

On the convergence of stochastic algorithms

Laurent Miclo
(based on works with
Marc Arnaudon)

Plan

- The problem of finding Fréchet means
- Simulated annealing
- Functional inequalities
- A stochastic algorithm finding Fréchet means
- Discrepancies between probability measures

- Proof of convergence

A priori hope: each of the 1h30 lectures will consist in approximately two of the above consecutive points.

1- Fréchet means

Consider (M, d) a compact metric space, endowed with a probability measure ν .

Fix $p \in [1, +\infty[$ and define the mapping

$$U : M \rightarrow \mathbb{R} \\ x \mapsto \int d^p(x, y) \nu(dy)$$

Denote \mathcal{M}_p the set of global minima of U .

An element x of \mathcal{M}_p is

called a p -mean associated
to ν .

For $p=2$, π is simply
called a mean, while for
 $p=1$, π is a median.

In the above definition the
compactness assumption
ensures that M_p is not
empty, but the notion of
 p -means can be extended
to general metric probability

spaces, under the condition that

\mathcal{D} admits a p -moment:

$\exists x \in M$ (equivalently $\forall x \in M$):

$$\int d^p(x, y) \nu(dy) < +\infty$$

When $M = \mathbb{R}^k$, with $k \in \mathbb{N}$,

endowed with the Euclidean

distance, one recovers the

usual notion of mean, or

median (for $k=1$).

But in general the p -mean

is not unique: consider for

instance the Haar measure ν
on the circle \mathbb{T} endowed
with the Riemannian distance
inherited from its usual embedding
in \mathbb{R}^2 :

$$\forall p \geq 1, \mathcal{M}_p = \mathbb{T}$$

A general problem is to find
 \mathcal{M}_p , from a sequence $(Y_n)_{n \in \mathbb{N}}$
of i.i.d. random variables
distributed according to ν
(in practice, the sample is

rather finite but large...).

In these lectures, we will only consider the Riemannian setting, where numerous applications are already taking place. For instance some infinite dimensional extensions of this framework are important in image processing: the distance between two images is obtained by associating some weights to natural infinitesimal transformations

enabling to pass from one to the other.

Other contexts are interesting too, e.g. weighted graphs: what is the mean size of a network, where the pre-distance between two sites is inversely proportional to the traffic between them?

So to end this introduction, let us recall the necessary notions from Riemannian geometry.

By definition, we are given on each tangent space of the compact manifold M a scalar product $\langle \cdot, \cdot \rangle$. It enables to define the gradient operator from $C^\infty(M) \rightarrow C^\infty(TM)$ by $\forall f \in C^\infty(M), \forall x \in M, \forall v \in T_x M,$

$$df_x(v) = \langle \nabla f, v \rangle_x$$

There is a natural measure μ associated to the Riemannian structure: in local coordinates

$$\text{if } \langle , \rangle_x = \sum_{k,l} g_{k,l}(x) dx^k dx^l,$$

$$\mu(dx) = \sqrt{\det(g(x))} dx^1 \dots dx^N$$

(with $N = \dim(M)$).

In the compact setting μ has finite weight, so we renormalize it to get a probability measure. The divergence operator is defined by duality on vector fields b via

$$\forall f \in C^\infty(M), \int \langle b, \nabla f \rangle d\mu = - \int \operatorname{div} b f d\mu$$

$$\begin{aligned} & \text{(in local coordinates, } \operatorname{div}(b) \\ & = \frac{1}{\sqrt{\det(g)}} \sum_k \partial_k (\sqrt{\det(g)} b^k). \end{aligned}$$

The Laplace-Beltrami operator is then given by

$$\begin{aligned} \Delta \cdot & = \operatorname{div}(\nabla \cdot) \\ & \left(= \frac{1}{\sqrt{\det(g)}} \sum_{k,l} \partial_k g^{k,l} \sqrt{\det(g)} \partial_l \cdot \right), \end{aligned}$$

which leads to the important integration by part formula:

$$\forall f, g \in C^\infty(M),$$

$$\int f \Delta g \, d\mu = \int g \Delta f \, d\mu$$

$$= - \int \langle \nabla f, \nabla g \rangle d\mu$$

The associated distance is given by

$$d(x, y) = \min_{\gamma: x \rightarrow y} \int_0^T \underbrace{\langle \dot{\gamma}_t, \dot{\gamma}_t \rangle}_{=: \|\dot{\gamma}_t\|^2} dt$$

where the minimum is over all smooth paths $\gamma: [0, T] \rightarrow M$ such that $\gamma(0) = x$, $\gamma(T) = y$.

A minimum path γ with $\|\dot{\gamma}_t\| \equiv 1$ is called a geodesic

with unitary speed. For x, y sufficiently close, it is unique and will be denoted $\gamma_{x,y}$.

It satisfies the Euler-Lagrange equation which enables to define for all times $t \in \mathbb{R}$.

2 - Simulated annealing

On a compact Riemannian manifold M , let be given U a smooth function. We would like to find the set \mathcal{M} of global minima of U .

The corresponding simulated annealing algorithm is a time-inhomogeneous Markov process $(X_t)_{t \geq 0}$, described in terms of its generators

$(L_{\beta_t})_{t \geq 0}$:

$\forall \beta \geq 0,$

$$L_{\beta} := \Delta \cdot - \beta \langle \nabla U, \nabla \cdot \rangle$$

and $\mathbb{R}_+ \ni t \mapsto \beta_t \in \mathbb{R}_+$

is a smooth evolution of the inverse of the temperature,

to be determined so that

for any neighborhood \mathcal{N} of \mathcal{M} , we have

$$\lim_{t \rightarrow +\infty} \mathbb{P}[X_t \in \mathcal{N}] = 1 \quad (1)$$

(so that to find with a good chance a point close to \mathcal{M} , it is sufficient to pick X_t for $t \geq 0$ large enough).

Rigorously, $X := (X_t)_{t \geq 0}$ is such that for any function $f \in C^\infty(M)$,

the process M^f defined by

$$\forall t \geq 0, M_t^f = f(X_t) - f(X_0) - \int_0^t L_{p_s} f(X_s) ds$$

is a martingale in the filtration generated by X .

Due to our regularity assumptions, once $\mathcal{L}(X_0)$ (the law of X_0)

is given, the above martingale problem uniquely determines $\mathcal{L}(X)$ on $C(\mathbb{R}_+, M)$.

There are probabilistic constructions of X which give a meaning to the description

$$dX_t = \sqrt{2} \sigma(X_t) dW_t - \beta_t \nabla U(X_t) dt$$

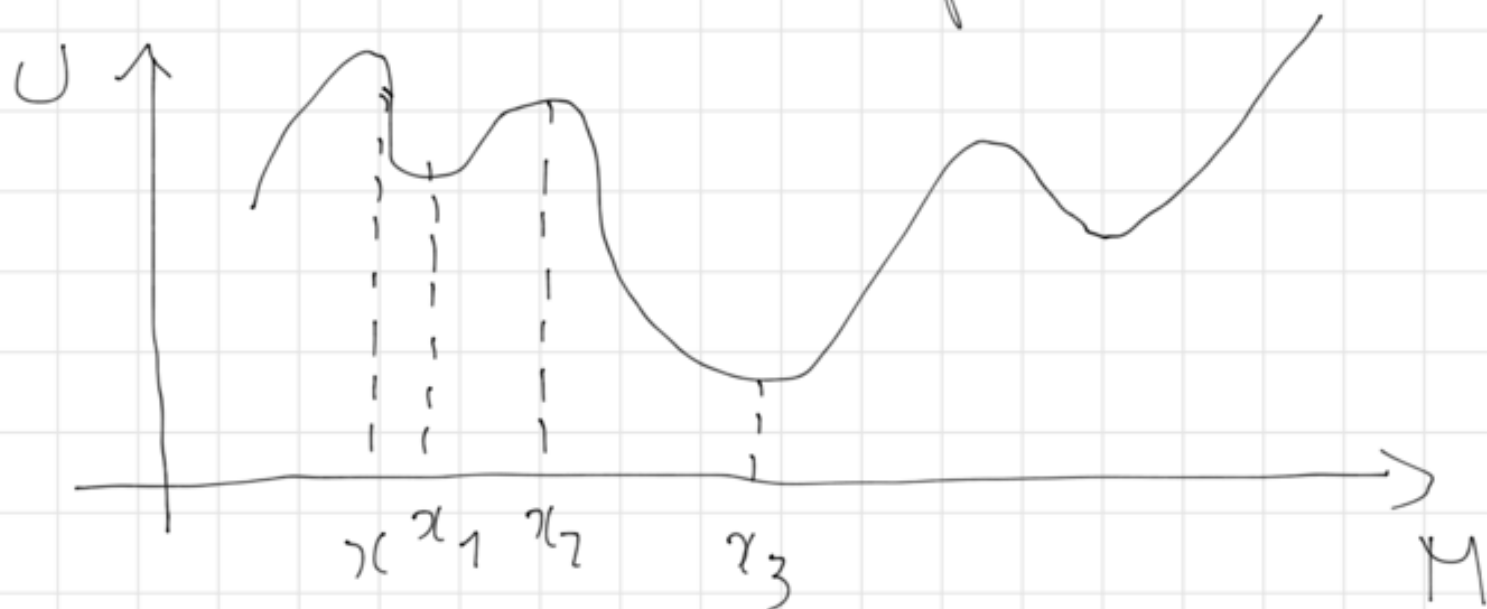
where for any $x \in M$

$$\sigma(x) : \mathbb{R}^N \rightarrow T_x M$$

is such that $\sigma(x) \sigma^*(x) = \text{id}$ (and where the matrix $\sigma(x)$ depends smoothly on x).

In the above formula, $(B_t)_{t \geq 0}$
is a standard Brownian motion
in \mathbb{R}^k .

We thus have the picture



At x , X_t "feels" the influence
of $-\beta_t \nabla U(X_t)$, which pushes it
toward the local minimum x_1 ,
as well as the perturbation

by the Brodian motion, which will help it to escape from the well whose bottom is x_1 , to pass through x_2 in direction of the global minimum x_3 .

To present nice evolutions

$\beta := (\beta_t)_{t \geq 0}$, we need to

introduce a critical constant

$b(U)$.

If $p: [0, 1] \rightarrow M$ is a continuous path, we denote

its elevation by

$$U(p) = \max_{t \in [0, 1]} U(p|t)$$

The minimal elevation between two points $x, y \in M$ is

$$U(x, y) = \min_{p \in \mathcal{P}_{x, y}} U(p)$$

set of continuous paths going from x to y .

Then we consider

$$b(U) := \max_{x, y \in M} U(x, y) - U(x) - U(y) + \min_M U$$

Fix $x_0 \in M$, this constant

is also equal to

$$b(U) = \max_{y \in M} U(\pi_0, y) - U(y)$$

namely, $b(U)$ is the highest

height of a well not containing π_0 :



For $\lambda \geq 1$ and $b > 0$, consider
the evolution of β given by

$$\forall t \geq 0, \quad \beta_t = b^{-1} \ln(1+t)$$

Holley, Kusuoka and Stroock [1989]

have proven the following result:

Theorem

If $b > b(\nu)$, then $\forall \mathcal{L}(X_0)$,

$$\lim_{t \rightarrow +\infty} \|\mathcal{L}(X_t) - \mu_{\beta_t}\|_{TV} = 0$$

If $b < b(\nu)$, then $\exists \mathcal{L}(X_0)$:

$$\lim_{t \rightarrow +\infty} \|\mathcal{L}(X_t) - \mu_{\beta_t}\|_{TV} \neq 0$$

where for any $\beta \geq 0$, μ_β is the Gibbs measure associated to the potential U and to the temperature β :

$$\mu_\beta(dx) = \frac{e^{-\beta U(x)} M(dx)}{\sum_{\beta} 1}$$

normalizing constant

Note that for $\beta \geq 0$ large, μ_β concentrates around \mathcal{M} :

For any neighborhood \mathcal{N} of \mathcal{M} ,

$$\lim_{\beta \rightarrow +\infty} \mu_\beta[\mathcal{N}] = 1$$

It follows that for the evolution of β corresponding to $\Delta \geq 1$ and $b > b(U)$, our goal (1) is fulfilled (recall the definition of the total variation norm between two probability measures μ and ν):

$$\|\mu - \nu\|_{TV} = 2 \sup_{\substack{A \subset M \\ \text{Borelian}}} \mu(A) - \nu(A)$$

$$= \sup_{\substack{f \in C^0(M) \\ \|f\|_\infty \leq 1}} \mu[f] - \nu[f]$$

Note that μ_β is reversible for

$$L_\beta : \forall f, g \in C^\infty(M),$$

$$\mu_\beta [f L_\beta g] = \mu_\beta [g L_\beta f]$$

Indeed, we have

$$L_\beta = e^{\beta U} \operatorname{div} (e^{-\beta U} \nabla \cdot)$$

so that

$$\mu_\beta [f L_\beta g] = \frac{1}{Z_\beta} \mu [f \operatorname{div} (e^{-\beta U} \nabla g)]$$

$$= -\frac{1}{Z_\beta} \mu [\langle \nabla f, e^{-\beta U} \nabla g \rangle]$$

$$= -\mu_\beta [\langle \nabla f, \nabla g \rangle]$$

In particular x_p is invariant for L_p , thus at any time $t \geq 0$, $\mathbb{E}[X_t]$ tends instantaneously to get closer to x_p . The above theorem states that if $(\beta_t)_{t \geq 0}$ is changing sufficiently slowly ($b > b(U)$), then $\mathbb{E}[X_t]$ succeeds to follow its target x_p and to get closer and closer to it in large times.

Remarks:

a) For $b > b(U)$, a.s. convergence

doesn't hold for X : it can be shown that the recurrence set of X is a.s. the connected component of $\{x : \cup_{\pi} |U| \leq \omega \text{ in } U + b\}$ containing \mathcal{M} , even if most of the time X_t is close to \mathcal{M} for large $t \geq 0$.

b) The constant $b(U)$ is not critical for (1), it is the constant $\tilde{b}(U)$ defined by

$$\tilde{b}(U) := \max_{y \in \pi} \min_{x_0 \in \mathcal{M}} U(x_0, y) - U(y)$$

A simple proof of ^{the first part of} the theorem consists in studying the evolution of a discrepancy between $m_f := \mathbb{I}(X_t)$ and μ_{p_t} , which is easier to differentiate than the total variation (absolute value is rather troublesome to manipulate).

In the above situation, the most appropriate choice is the relative entropy $E_{\text{rel}}(\mu_t, \mu_{p_t})$, but to avoid the investigation

of logarithmic Sobolev inequalities,
let us rather consider the functional

$$\forall t > 0,$$

$$I_t := \int \left(\frac{d\mu_t}{d\mu_{\beta_t}} - 1 \right)^2 d\mu_{\beta_t}$$

Radon-Nikodym density,

well-defined for $t > 0$ by

ellipticity of the generators

By Cauchy-Schwarz inequality,

$$\begin{aligned} \|\mu_t - \mu_{\beta_t}\|_{TV} &= \int \left| \frac{d\mu_t}{d\mu_{\beta_t}} - 1 \right| d\mu_{\beta_t} \\ &\leq \sqrt{I_t} \end{aligned}$$

it is thus sufficient to show that

$$\lim_{t \rightarrow +\infty} I_t = 0$$

To study the evolution of $t \mapsto I_t$, we must know how to differentiate m_t with respect to time. By the martingale problem, we have for any $f \in C(\mathbb{M})$,

$$\partial_t \mathbb{E}[f(X_t)] = \mathbb{E}[L_{P_t}[f](X_t)]$$

namely, in the weak sense

$$\partial [f \partial_t m_t] = \mu [L_{P_t}[f] m_t]$$

where the probability measure

w_t and its density $\frac{dw_t}{d\mu}$ are written in the same manner w_t (for $t > 0$, $w_t \in C^\infty(M)$).

It will be convenient to denote

$$f_t := \frac{dw_t}{d\mu_{\beta_t}}, \text{ for } t > 0.$$

We can now differentiate with respect to β , hence: for $t > 0$,

$$\dot{I}_t =$$

$$2 \int (\beta_t - 1) \left(\frac{\partial_t w_t}{\mu_{\beta_t}} - f_t \partial_t \ln \mu_{\beta_t} \right) d\mu_{\beta_t}$$

$$+ \int (\beta_t - 1)^2 \partial_t \ln \mu_{\beta_t} d\mu_{\beta_t}$$

$$= 2 \int (f_t - 1) \partial_t \mu_t d\mu$$

$$- \int \left((f_t - 1)^2 + 2(f_t - 1) \right) \partial_t \ln \mu_{\beta_t} d\mu_{\beta_t}$$

$$\leq 2 \int L_{\beta_t} (f_t - 1) f_t d\mu_{\beta_t}$$

$$+ \|\partial_t \ln \mu_{\beta_t}\|_{\infty} \left(\mathbb{I}_t + 2\sqrt{\mathbb{I}_t} \right)$$

$$= -2 \int \|\nabla f_t\|^2 d\mu_{\beta_t}$$

$$+ \|\partial_t \ln \mu_{\beta_t}\|_{\infty} \left(\mathbb{I}_t + 2\sqrt{\mathbb{I}_t} \right)$$

The last term is easy to evaluate:

$$\|\partial_t \ln \mu_{\beta_t}\|_{\infty} = \|\partial_{\beta} \ln \mu_{\beta}\|_{\infty} |\beta'_t|$$

and since for $x \in M$

$$\ln \mu_\beta(x) = -\beta U(x) + \ln \left(\int e^{-\beta U} d\mu \right)$$

$$\| \partial_\beta \ln \mu_\beta \|_\infty \leq \text{osc}(U)$$

$$:= \max_M U - \min_M U$$

Concerning the first term,
we need a spectral gap
estimate due to Holley, Kusuoka
and Stroock [1989]:

Proposition

| There exists a constant

C_M depending only on M , such that $\forall \beta \geq 0, \forall f \in C^\infty(M)$,

$$\mu_\beta[(\beta - 1)^2] \leq c(\beta) \mu_\beta[\|\nabla f\|^2]$$

with

$$c(\beta) = C_M (1 \vee (\beta \|\nabla U\|_\infty))^{5N-2} \exp(-b(U)\beta)$$

We will recall the proof of this Poincaré inequality in next section. For the time being let $J_t = \sqrt{I_t}$, we get for $t > 0$,

$$J_t' \leq -c_t J_t + K_t$$

with

$$c_t := c(\beta_t) - K_t/2$$

$$K_t := 2 \|U\|_\infty |\beta_t'|$$

By Gronwall's lemma, we conclude that

$$\lim_{t \rightarrow +\infty} J_t = 0$$

if $\int_0^{+\infty} c_t dt = +\infty$

$$\lim_{t \rightarrow +\infty} \frac{K_t}{c_t} = 0$$

This is satisfied if $b > b(U)$,
because then K_t is of order

$\frac{1}{T}$ and c_T of order $\frac{1}{T^{6/10}}$

(up to logarithmic correction).

3- functional inequalities

To avoid technicalities and to concentrate on ideas, we will deduce the wanted Poincaré inequality in the finite state space setting.

The underlying path approach is also due to Holley and Stroock [1988].

M is now a finite set endowed with a Markovian generator matrix $L := (L_{\alpha, \beta})_{\alpha, \beta \in M}$

Satisfying

$$\forall x \neq y, \quad |L(x, y)| \geq 0$$

$$\forall x, \quad \sum_y |L(x, y)| = 0$$

We assume irreducibility:

$$\forall x, y: \exists x = x_0, x_1, \dots, x_n = y$$

$$\text{with } |L(x_i, x_{i+1})| > 0 \quad \forall 0 \leq i < n \quad (2)$$

and reversibility: there exists μ probability measure such that

$$\forall x, y, \quad \mu(x) |L(x, y)| = \mu(y) |L(y, x)|$$

The operator L plays the role of the Laplacian.

Let U be a function on M and $\beta \geq 0$. The traditional Metropolis perturbation of L admitting the corresponding Gibbs measure $(\mu_\beta(x) \propto e^{-\beta U(x)} \mu(x))$ for reversible probability is defined via

$$\forall x, y, \quad L_\beta(x, y) = e^{-\beta(U(y) - U(x))_+} L(x, y)$$

The associated Dirichlet form is given by

$$\forall f \in \mathcal{F}(M) \quad (= \text{set of functions on } M),$$

$$\xi_{\beta}(f) = -\mu_{\beta}(f L_{\beta} f)$$

$$= \frac{1}{2} \sum_{x,y} \underbrace{\mu_{\beta}(x) L_{\beta}(x,y)}_{e^{-\beta(U(x) - U(y))} \mu_{\beta}(x) L(x,y)} (f(y) - f(x))^2$$

This replaces $\mu_{\beta}[\|\nabla f\|^2]$
in the continuous setting.

We would like to estimate

the best constant $c(\beta)$ such

that $\forall f \in \overline{\mathcal{F}}(M)$

$$\underbrace{\text{Var}(f, \mu_{\beta})}_{\mu_{\beta}[(f - \mu_{\beta}(f))^2]} \leq c(\beta) \xi_{\beta}(f)$$

Proposition

There exists a constant $C_L > 0$
such that

$$\forall \beta \geq 0, \quad |c(\beta)| \leq C_L e^{-\beta |U|}$$

where $b(U)$ is defined as
in the continuous setting, but
with the continuous paths
replaced by discrete paths
satisfying (2)

There is a matching lower
bound, but we will not use it.

For the proof, start with the alternative expression for the

variance: $\forall f \in \overline{\mathcal{F}}$,

$$\text{Var}(f, \mu_P) = \frac{1}{Z} \sum_{x,y} (f(y) - f(x))^2 \mu_P(x) \mu_P(y)$$

For given $x \neq y$, let

$\Gamma_{x,y} = (x = x_0, x_1, \dots, x_n = y)$ be

an admissible path with

$$U(\Gamma_{x,y}) = U(x,y)$$

Then we have

$$\begin{aligned} (f(y) - f(x))^2 &= (f(x_n) - f(x_{n-1}) \\ &+ f(x_{n-1}) - f(x_{n-2}) + \dots - f(x_0))^2 \end{aligned}$$

$$\leq n \sum_{e \in P_{x,y}} (df(e))^2$$

where $P_{x,y}$ is seen as the set of edges $\{(x_0, x_1), (x_1, x_2), \dots, (x_{n-1}, x_n)\}$ and where

$$df(e) = f(x'') - f(x') \text{ if } e = (x', x'').$$

Denote similarly

$$\mu_{\beta}(e) = \mu_{\beta}(x', L_{\beta}(x', x''))$$

Let N be the cardinality of \mathcal{H} . We have

$$\text{Var}(f, \mu_{\beta}) \leq \frac{N}{2} \sum_e (df(e))^2 \sum_{\substack{x,y \\ e \in P_{x,y}}} \mu_{\beta}(x) \mu_{\beta}(y)$$

$$= \frac{N}{Z} \sum_e (df(e))^2 M_\beta L_\beta(e) A_\beta(e)$$

where

$$A_\beta(e) := \frac{1}{M_\beta L_\beta(e)} \sum_{\substack{x, y \\ e \in P_{x, y}}} M_\beta(x) M_\beta(y)$$

$$= \frac{e^{-\beta(U(x') \vee U(x''))}}{\sum_{\substack{x', y \\ e \in P_{x', y}}} e^{-\beta(U(x') \vee U(y))}} \sum_{\substack{x', y \\ e \in P_{x', y}}} e^{-\beta(U(x') \vee U(y))}$$

with $e = (x', x'')$.

Using that

$$\sum_{\mathcal{B}} \underset{\uparrow}{\approx} e^{-\beta \min_M U}$$

up to upper and lower bounds independent of β

we get that

$$\max_{e \in P_{\eta, \gamma}} A_{\beta}(e) \approx e^{\beta(U(\eta, \gamma) - U(\mu) - U(\gamma) + \mu \eta)}$$

so

$$\max_e A_{\beta}(e) \approx e^{\beta b(U)}$$

which is what we wanted.

Note: it is sometimes more convenient to work with the relative entropy

$$Ent(\mu, \nu) = \begin{cases} \int \ln \frac{d\mu}{d\nu} d\nu, & \text{if } \mu \ll \nu \\ +\infty & , \text{ otherwise} \end{cases}$$

than with the quantity considered in the previous section.

In particular, we have Pinsker's inequality:

$$\| \mu - \nu \|_{TV} \leq \sqrt{2} \sqrt{\text{Ent}(\mu, \nu)}$$

For the previous approach to work, the Poincaré inequality

must be replaced by a logarithmic

Sobolev inequality: $\forall f \in \mathcal{F}, f \geq 0, f \neq 0$

$$\text{Ent}\left(\frac{f}{\int f d\mu}, \mu\right) \leq \tilde{C}(\beta) \Sigma_{\beta}(f)$$

It can be shown that the best constant $\tilde{C}(\beta)$ in this bound has at the exponential level a behavior similar to $r(\beta)$:

$$\exists \tilde{C}_L > 0: \forall \beta \geq 0$$

$$\tilde{C}(\beta) \leq \tilde{C}_L (1 + \beta) e^{6|\mu|\beta}$$

This is interesting on very large spaces, where the use of specific quantities (renormalization by the number of sites, e.g.) is required.

4- Another stochastic algorithm

We come back to the problem of finding the Fréchet mean of a probability measure ν on a compact Riemannian manifold M . We have at our disposal a sequence $(Y_n)_{n \in \mathbb{N}}$ of r.v. distributed according to ν .

We want to use them to construct an algorithm $(X_n)_{n \geq 0}$ finding $\underset{\mu}{\mu}$: for all neighborhood \mathcal{M}_2

\mathcal{N} of ω ,

$$\lim_{h \rightarrow +\infty} \mathbb{P}[X_h \in \mathcal{N}] = 1 \quad (3)$$

There is a proposition: in addition to the inverse temperature schedule $(\beta_r)_{r \geq 0}$, let be given $(\alpha_r)_{r \geq 0}$ an homogenization schedule, tuning the speed of use of the data.

More precisely, let $(N_r)_{r \geq 0}$ be a usual Poisson process (interjumps: independent

exponential variables of parameter 1).

Consider the time changed

process $N^{(\lambda)}$:

$$\forall t \geq 0, N_t^{(\lambda)} = N_{\int_0^t \lambda_s ds}$$

We assume that $\lambda_t \xrightarrow[t \rightarrow \infty]{} 0$

so the jumps of $N^{(\lambda)}$ occur more and more often. Call

them $0 < T_1 < T_2 < \dots < T_n < \dots$,

we have that $T_n \xrightarrow[n \rightarrow \infty]{} +\infty$

Between these times, X

evolves as a Brownian motion,
namely its generator is $\frac{\Delta}{2}$.

At time T_n , X jumps
from $X_{T_n^-}$ to $\gamma_{X_{T_n^-}, Y_n}$

with

$$\Delta_n = \alpha_{T_n} \beta_{T_n} d(X_{T_n^-}, Y_n)$$

Note that by construction

$\mathcal{L}(X_{T_n^-}) \ll \mu$, so the
event that $X_{T_n^-}$ is in the
sub-locus of Y_n is negligible.

Furthermore, α and β will

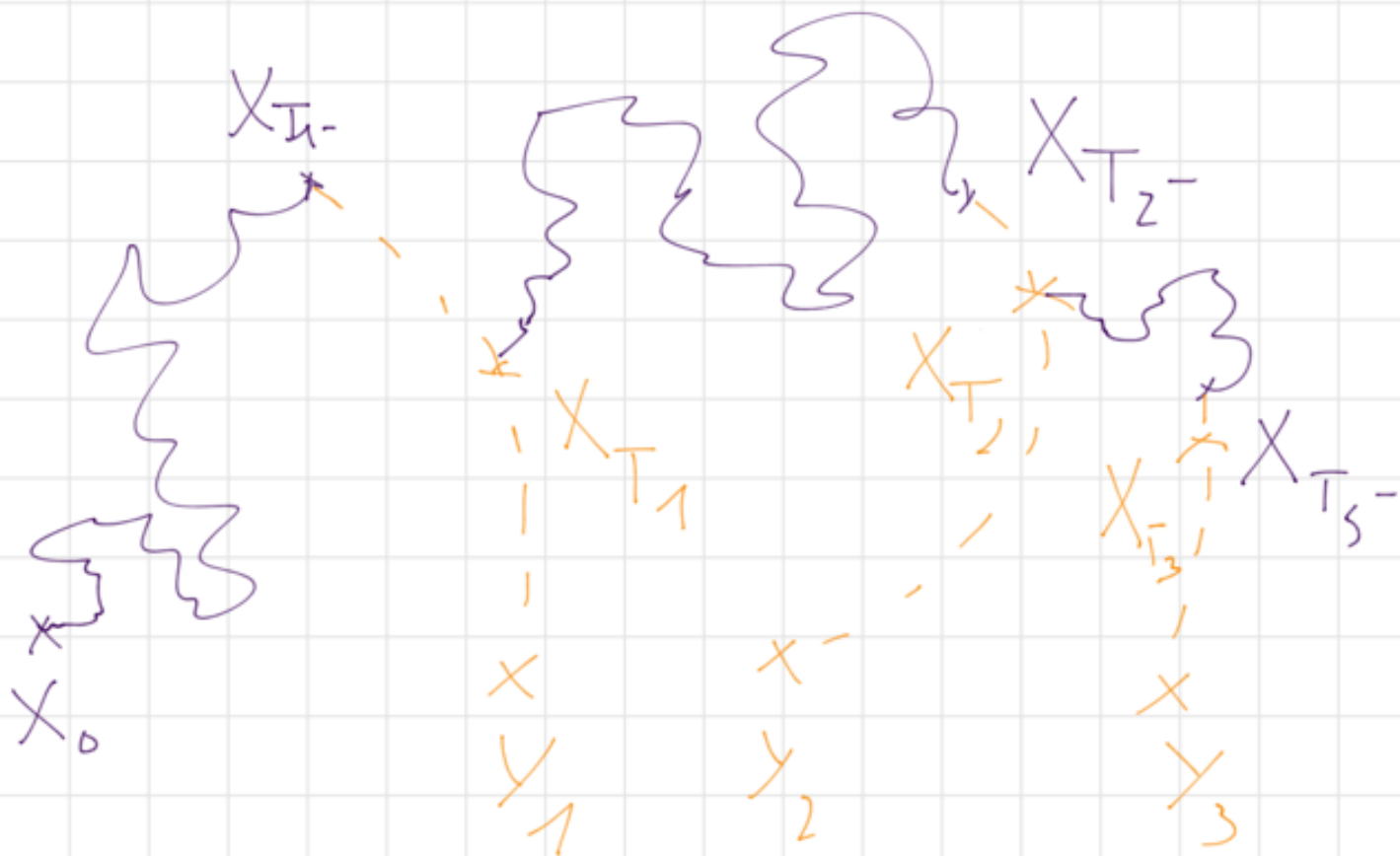
substly

$$\lim_{h \rightarrow 0} \alpha_t \beta_t = 0$$

thus X_{T_n} will end up being quite close

to $X_{T_n^-}$: as time goes on, there

will be a lot of small jumps:



Conjecture

Assume that ν admits a density (wrt μ) Hölder continuous of exponent $\alpha \in (0, 1]$. Then there exists $\bar{a} > 0$ depending on α and M

such that (3) is fulfilled if

$$\forall t \geq 0, \begin{cases} \alpha_t = (1+t)^{-\frac{1}{\alpha}} \\ \beta_t = b^{-1} \ln(1+t) \end{cases}$$

with $b > b(U)$

This result is true for the circle $\mathbb{T} = \mathbb{R} / 2\pi\mathbb{Z}$:

Theorem

The above conjecture is true
for $M = \mathbb{R}$ if $\tilde{a} = a$.

More generally, it is true
for p -means, if we take

$$\Delta_n = \frac{p}{2} \alpha_{T_n} \beta_{T_n} d^{p-1}(X_{T_n-1}, X_n)$$

and

$$\tilde{a} = \begin{cases} a & , \text{if } p=1 \text{ or } p \geq 2 \\ \min(a, p-1) & , \text{if } p \in (1, 2) \end{cases}$$

There is also a version
valid for general probability

measure ν : one needs to add some
time-inhomogeneous
noise to the sequence $(X_n)_{n \in \mathbb{N}}$.

More precisely, for $x \in M$ and
 $\kappa > 0$, consider on $T_x M$ the
probability measure $\bar{K}_{x, \kappa}$ whose
density with respect to the Lebesgue
measure is proportional to $(1 - \kappa \|v\|)_+$.

Denote by $K_{x, \kappa}$ the image of
 $\bar{K}_{x, \kappa}$ by the exponential mapping.

Construct the process $(Z_t)_{t \geq 0}$
as X but by replacing X_n

by a point sampled according to

$K_{\frac{1}{n}, K_Tn}(\cdot)$, where $(K_r)_{r \geq 0}$

is another schedule to tune.

If $M = \mathbb{T}$ and $p = 2$, we have:

Theorem

$$\lim_{t \rightarrow +\infty} P[Z_t \in \mathcal{N}] = 1$$

For any neighborhood \mathcal{N} of \mathcal{M} if

$$\forall t \geq 0 \quad \left\{ \begin{array}{l} \alpha_r = (1+t)^{-(2k+1)} \\ \beta_r = b^{-1} \ln(1+t) \\ \kappa_r = (1+t)^b \end{array} \right.$$

with $b > b/(\nu)$, $b > 0$.

How to choose k ?

The larger k is, the closer \checkmark will be to its approximation by its transport through the kernel $K_{\cdot, x_k}(\cdot)$ and the above convergence will be more precise. But the faster the algorithm uses the data $(Y_n)_{n \in \mathbb{N}}$. So in practice a trade-off has to be made between the number of r.v. Y_n at our disposal and the quality of the approximation of \mathcal{M} .

Heuristic of the convergence in the first theorem.

For $\alpha > 0$, $\beta \geq 0$ consider the generator defined by

$$\forall f \in C^\infty(M), \forall x \in M,$$

$$L_{\alpha, \beta}[f](x) = \frac{\Delta f(x)}{2} + \frac{1}{2} \int f(\gamma_{x,y}(s)) - f(x) \nu(dy)$$

with $\nu = \beta \alpha d(x, y)$.

At any time $t \geq 0$, $L_{\alpha, \beta}$ is the instantaneous generator of X .

We have, for fixed x, y, β ,

$$\lim_{\alpha \rightarrow 0_+} f(\gamma_{\alpha, y}(1)) - f(m)$$

$$= \beta \, d m, y | \langle \nabla f(m), \dot{\gamma}_{\alpha, y}(0) \rangle$$

for any $f \in C^\infty(M)$ and where

$$\dot{\gamma}_{\alpha, y}(0) = \left. \frac{d}{dt} \gamma_{\alpha, y}(t) \right|_{t=0} \in T_x M$$

Thus

$$\lim_{\alpha \rightarrow 0_+} L_{\alpha, \beta}[\gamma](m) = \frac{\Delta f(m)}{2} + \frac{\beta}{2} \langle F, \nabla f \rangle_x$$

where

$$F(m) = 2 \int d m, y \, \dot{\gamma}_{\alpha, y}(0) \nabla(dy)$$

$$= - \int \nabla_x d^2 m, y \, \nabla(dy)$$

$$= -\nabla U(x)$$

(skill with $p=2$).

We recover the generator of the simulated annealing algorithm associated to U , explaining why we expect that X will find x_0 in large time.

But more rigorous computations are needed to justify the choice of $(\alpha_k)_{k \geq 0}$. The same heuristic enables to understand

The convergence of $(Z_t)_{t \geq 0}$.

One drawback of these algorithms

is the computations of the

geodesic paths $\gamma_{x,y}$. In general

it is not easy, but this task

becomes simpler for symmetrical

spaces. For instance if M

is a sphere (of any dimension ≥ 1)

$$\gamma_{x,y}(\Delta) = \cos(\Delta)x + \sin(\Delta) \frac{y - \langle y, x \rangle x}{\sqrt{1 - \langle y, x \rangle^2}}$$

(For $|\langle y, x \rangle| \neq 1$).

Furthermore, still on the sphere, one does not need to entirely simulate trajectories of the Brownian motion: one can use the 1-dimensional Jacobi diffusion obtained by projection on the radius as well as a simple uniform sampling on subspheres. Similar remarks are valid for hori. Let us illustrate that with Π : to simulate $X_{T_{1-}}$, one sample

G a standard Gaussian distribution and reduce modulo 2π $X_0 + \sqrt{t} G$.

The above algorithm can be extended to treat more general cases, but at the expense of difficulties of implementation.

More precisely, consider $\rho: M \times M \rightarrow \mathbb{R}$ a continuous function and

$$U: x \mapsto \int \rho(x, y) \nu(dy)$$

In particular, ρ can be d or d^{\uparrow} , with $\rho > 0$, or the

length associated to a Finsler metric on M (still a compact Riemannian manifold).

Denote $(p(\delta, x, y))_{\delta > 0, x, y \in M}$ the heat kernel on M . We will use it to regularize the left argument of e .

So for $\epsilon > 0$, consider

$$p_{\delta}^{\epsilon}(x, y) = \int_M p(\delta, x, z) p(z, y) \mu(dz)$$

and

$$U_{\delta} : \mathcal{X} \rightarrow \int_M p_{\delta}(x, y) \nu(dy)$$

In addition to β and α , let
be given an approximation schedule $(\delta_k)_{k \geq 0}$
with values in $(0, +\infty)$.

A process X is constructed
as before, except that at
time T_n , the process jumps
from $X_{T_n^-}$ to

$$X_{T_n} := \phi_{\delta_{T_n}}(\alpha_{T_n} \beta_{T_n}, X_{T_n^-}, Y_n)$$

where for any $x, y \in M$,
 $(\phi_s(s, x, y))_{s \geq 0}$ is the flow
starting from x and associated

to the vector field

$$M \ni z \mapsto -\frac{1}{z} \nabla_z \rho_\delta(\cdot, y)$$

(hoping that we are able to construct this flow in practice).

Theorem

(3) is fulfilled if we take

$$\forall t \geq 0 \quad \left\{ \begin{array}{l} \alpha_t := 1/(1+t) \\ \beta_t := b \ln(1+t) \\ \delta_t := 1/\ln(2+t) \end{array} \right.$$

with $b > b(U)$

In the following considerations,
rather consider the Gibbs
measure

$$\mu_{\beta, \varepsilon}(dx) = \frac{e^{-\beta U_{\varepsilon}(x)}}{\int e^{-\beta U_{\varepsilon}(x)} dx} \mu(dx)$$

normalizing constant
and traditional estimates on
the heat kernel.

5 - Discrepancies between measures

At first thought, we would like to study the evolution of

$$t \mapsto D(\mathcal{L}(X_t), \Pi_t)$$

where Π_t is the instantaneous invariant measure at time t and where for any probability measures μ and ν

$$D(\mu, \nu) := \begin{cases} \int \left(\frac{d\mu}{d\nu} - 1 \right)^2 d\nu & \text{if } \mu \ll \nu \\ +\infty & \text{otherwise} \end{cases}$$

Unfortunately, the invariant measure $\bar{\mu}_{\alpha, \beta}$ associated to $L_{\alpha, \beta}$ is difficult to manage because of the combination of diffusion and jump parts. So we will rather investigate the evolution of

$$I: t \mapsto D(\mathcal{L}(X_t), \mu_t)$$

where μ_t is the Gibbs measure associated to the potential U and to the temperature β^{-1} .

At least, we have already

seen how to estimate its Poincaré constant.

What is missing is an evaluation of the discrepancy between μ_β and $\pi_{\alpha,\beta}$, for given $\alpha > 0$, $\beta > 0$.

Let $L_{\alpha,\beta}^*$ be the dual operator of $L_{\alpha,\beta}$ in $\mathcal{D}^2(\mu_\beta)$.

A crucial object is $L_{\alpha,\beta}^*[\mathbb{1}]$, the smaller it is, the closer μ_β is to $\pi_{\alpha,\beta}$. We have e.g.

$$L_{\alpha,\beta}^*[\mathbb{1}] = 0 \iff \mu_\beta = \pi_{\alpha,\beta}$$

Indeed, just use that
 $\forall f \in C^\infty(M)$,

$$\mu_\beta [L_{\alpha, \beta} [f]] = \mu_\beta [L_{\alpha, \beta}^* [1] f]$$

A priori $L_{\alpha, \beta}^* [1]$ is not out of reach, since it is constructed from $L_{\alpha, \beta}$ and μ_β which are given.

Thus $\forall f \in C^\infty(M)$, we compute via a straightforward integration by parts, that

$$L_{\alpha, \beta}^* [f] = \frac{1}{2} e^{\beta U} \Delta (e^{-\beta U} f) +$$

$$\frac{e^{\beta U}}{\alpha(1-\alpha\beta)} \int_{\text{ball}} \mathbb{1}_{\mathcal{O}(y, (1-\alpha\beta)\Pi)} T_{y, \frac{-\alpha\beta}{1-\alpha\beta}} [e^{-\beta U} f] v(dy) - \frac{f}{\alpha}$$

where $T_{y,s}$ is the operator defined by
 $\forall \varphi \in \Pi, T_{y,s} \varphi(x) = \int \varphi(x+y) s(dy)$

One has to be a little careful about the domain of $L_{\alpha,\beta}^*$, but as soon as v admits a continuous density, it is the space of functions whose second derivative in the distributional sense belongs to $\mathcal{D}^2(\mu)$.

We need next some computations,
based on not very exciting expansions,

to get:

Proposition:

Assume that the density of \sqrt{d} satisfies

$$\forall x, y \in \mathbb{T}, |D(x) - D(y)| \leq A d^a(x, y)$$

with $a \in (0, 1]$ and $A > 0$.

Then for $\beta \geq 1$ and $\alpha \in (0, 1/(2\beta^2))$,

$$\|L_{\alpha, \beta}^*\|_{\infty} \leq C(A) \max(\alpha \beta^4, \alpha^a \beta^{4a})$$

↑
constant depending only
on A

Similar estimates are also valid for $p \neq 2$ (curiously they are more involved for $1 < p < 2$), or in the case of a continuous p :

$$\forall \beta \geq 1, \forall \delta \in (0, 1], \forall \alpha \in (0, \delta^2 / (2\beta^2)),$$

$$\|L_{\alpha, \beta, \delta}^{\beta}[\mathbb{1}]\|_{\infty} \leq C \alpha \beta^4 \delta^{-4}$$

for an appropriate constant $C > 0$.

6 - End of the proof

Recall that for $t > 0$,

$$I_t := D / \mathcal{Z}(X_t), \mu_{\beta_t}$$

Its differentiation starts as in

the simulated annealing case:

$$I_t' \leq 2 \int L_{\alpha_r, \beta_r} [f_r^{-1}] f_r d\mu_{\beta_r}$$

$$+ 2 \|U\|_\infty |\beta_r'| (I_t + 2\sqrt{I_t'})$$

$$\leq (\text{diam}(M))^2$$

where $f_r = \frac{d\mathcal{Z}(X_t)}{d\mu_{\beta_t}}$

To go further, we write

$$\int L_{\alpha, \beta_t} [f_t - 1] f_t d\mu_{\beta_t} =$$

$$\int L_{\alpha, \beta_t} [f_t - 1] (f_t - 1) d\mu_{\beta_t}$$

$$+ \underbrace{\int L_{\alpha, \beta_t} [f_t - 1] d\mu_{\beta_t}}$$

$$= \int (f_t - 1) L_{\alpha, \beta_t} [1] d\mu_{\beta_t}$$

$$\leq \sqrt{\mathbb{I}_t} \|L_{\alpha, \beta_t}^* [1]\|_\infty$$

Let us decompose

$$L_{\alpha, \beta} = \underbrace{L_{\beta}} + R_{\alpha, \beta}$$

$$\frac{\Delta}{Z} - \beta \langle \nabla U, \nabla \cdot \rangle$$

We have already seen how to deal with the energy

$$\int (\beta_T - 1) L_{\beta_T} (\beta_T - 1) d\mu_{\beta_T}$$

via Poincaré inequalities at small temperature.

To treat the term

$$\int (\beta_T - 1) R_{\alpha_T, \beta_T} (\beta_T - 1) d\mu_{\beta_T}$$

more expansions are needed ...

in the end, we can find a constant $C'(A) > 0$ such that it is bounded above by

$$C'(A) \alpha_T \beta_T^3 \left(\int (\Delta f_T)^2 d\mu_{\beta_T} + \int (f_T^{-1})^2 d\mu_{\beta_T} \right) - \int (f_T^{-1}) L_{\beta_T} [f_T^{-1}] d\mu_{\beta_T}$$

Putting everything together, it appears that $\mathbb{I}_T = \sqrt{\mathbb{I}_T}$ satisfies

$$\mathbb{I}_T' \leq -\eta_T \mathbb{I}_T + \varepsilon_T$$

with

$$\eta_r := c_1(A) (\beta_r^{-3} \exp(-b(u)\beta_r) - \alpha_r \beta_r^3 - |\beta_r'|)$$

$$\xi_r := c_2(A) (\alpha_r \beta_r^4 + |\beta_r'|)$$

for two constants $c_1(A)$, $c_2(A)$ depending on A .

It follows that if

$$\lim_{r \rightarrow +\infty} \beta_r = +\infty$$

$$\int^{+\infty} \beta_r^{-3} e^{-b(u)\beta_r} dr = +\infty$$

$$\max(\alpha_r \beta_r^4, \alpha_r \beta_r^3, |\beta_r'|) \ll e^{-b(u)\beta_r} \beta_r$$

$r \rightarrow +\infty$

then $\lim_{r \rightarrow +\infty} J_r = 0$

This is in particular the case
for the explicit schedules α and β
given in the first theorem of
Section 4.