lambda/(L_f \lambda + 1)]$ to guarantee the stability of the Euler-Maruyama discretisation, and where we recall that L_f is the Lipschitz constant of ∇f . The values of γ within this range are subject to the a bias-variance trade-off. Precisely, large values of γ produce a fast-moving chain that convergences quickly and has low estimation variance, but potentially relatively high asymptotic bias. Conversely, small values of γ lead to low asymptotic bias, but produce a Markov chain that moves slowly and requires a large number of iterations to produce a stable estimate (such chains often also suffer from some additional bias from the transient or burn-in period). Because applications in imaging sciences involve high dimensionality and require moderately low computing times, as a general rule we recommend setting γ to a relatively large value. For example, in our experiments we use

$$\gamma \in [\lambda/5(L_f\lambda+1), \lambda/2(L_f\lambda+1)]$$
.

Observe that this range depends on the value of λ , which is also subject to a bias-variance tradeoff. Letting $\lambda \to 0$ to bring π^{λ} close to π reduces asymptotic bias, but forces $\gamma \to 0$ and consequently reduces significantly the efficiency of the chain. Conversely, increasing the value of λ accelerates the chain at the expense of some asymptotic bias. Based on our experience, and again with an emphasis on efficiency in high dimensional settings, we recommend using values of λ in the order of L_f^{-1} (there is no benefit in using larger values of λ because γ saturates at L_f^{-1}). In all our experiments we use $\lambda = 1/L_f$ and $\gamma \in [L_f^{-1}/10, L_f^{-1}/4]$ and obtain estimation errors of the order of 1%.

3.4. Connections to the proximal Metropolis-adjusted Langevin algorithm. We con-465 clude this section with a discussion of the connections between the proposed MYULA method 466 467 and the original proximal Metropolis-adjusted Langevin algorithm (Px-MALA) [36]. That algorithm is also based on a Euler-Maruvana approximation of a Langevin SDE targeting a 468 Moreau-Yoshide-type regularised approximation of π . However, unlike MYULA, that algo-469rithm uses this approximation as proposal mechanism to drive a Metropolis-Hastings (MH) 470471 algorithm targeting π (not the regularised approximation). The role of the MH is two-fold: it removes the asymptotic bias related to the approximations involved, and it provides a theoret-472 ical framework for Px-MALA by placing the scheme within the framework of MH algorithms 473(recall that many theoretical results regarding ULAs are very recent). However, as mentioned 474 475previously, the introduction of the MH step often slows down the algorithm, thus leading to higher estimation variance and longer chains (and potentially some bias from the chain's 476 initial transient regime). Of course, it also introduces a significant computational overhead 477 related to the computation of the MH acceptance ratio [36]. Another importance difference 478between MYULA and Px-MALA is that the latter uses the proximal operator of U, which 479is often unavailable and has to be approximated by using a forward-backward scheme based 480 on the decomposition U = f + q that we also use in this paper. This approximation error 481 is corrected in practice by the MH step, but it is not considered in the theoretical analysis 482 483 of the algorithm. Conversely, in MYULA this decomposition is explicit, both in the computational aspects of the method as well as in its theoretical analysis. Furthermore, the 484 485 theory for MYULA presented in this paper is significantly more complete than that currently available for Px-MALA and other MALAs. Finally, MYULA is also more robust and simple to implement than Px-MALA. For example, identifying suitable values of γ for MYULA is straightforward by using the guidelines described above, whereas setting γ for Px-MALA can be challenging and often requires using an adaptive MCMC approach based on a stochastic approximation scheme [36; 17].

4. Experimental results. In this section we illustrate the proposed methodology with four 491canonical imaging inverse problems related to image deconvolution and tomographic recon-492 struction with total-variation and ℓ_1 sparse priors. In the Bayesian setting these problems 493 are typically solved by MAP estimation, which delivers accurate solutions and can be com-494puted very efficiently by using proximal convex optimisation algorithm. Here we demonstrate 495MYULA by performing some advanced and challenging Bayesian analyses that are beyond the 496 scope of optimisation-based mathematical imaging methodologies. For example, in Section 497 4.1 we report two experiments where we use MYULA to perform Bayesian model choice for 498499 image deconvolution models, and where a novelty is that comparisons are performed intrinsically (i.e., without ground truth available) by computing the posterior probability of each 500model given the observed data. Following on from this, in Section 4.2 we report the two ad-501ditional experiments where we use MYULA to explore the posterior uncertainty about x and 502analyse specific aspects about the solutions delivered, particularly by computing simultaneous 503504 credible sets (joint Bayesian confidence sets).

Moreover, to assess the computational efficiency and the accuracy of MYULA we bench-505mark our estimations against the results of Px-MALA [36] targeting the exact posterior 506 $\pi(x) = p(x|y)$ (recall that this algorithm has no asymptotic estimation bias). We empha-507 sise at this point that we do not seek to compare explicitly and quantitatively the methods 508because: 1) MYULA and Px-MALA do not target the exact same stationary distribution; 2) 509 high-dimensional quantitative efficiency comparisons may depend strongly on the summary 510statistics used to define the efficiency metrics; and 3) results can often be marginally improved 511512by fine tuning the algorithm parameters (e.g., step sizes, burn-in periods, etc.). What our comparisons seek to demonstrate is that MYULA can deliver reliable approximate inferences 513with a computational cost that is often significantly lower than Px-MALA, and more impor-514515tantly, that it provides a general, robust, and theoretically sound computational framework for performing advanced Bayesian analyses for imaging problems. In all experiments the Lip-516 517 chitz constant L_f required to set the value γ was computed explicitly. Experiments were 518 conducted on a Apple Macbook Pro computer running MATLAB 2015.

exp:BMS

519

4.1. Bayesian model selection.

4.1.1. Bayesian analysis and computation. Most mathematical imaging problems can be 520solved with a range of alternative models. Currently, the predominant approach to select the 521best model for a specific problem is to compare their estimations against ground truth. For 522example, given K alternative Bayesian models $\mathcal{M}_1, \ldots, \mathcal{M}_K$, practitioners often benchmark 523models by artificially degrading a set of test images, computing the MAP estimator for each 524525model and image, and then measuring estimation error with respect to the truth. The model with the best overall performance is then used in applications to analyse real data. Of course 526527 this approach to model selection has some limitations: 1) it relies strongly on test data that

sec:experiments

528 may not be representative of the unknown, and 2) conclusions can depend on the estimation 529 error metrics used.

An advantage of formulating inverse problems within the Bayesian framework is that, in addition to strategies to perform point estimation, this formalism also provides theory to compare models objectively and intrinsically, and hence perform model selection in the absence of ground truth. Precisely, K alternative Bayesian models are compared through their marginal posterior probabilities

<u>margPost35</u> (19) $p(\mathcal{M}_j|y) = \frac{p(y|\mathcal{M}_j)K^{-1}}{\sum_{k=1}^{K} p(y|\mathcal{M}_k)K^{-1}}, \quad j = \{1, \dots, K\},$

where for objectiveness here we use an uniform prior on the auxiliary variable j indexing the models, $p(y|\mathcal{M}_j)$ is the marginal likelihood

-margLike3 (20)
$$p(y|\mathcal{M}_j) = \int p(x,y|\mathcal{M}_j) \mathrm{d}x, \quad j = \{1,\dots,K\},$$

measuring model-fit-to-data and $p(y, x|\mathcal{M}_j)$ is the joint probability density associated with \mathcal{M}_j (see Appendix B for details regarding the case of improper priors). Following Bayesian decision theory, to perform model selection we simply chose the model with the highest posterior probability (this is equivalent to performing MAP estimation on the model index j):

543
$$\mathcal{M}^* = \underset{j \in \{1, \dots, K\}}{\operatorname{arg\,max}} p(\mathcal{M}_j | y).$$

ssec:exp1

558

From a computation viewpoint, performing Bayesian model selection for imaging problems is challenging because it requires evaluating the likelihoods $p(y|\mathcal{M}_j)$ up to a proportionality constant, or equivalently the Bayes factors $p(y|\mathcal{M}_j)/p(y|\mathcal{M}_i)$ for $i, j \in \{1, \dots, K\}$ (see Appendix C.2 for details regarding the case of improper priors). Here we perform this computation by Monte Carlo integration. Precisely, given n samples X_1^M, \dots, X_n^M from $p(x|y, \mathcal{M}_j)$, we approximate the marginal likelihood of model \mathcal{M}_j by using the truncated harmonic mean estimator [41]

$$\overline{\text{nonicEstimators}} \quad (21) \qquad \qquad p(y|\mathcal{M}_j) \approx \left(\sum_{k=1}^n \frac{\mathbbm{1}_{\mathsf{A}^\star}(X_k^M)}{p(X_k^M, y|\mathcal{M}_j)}\right)^{-1} \operatorname{Vol}(\mathsf{A}^\star) \,, \quad j = \{1, 2, K\}$$

where for all $x, y, p(x, y|\mathcal{M}_j)$ is joint density of \mathcal{M}_j and $\mathsf{A}^* = \bigcup_{j=1}^K \mathcal{C}^*_{j,\alpha}$ is the union of highest posterior density regions (24) of each model at level $(1 - \alpha)$ (see Section 4.2 for details about HPD regions). In our experiments we use the samples to calibrate each $\mathcal{C}^*_{j,\alpha}$ for $\alpha = 0.8$. Notice that it is not necessary to compute Vol(A^*) to calculate (21) because the normalisation is retrieved via $\sum_{j=1}^K p(\mathcal{M}_j|y) = 1$. See Appendix C for more details about this estimator and its use to compute the Bayes factors.

4.1.2. Experiment 1: Image deconvolution with total-variation prior.

Experiment setup. To illustrate the Bayesian model selection approach we consider an 559image deconvolution problem with three alternative models related to three different blur 560 operators. The goal of image deconvolution is to recover a high-resolution image $x \in \mathbb{R}^n$ 561from a blurred and noisy observation y = Hx + w, where H is a circulant blurring matrix 562563 and $w \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$. This inverse problem is ill-conditioned, a difficulty that Bayesian image deconvolution methods address by exploiting the prior knowledge available. For this first 564experiment we consider three alternative models involving three different blur operators H_1 , 565 H_2 , and H_3 . With regards to the prior, we use the popular total-variation prior that promotes 566 regularity by using the pseudo-norm $TV(x) = \|\nabla_d x\|_{1-2}$, where $\|\cdot\|_{1-2}$ is the composite $\ell_1 - \ell_2$ 567 norm and ∇_d is the two-dimensional discrete gradient operator. The posterior distribution 568p(x|y) for the models is given by 569

onvolution
$$\overline{\mathbf{5V0}}$$
 (22)

lecc

$$\mathcal{M}_j: \quad \pi(x) \propto \exp\left[-(\|y - H_j x\|^2/2\sigma^2) - \beta T V(x)\right]$$

with fixed hyper-parameters $\sigma > 0$ and $\beta > 0$ set manually by an expert. This density is logconcave and MAP estimation can be performed efficiently by proximal convex optimisation.

Figure 4 presents an experiment with the Boat test image of size $d = 256 \times 256$ pixels. Figure 4(a) shows a blurred and noisy observation y, generated by using a 5 × 5 uniform blur and Gaussian noise with $\sigma = 0.47$, related to a blurred signal-to-noise ratio of 40dB. Moreover, Figures 4(b)-(d) show the MAP estimates associated with three alternative instances of model involving the following blur operators:

 $578 \\ 579$

580

• \mathcal{M}_1 : H_1 is the correct 5×5 uniform blur operator.

- \mathcal{M}_2 : H_2 is a mildly misspecified 6×6 uniform blur operator.
- \mathcal{M}_3 : H_3 is a strongly misspecified 7×7 uniform blur operator.

581 (All models share the same hyper-parameter values $\sigma = 0.47$ and $\beta = 0.03$ selected manually to produce good image deconvolution results.) We observe in Figure 4 that models \mathcal{M}_1 and 582 \mathcal{M}_2 have produced sharp images with fine detail, whereas \mathcal{M}_3 is clearly misspecified. In terms 583 of estimation performance with respect to the truth, as expected the estimate of Figure 4(c)584corresponding to model \mathcal{M}_1 achieves the highest peak signal-to-noise-ratio (PSNR) of 33.8dB, 585 \mathcal{M}_2 scores 33.4dB, and \mathcal{M}_3 scores 13.4dB. Finally, computing the estimates displayed in Figure 586 4 with a forward-backward optimisation algorithm [17], which is algorithmically similar to 587 MYULA, required approximately 1000 iterations and 30 seconds per model¹. 588

Model selection in the absence of ground truth. We now demonstrate the Bayesian approach 589to perform model selection intrinsically. Precisely, we ran 10^5 iterations of MYULA with 590the specific blur operators corresponding to \mathcal{M}_1 , \mathcal{M}_2 , and \mathcal{M}_3 . For this experiment we implemented MYULA with $f(x) = ||y - H_j x||^2 / 2\sigma^2$ and $g(x) = \beta TV(x)$, with fixed algorithm parameters $\lambda = L_f^{-1} = 0.45$ and $\gamma = L_f^{-1}/5 = 0.1$, and by using Chambolle's algorithm 591592593 [7] to evaluate the proximal operator of the TV-norm. Computing these samples required 594 approximately 30 minutes per model. Following on from this, we used the samples to calibrate 595the high-posterior-density regions \mathcal{C}_{i}^{\star} of each model at level 20%, and then computed the Bayes 596 factors between the models by using (21) (see C.1 for details). 597

¹The computation of the MAP estimates with the SALSA convex optimisation algorithm [1], which is faster than the forward-backward splitting algorithm, required 2 seconds per model.



Figure 4. Deconvolution experiment - Boat test image $(256 \times 256 \text{ pixels})$: (a) Blurred and noisy image y, (b)-(d) MAP estimators corresponding to models \mathcal{M}_1 , \mathcal{M}_2 , and \mathcal{M}_3 .

By applying this procedure we obtained that \mathcal{M}_1 has the highest posterior probability $p(\mathcal{M}_1|y) = 0.964$, followed by $p(\mathcal{M}_2|y) = 0.036$ and $p(\mathcal{M}_3|y) < 0.001$ (the values of the Bayes factors for this experiment are $\hat{B}_{1,2}(y) = 26.8$ and $\hat{B}_{1,3}(y) > 10^3$). These results, which have been computing without using any form of ground truth, are in agreement with the PSNR values calculated by using the true image and provide strong evidence in favour of model \mathcal{M}_1 . They also confirm the good performance of the Bayesian model selection technique.

Comparison with proximal MALA. We conclude this first experiment by benchmarking our 604 estimations against Px-MALA, which targets (22) exactly. Precisely, we recalculated the 605models' posterior probabilities (31) with Px-MALA and obtained that $p(\mathcal{M}_1|y) = 0.962$, 606 $p(\mathcal{M}_2|y) = 0.038$, and $p(\mathcal{M}_3|y) < 0.001$, indicating that the MYULA estimate has an ap-607 proximation error of the order of 0.5% (to obtain accurate estimates for Px-MALA we used 608 $n = 10^7$ iterations with an adaptive time-step targeting an average acceptance rate of order 609 45%). Moreover, comparing the chains generated with MYULA and Px-MALA revealed that 610611 MYULA is significantly more computationally efficient than Px-MALA. For illustration, Fig.

This manuscript is for review purposes only.

FibBoat1



Figure 5. MYULA and Px-MALA comparison: (a) Convergence of the chains to the typical set of (22) under model \mathcal{M}_1 (logarithmic scale), (b) chain autocorrelation function (ACF).

5(a) shows the transient regimes of the MYULA and Px-MALA chains related \mathcal{M}_1 , where 612 starting from a common initial condition the chains converge to the posterior typical set² of 613 p(x|y) (to improve visibility this is displayed in logarithmic scale). Observe that MYULA 614requires around 10^2 iterations to navigate the parameter space and reach the typical set, 615whereas Px-MALA requires 10^4 iterations. Furthermore, to compare the efficiency of the 616 chains in stationarity, Fig. 5(b) shows the autocorrelation function of the chains generated 617 by MYULA and Px-MALA. To highlight the efficiency of MYULA we have used the chains' 618 slowest component³ as summary statistic. Again, observe that MYULA is clearly significantly 619 more efficient than Px-MALA. From a practitioner's viewpoint, this efficiency advantage is 620 further accentuated by the fact that MYULA iterations are almost twice less computationally 621 622 expensive than Px-MALA iterations, which include the MH step.

eriment-2:-image

4.1.3. Experiment 2: Image deconvolution with wavelet frame.

Experiment setup. The second model selection experiment we consider involves three alternative image deconvolution models with different priors. This experiment is more challenging than the previous one because priors operate indirectly on y through x. We consider three models of the form

rolutionL1wave
$$8$$
 (23)

nν

$$\mathcal{M}_j: \quad p(x|y) \propto \exp\left[-(\|y - Hx\|^2/2\sigma^2) - \beta_j\|\Psi_j x\|_1\right]$$

629 where Ψ_i is a model dependent frame:

630 • \mathcal{M}_1 : Ψ_1 is a redundant Haar frame with 6-level, and $\beta_1 = 0.02$ is selected automatically 631 by using a hierarchical Bayesian method [38],

²In stationarity, x|y is with very high probability in the neighbourhood of the (d-1)-dimensional shell $\{x : U(x) = \mathbb{E}[U(x)|y]\}$, see [37]

³The chain's slowest component was identified by doing an approximate singular value decomposition of the chain's covariance matrix and then projecting the samples on the dominant eigenvector.

FibBoat2

632 633 • \mathcal{M}_2 : Ψ_2 is a redundant Haar frame with 3-level, and $\beta_2 = 0.02$ is selected automatically by using a hierarchical Bayesian method [38],

634 635 • \mathcal{M}_3 : Ψ_3 is a redundant Haar frame with 3-level, and $\beta_3 = 0.003$ is selected automatically by using the L-curve method [20].

To make the selection problem even more challenging, in this experiment we use a higher noise level $\sigma = 1.76$, related to a blurred signal-to-noise ratio of 30dB. We note that (23) is logconcave and MAP estimation can be performed efficiently by proximal convex optimisation.

Fig. 6 presents an experiment with the Flinstones test image of size $d = 256 \times 256$ 639 pixels. Fig. 4(a) shows the blurred and noisy observation y used in this experiment, which we 640generated by using a 5 \times 5 uniform blur and $\sigma = 1.76$, and Fig. 6(b)-(d) show the MAP esti-641 mates obtained with \mathcal{M}_1 , \mathcal{M}_2 , and \mathcal{M}_3 by convex optimisation (we used a forward-backward 642 splitting algorithm [17] that is algorithmically similar to MYULA, and which required approx-643 imately 2500 iterations and 2 minutes per model⁴). We observe in Figure 4 that models \mathcal{M}_1 644 and \mathcal{M}_2 have produced sharp images with fine detail, whereas \mathcal{M}_3 is misspecified. In terms of 645 estimation performance with respect to the truth, the estimate of Figure 6(c) corresponding 646 to model \mathcal{M}_2 achieves the highest peak signal-to-noise-ratio (PSNR) of 20.8dB, \mathcal{M}_1 scores 647 20.6dB, and \mathcal{M}_3 scores 11.6dB. 648

649 Model selection in the absence of ground truth. Similarly to the previous experiment, we used MYULA to perform Bayesian model selection intrinsically. Precisely, we used MYULA 650 to generate three sets of $n = 10^5$ samples X_1^M, \ldots, X_n^M approximately distributed according to (23) with the parameters corresponding to $\mathcal{M}_1, \mathcal{M}_2$, and \mathcal{M}_3 . For this experiment we 651 652implemented MYULA with $f(x) = \|y - Hx\|^2/2\sigma^2$ and $g(x) = \beta_j \|\Psi_j x\|_1$, with fixed algorithm 653 parameters $\lambda = L_f^{-1} = 4.5$ and $\gamma = L_f^{-1}/5 = 0.9$. Computing these samples required 50 minutes per model. Following on from this, we used the samples to calibrate the high-posterior-654 655 density regions \mathcal{C}_{j}^{\star} of each model at level 20%, and then computed the Bayes factors between 656 657 the models by using (21) (see C.1 for details).

By applying this procedure we obtained that \mathcal{M}_2 has the highest posterior probability $p(\mathcal{M}_2|y) = 0.42$, followed by $p(\mathcal{M}_1|y) = 0.32$ and $p(\mathcal{M}_3|y) = 0.26$ (the values of the Bayes factors for this experiment are $\hat{B}_{2,1}(y) = 1.31$ and $\hat{B}_{2,3}(y) = 1.62$). Note that these results, which have been computing without using any form of ground truth, are in agreement with the PSNR values calculated by using the true image and indicate that \mathcal{M}_2 is the most appropriate model for data y.

Comparison with proximal MALA. Again, we conclude our second experiment by bench-664 665 marking our estimations against Px-MALA, which targets (23) exactly. Precisely, we recalculated the models' posterior probabilities (31) with Px-MALA and obtained that $p(y|\mathcal{M}_1) =$ 6660.41, $p(y|\mathcal{M}_2) = 0.33$, and $p(y|\mathcal{M}_3) = 0.26$, indicating that the MYULA estimate has an 667 approximation error of the order of 0.5% (to obtain accurate estimates for Px-MALA we used 668 $n = 10^7$ iterations with an adaptive time-step targeting an average acceptance rate of order 669 670 45%). Moreover, efficiency analyses indicate that in this case MYULA is approximately an order of magnitude more efficient per iteration than Px-MALA, with an additional advantage 671 in terms of time-normalised computational efficiency because of a lower computational cost 672

⁴The computation of the MAP estimates with the SALSA convex optimisation algorithm [1], which is faster than the forward-backward splitting algorithm, required 4 seconds per model.

FibFlin



Figure 6. Deconvolution experiment - Flinstones test image $(256 \times 256 \text{ pixels})$: (a) Blurred and noisy image y, (b)-(d) MAP estimators corresponding to models \mathcal{M}_1 , \mathcal{M}_2 , and \mathcal{M}_3 .

per iteration.

673 exp:BUQ

674

4.2. Bayesian uncertainty quantification via posterior credible sets.

675 **4.2.1. Bayesian analysis and computation.** As mentioned earlier, point estimators such 676 as \hat{x}_{MAP} deliver accurate results but do not provide information about the posterior un-677 certainty of x. Given the uncertainty that is inherent to ill-posed and ill-conditioned inverse 678 problems, it would be highly desirable to complement point estimators with posterior credibil-679 ity sets that indicate the region of the parameter space where most of the posterior probability 680 mass of x lies. This is formalised in the Bayesian decision theory framework by computing 681 credible regions [40]. A set C_{α} is a posterior credible region with confidence level $(1 - \alpha)$ if

682
$$\mathbb{P}\left[x \in \mathcal{C}_{\alpha}|y\right] = 1 - \alpha$$

It is easy to check that for any $\alpha \in (0, 1)$ there are infinitely many regions of the parameter space that verify this property. Among all possible regions, the so-called *highest posterior*

This manuscript is for review purposes only.

685 density (HPD) region has minimum volume [40], and is given by

HRDG
$$(24)$$

$$\mathcal{C}^{\star}_{lpha} = \{oldsymbol{x}: U(x) \leq \eta_{lpha}\}$$

687 with $\eta_{\alpha} \in \mathbb{R}$ chosen such that $\int_{\mathcal{C}_{\alpha}^{\star}} p(x|y) dx = 1 - \alpha$ holds. This joint credible set has the 688 important advantage that it can be enumerated by simply specifying the scalar value η_{α} .

From a computation viewpoint, calculating credible sets for images is very challenging because it requires solving very high-dimensional integrals of the form $\int_{\mathcal{C}_{\alpha}^{\star}} p(x|y) dx$. In this work, we use MYULA to approximate these integrals.

4.2.2. Experiment 3: Tomographic image reconstruction.

Experiment setup. The third experiment we consider is a tomographic image reconstruction 693 problem with a total-variation prior. The goal is to recover the image $x \in \mathbb{R}^n$ from an 694 incomplete and noisy set of Fourier measurements y = AFx + w, where F is the discrete 695 Fourier transform operator, A is a tomographic sampling mask, and $w \sim \mathcal{N}(0, \sigma^2 I_n)$. This 696 inverse problem is ill-posed, resulting in significant uncertainty about the true value of x. 697 Similarly to Experiment 1, in this experiment we regularise the problem and reduce the 698 uncertainty about x by using a total-variation prior promoting piecewise regular images. The 699 700 resulting posterior p(x|y) is

tomographild
$$(25)$$

$$\pi(x) \propto \exp\left[-\|y - AFx\|^2/2\sigma^2 - \beta TV(x)\right].$$

with fixed hyper-parameters $\sigma > 0$ and $\beta > 0$ set manually by an expert. We note that this density is log-concave and MAP estimation can be performed efficiently by proximal convex optimisation.

Figure 7 presents an experiment with the Shepp-Logan phantom magnetic resonance image (MRI) of size $d = 128 \times 128$ pixels presented in Figure 7(a). Figure 7(b) shows a noisy tomographic measurement y of this image, contaminated with Gaussian noise with $\sigma = 7 \times 10^{-2}$ (to improve visibility Figure 7(b) shows the amplitude of the Fourier coefficients in logarithmic scale, with black regions representing unobserved coefficients). Notice from Figure 7(b) that only 15% of the original Fourier coefficients are observed. Moreover, Figure 7(c) shows the Bayesian estimate \hat{x}_{MAP} associated with (25) with hyper-parameter value $\beta = 5$.

Bayesian uncertainty analysis. We now conduct a simple Bayesian uncertainty analysis to 712illustrate how posterior credible sets can inform decision-making. For illustration, suppose 713 714that the structure highlighted in red in Figure 7(c) is relevant from a clinical viewpoint because it provides important information for diagnosis or treatment related decision-making. 715Also, suppose that we first observe this structure in the Bayesian estimate \hat{x}_{MAP} and that, 716 following on from this, we wish to explore the posterior uncertainty about x to learn more 717about the structure. In particular, here we conduct a simple analysis to show that there is 718lack of confidence regarding the presence of this structure in the true image (i.e., the structure 719 could be an artefact). Precisely, this is achieved by computing the HDP credible region $\mathcal{C}^{\star}_{\alpha}$ 720 and showing that it includes solutions that are essentially equivalent to \hat{x}_{MAP} except for the 721 fact that they do not have the structure of interest. 722

As alternative solution or "counter example" of \hat{x}_{MAP} , consider the image x_{\dagger} displayed in Figure 7(d). This image is equivalent to \hat{x}_{MAP} except for the fact that the structure of interest



Figure 7. Tomography experiment: (a) Shepp-Logan phantom image $(128 \times 128 \text{ pixels})$, (b) tomographic observation y (amplitude of Fourier coefficients in logarithmic scale), (c) MAP estimator \hat{x}_{MAP} , (d) counter example image x_{\dagger} .

FigMRI1

has been removed (we generated this image by modifying \hat{x}_{MAP} by applying a segmentationinpainting process to replace the structure with the surrounding intensity level). Of course, clinicians observing x_{\dagger} images would potentially arrive to significantly different conclusions about the diagnosis or the treatment required. This test image scores $U(x_{\dagger}) = 1.27 \times 10^4$.

To determine if x_{\dagger} belongs to $\mathcal{C}^{\star}_{\alpha}$ we used MYULA to generate $n = 10^5$ samples from 729 (25), and calculated the HPD threshold η_{α} by estimating the $(1 - \alpha)$ -quantile of U(x) (we 730implemented the algorithm with $f(x) = ||y - AFx||^2/2\sigma^2$ and $g(x) = \beta TV(x)$, with fixed parameters $\lambda = L_f^{-1} = 1 \times 10^{-4}$ and $\gamma_k = L_f^{-1}/10 = 10^{-5}$, and by using Chambolle's algorithm [7] to evaluate the proximal operator of the TV-norm). Fig. 8(a) shows the threshold values η_{α} 731732 733 for a range of values of $\alpha \in [0.01, 0.99]$. Observe that $U(x_{\dagger}) = 1.27 \times 10^4$ is significantly lower 734 than the values displayed in Fig. 8(a), indicating that the counter example image x_{\dagger} belongs 735to set of likely solutions to the inverse problem (e.g., at level 90% $\eta_{0.10} = 2.34 \times 10^4$ hence 736 $x_{\dagger} \in \mathcal{C}_{0,10}^{\star}$, for information $U(\hat{x}_{MAP}) = 1.21 \times 10^4$). Based on this we conclude that, with the 737

current number of observations and noise level, it is not possible to assert confidently that the

- 739 structure considered is present in the true image. Consequently, we would recommend that
- ⁷⁴⁰ this data is not used as primary evidence to support decision-making about this structure.
- Generating the Monte Carlo samples and computing the HPD threshold values required 15minutes.
- *Comparison with proximal MALA*. We conclude this experiment by benchmarking our es-743 timations against Px-MALA, which targets (25) exactly (to obtain accurate estimates for 744Px-MALA we use $n = 10^7$ iterations with an adaptive time-step targeting an average accep-745tance rate of order 45%). The HPD threshold values η_{α} obtained with Px-MALA are reported 746 747 in Fig. 8(a), notice the approximation error of order of 3% due to MYULA's estimation bias (this does not affect the conclusions of the experiment). With regards to computational per-748 formance, an efficiency analysis of the two algorithms indicates that for this model MYULA is 749 approximately two orders of magnitude more efficient than Px-MALA in terms of integrated 750autocorrelation time (for illustration Fig. 8(b) compares the autocorrelation functions for 751 slowest component⁵ of the MYULA and Px-MALA chains. 752
 - 2.36 Px-MALA MYULA 0.9 2.34 0.8 2.32 0.7 2.3 0.6 జ^ర 2.28 Ŭ 0.5 0.4 2.26 0.3 2.24 0.2 2.22 x-MAl 0.1 MYULA 2.2 L 0 0 0.2 0.8 0.4 0.6 500 1500 2000 2500 1000 $1-\alpha$ lag (a) (b)

Figure 8. Tomography experiment: (a) HDP region thresholds η_{α} for MYULA and Px-MALA, (b) chain autocorrelation functions for MYULA and Px-MALA.

FigMRI3

4.2.3. Experiment 4: Sparse image deconvolution with an ℓ_1 prior.

Experiment setup. The fourth experiment we consider is a sparse image deconvolution problem with a Laplace or ℓ_1 prior. Again, we aim to recover $x \in \mathbb{R}^n$ from y = Hx + w, where *H* is a circulant blurring matrix and $w \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$. We expect sparse solutions and use a Laplace prior related to the ℓ_1 norm of x. The resulting posterior p(x|y) is

deconvolutions (26)
$$\pi(x) \propto \exp\left[-\|y - Hx\|^2/2\sigma^2 - \beta \|x\|_1\right].$$

ssec:exp4

753

with fixed hyper-parameters $\sigma > 0$ and $\beta > 0$ set manually by an expert. Similarly to the previous experiments, we notice that this density is log-concave and MAP estimation can be

⁵Again, the chain's slowest component was identified by doing an approximate singular value decomposition of the chain's covariance matrix and then projecting the samples on the dominant eigenvector.

761 performed efficiently by proximal convex optimisation.

Figure 9 presents an experiment with a microscopy dataset of [49] related to high-resolution 762 live cell imaging. Figure 9(a) shows an observation y of field of size $4\mu m \times 4\mu m$ containing 763 100 molecules. This low-resolution observation has been acquired with an instrument specific 764 765 point-spread-function of size 16×16 pixels and a blurred signal-to-noise ratio of 20dB (see [49] for more details). Figure 9(b) shows the Bayesian estimate \hat{x}_{MAP} associated with (26) 766 with hyper-parameter value $\alpha = 0.01$ (notice that \hat{x}_{MAP} is displayed in logarithmic scale to 767 improve visibility). Computing this estimate with a forward-backward splitting optimisation 768 algorithm, which is algorithmically similar to MYULA, required approximately 5 minutes⁶. 769

770 Bayesian uncertainty analysis. As second example of Bayesian uncertainty quantification, we use $\mathcal{C}^{\star}_{\alpha}$ to examine the uncertainty about the position of the group of molecules highlighted 771 in red in Fig. 9, which we assume to be relevant for an application considered. Precisely, we 772 used $n = 10^5$ samples generated with MYULA to compute $\mathcal{C}^{\star}_{\alpha}$ with $\alpha = 0.01$ related to the 773 99% confidence level, and obtained the threshold value $\eta_{0.01} = 9.69 \times 10^4$. Following on from 774 this, to explore $\mathcal{C}_{0.01}^{\star}$ to quantify the uncertainty about the exact position of the molecules, 775 we generated several surrogate test images by modifying \hat{x}_{MAP} by displacing the molecules in 776 different directions until these surrogates exit $C_{0.01}^{\star}$ (similarly to the previous experiment, the 777 resulting empty space was filled by inpainting). Figure 9(c) shows the posterior uncertainty 778 of the molecule positions (note that for visibility the figure focuses on the region of interest). 779 This analysis reveals that the uncertainty at level 99% is of the order of ± 5 pixels vertically 780 and ± 8 pixels horizontally, corresponding to $\pm 78nm$ and $\pm 125nm$. It is worth mentioning 781that these results are in close in agreement with the experimental precision results reported in 782 [49], which identified an average precision of the order of 80nm for the one hundred molecules. 783 Comparison with proximal MALA. Again, we conclude the experiment by benchmarking our 784 estimations against Px-MALA, which targets (26) exactly (to obtain accurate estimates for 785Px-MALA we use $n = 2 \times 10^7$ iterations with an adaptive step-size targeting an acceptance rate 786 of the order of 45%). Figure 9(d) compares the estimations of the threshold values η_{α} obtained 787 with MYULA and Px-MALA for different values of α , indicating that the approximation errors 788 of MYULA are of the order of 0.1%. Moreover, performance analyses based on the chains 789 790 generated with each algorithm indicate that in this case MYULA is approximately one order 791 of magnitude more computationally efficient than Px-MALA.

sec:conclusion 792

5. Discussion and conclusion. This paper presented a new and general proximal MCMC 793 methodology to perform Bayesian computation in log-concave models, with a focus on enabling advanced Bayesian analyses for imaging inverse problems that are convex and not 794smooth, and currently solved mainly by convex optimisation. The methodology is based on a 795 Moreau-Yoshida-type regularised approximation of the target density that is by construction 796 is log-concave and Lipchitz continuously differentiable, and which can be addressed efficiently 797 by using an unadjusted Langevin MCMC algorithm. We provided a detailed theoretical anal-798 ysis of this scheme, including asymptotic as well as non-asymptotic convergence results, and 799 bounds on the convergence rate of the chains with explicit dependence on model dimension. In 800 addition to being highly computational efficient and having a strong theoretical underpinning, 801

⁶The computation of the MAP estimate with the SALSA [1] convex optimisation algorithm, which is faster than the forward-backward splitting algorithm, required 2.3 seconds.



Figure 9. Microscopy experiment: (a) Blurred image y (256 × 256 pixels, $4\mu m \times 4\mu m$)), (b) MAP estimate \hat{x}_{MAP} (logarithmic scale), (c) molecule position uncertainty quantification (vertical: $\pm 78nm$, horizontal $\pm 125nm$), (d) HDP region thresholds η_{α} for MYULA and Px-MALA.

this new methodology is general and can be applied straightforwardly to most problems solved by proximal optimisation, particularly all problems solved by using forward-backward splitting techniques. The proposed methodology was finally demonstrated with four experiments related to image deconvolution and tomographic reconstruction with total-variation and ℓ_1 priors, where we conducted a range of challenging analyses related to model comparison and uncertainty quantification, and where we reported estimation accuracy and computational efficiency comparisons with the proximal Metropolis-adjusted Langevin algorithm.

Furthremore, observe that the regularisation strategy used in this paper, based on the Moreau-Yoshida envelope, can also be applied straightforwardly to the Hamiltonian Monte

This manuscript is for review purposes only.

FigMicro

Carlo algorithm [17]. The theoretical and empirical properties of this algorithm are currently
under investigation and will be reported separately.

Finally, it is worth mentioning that MYULA can also be applied to some models that are 813 not log-concave, for example multi-modal models where the smooth term f is not convex. 814 815 In this case, the non-smooth term q must be associated with a proper prior to guarantee that the approximation π^{λ} is proper. It is possible to derive non-asymptotic convergence 816 results for these models; however, unlike the log-concave case, here the dependence w.r.t. 817 to the dimension of the model is difficult to analyse (see [45] for details). The analysis of 818 the performance of MYULA in non-convex settings is an important perspective for future 819 work. Also, another important perspective for future work is to investigate ways in which 820 the proposed regularisation approach, combined with an appropriate convex relaxation, could 821 enable Langevin and Hamiltonian sampling in high-dimensional spaces that are discrete. 822

6. Acknowledgements. Part of this work was conducted when Marcelo Pereyra held a Marie Curie Intra-European Research Fellowship for Career Development at the School of Mathematics of the University of Bristol, and a Visiting Scholar at the School of Mathematical and Computer Sciences of Heriot-Watt University.

References.

827

829 830

832

834 835

837

838

840

841 842

844

846

847

848

850 851

853

Figueiredo20118

Atchade20161

iaux:guillin:20083

aario:laine:20146

Bonettini2013

Candes:2013

Chambosle5

Chan20349

Combettes20112

- M. AFONSO, J. M. BIOUCAS-DIAS, AND M. A. T. FIGUEIREDO, An augmented Lagrangian approach to the constrained optimization formulation of imaging inverse problems, IEEE. Trans. on Image Process., 20 (2011), pp. 681–695.
- Y. F. ATCHADÉ, A Moreau-Yosida approximation scheme for a class of high-dimensional posterior distributions, ArXiv e-prints, (2015), arXiv:1505.07072.
- [3] D. BAKRY, F. BARTHE, P. CATTIAUX, AND A. GUILLIN, A simple proof of the Poincaré inequality for a large class of probability measures., Electronic Communications in Probability [electronic only], 13 (2008), pp. 60–66, http://eudml.org/doc/225690.
- [4] J. M. BARDSLEY, A. SOLONEN, H. HAARIO, AND M. LAINE, Randomize-then-optimize: A method for sampling from posterior distributions in nonlinear inverse problems, SIAM Journal on Scientific Computing, 36 (2014), pp. A1895–A1910.
- [5] S. BONETTINI, A. CORNELIO, AND M. PRATO, A new semiblind deconvolution approach for Fourier-based image restoration: An application in astronomy, SIAM J. Imaging Sci., 6 (2013), pp. 1736–1757, doi:10.1137/120873169, http://dx.doi.org/10.1137/120873169, arXiv:http://dx.doi.org/10.1137/120873169.
- [6] E. J. CANDÈS, Y. C. ELDAR, T. STROHMER, AND V. VORONINSKI, Phase retrieval via matrix completion, SIAM J. Imaging Sci., 6 (2013), pp. 199–225.
- [7] A. CHAMBOLLE, An algorithm for total variation minimization and applications, Journal of Mathematical Imaging and Vision, 20 (2004), pp. 89–97, doi:10.1023/B:JMIV.0000011325.36760.1e, http://dx.doi.org/10.1023/B:JMIV.0000011325.36760.1e.
- [8] R. H. CHAN, J. YANG, AND X. YUAN, Alternating direction method for image inpainting in wavelet domains, SIAM J. Imaging Sci., 4 (2011), pp. 807–826, doi:10.1137/100807247, http://dx.doi.org/10.1137/100807247, arXiv:http://dx.doi.org/10.1137/100807247.
 - [9] P. L. COMBETTES AND J.-C. PESQUET, Fixed-Point Algorithms for Inverse Problems in Science and Engineering, Springer New York, New York, NY, 2011, ch. Proximal

854		Splitting Methods in Signal Processing, pp. 185–212.
pesquet:20115	[10]	P. L. COMBETTES AND JC. PESQUET, Proximal splitting methods in signal process-
856		ing, in Fixed-Point Algorithms for Inverse Problems in Science and Engineering, H. H.
857		Bauschke, R. S. Burachik, P. L. Combettes, V. Elser, D. R. Luke, and H. Wolkowicz,
858		eds., Springer New-York, 2011, pp. 185–212.
dalalyan:20149	[11]	A. S. DALALYAN, Theoretical guarantees for approximate sampling from smooth and log-
860		concave densities, Journal of the Royal Statistical Society: Series B (Statistical Methodol-
861		ogy), (2016), pp. n/a–n/a, doi:10.1111/rssb.12183, http://dx.doi.org/10.1111/rssb.12183.
Donoho20062	[12]	D. L. DONOHO, Compressed sensing, IEEE Trans. Inf. Theory, 52 (2006), pp. 1289–1306.
s:moulines:20163	[13]	A. DURMUS AND E. MOULINES, Non-asymptotic convergence analysis for the unadjusted
864		langevin algorithm, Accepted for publication in Ann. Appl. Probab. 1507.05021, arXiv,
865		July 2015, http://arxiv.org/pdf/1507.05021v1.pdf.
eberle:20366	[14]	A. EBERLE, Reflection couplings and contraction rates for diffusions, Probab. Theory
867		Related Fields, (2015), pp. 1–36, doi:10.1007/s00440-015-0673-1, http://dx.doi.org/10.
868		1007/s00440-015-0673-1.
ermak:19869	[15]	D. L. ERMAK, A computer simulation of charged particles in solution. i. technique and
870		equilibrium properties, The Journal of Chemical Physics, 62 (1975), pp. 4189–4196.
zano:levan:2001	[16]	M. FLORENZANO AND C. LE VAN, Finite dimensional convexity and optimization, vol. 13
872		of Studies in Economic Theory, Springer-Verlag, Berlin, 2001, doi:10.1007/978-3-642-
873		$56522\mathchar`eq$, http://dx.doi.org/10.1007/978-3-642-56522-9. In cooperation with Pascal Gour-
874		del.
Green20155	[17]	P. J. GREEN, K. ŁATUSZYŃSKI, M. PEREYRA, AND C. P. ROBERT, Bayesian compu-
876		tation: a summary of the current state, and samples backwards and forwards, Statistics
877		and Computing, 25 (2015), pp. 835–862, doi:10.1007/s11222-015-9574-5, http://dx.doi.
878		org/10.1007/s11222-015-9574-5.
grenander:19839	[18]	U. GRENANDER, Tutorial in pattern theory. Division of Applied Mathematics, Brown
880		University, Providence, 1983.
nder:miller:19941	[19]	U. GRENANDER AND M. I. MILLER, Representations of knowledge in complex sys-
882		<i>tems</i> , J. Roy. Statist. Soc. Ser. B, 56 (1994), pp. 549–603, http://links.jstor.org/sici?
883		sici=0035-9246(1994)56:4(549:ROKICS)2.0.CO;2-2&origin=MSN. With discussion and
884	[0.0]	a reply by the authors.
Hanke19985	[20]	M. HANKE AND C. HANSEN, Regularization methods for large-scale problems, Surv.
886		Math. Ind., 3 (1993), pp. 253–315.
Haro:20187	[21]	G. HARO, A. BUADES, AND JM. MOREL, Photographing paintings by image fusion,
888		SIAM J. Imaging Sci., 5 (2012), pp. 1055–1087, doi:10.1137/120873923, http://dx.doi.
889		org/10.1137/120873923, arXiv:http://dx.doi.org/10.1137/120873923.
somersalo:20050	[22]	J. KAIPIO AND E. SOMERSALO, Statistical and Computational Inverse Problems,
891	[00]	Springer, New-York, 2005.
khasminskii:19602	[23]	K. Z. KHAS MINSKII, Ergodic properties of recurrent diffusion processes and stabilization
893		of the solution to the cauchy problem for parabolic equations, Theory of Probability &
894		Its Applications, 5 (1960), pp. 179–196, doi:10.1137/1105016, http://dx.doi.org/10.1137/
895	[0.4]	1105016, arAiv:http://dx.doi.org/10.113//1105016.
Lebrun:20136	[24]	M. LEBRUN, A. BUADES, AND J. M. MOREL, A nonlocal Bayesian image denoising
897		<i>algorithm</i> , SIAM J. Imaging Sci., 6 (2013), pp. 1665–1688.

- [25] F. LUCKA, Fast markov chain monte carlo sampling for sparse bayesian inference in highdimensional inverse problems using l1-type priors, Inverse Problems, 28 (2012), p. 125012.
- [26] F. LUCKA, Fast gibbs sampling for high-dimensional bayesian inversion, Inverse Problems, 32 (2016), p. 115019, http://stacks.iop.org/0266-5611/32/i=11/a=115019.
- [27] M. LUSTIG, D. DONOHO, AND J. M. PAULY, Sparse mri: The application of compressed sensing for rapid mr imaging, Magnetic Resonance in Medicine, 58 (2007), pp. 1182–1195.
- [28] J.-M. MARIN AND C. ROBERT, Bayesian core: a practical approach to computational Bayesian statistics, Springer Science & Business Media, 2007.
- [29] S. MEYN AND R. TWEEDIE, Markov Chains and Stochastic Stability, Cambridge University Press, New York, NY, USA, 2nd ed., 2009.
- [30] V. I. MORGENSHTERN AND E. J. CANDÈS, Super-resolution of positive sources: The discrete setup, SIAM J. Imaging Sci., 9 (2016), pp. 412–444.
- neal:19920[31]R. M. NEAL, Bayesian learning via stochastic dynamics, in Advances in Neural Infor-
mation Processing Systems 5, [NIPS Conference], San Francisco, CA, USA, 1993, Mor-
gan Kaufmann Publishers Inc., pp. 475–482, http://dl.acm.org/citation.cfm?id=645753.913667903.
 - [32] N. PARIKH AND S. BOYD, *Proximal algorithms*, Foundations and Trends in Optimization, 1 (2013), pp. 123–231.
 - [33] G. PARISI, Correlation functions and computer simulations, Nuclear Physics B, 180 (1981), pp. 378–384.
 - [34] T. PARK AND G. CASELLA, The Bayesian lasso, J. Amer. Statist. Assoc., 103 (2008), pp. 681–686, doi:10.1198/016214508000000337, http://dx.doi.org/10.1198/ 016214508000000337.
 - [35] M. PEREYRA, Maximum-a-posteriori estimation with Bayesian confidence regions, SIAM 922 J. Imaging Sci. to appear.
 - [36] M. PEREYRA, Proximal Markov chain Monte Carlo algorithms, Statistics and Computing, (2015). open access paper, http://dx.doi.org/10.1007/s11222-015-9567-4.
 - [37] M. PEREYRA, Maximum-a-posteriori estimation with Bayesian confidence regions, ArXiv e-prints, (2016), arXiv:1602.08590.
 - [38] M. PEREYRA, J. M. BIOUCAS-DIAS, AND M. A. T. FIGUEIREDO, Maximum-aposteriori estimation with unknown regularisation parameters, in Proc. European Signal Proc. Conf. (EUSIPCO), Nice, France, Sep. 2015., Aug 2015, pp. 230–234.
 - [39] M. PEREYRA, P. SCHNITER, E. CHOUZENOUX, J.-C. PESQUET, J.-Y. TOURNERET, A. HERO, AND S. MCLAUGHLIN, A survey of stochastic simulation and optimization methods in signal processing, IEEE. J. Selected Topics in Signal Process., 10 (2016), pp. 224–241.
 - [40] C. P. ROBERT, The Bayesian Choice (second edition), Springer Verlag, New-York, 2001.
 - [41] C. P. ROBERT AND D. WRAITH, Computational methods for Bayesian model choice, AIP Conf. Proc., 1193 (2009), pp. 251–262.
 - [42] G. O. ROBERTS AND R. L. TWEEDIE, Exponential convergence of Langevin distributions and their discrete approximations, Bernoulli, 2 (1996), pp. 341–363, doi:10.2307/3318418, http://dx.doi.org/10.2307/3318418.
 - [43] R. T. ROCKAFELLAR AND R. J.-B. WETS, *Variational analysis*, vol. 317 of Grundlehren der Mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences],

30

lucka:20128

lucka:2096

Lustig:20072

rt:2007:bayesian

n:tweedie:2009

Veniamin2096

Parikh20134

parisi:19816

ark:casella:2008a8

Pereyra2016B1

Pereyra2018

Pereyra:20168

EUSIPCO'20187

Pereyra2016

899

903

905

907

909

915

917

919

920

924

926

928

929

931

932

933

936

938

939

941

cprbayes

Robert2009A98

e-Langevin:19967

afellar:wets:1998

Springer-Verlag, Berlin, 1998, doi:10.1007/978-3-642-02431-3, http://dx.doi.org/10.1007/978-3-642-02431-3.

942

943

945 946

947

949

950

952

954

955

957

958 959

961

964

965

Zhu20920

ll:friedman:19784

alay:tubaro:19948

Vollmer20151

solonen:al:2016

Zhengming:20926

roof-crefpr-meas

fun-convex-gene

q:max

- [44] P. J. ROSSKY, J. D. DOLL, AND H. L. FRIEDMAN, Brownian dynamics as smart Monte Carlo simulation, The Journal of Chemical Physics, 69 (1978), pp. 4628–4633, doi:http://dx.doi.org/10.1063/1.436415, http://scitation.aip.org/ content/aip/journal/jcp/69/10/10.1063/1.436415.
- [45] D. TALAY AND L. TUBARO, Expansion of the global error for numerical schemes solving stochastic differential equations, Stochastic Anal. Appl., 8 (1990), pp. 483–509 (1991), doi:10.1080/07362999008809220, http://dx.doi.org/10.1080/07362999008809220.
- [46] S. J. VOLLMER AND K. C. ZYGALAKIS, (Non-) asymptotic properties of Stochastic Gradient Langevin Dynamics, ArXiv e-prints, (2015), arXiv:1501.00438.
- [47] Z. WANG, J. M. BARDSLEY, A. SOLONEN, T. CUI, AND Y. M. MARZOUK, Bayesian inverse problems with l_1 priors: a randomize-then-optimize approach, arXiv preprint arXiv:1607.01904, (2016).
- [48] Z. XING, M. ZHOU, A. CASTRODAD, G. SAPIRO, AND L. CARIN, Dictionary learning for noisy and incomplete hyperspectral images, SIAM J. Imaging Sci., 5 (2012), pp. 33–56, doi:10.1137/110837486, http://dx.doi.org/10.1137/110837486, arXiv:http://dx.doi.org/10.1137/110837486.
- [49] L. ZHU, W. ZHANG, D. ELNATAN, AND B. HUANG, Faster STORM using compressed sensing, Nat. Meth., 9 (2012), pp. 721–723.

Appendix A. Proof of Proposition 1. We preface the proof by a Lemma.

Lemma 4. Let $g : \mathbb{R}^d \to (-\infty, +\infty]$ be a lower bounded, l.s.c convex function satisfying $0 < \int_{\mathbb{R}^d} e^{-g(y)} dy < +\infty$. Then there exists $x_g \in \mathbb{R}^d$, $R_g, \rho_g > 0$ such that for all $x \in \mathbb{R}^d$, $x \notin B(x_g, R_g), g(x) - g(x_g) \ge \rho_g ||x - x_g||$.

Proof. The proof is a simple extension of the one of [3, Theorem 2.2.2], where g is assumed to be continuously differentiable.

We first show that g is finite on a non-empty open set of \mathbb{R}^d . Note since $\int_{\mathbb{R}^d} e^{-g(y)} dy > 0$, the set $\{g < \infty\}$ can not be contained in a k-dimensional hyperplane, for $k \in \{0, \dots, d-1\}$. Then, there exists d+1 points $\{v_i\}_{0 \le i \le d} \subset \{g < \infty\}$ such that the vectors $\{v_i - v_0\}_{1 \le i \le d}$ are linearly independent. Denote by $co(v_0, \dots, v_d)$ the convex hull of $\{v_i\}_{0 \le i \le d}$ defined by

$$\operatorname{co}(\mathbf{v}_0,\cdots,\mathbf{v}_d) = \left\{ \sum_{i=0}^d \alpha_i \mathbf{v}_i \mid \sum_{i=0}^d \alpha_i = 1, \forall i \in \{0,\cdots,d\}, \ \alpha_i \ge 0 \right\} .$$

968 Since g is convex and $co(v_0, \dots, v_d) \subset \{g < \infty\}$, we have

$$\underbrace{\text{conv'hull6}}_{y \in \text{co}(\mathbf{v}_0, \cdots, \mathbf{v}_d)} |\mathbf{g}(y)| \le M_{\text{co}} = \max_{i \in \{0, \cdots, d\}} \{ |\mathbf{g}(\mathbf{v}_i)| \} .$$

970 It follows from $\{v_i\}_{0 \le i \le d} \subset \{g < \infty\}$ and g is lower bounded that M_{co} is finite. Finally by 971 [16, Lemma 1.2.1], $co(v_0, \cdots, v_d)$ has non empty interior.

Consider now the set $\{g \leq M_{co} + 1\}$. We prove by contradiction that it is a bounded subset of \mathbb{R}^d . Assume that for all $R \geq 0$, there exists $x_R \in \{g \leq M_{co} + 1\}$ and $x_R \notin B(v_0, R)$. Then since $\{g \leq M_{co} + 1\}$ is convex, it contains the convex hull of $\{v_0, \dots, v_d, x_R\}$. Since 975 $\operatorname{co}(v_0, \dots, v_d)$ has non empty interior, the volume of $\operatorname{co}(v_0, \dots, v_d, x_R)$ grows at least linearly 976 in R and the volume corresponding to $\{g \leq M_{co} + 1\}$ is infinite taking the limit as R goes to 977 ∞ . On the other hand, by assumption and since $\{v_0, \dots, v_d, x_R\} \subset \{g \leq M_{co} + 1\}$, we have 978 using the Markov inequality

979
$$\operatorname{Vol}(\{g \le M_{\rm co} + 1\}) \le e^{M_{\rm co} + 1} \int_{\{g \le M_{\rm co} + 1\}} e^{-g(y)} \mathrm{d}y < +\infty$$

which leads to a contradiction. Then there exists $R_{g} \geq 0$, such that $\{g \leq M_{co}+1\} \subset B(v_{0}, R_{g})$. For all $x \notin B(v_{0}, R_{g})$, consider $y = R_{g}(x - v_{0}) ||x - v_{0}||^{-1} + v_{0}$. Note that $y \notin \{g \leq M_{co}+1\}$, so $g(y) \geq M_{co} + 1$. Now using the convexity of g, we have for all $x \notin B(v_{0}, R_{g})$,

983
$$M_{\rm co} + 1 \le g(y) \le R_{\rm g} \|x - v_0\|^{-1} (g(x) - g(v_0)) + g(v_0) .$$

984 Since $g(v_0) \leq M_{co}$, we get

985
$$(g(x) - g(v_0)) \ge R_g^{-1} ||x - v_0||$$

986 and the proof is concluded setting
$$x_{\rm g} = v_0$$
.

Proof of Proposition 1. a) We first assume that H 2-(i) holds. By (6), $U \ge U^{\lambda}$ and therefore $0 < \int_{\mathbb{R}^d} e^{-U(y)} dy < \int_{\mathbb{R}^d} e^{-U^{\lambda}(y)} dy$. We now prove $e^{-g^{\lambda}}$ is integrable with respect to the Lebesgue measure, which implies $y \mapsto e^{-U^{\lambda}(y)}$ is integrable as well since f is assumed to be lower bounded. By H1 and Lemma 4, there exist $\rho_g > 0$, $x_g \in \mathbb{R}^d$ and $M_1 \in \mathbb{R}$ such that for all $x \in \mathbb{R}^d$, $g(x) - g(x_g) \ge M_1 + \rho_g ||x - x_g||$. Thus, for all $x \in \mathbb{R}^d$, we have by (6) and (9)

992
$$g^{\lambda}(x) - g(x_g) = g(\operatorname{prox}_g^{\lambda}(x)) - g(x_g) + (2\lambda)^{-1} \left\| x - \operatorname{prox}_g^{\lambda}(x) \right\|^2$$

993

$$\geq M_1 + \rho_g \left\| \operatorname{prox}_g^{\lambda}(x) - x_g \right\| + (2\lambda)^{-1} \left\| x - \operatorname{prox}_g^{\lambda}(x) \right\|^2$$

te-measure-M¥94 995

1001

(28)
$$\geq M_1 + \inf_{y \in \mathbb{R}^d} \{ \rho_g \| y - x_g \| + (2\lambda)^{-1} \| x - y \|^2 \} \geq M_1 + h^{\lambda}(x) ,$$

where $h^{\lambda}(x)$ is the λ -Moreau Yosida envelope of $h(x) = \rho_g ||x - x_g||$. By [32, Section 6.5.1], the proximal operator associated with the norm is the block soft thresholding given for all $\lambda > 0$ and $x \in \mathbb{R}^d \setminus \{0\}$ by $\operatorname{prox}_{h}^{\lambda}(x) = \max(0, 1 - \lambda/||x||)x$ and $\operatorname{prox}_{h}^{\lambda}(0) = 0$. Therefore using again (6), it follows that there exists $M_2 \in \mathbb{R}$ such that for all $x \in \mathbb{R}^d$,

 $\mathbf{h}^{\lambda}(x) \ge \rho_a \|x - x_a\| + M_2 \; .$

Combining this inequality with (28) concludes the proof.

1002 We now assume that H2-(ii) holds. First, we show that for all $\lambda > 0$

$$\boxed{-\operatorname{eq:unifprox}(3)} \quad (29) \qquad \qquad \sup_{x \in \mathbb{R}^d} \{g(x) - g^{\lambda}(x)\} \le \lambda \left\|g\right\|_{\operatorname{Lip}}^2 / 2 ,$$

1004 Indeed if this inequality holds, then for all $x \in \mathbb{R}^d$, we have

1005
$$f(x) + g(x) - \lambda \|g\|_{\text{Lip}}^2 / 2 \le f(x) + g^{\lambda}(x) .$$

1006 Therefore by assumption

1007
$$\int_{\mathbb{R}^d} e^{-U^{\lambda}(x)} dx \le e^{\lambda \|g\|_{\text{Lip}}^2/2} \int_{\mathbb{R}^d} e^{-U(x)} dx < +\infty .$$

1008 We now prove (29). Using that g is Lipschitz, we have by (6), for all $x \in \mathbb{R}^d$

$$1009 \quad g(x) - g^{\lambda}(x) = g(x) - \inf_{y \in \mathbb{R}^d} \left\{ g(y) + (2\lambda)^{-1} \|x - y\|^2 \right\} = \sup_{y \in \mathbb{R}^d} \left\{ g(x) - g(y) - (2\lambda)^{-1} \|x - y\|^2 \right\}$$

$$1010 \\ 1011 \qquad \qquad \leq \sup_{y \in \mathbb{R}^d} \left\{ \|g\|_{\operatorname{Lip}} \|x - y\| - (2\lambda)^{-1} \|x - y\|^2 \right\} \leq \lambda \|g\|_{\operatorname{Lip}}^2 / 2 ,$$

1012 where we have used that the maximum of $u \mapsto au - bu^2$, for $a, b \ge 0$, is given by $a^2/(4b)$.

1013 b) This point is a straightforward consequence of (8) and (7).

1014 c) Since π has also a density with respect to the Lebesgue measure and $U^{\lambda}(x) \leq U(x)$ for all 1015 $x \in \mathbb{R}^d$, we have for all $\lambda > 0$

$$\|\pi^{\lambda} - \pi\|_{\mathrm{TV}} = \int_{\mathbb{R}^d} \left|\pi^{\lambda}(x) - \pi(x)\right| \,\mathrm{d}x \le 2A_{\lambda} ,$$

1017 where $A_{\lambda} = \int_{\mathbb{R}^d} \{1 - e^{g^{\lambda}(x) - g(x)}\} \pi^{\lambda}(x) dx = 1 - \left\{\int_{\mathbb{R}^d} e^{-U^{\lambda}(x)} dx\right\}^{-1} \int_{\mathbb{R}^d} e^{-U(x)} dx$. By (10), for 1018 all $x \in \mathbb{R}^d$, we get $\lim_{\lambda \downarrow 0} \uparrow U^{\lambda}(x) = U(x)$. We conclude by applying the monotone convergence 1019 theorem.

1020 d) Using that for all $x \in \mathbb{R}^d$, $g^{\lambda}(x) \leq g(x)$ and $1 - e^{-u} \leq u$ for all $u \geq 0$, (30) shows that

1021
$$\|\pi^{\lambda} - \pi\|_{\mathrm{TV}} \le 2 \int_{\mathbb{R}^d} \{g(x) - g^{\lambda}(x)\} \pi^{\lambda}(x) \mathrm{d}x \; .$$

1022 Then the proof follows from (29).

proper-imp priors**Appendix B. Model selection using improper priors.** Model selection using improper102410241025priors can lead to tedious considerations [40]. Indeed, in that case the joint density of each1026model is not defined. However, this difficulty can be avoided when the considered models share1026the same improper prior distribution see [28]. Let $\mathcal{M}_1, \ldots, \mathcal{M}_K$ be K alternative Bayesian1027models having the same improper distribution with density $\tilde{p}(x)$ on \mathbb{R}^d and associated to the1028family of likelihood functions $p_i(y|x)$ such that for all $i \in \{1, \ldots, K\}, \int_{\mathbb{R}^d} p_i(y|x)\tilde{p}(x)dx < +\infty$.1029The marginal posterior probabilities of $\mathcal{M}_1, \ldots, \mathcal{M}_K$ are then defined by

$$\boxed{-\operatorname{margPdst30}} \quad (31) \qquad \qquad \tilde{p}(\mathcal{M}_j|y) = \frac{\tilde{p}(y|\mathcal{M}_j)K^{-1}}{\sum_{k=1}^{K} \tilde{p}(y|\mathcal{M}_k)K^{-1}}, \quad j \in \{1, \dots, K\},$$

1031 where for all $j \in \{1, ..., K\}$,

$$\tilde{p}(y|\mathcal{M}_j) = \int_{\mathbb{R}^d} p_i(y|x)\tilde{p}(x)\mathrm{d}x$$

HME3

1032

oı

Appendix C. Truncated harmonic mean estimator.

C.1. Case of proper prior distributions. Consider a positive probability density p on $\mathbb{R}^d \times \mathbb{R}^m$ for $d, m \in \mathbb{N}^*$ of the form: $p(x, y) = f(x, y) / \int_{\mathbb{R}^d \times \mathbb{R}^m} f(z, w) dz dw$. Assume that f is 1035 known but not the normalization constant of p. Here p plays the role of a joint distribution 1036 of the data and the parameters. It can be defined if we take a proper prior distribution for 1037 1038 the parameters. Define for any bounded Borel set $A \in \mathcal{B}(\mathbb{R}^d)$

-harmonicmebn39

harmonic meblo43

onicmear

case-proper-prior

1034

1040 1041

$$I_{\mathcal{A}}(f,y) = \int_{\mathbb{R}^d} \mathbb{1}_{\mathcal{A}}(x) \frac{p(x|y)}{f(x,y)} dx$$
$$= \int_{\mathbb{R}^d} \mathbb{1}_{\mathcal{A}}(x) \frac{p(x|y)}{p(x,y)} dx \Big/ \int_{\mathbb{R}^d \times \mathbb{R}^m} f(z,w) dz dw$$

Since p(x|y) = p(x,y)/p(y), the following identity holds 1042

(33)
$$p(y) = \operatorname{Vol}(A) \left\{ I_A(f, y) \int_{\mathbb{R}^d \times \mathbb{R}^m} f(z, w) \mathrm{d} z \mathrm{d} w \right\}^{-1}$$

For all $y \in \mathbb{R}^m$ and $A \in \mathcal{B}(\mathbb{R}^d)$, we define the truncated harmonic mean estimator of $I_A(f, y)$ 1044 1045bv

$$\hat{I}_{A}(f,y) = \sum_{k=1}^{n} \frac{\mathbb{1}_{A}(X_{k})}{f(X_{k},y)} ,$$

where $(X_k)_{k>1}$ is an ergodic Markov chain targeting p(x|y) to ensure that the defined estimator 1047 almost surely converges to $I_A(f, y)$ given by (32). 1048

Let p_1, p_2 be two positive distributions on $\mathbb{R}^d \times \mathbb{R}^m$, associated with their two unormalized 1049 versions $f_1, f_2 : \mathbb{R}^d \times \mathbb{R}^m \to \mathbb{R}_+$. We aim to estimate $p_1(y)/p_2(y)$. By (33), we have 1050

1051
$$\frac{p_1(y)}{p_2(y)} = \frac{\int_{\mathbb{R}^d \times \mathbb{R}^m} f_2(z, w) dz dw}{\int_{\mathbb{R}^d \times \mathbb{R}^m} f_1(z, w) dz dw} \frac{I_A(f_2, y)}{I_A(f_1, y)}$$

Using (34), we estimate this ratio by 1052

1053
$$\frac{p_1(y)}{p_2(y)} \approx \hat{B}_{1,2}(y) = \frac{\int_{\mathbb{R}^d \times \mathbb{R}^m} f_2(z, w) dz dw}{\int_{\mathbb{R}^d \times \mathbb{R}^m} f_1(z, w) dz dw} \frac{\hat{I}_A(f_2, y)}{\hat{I}_A(f_1, y)}$$

However, we need to compute the ratio $\int_{\mathbb{R}^d \times \mathbb{R}^m} f_2(z, w) dz dw / \int_{\mathbb{R}^d \times \mathbb{R}^m} f_1(z, w) dz dw$. 1054

Assume that for $i = 1, 2, f_i(x, y) = h_i(x, y)g_i(x)$, for some measurable functions h_i : 1055 $\mathbb{R}^d \times \mathbb{R}^m \to \mathbb{R}^*_+, g_i : \mathbb{R}^d \to \mathbb{R}^*_+$ such that $\int_{\mathbb{R}^m} h_i(x, y) dy$ does not depend on x. Note that this 1056assumption holds in Section 4.1.3. We distinguish two cases: 1057

1. If for $i = 1, 2, g_i$ is integrable, we get 1058

1059
$$\hat{B}_{1,2}(y) = \frac{\int_{\mathbb{R}^d} g_2(z) dz}{\int_{\mathbb{R}^d} g_1(z) dz} \frac{\hat{I}_{A}(f_2)}{\hat{I}_{A}(f_1)}$$

In the case where the ratio $\int_{\mathbb{R}^d} g_2(z) dz / \int_{\mathbb{R}^d} g_1(z) dz$ is unknown, such as with the priors 1060 considered in the experiment reported in Section 4.1.3, we use a Monte Carlo algorithm 1061

(32)

1062	such as MYULA or Px-MALA to compute it. Observe that this computation can be
1063	performed offline when the ratio does not depend on the value of y .
1064	2. If there exists a function $g: \mathbb{R}^d \to \mathbb{R}^*_+$ and two real numbers $\lambda_1, \lambda_2 > 0$ such that for
1065	$i = 1, 2, g_i(x) = g(\lambda_i x)$ for all $x \in \mathbb{R}^d$, we get for all $R > 0$
1066	
1067	$\int_{\mathbb{R}^d \times \mathbb{R}^m} \mathbb{1}_{\mathrm{B}(0,R)} f_2(z,w) \mathrm{d}z \mathrm{d}w / \int_{\mathbb{R}^d \times \mathbb{R}^m} \mathbb{1}_{\mathrm{B}(0,\lambda_1 \lambda_2^{-1} R)} f_1(z,w) \mathrm{d}z \mathrm{d}w$
1068 1069	$= \int_{\mathbb{R}^d} \mathbb{1}_{\mathrm{B}(0,R)} g_2(z) \mathrm{d}z / \int_{\mathbb{R}^d} \mathbb{1}_{\mathrm{B}(0,\lambda_1\lambda_2^{-1}R)} g_1(z) \mathrm{d}z = (\lambda_1/\lambda_2)^d \ .$
1070	Since for all $a > 0$ and $i = 1, 2$,

1071
$$\int_{\mathbb{R}^d \times \mathbb{R}^m} f_i(z, w) dz dw = \lim_{R \to +\infty} \int_{\mathbb{R}^d \times \mathbb{R}^m} \mathbb{1}_{\mathrm{B}(0, aR)} f_i(z, w) dz dw ,$$
1072 we get

ълат

we get

NANZTIT A

1073

1074

1075

proper-imp prior

rginal improper76

107

$$\int_{\mathbb{R}^d \times \mathbb{R}^m} f_2(z, w) \mathrm{d}z \mathrm{d}w \bigg/ \int_{\mathbb{R}^d \times \mathbb{R}^m} f_1(z, w) \mathrm{d}z \mathrm{d}w = (\lambda_1/\lambda_2)^d$$

C.2. Case of improper prior distributions. Let $f : \mathbb{R}^d \times \mathbb{R}^m \to \mathbb{R}_+$ such that for all $y \in \mathbb{R}^m$,

(35)
$$\tilde{p}(y) = \int_{\mathbb{R}^d} f(x, y) \mathrm{d}x < +\infty \; .$$

Here, f plays the role of an improper joint density of the data and the parameters as the prior 1077distribution is improper. This setting corresponds to Section 4.1.2. Define for all $y \in \mathbb{R}^m$ 1078 the conditional distribution on $\mathbb{R}^d \times \mathbb{R}^m$ by $p(x|y) = f(x,y)/\tilde{p}(y)$, where \tilde{p} is defined by (35). 1079 Then, define for any bounded Borel set $A \in \mathcal{B}(\mathbb{R}^d)$ 1080

$$I_{\mathcal{A}}(f,y) = \int_{\mathbb{R}^d} \mathbb{1}_{\mathcal{A}}(x) \frac{p(x|y)}{f(x,y)} dx .$$

Then by (35), we get 1082

(37)

mean improper83

$$\tilde{p}(y) = \operatorname{Vol}(A)/I_A(f, y)$$

For all $y \in \mathbb{R}^m$ and $A \in \mathcal{B}(\mathbb{R}^d)$, we define the truncated harmonic mean estimator of $I_A(f, y)$ 1084 1085as in Appendix C.1 by (34).

Let now $f_1, f_2 : \mathbb{R}^d \times \mathbb{R}^m \to \mathbb{R}_+$, satisfying for all i = 1, 2 and $y \in \mathbb{R}^m$, $\tilde{p}_i(y) =$ 1086 $\int_{\mathbb{R}^d} f_i(x,y) dx < +\infty$. We aim to estimate $\tilde{p}_1(y) / \tilde{p}_2(y)$. But by (37), we have 1087

1088
$$\frac{\tilde{p}_1(y)}{\tilde{p}_2(y)} = \frac{I_{\rm A}(f_2, y)}{I_{\rm A}(f_1, y)}$$

Using (36) and (34), we estimate this ratio by 1089

1090
$$\frac{\tilde{p}_1(y)}{\tilde{p}_2(y)} \approx \hat{B}_{1,2}(y) = \frac{I_A(f_2, y)}{\hat{I}_A(f_1, y)}$$