

## TP 3 : Hasting-Metropolis and Gibbs samplers

### Exercise 1: Hasting-Metropolis within Gibbs – Stochastic Approximation EM

We observe a group of  $N$  (independent) individuals. For the  $i$ -th individual, we have  $k_i$  measurements  $y_{i,j} \in \mathbb{R}$ . In studies on the progression of diseases, measurements  $y_{i,j}$  can be measures of weight, volume of brain structures, protein concentration, tumoral score, etc. over time. We assume that each measurement  $y_{i,j}$  are independent and obtained at time  $t_{i,j}$  with  $t_{i,1} < \dots < t_{i,k_i}$ .

#### 1.A – A population model for longitudinal data

We wish to model an average progression as well as individual-specific progressions of the measurements from the observations  $(y_{i,j})_{1 \leq i \leq N, 1 \leq j \leq k_i}$ . To do that, we consider a hierarchical model defined as follows.

- i. We assume that the average trajectory is the straight line which goes through the point  $p_0$  at time  $t_0$  with velocity  $v_0$

$$d(t) := p_0 + v_0(t - t_0)$$

where

$$p_0 \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\overline{p_0}, \sigma_{p_0}^2) \quad ; \quad t_0 \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\overline{t_0}, \sigma_{t_0}^2) \quad ; \quad v_0 \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\overline{v_0}, \sigma_{v_0}^2)$$

and  $\sigma_{p_0}, \sigma_{t_0}, \sigma_{v_0}$  are fixed variance parameters. While we consider straight lines, we can also fix  $\overline{p_0}$ .

- ii. For the  $i$ -th individual, we assume a trajectory of progression of the form

$$d_i(t) := d(\alpha_i(t - t_0 - \tau_i) + t_0).$$

The trajectory of the  $i$ -th individual corresponds to an affine reparametrization of the average trajectory. This affine reparametrization, given by  $t \mapsto \alpha_i(t - t_0 - \tau_i) + t_0$ , allows to characterize changes in speed and delay in the progression of the  $i$ -th individual with respect to the average trajectory. Moreover, we assume that for all  $i$ -th individual measurements

$$\begin{cases} y_{i,j} = d_i(t_{i,j}) + \varepsilon_{i,j} & \text{where } \varepsilon_{i,j} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2) \\ \alpha_i = \exp(\xi_i) & \text{where } \xi_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_\xi^2) \\ \tau_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_\tau^2) \end{cases}.$$

The parameters of the model are  $\theta = (\overline{t_0}, \overline{v_0}, \sigma_\xi, \sigma_\tau, \sigma)$ . For  $1 \leq i \leq N$ ,  $z_i = (\alpha_i, \tau_i)$  are random variables called *random effects* and  $z_{pop} = (t_0, v_0)$  are called *fixed effects*. The fixed effects are used to model the group progression whereas random effects model individual progressions. Likewise, we define  $\theta_i = (\sigma_\xi, \sigma_\tau, \sigma)$  and  $\theta_{pop} = (\overline{t_0}, \overline{v_0})$ .

We consider a bayesian framework and assume the following *a priori* on the parameters  $\theta$  :

$$\overline{t_0} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\overline{\overline{t_0}}, s_{\overline{t_0}}^2) \quad ; \quad \overline{v_0} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\overline{\overline{v_0}}, s_{\overline{v_0}}^2)$$

$$\sigma_\xi \stackrel{\text{i.i.d.}}{\sim} \mathcal{W}^{-1}(v_\xi, m_\xi) \quad ; \quad \sigma_\tau \stackrel{\text{i.i.d.}}{\sim} \mathcal{W}^{-1}(v_\tau, m_\tau) \quad ; \quad \sigma \stackrel{\text{i.i.d.}}{\sim} \mathcal{W}^{-1}(v, m).$$

where  $\mathcal{W}^{-1}(v, m)$  is the inverse-Wishart distribution :

$$f_{\mathcal{W}^{-1}}(\sigma) = \frac{1}{\Gamma\left(\frac{m_\sigma}{2}\right)} \left( \frac{v}{\sigma\sqrt{2}} \exp\left(-\frac{v^2}{2\sigma^2}\right) \right)^{m_\sigma}.$$

1. Write the complete log-likelihood of the previous model  $\log q(y, z, \theta)$  and show that the proposed model belongs to the curved exponential family.
2. Generate synthetic data from the model by taking some reasonable values for the parameters.

### 1.B – HMwG-SAEM – Hasting-Metropolis within Gibbs sampler

In order to estimate – by a *maximum a posteriori* for instance – the parameters of this statistical model, we will use the SAEM – Stochastic Approximation EM algorithm. However, this algorithm requires that we are able to sample from the *a posteriori* distribution, see algorithm 2.

If we consider  $\pi$ , a density defined on an open set  $\mathcal{U}$  of  $\mathbb{R}^n$  ( $d \geq 2$ ) and if we denote, for  $i \in \llbracket 1, n \rrbracket$ ,  $\pi_i$  the  $i$ th full conditional of  $\pi$  given by :

$$\pi_i(z_i \mid z_{-i}) \propto \pi(z)$$

where  $z_{-i} = \{z_1, \dots, z_{i-1}, z_{i+1}, \dots, z_n\}$  we recall that the classical Gibbs sampler writes as follows :

---

#### Algorithm 1: Gibbs Sampler

---

```

1 Given  $z^{(k)} = (z_1^{(k)}, \dots, z_n^{(k)})$ 
2 for  $i = 1$  to  $n$  do
3    $z_i^{(k+1)} \sim \pi_i(z_i \mid z_1^{(k+1)}, \dots, z_{i-1}^{(k+1)}, z_{i+1}^{(k)}, \dots, z_n^{(k)})$     (*)
4 end
```

---

When direct sampling from the full conditionals is not possible, the step (\*) is often replaced with a Metropolis-Hastings step. The resulting MCMC algorithm is called *hybrid Gibbs sampler* or *Metropolis-Hastings within Gibbs sampler*.

3. Propose a Metropolis-Hastings within Gibbs sampler to sample from the a posteriori distribution of  $z_i = (\xi_i, \tau_i)$ .
4. Likewise, propose a HMwG sampler for the a posteriori distribution of  $z_{pop} = (t_0, v_0)$ .
5. Compute the optimal parameters

$$\theta^{(k)} = \underset{\theta \in \Theta}{\operatorname{argmax}} \{ -\Phi(\theta) + \langle S_k \mid \Psi(\theta) \rangle \}$$

and implement the HMwG-SAEM in order to find the MAP. In particular, we assume that the MAP exists. Use the question 2 to check your algorithm.

For step-sizes  $\varepsilon_k$  we can choose a parameter  $N_b$  – burn-in parameter – and define

$$\forall k \in \mathbb{N}, \quad \begin{cases} 1 & \text{if } 1 \leq k \leq N_b \\ (k - N_b)^{-\alpha} & \text{otherwise} \end{cases}$$

where  $\alpha \in [\frac{1}{2}, 1[$  is necessary to ensure the convergence of the MCMC-SAEM. See [AKT10, AK15].

**Remark :** Contrary to Bayesian inference, where burn-in traditionally refers to a certain amount of samples which are discarded, here the term burn-in refers to memoryless approximation steps. In other words, during the burn-in phase, the information contained in  $z^{(k)}$  is not used in the approximation of the sufficient statistics. In practice, the burn-in period is chosen to be half of the maximum number of iterations.

**Algorithm 2:** MCMC-SAEM (for curved exponential family)

```

1 Given data  $y$  and initial guess  $\theta^{(0)}$ 
2 Initialization :  $z \leftarrow 0, S \leftarrow 0, \theta \leftarrow \theta^{(0)}$  and step-sizes  $(\varepsilon_k)_{k \geq 0}$ .
3 repeat
4   | Simulation using a Metropolis-Hastings-within Gibbs sampler :  $z^{(k)} \stackrel{\text{i.i.d.}}{\sim} \pi_{y, \theta^{(k-1)}}(z^{(k-1)}, \cdot)$ 
5   | Stochastic approximation :  $S \leftarrow S + \varepsilon_k (S(y, z^{(k)}) - S)$ ,
6   | Maximization :  $\theta \leftarrow \underset{\theta \in \Theta}{\operatorname{argmax}} \{-\Phi(\theta) + \langle S \mid \Psi(\theta) \rangle\}$ 
7 until convergence;
```

We can improve the sampling step for big dataset by considering a Block HMwG sampler instead of a "one-at-a-time" as described above HMwG sampler. In the Block version, each Metropolis-Hastings step of the algorithm consists in a multivariate symmetric random walk. Then, the Block MHwG sampler updates simultaneously block (or sets) of latent variables given the others.

6. Explain what are the advantages of a Block Gibbs sampler over a "one-at-a-time" Gibbs sampler for our model.
7. Implement a Block HMwG sampler by choosing a block for the fixed effects and a block by individuals.

The model studied in this exercise is a very simplified version of the model proposed by Jean-Baptiste Schiratti in his Phd-Thesis. For more details, you can refer to [SACD15, Sch16].

## Exercise 2: Multiplicative Hasting-Metropolis

Let  $f$  be a density function on  $] -1, 1[$ . We consider the *multiplicative Hasting-Metropolis algorithm* defined as follows.

Let  $X$  be the current state of the Markov chain.

- (i) We sample  $\varepsilon$  from the probability density function  $f$  and a random variable  $\mathcal{B}$  from the Bernoulli distribution with parameter  $\frac{1}{2}$ .
- (ii) If  $\mathcal{B} = 1$ , we set  $Y = \varepsilon X$ . Otherwise, we set  $Y = \frac{X}{\varepsilon}$ . Then, we accept the candidate  $Y$  with a probability given by  $\alpha(X, Y)$ , the usual Hasting-Metropolis acceptance ratio.

1. Determine the density of the jumping distribution  $Y \sim q(X, Y)$ .
2. Compute the acceptance ratio  $\alpha$  so that the chain has a given distribution  $\pi$  as invariant distribution.
3. Implement this sampler for two different target distributions : the first one being a distribution from which we can sample using the inverse transform method and the second one is of your choice.
4. Evaluate, in each case, the match of your samples with the true distribution.

### Exercise 3: Data augmentation

Let  $f : (x, y) \in \mathbb{R}^p \times \mathbb{R}^q \mapsto f(x, y) \in \mathbb{R}^+$  be a density with respect to the Lebesgue measure on  $\mathbb{R}^{p+q}$ .  
 Les define

$$f_X(x) := \int f(x, y) dy ; \quad f_Y(y) = \int f(x, y) dx ;$$

and

$$\forall y \in \mathbf{Y} := \{y \in \mathbb{R}^q \mid f_Y(y) > 0\}, \quad f_{X|Y}(x, y) := \frac{f(x, y)}{f_Y(y)} ;$$

$$\forall x \in \mathbf{X} := \{x \in \mathbb{R}^p \mid f_X(x) > 0\}, \quad f_{X|Y}(x, y) := \frac{f(x, y)}{f_Y(y)} .$$

We define a bivariate chain  $\{(X_n, Y_n), n \geq 0\}$  as in the following algorithm.

#### Algorithm 3: Data augmentation

```

1 Given  $(X_0, Y_0) \in \mathbb{R}^p \times \mathbb{R}^q$  and  $N \in \mathbb{N}$ 
2 for  $n = 1$  to  $N$  do
3    $X_n \sim f_{X|Y}(\cdot, Y_{n-1})$ 
4    $Y_n \sim f_{Y|X}(X_n, \cdot)$ 
5 end
6 return  $\{(X_n, Y_n), 0 \leq n \leq N\}$ 
```

1. Show that the bivariate process  $\{(X_n, Y_n), n \geq 0\}$  is a Markov chain. Give the expression of its transition kernel as a function of the quantities defined above.

2. Show that  $\{Y_n, n \geq 0\}$  is a Markov chain : give the expression of its transition kernel and prove that  $f_Y(y) dy$  is invariant for this kernel.

Hereafter, we consider the case when

$$f(x, y) = \frac{4}{\sqrt{2\pi}} y^{\frac{3}{2}} \exp \left[ -y \left( \frac{x^2}{2} + 2 \right) \right] \mathbb{1}_{\mathbb{R}^+}(y)$$

3. Describe a Gibbs algorithm to approximate the distribution on  $\mathbb{R} \times \mathbb{R}$  with density  $f$ .

We can use a gamma distribution sampler : `numpy.random.gamma` in python or `gamrnd` from the Statistics and Machine Learning Toolbox in Matlab. We can also find a Matlab toolbox-free Gamma Generator in the Handbook of Monte Carlo Methods [KTB13] :  
<https://people.smp.uq.edu.au/DirkKroese/montecarlohandbook/probdist/>.

4. Let  $H$  be a bounded function on  $\mathbb{R}$ . Explain how to approximate

$$\int_{\mathbb{R}} \frac{H(x)}{(4+x^2)^{\frac{5}{2}}} dx$$

from the output  $\{(X_n, Y_n), 0 \leq n \leq N\}$  of this Gibbs sampler.

## References

- [AK15] Stéphanie Allasonnière and Estelle Kuhn. Convergent stochastic expectation maximization algorithm with efficient sampling in high dimension. application to deformable template model estimation. *Computational Statistics & Data Analysis*, 91:4 – 19, 2015.
- [AKT10] Stéphanie Allasonnière, Estelle Kuhn, and Alain Trouvé. Construction of bayesian deformable models via a stochastic approximation algorithm: A convergence study. *Bernoulli*, 16(3):641–678, 08 2010.
- [Dut12] Somak Dutta. Multiplicative random walk metropolis-hastings on the real line. *Sankhya B*, 74(2):315–342, 2012.
- [KTB13] D.P. Kroese, T. Taimre, and Z.I. Botev. *Handbook of Monte Carlo Methods*. Wiley Series in Probability and Statistics. Wiley, 2013.
- [SACD15] Jean-Baptiste Schiratti, Stéphanie Allasonnière, Olivier Colliot, and Stanley Durrleman. Learning spatiotemporal trajectories from manifold-valued longitudinal data. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28*, pages 2404–2412. Curran Associates, Inc., 2015.
- [Sch16] Jean-Baptiste Schiratti. *Models and algorithms to learn spatiotemporal changes from longitudinal manifold-valued observations*. PhD thesis, École polytechnique, 2016.