

Sampling from log-concave density

Alain Durmus, Eric Moulines, Marcelo Pereyra

Telecom ParisTech, Ecole Polytechnique, Bristol University

1 Motivation

2 Framework

3 Sampling from strongly log-concave density

4 Sampling from log-concave density

5 Non-smooth potentials

6 Numerical illustrations

7 Conclusion

Introduction

- Sampling distribution over high-dimensional state-space has recently attracted a lot of research efforts in computational statistics and machine learning community...
- **Applications** (non-exhaustive)
 - 1 Bayesian inference for high-dimensional models and Bayesian non parametrics
 - 2 Bayesian linear inverse problems (typically function space problems converted to high-dimensional problem by Galerkin method)
 - 3 Aggregation of estimators and experts
- Most of the sampling techniques known so far **do not scale** to high-dimension... Challenges are numerous in this area...

Bayesian setting (I)

- In a Bayesian setting, a parameter $\beta \in \mathbb{R}^d$ is embedded with a **prior distribution** ξ and the observations are given by a **probabilistic model**:

$$Y \sim \ell(\cdot|\beta)$$

The inference is then based on the **posterior distribution**:

$$\pi(d\beta|Y) = \frac{\xi(d\beta)\ell(Y|\beta)}{\int \ell(Y|u)\xi(du)}.$$

In most cases the normalizing constant is **not tractable**:

$$\pi(d\beta|Y) \propto \xi(d\beta)\ell(Y|\beta).$$

Bayesian setting (II)

Bayesian decision theory relies on computing expectations:

$$\int_{\mathbb{R}^d} f(\boldsymbol{\beta}) \ell(Y|\boldsymbol{\beta}) \xi(d\boldsymbol{\beta})$$

Generic problem: estimation of an expectation $\mathbb{E}_{\pi}[f]$, where

- π is known up to a multiplicative factor ;
- we do not know how to sample from π (no basic Monte Carlo estimator);

Examples: Logistic and probit regression

- **Likelihood:** Binary regression set-up in which the binary observations (responses) (Y_1, \dots, Y_n) are conditionally independent Bernoulli random variables with success probability $F(\beta^T X_i)$, where
 - 1 X_i is a d dimensional vector of known covariates,
 - 2 β is a d dimensional vector of unknown regression coefficient
 - 3 F is a distribution function.
- Two important special cases:
 - 1 **probit regression:** F is the standard normal distribution function,
 - 2 **logistic regression:** F is the standard logistic distribution function,
 $F(t) = e^t / (1 + e^t)$.

Examples: Logistic and probit regression

- The posterior density distribution of β is given by Bayes' rule, up to a proportionality constant by $\pi(\beta|(Y, X)) \propto \exp(-U(\beta))$, where the potential $U(\beta)$ is given by

$$U(\beta) = - \sum_{i=1}^p \{Y_i \log F(\beta^T X_i) + (1 - Y_i) \log(1 - F(\beta^T X_i))\} + g(\beta),$$

where g is the log-density of the prior distribution.

- Two important cases:

- Gaussian prior: $g(\beta) = -(1/2)\beta^T \Sigma_{\beta} \beta$, ridge regression.
- Laplace prior: $g(\beta) = -\lambda \sum_{k=1}^d |\beta_k|$, lasso regression.

New challenges

Problem the number of predictor variables d is **large** (10^4 and up).

Examples

- text categorization,
- genomics and proteomics (gene expression analysis), ,
- other data mining tasks (recommendations, longitudinal clinical trials, ..).

Data Augmentation

- The most popular algorithms for Bayesian inference in **ridge** binary regression models are based on **data augmentation**:
 - 1 probit link: Albert and Chib (1993).
 - 2 logistic link: Polya-Gamma sampler, Polsson and Scott (2012)... !
- Bayesian lexicon:
 - **Data Augmentation** instead on sampling $\pi(\beta|(Y, X))$ sample $\pi(\beta, W|(Y, X))$ and marginalize W .
 - Typical application of the Gibbs sampler: sample in turn $\pi(\beta|W, Y, X)$ and $\pi(W|\beta, X, Y)$
 - The choice of the DA should make these two steps reasonably easy...

Data Augmentation algorithms

- These two algorithms have been shown to be uniformly geometrically ergodic, **BUT** the constants depends highly on the dimension.
- The algorithms are very demanding in terms of computational ressources...
 - applicable only when is d small 10 to moderate 100 but certainly not when d is large (10^4 or more).
 - convergence time prohibitive as soon as $d \geq 10^2$.

A daunting problem ?

- In the case of the **ridge** regression, the potential $\beta \mapsto U(\beta)$ is **smooth, strongly convex**
- In the case of the **lasso** regression, the potential $\beta \mapsto U(\beta)$ is **non-smooth** but still **convex...**
- A wealth of reasonably fast optimisation algorithms are available to solve this problem in high-dimension...

- 1 Motivation
- 2 Framework**
- 3 Sampling from strongly log-concave density
- 4 Sampling from log-concave density
- 5 Non-smooth potentials
- 6 Numerical illustrations
- 7 Conclusion

Framework

- Denote by π a target density w.r.t. the Lebesgue measure on \mathbb{R}^d , known up to a normalisation factor

$$x \mapsto e^{-U(x)} / \int_{\mathbb{R}^d} e^{-U(y)} dy ,$$

Implicitly, $d \gg 1$.

- Assumption:** U is L -smooth : continuously differentiable and there exists a constant L such that for all $x, y \in \mathbb{R}^d$,

$$\|\nabla U(x) - \nabla U(y)\| \leq L\|x - y\| .$$

Langevin diffusion

■ Langevin SDE:

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

where $(B_t)_{t \geq 0}$ is a d -dimensional Brownian Motion.

- Denote by $(P_t)_{t \geq 0}$ the semigroup of the diffusion,
 $P_t(x, A) = \mathbb{E}_x [Y_t \in A]$.
- $(P_t)_{t \geq 0}$ is
 - aperiodic, strong Feller (all compact sets are small).
 - reversible w.r.t. to π (admits π as its unique invariant distribution).
- $\pi \propto e^{-U}$ is **reversible** \rightsquigarrow the unique **invariant probability** measure.
For all $x \in \mathbb{R}^d$, measurable and bounded functions $f : \mathbb{R}^d \rightarrow \mathbb{R}$

$$\lim_{t \rightarrow +\infty} P_t f(x) = \lim_{t \rightarrow +\infty} \mathbb{E}_x [f(Y_t)] = \int_{\mathbb{R}^d} f(y) d\pi(y) .$$

Discretized Langevin diffusion

- **Idea:** Sample the diffusion paths, using for example the **Euler-Maruyama (EM)** scheme:

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

where

- $(Z_k)_{k \geq 1}$ is i.i.d. $\mathcal{N}(0, I_d)$
 - $(\gamma_k)_{k \geq 1}$ is a sequence of stepsizes, which can either be held constant or be chosen to decrease to 0 at a certain rate.
- Closely related to the **gradient algorithm**.

Discretized Langevin diffusion: constant stepsize

- When $\gamma_k = \gamma$, then $(X_k)_{k \geq 1}$ is an **homogeneous Markov chain** with Markov kernel R_γ
- Under some appropriate conditions, this Markov chain is irreducible, positive recurrent \leadsto unique invariant distribution π_γ .
- **Problem:** $\pi_\gamma \neq \pi$.

- When $(\gamma_k)_{k \geq 1}$ is nonincreasing and non constant $(X_k)_{k \geq 1}$ is an **inhomogeneous Markov chain** associated with the sequence of Markov kernel $(R_{\gamma_k})_{k \geq 1}$
- Denote by $\delta_x Q_{\gamma}^p$ the law of X_p started at x .
- Reminder: the diffusion converges to the target distribution π .
- Question: since the EM discretization approximates the diffusion, can it be used to sample from π :
 - Is $\delta_x Q_{\gamma}^p$ closed to π and in which sense ?
 - Can we have some theoretical guarantees ? In particular what is the dependence on the dimension d ?

Metric on probability spaces

Definition

For μ, ν two probabilities measure on \mathbb{R}^d , define

$$\|\mu - \nu\|_{\text{TV}} = \sup_{|f| \leq 1} |\mathbb{E}_\mu[f] - \mathbb{E}_\nu[f]| .$$

$$W_1(\mu, \nu) = \sup_{\|f\|_{\text{Lip}} \leq 1} |\mathbb{E}_\mu[f] - \mathbb{E}_\nu[f]| .$$

- 1 Motivation
- 2 Framework
- 3 Sampling from strongly log-concave density**
- 4 Sampling from log-concave density
- 5 Non-smooth potentials
- 6 Numerical illustrations
- 7 Conclusion

Wasserstein distance convergence

- We assume in this part that U is strongly convex: there exist and $m > 0$, such that for all $x, y \in \mathbb{R}^d$,

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

- For all $x \in \mathbb{R}^d$,

$$W_1^2(\delta_x P_t, \pi) \leq e^{-mt} \int_{\mathbb{R}^d} \|y - x\|^2 \pi(dy) .$$

- Assume U is L -smooth and strongly convex. Let $(\gamma_k)_{k \geq 1}$ be a nonincreasing sequence with $\gamma_1 \leq 1/(m + L)$.
- Then for all $x \in \mathbb{R}^d$ and $n \geq 1$,

$$W_1^2(\delta_x Q_\gamma^n, \pi) \leq u_n^{(1)}(\gamma) \int_{\mathbb{R}^d} \|y - x\|^2 \pi(dy) + u_n^{(2)}(\gamma),$$

where $(u_n^{(1)}(\gamma), u_n^{(2)}(\gamma))_{n \geq 1}$ are explicit.

- We have if $\lim_{k \rightarrow +\infty} \gamma_k = 0$ and $\lim_{k \rightarrow +\infty} \Gamma_k = +\infty$,

$$\lim_{n \rightarrow +\infty} W_1^2(\delta_x Q_\gamma^n, \pi) = 0,$$

with explicit convergence.

- Order of convergence for decreasing stepsize.

	$\alpha \in (0, 1)$	$\alpha = 1$
Order of convergence	$\mathcal{O}(n^{-\alpha})$	$\mathcal{O}(n^{-1})$

Table : Order of convergence of $W_2(\delta_x Q_\gamma^n, \pi)$ for $\gamma_k = \gamma_1 k^{-\alpha}$

When $(\gamma_k)_{k \geq 1}$ is constant:

- We optimize γ and n to get $W_1(\delta_x Q_\gamma^n, \pi) \leq \epsilon$. In particular, we need

$$n = \mathcal{O}(d\epsilon^{-2}) .$$

- If the number of iterations n is fixed, we can optimize γ and we find a bound in $W_1(\delta_x Q_\gamma^n, \pi) \leq \mathcal{O}(\sqrt{n})$.

- To improve the bound, we make a regularity assumption on U : The potential U is three times continuously differentiable and there exists \tilde{L} such that for all $x, y \in \mathbb{R}^d$:

$$\|\nabla^2 U(x) - \nabla^2 U(y)\| \leq \tilde{L} \|x - y\| .$$

- Assume U is strongly convex L -smooth and satisfies the condition above. Let $(\gamma_k)_{k \geq 1}$ be a nonincreasing sequence with $\gamma_1 \leq 1/(m + L)$.
- Then for all $x \in \mathbb{R}^d$ and $n \geq 1$,

$$W_1^2(\delta_x Q_\gamma^n, \pi) \leq u_n^{(3)}(\gamma) \int_{\mathbb{R}^d} \|y - x\|^2 \pi(dy) + u_n^{(4)}(\gamma) ,$$

where $(u_n^{(3)}(\gamma), u_n^{(4)}(\gamma))_{n \geq 1}$ are explicit.

- Order of convergence for decreasing stepsize.

	$\alpha \in (0, 1)$	$\alpha = 1$
Order of convergence	$\mathcal{O}(n^{-2\alpha})$	$\mathcal{O}(n^{-2})$

Table : Order of convergence of $W_1(\delta_x Q_\gamma^n, \pi)$ for $\gamma_k = \gamma_1 k^{-\alpha}$

When $(\gamma_k)_{k \geq 1}$ is constant:

- We optimize γ and n to get $W_1(\delta_x Q_\gamma^n, \pi) \leq \epsilon$. In particular, we need

$$n = \mathcal{O}(\sqrt{d}\epsilon^{-1}).$$

- If the number of iterations n is fixed, we can optimize γ and we find a bound in $W_1(\delta_x Q_\gamma^n, \pi) \leq \mathcal{O}(n^{-1})$.

- 1 Motivation
- 2 Framework
- 3 Sampling from strongly log-concave density
- 4 Sampling from log-concave density**
- 5 Non-smooth potentials
- 6 Numerical illustrations
- 7 Conclusion

Convergence of the Euler discretization

- If we assume that U is **convex**, L -smooth.
- Explicit bound for $\|\delta_x Q_\gamma^p - \pi\|_{\text{TV}}$.
- If $\lim_{\gamma_k \rightarrow +\infty} \gamma_k = 0$, and $\sum_k \gamma_k = +\infty$ then

$$\lim_{p \rightarrow +\infty} \|\delta_x Q_\gamma^p - \pi\|_{\text{TV}} = 0 .$$

- Computable bounds for the convergence.

Target precision ϵ : the convex case

- For constant stepsizes, We can optimize γ and p to get

$$\|\delta_x Q_\gamma^p - \pi\|_{\text{TV}} \leq \epsilon .$$

- | | d | ϵ | L |
|----------|-----------------------|--|-----------------------|
| γ | $\mathcal{O}(d^{-4})$ | $\mathcal{O}(\epsilon^2 / \log(\epsilon^{-1}))$ | $\mathcal{O}(L^{-2})$ |
| p | $\mathcal{O}(d^7)$ | $\mathcal{O}(\epsilon^{-2} \log^2(\epsilon^{-1}))$ | $\mathcal{O}(L^2)$ |

- We can also at fixed iteration, optimize the stepsize γ .
- Dependence on the dimensions: comes from the fact that the convergence of the diffusion in the convex case also depends on the dimension.

- 1 Motivation
- 2 Framework
- 3 Sampling from strongly log-concave density
- 4 Sampling from log-concave density
- 5 Non-smooth potentials**
- 6 Numerical illustrations
- 7 Conclusion

Non-smooth potentials

The target distribution has a density π with respect to the Lebesgue measure on \mathbb{R}^d of the form $x \mapsto e^{-U(x)} / \int_{\mathbb{R}^d} e^{-U(y)} dy$ where $U = f + g$, with $f : \mathbb{R}^d \rightarrow \mathbb{R}$ and $g : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ are two lower bounded, convex functions satisfying:

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

$$\|\nabla f(x) - \nabla f(y)\| \leq L_f \|x - y\| .$$

- 2 g is lower semi-continuous and $\int_{\mathbb{R}^d} e^{-g(y)} dy \in (0, +\infty)$.

Moreau-Yosida regularization

- Let $h : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ be a l.s.c convex function and $\lambda > 0$. The λ -Moreau-Yosida envelope $h^\lambda : \mathbb{R}^d \rightarrow \mathbb{R}$ and the proximal operator $\text{prox}_h^\lambda : \mathbb{R}^d \rightarrow \mathbb{R}^d$ associated with h are defined for all $x \in \mathbb{R}^d$ by

$$h^\lambda(x) = \inf_{y \in \mathbb{R}^d} \left\{ h(y) + (2\lambda)^{-1} \|x - y\|^2 \right\} \leq h(x) .$$

- For every $x \in \mathbb{R}^d$, the minimum is achieved at a unique point, $\text{prox}_h^\lambda(x)$, which is characterized by the inclusion

$$x - \text{prox}_h^\lambda(x) \in \gamma \partial h(\text{prox}_h^\lambda(x)) .$$

- The **Moreau-Yosida envelope** is a regularized version of g , which approximates g from below.

Properties of proximal operators

- As $\lambda \downarrow 0$, converges h^λ converges pointwise h , i.e. for all $x \in \mathbb{R}^d$,

$$h^\lambda(x) \uparrow h(x), \quad \text{as } \lambda \downarrow 0.$$

- The function h^λ is convex and continuously differentiable

$$\nabla h^\lambda(x) = \lambda^{-1}(x - \text{prox}_h^\lambda(x)).$$

- The proximal operator is a monotone operator, for all $x, y \in \mathbb{R}^d$,

$$\langle \text{prox}_h^\lambda(x) - \text{prox}_h^\lambda(y), x - y \rangle \geq 0,$$

which implies that the Moreau-Yosida envelope is **L -smooth**:

$$\|\nabla h^\lambda(x) - \nabla h^\lambda(y)\| \leq \lambda^{-1} \|x - y\|, \text{ for all } x, y \in \mathbb{R}^d.$$

MY regularized potential

- If g is not differentiable, but the proximal operator associated with g is available, its λ -Moreau Yosida envelope g^λ can be considered.
- This leads to the approximation of the potential $U^\lambda : \mathbb{R}^d \rightarrow \mathbb{R}$ defined for all $x \in \mathbb{R}^d$ by

$$U^\lambda(x) = f(x) + g^\lambda(x) .$$

Theorem

Under (H), for all $\lambda > 0$, $0 < \int_{\mathbb{R}^d} e^{-U^\lambda(y)} dy < +\infty$.

Some approximation results

Theorem

Assume (H).

- 1 Then, $\lim_{\lambda \rightarrow 0} \|\pi^\lambda - \pi\|_{\text{TV}} = 0$.
- 2 Assume in addition that g is Lipschitz. Then for all $\lambda > 0$,

$$\|\pi^\lambda - \pi\|_{\text{TV}} \leq \lambda \|g\|_{\text{Lip}}^2 .$$

- 3 If $g = \iota_{\mathcal{K}}$ where \mathcal{K} is a convex body of \mathbb{R}^d . Then for all $\lambda > 0$ we have

$$\|\pi^\lambda - \pi\|_{\text{TV}} \leq 2 (1 + D(\mathcal{K}, \lambda))^{-1} ,$$

where $D(\mathcal{K}, \lambda)$ is explicit in the proof, and is of order $\mathcal{O}(\lambda^{-1})$ as λ goes to 0.

The MYULA algorithm-I

Given a regularization parameter $\lambda > 0$ and a sequence of stepsizes $\{\gamma_k, k \in \mathbb{N}^*\}$, the algorithm produces the Markov chain $\{X_k^M, k \in \mathbb{N}\}$:
for all $k \geq 0$,

$$X_{k+1}^M = X_k^M - \gamma_{k+1} \{ \nabla f(X_k^M) + \lambda^{-1} (X_k^M - \text{prox}_g^\lambda(X_k^M)) \} + \sqrt{2\gamma_{k+1}} Z_{k+1},$$

where $\{Z_k, k \in \mathbb{N}^*\}$ is a sequence of i.i.d. d -dimensional standard Gaussian random variables.

The MYULA algorithm-II

- The ULA target the smoothed distribution π^λ .
- To compute the expectation of a function $h : \mathbb{R}^d \rightarrow \mathbb{R}$ under π from $\{X_k^M ; 0 \leq k \leq n\}$, an importance sampling step is used to correct the regularization.
- This step amounts to approximate $\int_{\mathbb{R}^d} h(x)\pi(x)dx$ by the weighted sum

$$S_n^h = \sum_{k=0}^n \omega_{k,n} h(X_k) , \text{ with } \omega_{k,n} = \left\{ \sum_{k=0}^n \gamma_k e^{\bar{g}^\lambda(X_k^M)} \right\}^{-1} \gamma_k e^{\bar{g}^\lambda(X_k^M)} ,$$

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

$$\bar{g}^\lambda(x) = g^\lambda(x) - g(x) = g(\text{prox}_g^\lambda(x)) - g(x) + (2\lambda)^{-1} \|x - \text{prox}_g^\lambda(x)\|^2 .$$

- 1 Motivation
- 2 Framework
- 3 Sampling from strongly log-concave density
- 4 Sampling from log-concave density
- 5 Non-smooth potentials
- 6 Numerical illustrations**
- 7 Conclusion

Image deconvolution

- **Objective** recover an original image $\mathbf{x} \in \mathbb{R}^n$ from a blurred and noisy observed image $\mathbf{y} \in \mathbb{R}^n$ related to \mathbf{x} by the linear observation model $\mathbf{y} = H\mathbf{x} + \mathbf{w}$, where H is a linear operator representing the blur point spread function and \mathbf{w} is a Gaussian vector with zero-mean and covariance matrix $\sigma^2 \mathbf{I}_n$.
- This inverse problem is usually ill-posed or ill-conditioned: exploits prior knowledge about \mathbf{x} .
- One of the most widely used image prior for deconvolution problems is the improper total-variation norm prior, $\pi(\mathbf{x}) \propto \exp(-\alpha \|\nabla_d \mathbf{x}\|_1)$, where ∇_d denotes the discrete gradient operator that computes the vertical and horizontal differences between neighbour pixels.

$$\pi(\mathbf{x}|\mathbf{y}) \propto \exp \left[-\|\mathbf{y} - H\mathbf{x}\|^2 / 2\sigma^2 - \alpha \|\nabla_d \mathbf{x}\|_1 \right].$$

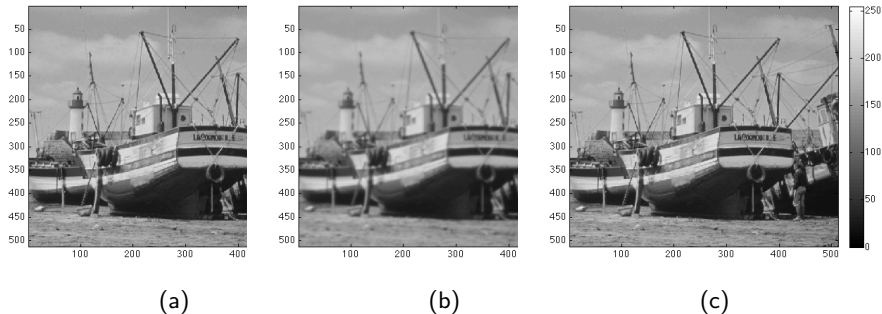


Figure : (a) Original Boat image (256×256 pixels), (b) Blurred image, (c) MAP estimate.

Credibility intervals

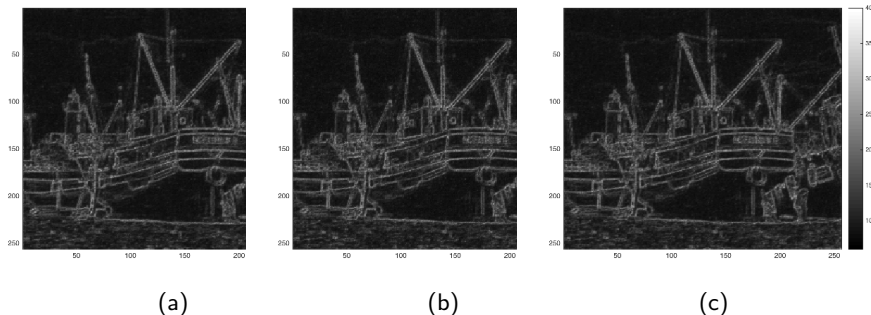


Figure : (a) Pixel-wise 90% credibility intervals computed with proximal MALA (computing time 35 hours), (b) Approximate intervals estimated with MYULA using $\lambda = 0.01$ (computing time 3.5 hours), (c) Approximate intervals estimated with MYULA using $\lambda = 0.1$ (computing time 20 minutes).

- 1 Motivation
- 2 Framework
- 3 Sampling from strongly log-concave density
- 4 Sampling from log-concave density
- 5 Non-smooth potentials
- 6 Numerical illustrations
- 7 Conclusion**

What's next ?

- **Extension of this work**
 - Richardson-Romberg interpolation: debiasing for smooth functionals with non-asymptotic bounds on the MSE.
 - Langevin meets Gibbs: ULA within Gibbs.
 - detailed comparison with MALA
- Thank you for your attention.