# Gaussian Processes to speed up MCMC with automatic exploratory-exploitation effect

**Alessio Benavoli, Jason Wyse, Arthur White**, SCSS, Trinity College Dublin, Ireland

## Introduction

Probabilistic programming languages rely on general-purpose automatic MC sampling techniques for computing accurate approximations of the posterior. These techniques require repeated evaluations of the likelihood function, which can be computationally demanding in some applications. We consider the problem of sampling from a posterior

$$\pi(\theta|D) \propto p(D|\theta)p(\theta),$$

where $D$ denotes data and $\theta \in \Theta$ is a vector of unknown parameters, in the case where the likelihood $p(D|\theta)$ is costly to evaluate. We discuss two-stage algorithms.

## Adaptive Metropolis-Hastings

**Stage 1:** adaptive Metropolis-Hastings (MH) algorithm which employs an adaptively tuned Gaussian Process (GP) surrogate model at the first stage to filter out poor proposals.

**Stage 2:** If a proposal is not filtered out, at the second stage a full (expensive) log-likelihood evaluation is carried out and used to decide whether it is accepted as the next state. Introduction of the 1st stage saves computation on poor proposals.

**Key contribution** is in the form of the acceptance probability in the first stage obtained by marginalising out the GP function. This makes the acceptance ratio dependent on the variance of the GP, which naturally results in an exploration-exploitation trade-off similar to the one of Bayesian Optimisation

## MALA

The second algorithm is a two-stage form of Metropolis adjusted Langevin algorithm (MALA). Here, we use GP as a surrogate for the log-likelihood function again, but in this case the GP is also used to approximate the gradient required for MALA updating, using a well known result that the gradient of a GP is also a GP. **See papers for details.**

The approximation we use is

$$\mathrm{LL}(D|\theta) := \ln p(D|\theta) \approx \widetilde{\mathrm{LL}}_t(D|\theta) \sim \mathrm{GP}(\mu(\theta|\mathscr{I}_t), k(\theta,\theta^*|\mathscr{I}_t)) \quad (1)$$

where $\mathscr{I}_t$ denotes the set of $t$ full evaluations of the log-likelihood by the current iteration, and $\theta^*$ collectively denotes the parameter values at which these evaluations were made.

Notation and explicit definition of $\mathscr{I}_t$:
- $\mathscr{S}_k$ denotes the points sampled up to the iteration $k$ of the algorithm;
- $\theta^{(k)}$ denotes the most recent element in $\mathscr{S}_k$ and $\theta^*$ denotes the proposed state
- $\mathscr{I}_t = \{(\theta^{(i)}, \mathrm{LL}(D|\theta^{(i)})): \ i = 1, \ldots, t\}$ denotes the $t \le k$ exact likelihood evaluations performed up to iteration $k$.

We use a noise free GP as a surrogate model for the log-likelihood and denote by $\mathrm{GP}_k(\mu(\theta|\mathscr{I}_t), k(\theta,\theta|\mathscr{I}_t))$ the posterior GP at the iteration $k$ conditioned on the collection $\mathscr{I}_t$. We use $\widetilde{\mathrm{LL}}_k(\theta)$ to denote the GP-distributed log-likelihood. We choose the parameters of the GP to satisfy the following exact interpolation property.

**Assumption 1** *The prior mean and covariance function of the GP are selected to guarantee exact interpolation:*

$$\mu(\theta^{(i)}|\mathscr{I}_t) = \mathrm{LL}(D|\theta^{(i)}), \qquad k(\theta^{(i)}, \theta|\mathscr{I}_t) = 0,$$

*for all $\theta^{(i)}$ with a corresponding entry in $\mathscr{I}_t$ and $\theta \in \Theta$.*

The two stages of the MH algorithm are as follows.

### Stage 1

Use the predictive posterior GP (conditioned on the collection $\mathscr{I}_t$) to approximate the log-likelihood. Define the first stage acceptance probability:

$$\tilde{\alpha}^{(1)}(\theta^{(k)}, \theta^*) = 1 \wedge \frac{\exp(\widetilde{\mathrm{LL}}_t(\theta^*))p(\theta^*)q(\theta^{(k)}|\theta^*)}{\exp(\widetilde{\mathrm{LL}}_t(\theta^{(k)}))p(\theta^{(k)})q(\theta^*|\theta^{(k)})} \quad (2)$$

where $\widetilde{\mathrm{LL}}_t(\cdot) \sim \mathrm{GP}_k(\mu(\theta|\mathscr{I}_t), k(\theta,\theta|\mathscr{I}_t))$ and we use the shorthand notation $a \wedge b = \min(a, b)$. Note that, because of the exact interpolation property in Assumption 1, it results that $\widetilde{\mathrm{LL}}_t(\theta^{(k)}) = \mathrm{LL}_t(\theta^{(k)})$.

The acceptance probability $\tilde{\alpha}^{(1)}(\theta^{(k)}, \theta^*)$ (respectively, $\tilde{\alpha}^{(1)}(\theta^*, \theta^{(k)})$) depends on $\widetilde{\mathrm{LL}}_t(\theta^*) - \widetilde{\mathrm{LL}}_t(\theta^{(k)})$ (respectively, $-\widetilde{\mathrm{LL}}_t(\theta^*) + \widetilde{\mathrm{LL}}_t(\theta^{(k)})$ ) which is GP distributed.

**Proposition 1** *The distribution of $e^{\widetilde{\mathrm{LL}}_t(\theta)}$ is Lognormal $(\mu(\theta|\mathscr{I}_t), k(\theta,\theta|\mathscr{I}_t))$, and its mean is*

$$e^{\mu(\theta|\mathscr{I}_t)+\frac{1}{2}k(\theta,\theta|\mathscr{I}_t)}. \quad (3)$$

By exploiting Proposition 1, we remove the dependence of the acceptance probability on $\widetilde{\mathrm{LL}}$ in (2) resulting in the acceptance probability:

$$\alpha^{(1)}(\theta^{(k)}, \theta^*) = 1 \wedge \frac{e^{\mu(\theta^*|\mathscr{I}_t)+\frac{1}{2}k(\theta^*,\theta^*|\mathscr{I}_t)}p(\theta^*)q(\theta^{(k)}|\theta^*)}{e^{\mu(\theta^{(k)}|\mathscr{I}_t)}p(\theta^{(k)})q(\theta^*|\theta^{(k)})}, \quad (4)$$

where $\mu(\theta^{(k)}|\mathscr{I}_t) = \mathrm{LL}(\theta^{(k)})$ is the exact log-likelihood (by Assumption 1).

Given (4), in Stage 1, we accept $\theta^*$ with probability $\alpha^{(1)}(\theta^{(k)}, \theta^*)$, otherwise $\theta^{(k+1)} = \theta^{(k)}$. This defines the following transition kernel at Stage 1:

$$Q_k^*(A|\theta^{(k)}) = \int_A \alpha^{(1)}(\theta^{(k)}, \theta^*)q(\theta^*|\theta^{(k)})d\theta^* + I_A(\theta)\int_\Theta (1-\alpha^{(1)}(\theta^{(k)}, \theta^*))d\theta^*. \quad (5)$$

One can show that the above transition kernel satisfies the detailed balance property for the approximated target distribution $e^{\mu(\theta|\mathscr{I}_t)+\frac{1}{2}k(\theta,\theta|\mathscr{I}_t)}p(\theta)$.

### Stage 2.

We perform another MH acceptance step, evaluating the exact log-likelihood. Let $\theta^*$ denote a point sampled from $q_k^*(\theta^*|\theta^{(k)}) := Q_k^*(d\theta^*|\theta^{(k)})$. Note that, $\theta^*$ is either equal to the point $\theta^*$ sampled at Stage 1 or to $\theta^{(k)}$ if $\theta^*$ was rejected at Stage 1. So, with probability

$$\alpha^{(2)}(\theta^{(k)}, \theta^*) = 1 \wedge \frac{\exp(\mathrm{LL}(D|\theta^*))p(\theta^*)q_k^*(\theta^{(k)}|\theta^*)}{\exp(\mathrm{LL}(D|\theta^{(k)}))p(\theta^{(k)})q_k^*(\theta^*|\theta^{(k)})}$$

$$= 1 \wedge \frac{\exp(\mathrm{LL}(D|\theta^*))p(\theta^*)q(\theta^{(k)}|\theta^*)\alpha^{(1)}(\theta^*, \theta^{(k)})}{\exp(\mathrm{LL}(D|\theta^{(k)}))p(\theta^{(k)})q(\theta^*|\theta^{(k)})\alpha^{(1)}(\theta^{(k)}, \theta^*)}, \quad (6)$$
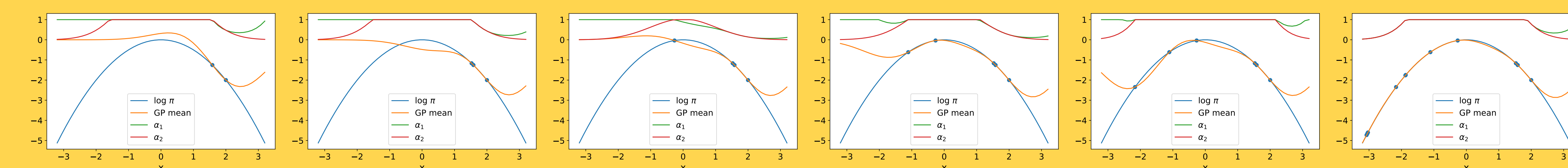
we accept $\theta^*$, otherwise $\theta^{(k+1)} = \theta^{(k)}$.

### Convergence

$$\limsup_{\substack{t\to\infty \\ \theta\in\Theta}} ||P_t(\cdot|\theta) - P_{t-1}(\cdot|\theta)|| = 0 \quad \text{in probability.} \quad (7)$$

## Toy Example

We consider a 1D case with $\pi(\theta) \propto e^{-\frac{x^2}{2}}$. It can be noticed how $\alpha_1$ converges to $\alpha_2$ at the increase of the log-likelihood evaluations.



## Simulations

We consider five target distributions.

**T1:** 2D posterior of the parameters $a, b$ of the banana shape distribution (true value set to $a = 0.2, b = 2$);

**T2:** 3D posterior of the parameters $a, b, \sigma$ of the nonlinear regression model $y = a\frac{x}{x+b} + \varepsilon$, $\varepsilon \sim N(0, \sigma^2)$ (true value set to $a = 0.14, b = 50, \sigma = 0.1$).

**T3:** 3D posterior of the parameters $\ell_1, \ell_2, \sigma^2$ of the SE kernel for a GP-classifier.

**T4:** 4D posterior of the parameters $\beta, \gamma, \sigma_1, \sigma_2$ of a Susceptible, Infected, Recovery (SIR) model.

**T5:** 5D posterior of the parameters $\beta_0, \ldots, \beta_4$ of a parametric logistic regression problem.

We compare our two-stage algorithm with the standard implementations of MH and MALA. For each target problem and in each simulation, we generate 2500 samples (500 for burnin). We have deliberately selected a small number of samples to show that our approach converges quickly, which is important in computationally expensive applications.

|    |         | AR   | ESS | ESJD  | Eval% | SD    |
|----|---------|------|-----|-------|-------|-------|
| T1 | MH      | 0.37 | 90  | 0.13  | 100   | 0.02  |
|    | GP-MH   | 0.36 | 113 | 0.13  | **41** | 0.02  |
|    | MALA    | 0.26 | 73  | 0.2   | 100   | 0.03  |
|    | GP-MALA | 0.26 | 75  | 0.2   | **35** | 0.02  |
| T2 | MH      | 0.28 | 138 | 32.6  | 100   | 339   |
|    | GP-MH   | 0.27 | 133 | 31    | **39** | 339   |
|    | MALA    | 0.26 | 220 | 51    | 100   | 316   |
|    | GP-MALA | 0.21 | 147 | 29    | **43** | 255   |
| T3 | MH      | 0.42 | 137 | 0.44  | 100   | 4.1   |
|    | GP-MH   | 0.42 | 135 | 0.38  | **42** | 3.5   |
|    | MALA    | 0.44 | 133 | 0.48  | 100   | 4.1   |
|    | GP-MALA | 0.43 | 134 | 0.45  | **45** | 3.5   |
| T4 | MH      | 0.1  | 51  | 0.003 | 100   | 0.009 |
|    | GP-MH   | 0.1  | 45  | 0.003 | **15** | 0.009 |
| T5 | MH      | 0.29 | 98  | 0.002 | 100   | 0.006 |
|    | GP-MH   | 0.29 | 102 | 0.002 | **35** | 0.006 |
|    | MALA    | 0.67 | 339 | 0.009 | 100   | 0.006 |
|    | GP-MALA | 0.67 | 368 | 0.009 | **68** | 0.006 |