

# Dynamic Models with Latent Variables

Jean-Michel Zakoian

CREST-ENSAE

Bayesian and simulated methods

1 Simulation algorithms

2 Examples

## Simulation-based method

- In recent years, inference methods based on simulations have developed tremendously.
- Such methods are particularly appropriate when the likelihood cannot be used directly.
- One popular method is the **Monte Carlo Markov Chain (MCMC)** technique.

# Bayesian inference and MCMC

The *Bayesian* approach combines information brought by the data with *a priori* ideas on the parameters.

Such ideas are represented by probability distributions, called *prior* distributions, where the parameters have the status of random variables.

The goal of the inference is to obtain *posterior distributions* of the parameters, which are deduced, thanks to the Bayes formula, from the *a priori* distribution and the law of the observations conditional on the parameters.

## Notations

Parameter :  $\theta \in \mathbb{R}^d$ .

Observations :  $\mathbf{X}$ .

*A priori* density of the parameter :  $\pi$ .

*A posteriori* density of the parameter :  $\pi(\cdot | \mathbf{X})$

Likelihood (the density of the observations conditional to  $\theta$ ) :  $f(\cdot | \theta)$ .

By the Bayes formula we have :

$$\pi(\theta | \mathbf{X}) = \frac{f(\mathbf{X} | \theta)\pi(\theta)}{\int f(\mathbf{X} | \vartheta)\pi(\vartheta)d\vartheta} \propto f(\mathbf{X} | \theta)\pi(\theta).$$

## Bayesian estimator of $g(\theta)$

Obtained by minimizing the *a posteriori cost*

$$\int L(\delta, \vartheta) \pi(\vartheta | \mathbf{X}) d\vartheta,$$

where  $L$  is a positive **cost function**.

For the **quadratic cost**,  $L(\delta, \theta) = \|g(\theta) - \delta\|^2$ , the Bayesian estimator is (if the integral exists),

$$\hat{\delta} = E\{g(\theta) | \mathbf{X}\} = \int g(\vartheta) \pi(\vartheta | \mathbf{X}) d\vartheta.$$

This expression has little practical interest because in general,

- (i) the *posterior* distribution is not known explicitly,
- (ii) the integral is difficult to compute. This problem can be handled by *numerical methods*, or by **simulations based methods**.

## Simulations based methods

Suppose that a sequence  $\{\theta^{(i)}(\mathbf{X})\}_{i \geq 1}$  of variables with density  $\pi(\cdot | \mathbf{X})$  is available.

Assuming that the **ergodic theorem** can be applied, conditionally to  $\mathbf{X}$ ,

$$\hat{\delta}_N = \frac{1}{N} \sum_{i=1}^N g\{\theta^{(i)}(\mathbf{X})\} \rightarrow E\{g(\theta) | \mathbf{X}\}, \quad p.s. \text{ as } N \rightarrow \infty.$$

# MCMC

The MCMC approach is a set of techniques allowing to generate such sequences  $\{\theta^{(i)}(\mathbf{X})\}_{i \geq 1}$ .

The basic idea of the MCMC method is to approximate a probability law  $\mathbb{P}$  using a Markov chain  $(\theta^{(i)})$ , which is irreducible, aperiodic, and admits  $\mathbb{P}$  as invariant law.

The target law is simulated using an initial value  $\theta^{(0)}$ , and the transition probabilities of the chain to generate  $\theta^{(1)}, \dots, \theta^{(N)}$ .

For  $N$  sufficiently large,  $\theta^{(N)}$  can be considered as a realization drawn from the law  $\mathbb{P}$ .

Two methods for constructing such chains are *the Metropolis-Hastings algorithm* and *the Gibbs sampling*.



- 1 Simulation algorithms
  - Accept-reject method
  - Metropolis-Hastings algorithm
  - Gibbs sampling
- 2 Examples

## Principle

We wish to simulate the density  $\pi : \theta \rightarrow \pi(\theta)$  of some rv  $\theta$ .

Finding an **explicit formula** for  $F^{-1}(U)$ , where  $F$  is the cdf (cumulative distribution function) of  $\theta$  and  $U \sim \mathcal{U}[0,1]$  is **not always possible**.

Moreover, even if it is, there may be alternative methods for generating a rv distributed as  $F$  that are **more efficient** than the inverse transform method.

**Basic idea** : find an alternative probability distribution  $G$ , with density function  $g(\theta)$ , from which we already have an efficient simulation algorithm, but also such that the function  $g(\theta)$  is "close" to  $f(\theta)$ .

## Accept-reject algorithm

Suppose that for all  $\theta$

$$\pi(\theta) \leq cg(\theta),$$

where  $c$  is a constant ( $c \geq 1$ ) and  $g$  is a density which can be easily simulated.

### Algorithm :

- ① Generate  $Z \sim g$  and generate  $U \sim \mathcal{U}_{[0, cg(Z)]}$ .
- ② Accept  $Z$  if  $U < \pi(Z)$ , otherwise reject and go back to 1.

### Equivalently :

- ① Generate  $Z \sim g$  and  $U \sim \mathcal{U}_{[0,1]}$ .
- ② Accept  $Z$  if  $U < \frac{\pi(Z)}{cg(Z)}$ , otherwise reject and go back to 1.

## Proof that the algorithm works

$Z^*$  : the variable generated by the algorithm.

For all  $z^*$ ,

$$P[Z^* < z^*] = P[Z < z^* \mid U < \pi(Z)] = \frac{P[Z < z^*, U < \pi(Z)]}{P[U < \pi(Z)]}.$$

Moreover,

$$P[Z < z^*, U < \pi(Z)] = \int_{z < z^*} g(z) \frac{1}{cg(z)} \int_{u < \pi(z)} du dz = \frac{1}{c} \int_{z < z^*} \pi(z) dz.$$

and, by taking  $z^* = \infty$ ,

$$P[U < \pi(Z)] = \frac{1}{c} \int \pi(z) dz = \frac{1}{c}.$$

Finally,

$$P[Z^* < z^*] = \int_{z < z^*} \pi(z) dz.$$

## Remarks

- By iterating the procedure we get an iid sample distributed as  $\pi$ .
- The probability of accepting a simulation of  $Z$  is  $1/c$  : the constant  $c$  has to be chosen as small as possible to minimize the computation time.
- $c$  depends on the tails of the laws  $\pi$  and  $g$  : in particular, the ratio  $\pi/g$  has to be bounded (it is not possible to simulate a Cauchy from a Gaussian, but the converse is possible).
- If  $c$  is chosen too small so that the condition  $\pi(\theta) \leq cg(\theta)$  is not satisfied for some values of  $\theta$ , the simulated density is not  $\pi$  but rather

$$\pi^*(\theta) \propto \min\{\pi(\theta), cg(\theta)\}.$$

## Simulation of an a posteriori distribution

We wish to simulate the *a posteriori* distribution

$$\pi(\theta \mid \mathbf{X}) = h(\mathbf{X})\pi^*(\theta, \mathbf{X})$$

where  $\pi^*$  is a known function but  $h$  has a complicated non explicit form.

If there exists a density  $g^*(\cdot \mid \mathbf{X})$  such that

$$\pi^*(\theta, \mathbf{X}) \leq cg^*(\theta \mid \mathbf{X}),$$

the method can be adapted as follows :

- ① Generate  $\theta$  distributed as  $g^*(\cdot \mid \mathbf{X})$  and  $U \sim \mathcal{U}_{[0,1]}$ .
- ② Accept  $\theta$  if  $U < \frac{\pi^*(\theta, \mathbf{X})}{cg^*(\theta \mid \mathbf{X})}$ , otherwise reject and go back to 1.

# Metropolis-Hastings (MH) algorithm

Metropolis et al. (1953) [discrete state space], Hastings (1970) [statistical setting].

The MH algorithm allows to generate a density on which we have little information.

To simulate the *posterior* density  $\pi(\cdot | \mathbf{X})$ , the MH algorithm allows to generate a sequence of variables  $\{\theta^{(i)}\}_{i \geq 1} = \{\theta^{(i)}(\mathbf{X})\}_{i \geq 1}$  which, conditionally to  $\mathbf{X}$ , is a **Markov chain** with **stationary distribution**  $\pi(\cdot | \mathbf{X})$ .

## Markov chains

$(X_t)$  is an homogenous Markov chain on  $(E, \mathcal{E})$  if :

$$\forall x \in E, \forall B \in \mathcal{E}, \forall s, t \in \mathbb{N},$$

$$P(X_{s+t} \in B \mid X_r, r < s; X_s = x) = P(X_{s+t} \in B \mid X_s = x) := P^t(x, B).$$

The mapping  $P: E \times \mathcal{E} \rightarrow [0, 1]$  is called *transition kernel* and it satisfies :

- (i)  $\forall B \in \mathcal{E}$ , the function  $P(\cdot, B)$  is measurable;
- (ii)  $\forall x \in E$ , the function  $P(x, \cdot)$  is a probability over  $(E, \mathcal{E})$ .

The law of  $(X_t)$  is characterized by an initial probability  $\mu$  and a transition kernel  $P$ .



## Markov chains

Under certain conditions, there exists a probability  $\pi$  such that  
 $\forall x \in E, \forall B \in \mathcal{E}$ ,

$$P^t(x, B) \rightarrow \pi(B), \quad \text{as } t \rightarrow \infty.$$

$\pi$  is called **invariant probability** and it satisfies :

$$\forall B \in \mathcal{E}, \quad \pi(B) = \int P(x, B) \pi(dx).$$

Under additional conditions, the chain is **ergodic** and we have, a.s.

$$\frac{1}{n} \sum_{t=1}^n g(X_t) \rightarrow \int g d\pi, \quad \text{as } n \rightarrow \infty$$

for any function  $g$  which is with respect to  $\pi$ .

# MH algorithm

The algorithm is based on a transition kernel of the form

$$Q(\theta, d\theta') = q(\theta, \theta')\lambda(d\theta'), \quad \theta \in \mathbb{R}^d.$$

The transition densities,  $q(\theta, \cdot)$ , also called *jump densities*, have to be chosen in order to explore the parameter space. They may depend on  $\mathbf{X}$  and are used at each step of the algorithm to simulate a new possible value and to decide whether this value will be kept or rejected.

## MH algorithm

- ① Choose an initial value  $\theta^{(0)}$  such that  $\pi(\theta^{(0)} | \mathbf{X}) > 0$ .
- ② For  $i=1,2,\dots$ 
  - ① Generate  $\theta^*$  distributed as  $Q(\theta^{(i-1)}, \cdot)$ .
  - ② Compute the ratio

$$r = \frac{\pi(\theta^* | \mathbf{X})}{\pi(\theta^{(i-1)} | \mathbf{X})} \frac{q(\theta^{(i-1)}, \theta^*)}{q(\theta^*, \theta^{(i-1)})}.$$

- ③ Take

$$\theta^{(i)} = \begin{cases} \theta^* & \text{with probability } \min(r, 1) \\ \theta^{(i-1)} & \text{with probability } 1 - \min(r, 1). \end{cases}$$

**Remark :** the condition  $\pi(\theta^{(0)} | \mathbf{X}) > 0$  ensures that the ratio is well defined for all  $i$ .

## Remarks

- If  $q(x, y) = q(y, x)$  for all  $(x, y)$ , the acceptance rule does not depend on  $q$ . The probability of acceptance of  $\theta^*$  reduces to

$$\min\left(\frac{\pi(\theta^* | \mathbf{X})}{\pi(\theta^{(i-1)} | \mathbf{X})}, 1\right)$$

and  $\theta^*$  will be always accepted if jumping from  $\theta^{(i-1)}$  to  $\theta^*$  increases the *a posteriori* density; otherwise,  $\theta^*$  is accepted with a probability equal to the ratio  $r$  of the *posterior* densities.

- In the general case,  $\theta^*$  is accepted with probability 1 only if the ratio of *posterior* densities is greater than the ration of densities of jump, from  $\theta^{(i-1)}$  to  $\theta^*$  and from  $\theta^*$  to  $\theta^{(i-1)}$ .
- **$r$  does not depend on normalizing constants** involved in the conditional densities. In particular, it is sufficient to know  $\pi(\theta^* | \mathbf{X})$  up to a multiplicative constant to use this algorithm. Thus, it suffices to evaluate  $f(\mathbf{X} | \cdot)\pi(\cdot)$  in  $\theta^*$  and  $\theta^{(i-1)}$ .

## Examples of transition kernels

(i) *Kernel associated with the random walk* : at step  $i$ , the value  $\theta^*$  is generated using the model

$$\theta^* = \theta^{(i-1)} + \epsilon$$

where  $\epsilon$  is a centered variable with density  $f$ , independent of  $\theta^{(i-1)}$ .  
Thus

$$q(\theta, \theta') = f(\theta' - \theta).$$

If  $f$  is symmetric around 0, the kernel is symmetric.

Note that the sequence  $(\theta^{(i)})$  is not a random walk, due to repetitions  $(\theta^{(i)} = \theta^{(i-1)})$  which are non independent from  $\theta^{(i-1)}$ .

## Examples of transition kernels

(ii) *Independent Kernel* : the values  $\theta^*$  are drawn, independently of  $\theta^{(i-1)}$ , from a density  $f$  :

$$q(\theta, \theta') = f(\theta')$$

and the probability of acceptance of the value  $\theta^*$  is

$$\min\left(\frac{\omega(\theta^*)}{\omega(\theta^{(i-1)})}, 1\right)$$

where  $\omega(\theta) = \pi(\theta \mid \mathbf{X})/f(\theta)$ .

The function  $\omega$  can be interpreted as an importance function used to simulate  $\pi$  from simulations of the law  $f$ .

Low-weight candidates will be seldom accepted ; conversely, large-weight candidates will generally be chosen repeatedly.

## Examples of transition kernels

(iii) *Kernel associated with an accept-reject algorithm* : particular case of independent kernel where  $f = \pi(\cdot | \mathbf{X})$  is simulated by the accept-reject method.

If the condition  $\pi(\theta | \mathbf{X}) \leq cg(\theta)$  is in failure, the density generated for the  $\theta^*$  will be

$$\pi^*(\theta^* | \mathbf{X}) \propto \min\{\pi(\theta^* | \mathbf{X}), cg(\theta^*)\}.$$

Letting  $C = \{\theta | \pi(\theta | \mathbf{X}) \leq cg(\theta)\}$ , the acceptance probability of the value  $\theta^*$  is

$$\begin{cases} 1, & \text{if } \theta^{(i-1)} \in C \\ \frac{cg(\theta^{(i-1)})}{\pi(\theta^{(i-1)} | \mathbf{X})}, & \text{if } \theta^{(i-1)} \notin C, \theta^* \in C, \\ \min\left(\frac{\pi(\theta^* | \mathbf{X})g(\theta^{(i-1)})}{\pi(\theta^{(i-1)} | \mathbf{X})g(\theta^*)}, 1\right), & \text{if } \theta^{(i-1)} \notin C, \theta^* \notin C. \end{cases}$$

## Examples of transition kernels

- Some proposals are rejected when  $\theta^{(i-1)} \notin C$ . This value is thus repeated, which compensates the deficiency, due to the failure of the bound  $cg$ , in this region.  
Dependence is thus introduced to alleviate this problem.
- If the deficiency never holds, the values accepted by the accept-reject algorithm are also accepted by the global algorithm and the resulting sample is iid.



## Examples of transition kernels

(iv) *Kernel associated with an AR :*

$$\theta^* = a + b(\theta^{(i-1)} - a) + \epsilon$$

where  $\epsilon$  is a variable with density  $f$ , independent from  $\theta^{(i-1)}$ . Thus

$$q(\theta, \theta') = f(\theta' - a - b(\theta - a)).$$

The choice of  $b < 0$  allows to introduce negative correlations between the successive value of the algorithm and to explore faster the support of  $\pi$ .

## Markov chain property

Conditionally to  $\mathbf{X}$ , the sequence  $(\theta^{(i)}(\mathbf{X}))_{i \geq 0}$  is a Markov chain on  $\mathbb{R}^d$ .

It can be shown that, under certain conditions on the transition kernel, this MC admits  $\pi(\cdot | \mathbf{X})$  as an invariant probability measure and is **ergodic** :

$$\hat{\delta}_N = \frac{1}{N} \sum_{i=1}^N g\{\theta^{(i)}(\mathbf{X})\} \rightarrow E_{\pi}\{g(\theta) | \mathbf{X}\}, \quad p.s. \text{ as } N \rightarrow \infty.$$

## Example : STAR model

$$Y_t = a_0 + a_1 Y_{t-1} + \frac{(b_0 - a_0) + (b_1 - a_1) Y_{t-1}}{1 + \exp\{-\gamma(Y_{t-1} - c)\}} + \epsilon_t, \quad \epsilon_t \sim \text{IID}(0, \sigma^2), \quad \gamma > 0.$$

► Gibbs estimation

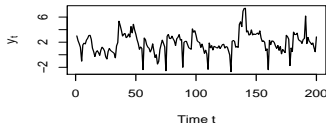
$E(Y_t | Y_{t-1} = y)$  varies continuously from  $a_0 + a_1 y$  (when  $y \rightarrow -\infty$ ) to  $b_0 + b_1 y$  (when  $y \rightarrow +\infty$ ).

Simulation of

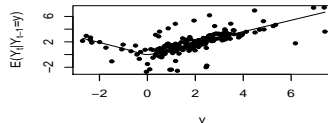
$$Y_t = -0.95 Y_{t-1} + \frac{1.85 Y_{t-1}}{1 + e^{-5 Y_{t-1}}} + \epsilon_t$$

where  $\epsilon_t \sim 0.9\mathcal{N}(0, 0.5^2) + 0.05\mathcal{N}(3, 1) + 0.05\mathcal{N}(-3, 1)$ .

(a) Simulation ( $Y_t$ ) of a STAR



(b) Points ( $Y_{t-1}, Y_t$ ) and regression function



## Example : STAR model

**Independent prior laws :**  $c$  fixed to 0,  $a_0$ ,  $a_1$ ,  $b_0$ , and  $b_1$  follow a  $\mathcal{N}(0, 1)$ ,  $\gamma \sim \mathcal{E}(1)$ , and  $\sigma^2 \sim IG(1, 1)^*$ .

We thus have, for  $\theta = (a_0, a_1, b_0, b_1, \gamma, \sigma^2)$ ,

$$\begin{aligned}\pi(\theta | \mathbf{X}) &\propto f(\mathbf{X} | \theta) \pi(\theta) \\ &= \prod_{t=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{1}{2\sigma^2} (Y_t - m_{t-1})^2 \right\} \\ &\quad \times \frac{1}{(2\pi)^2} e^{-(a_0^2 + a_1^2 + b_0^2 + b_1^2)/2} e^{-\gamma} \frac{e^{-1/\sigma^2}}{(\sigma^2)^2}\end{aligned}$$

---

\*. the inverse gamma distribution  $IG(a, b)$  has density

$$f(x) = \frac{b^a}{\Gamma(a)} \frac{e^{-b/x}}{x^{a+1}} 1_{[0, +\infty[}(x).$$

## Example : STAR model

The MH algorithm is used with transition kernel derived from the random walk,

$$Q(x, y) \sim \mathcal{N}(x, \tau I_6).$$

The probability of acceptance of  $\theta^*$  is thus  $\min\left(\frac{\pi(\theta^*|\mathbf{X})}{\pi(\theta^{(i-1)}|\mathbf{X})}, 1\right)$ .

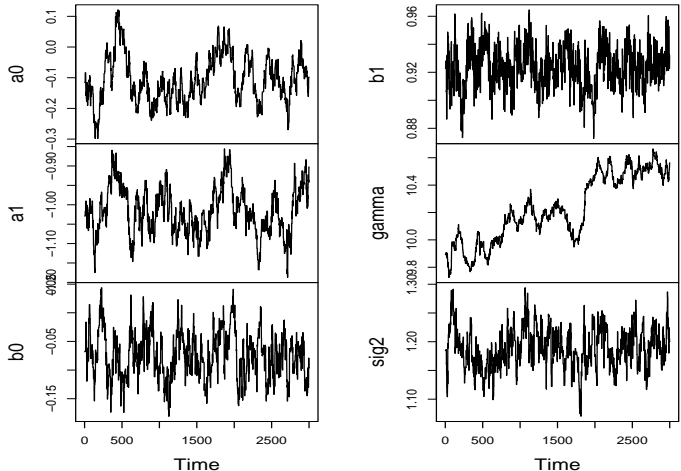
In practice the parameter  $\tau$  is very important for the performance of the algorithm :

- large values of  $\tau$  entail small rates of acceptance ( $\theta^*$  is likely to fall in a region with low density  $\pi(B|\mathbf{X})$ , and thus is likely to be rejected).
- Small values of  $\tau$  have high rates of acceptance and correspond to small moves of the MC.

For  $\tau = 0.02$  and  $n = 3000$  the acceptance rate for  $\theta^*$  is 0.37 but the MC has not converged.

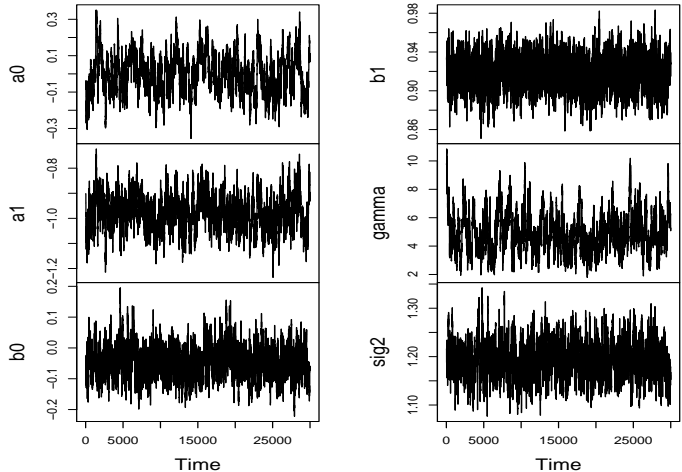
## Example : STAR model

### Markov chain simulated by Metropolis



## Example : STAR model

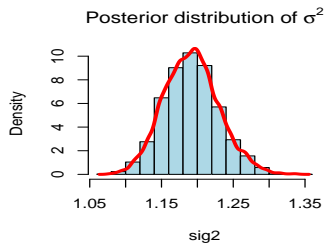
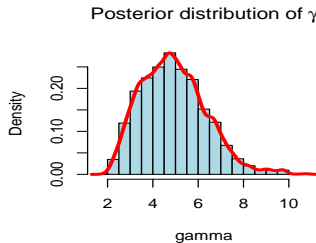
**Markov chain simulated by Metropolis**



## Example : STAR model

True parameter values :

$$a_0 = b_0 = c = 0, \quad a_1 = -0.95, \quad b_1 = 0.9, \quad \gamma = 5, \quad \sigma^2 = 1.225.$$





## Gibbs sampling (Geman and Geman (1984), Gelfand and Smith (1990))

An alternative method for generating a MC having a specified invariant distribution.

However, the algorithm

- is not based on an acceptance/rejection mechanism : all simulated values are accepted ;
- only applies to multivariate distributions ;
- requires more information than MH on the target law.

$\theta_{[-j]} \in \mathbb{R}^{d-1}$  : parameter vector **without**  $\theta_j$ .

We assume that the conditional densities

$$\pi^j(\theta_j | \mathbf{X}, \theta_{[-j]}), \quad j = 1, \dots, d$$

**are available and can be simulated.**

# Gibbs sampler

- ① Choose initial values  $\theta_j^{(0)}$ ,  $j=2,\dots,d$ .
- ② For  $i=1,2,\dots$ 
  - ① Generate  $\theta_1^{(i)}$  in the law  $\pi^1(\cdot | \mathbf{X}, \theta_{[-1]}^{(i-1)})$
  - ② For  $\ell=2,\dots,d$ , generate  $\theta_\ell^{(i)}$  in the law

$$\pi^\ell(\cdot | \mathbf{X}, \theta_1^{(i)}, \dots, \theta_{\ell-1}^{(i)}, \theta_{\ell+1}^{(i-1)}, \dots, \theta_d^{(i-1)})$$

- ③ Take  $\theta^{(i)} = (\theta_1^{(i)}, \dots, \theta_d^{(i)})$ .

## Remarks

- The simulated laws are those of the coordinates of  $\theta$ , conditional on the other coordinates.  
The components can be multivariate : an extension consists in simulating blocks of  $\theta$ , provided that their conditional distributions given the other blocks be available.
- It can be interesting to interpret the target law  $\pi(\theta | \mathbf{X})$ , as the marginal law of a higher-dimensional vector whose conditional distributions are easier to simulate than those of  $\theta$  (data augmentation technique).

## Why does it work ?

Let  $\pi_i$  the marginal densities (given  $\mathbf{X}$ ) of the  $\theta_i$ 's. We introduce the *positivity* assumption :

$$\mathbf{H}: \quad \forall \theta, \quad \prod_{i=1}^d \pi_i(\theta_i | \mathbf{X}) > 0 \quad \implies \quad \pi(\theta | \mathbf{X}) > 0,$$

Theorem (Hammersley and Clifford (1970))

*Under Assumption  $\mathbf{H}$ , for all  $\theta, \theta'$  such that  $\pi(\theta | \mathbf{X}) > 0$  and  $\pi(\theta' | \mathbf{X}) > 0$ ,*

$$\frac{\pi(\theta | \mathbf{X})}{\pi(\theta' | \mathbf{X})} = \prod_{i=1}^d \frac{\pi^i(\theta_i | \mathbf{X}, \theta_1, \dots, \theta_{i-1} \theta'_{i+1} \dots \theta'_d)}{\pi^i(\theta'_i | \mathbf{X}, \theta_1, \dots, \theta_{i-1} \theta'_{i+1} \dots \theta'_d)}.$$

## Proof :

For any  $\theta$  such that  $\pi(\theta | \mathbf{X}) > 0$ ,

$$\pi(\theta | \mathbf{X}) = \pi^d(\theta_d | \mathbf{X}, \theta_1, \dots, \theta_{d-1}) \pi_{[-d]}(\theta_1, \dots, \theta_{d-1} | \mathbf{X})$$

where  $\pi_{[-i]}(\cdot | \mathbf{X})$  is the density of the vector  $\theta$  without  $\theta_i$ .

Moreover, for all  $\theta'$

$$\pi(\theta_1, \dots, \theta_{d-1}, \theta'_d | \mathbf{X}) = \pi^d(\theta'_d | \mathbf{X}, \theta_1, \dots, \theta_{d-1}) \pi_{[-d]}(\theta_1, \dots, \theta_{d-1} | \mathbf{X}).$$

It follows that

$$\pi(\theta | \mathbf{X}) = \frac{\pi^d(\theta_d | \mathbf{X}, \theta_1, \dots, \theta_{d-1})}{\pi^d(\theta'_d | \mathbf{X}, \theta_1, \dots, \theta_{d-1})} \pi(\theta_1, \dots, \theta_{d-1}, \theta'_d | \mathbf{X}).$$

## Proof (continued)

Similarly

$$\begin{aligned} & \pi(\theta_1, \dots, \theta_{d-1}, \theta'_d | \mathbf{X}) \\ = & \pi^{d-1}(\theta_{d-1} | \mathbf{X}, \theta_1, \dots, \theta_{d-2}, \theta'_d) \pi_{[-(d-1)]}(\theta_1, \dots, \theta_{d-2}, \theta'_d | \mathbf{X}), \\ & \pi(\theta_1, \dots, \theta'_{d-1}, \theta'_d | \mathbf{X}) \\ = & \pi^{d-1}(\theta'_{d-1} | \mathbf{X}, \theta_1, \dots, \theta_{d-2}, \theta'_d) \pi_{[-(d-1)]}(\theta_1, \dots, \theta_{d-2}, \theta'_d | \mathbf{X}), \end{aligned}$$

hence

$$\begin{aligned} \pi(\theta | \mathbf{X}) = & \frac{\pi^d(\theta_d | \mathbf{X}, \theta_1, \dots, \theta_{d-1})}{\pi^d(\theta'_d | \mathbf{X}, \theta_1, \dots, \theta_{d-1})} \frac{\pi^{d-1}(\theta_{d-1} | \mathbf{X}, \theta_1, \dots, \theta_{d-2}, \theta'_d)}{\pi^{d-1}(\theta'_{d-1} | \mathbf{X}, \theta_1, \dots, \theta_{d-2}, \theta'_d)} \\ & \times \pi(\theta_1, \dots, \theta_{d-2}, \theta'_{d-1}, \theta'_d | \mathbf{X}). \end{aligned}$$

Continuing this way the replacement of the coordinates of  $\theta$  by those of  $\theta'$ , we get the announced formula.

## Proof (continued)

Finally, we check the positivity of the denominator in this formula.  
We have

$$\begin{aligned} & \pi(\theta \mid \mathbf{X}) > 0, \quad \pi(\theta' \mid \mathbf{X}) > 0 \\ \Rightarrow & \pi_i(\theta_i \mid \mathbf{X}) > 0, \quad i = 1, \dots, d-1 \quad \text{and} \quad \pi_d(\theta'_d \mid \mathbf{X}) > 0 \\ \Rightarrow & \pi(\theta_1, \dots, \theta_{d-1}, \theta'_d \mid \mathbf{X}) > 0, \quad \text{by assumption } \mathbf{H} \\ \Rightarrow & \pi^d(\theta'_d \mid \mathbf{X}, \theta_1, \dots, \theta_{d-1}) > 0. \end{aligned}$$

The other terms of the product are treated similarly.

## Remarks

- ① This relation allows to get, theoretically, the joint density from the conditional densities (by integrating the inverses of each side of the equality w.r.t.  $\theta'$ ) :

$$\pi(\theta | \mathbf{X}) = \left( \int \prod_{i=1}^d \frac{\pi^i(\theta'_i | \mathbf{X}, \theta_1, \dots, \theta_{i-1}, \theta'_{i+1}, \dots, \theta'_d)}{\pi^i(\theta_i | \mathbf{X}, \theta_1, \dots, \theta_{i-1}, \theta'_{i+1}, \dots, \theta'_d)} d\theta' \right)^{-1}.$$

- ② It can be shown that under Assumption **H**,  $(\theta^{(i)})$  is an **ergodic Markov chain** admitting  $\pi(\cdot | \mathbf{X})$  as invariant probability measure.



## Hybrid algorithms

It can be worth **combining the Gibbs and MH algorithms** :

- the MH algorithm does not take into account possible available information on the conditional laws. Moreover, it is not well suited for hierarchical structures (such as the SV model).
- the Gibbs algorithm cannot be used alone if certain conditional densities are not available in closed form or cannot be easily simulated.
- Even if the Gibbs algorithm can be implemented, the convergence to the stationary distribution of the MC can be very slow because the components are modified one by one.

## 1 Simulation algorithms

## 2 Examples

- AR(1) model estimation by the Gibbs algorithm
- STAR(1) model estimation by an hybrid algorithm

## AR(1) model

$$X_t = \omega + \beta X_{t-1} + \sigma v_t, \quad (v_t) \text{ iid } \mathcal{N}(0, 1), \quad \sigma > 0.$$

**Aim** : generate a MC with stationary *a posteriori* distribution,  $\pi(\theta | \mathbf{X})$ , where  $\theta = (\omega, \beta, \sigma)'$ .

**Prior laws** :  $(\omega, \beta)$  and  $\sigma^2$  are *a priori* independent with

$$(\omega, \beta) \sim \mathcal{N}((\omega^0, \beta^0), \Sigma^0), \quad \sigma^2 \sim IG(a, b)$$

$a > 0, b > 0, \omega^0 \in \mathbb{R}, \beta^0 \in \mathbb{R}, \Sigma^0$  positive definite.

The *posterior conditional laws* of  $(\omega, \beta)$  and  $\sigma$  are denoted

$$\pi^1((\omega, \beta) | \mathbf{X}, \sigma) \quad \text{and} \quad \pi^2(\sigma^2 | \mathbf{X}, \omega, \beta).$$

## Notations

Likelihood conditional to a fixed initial value  $X_0$  :

$$f(\mathbf{X} | \theta) \propto \exp \left\{ \frac{-1}{2\sigma^2} \sum_{t=1}^n (X_t - \omega - \beta X_{t-1})^2 \right\}.$$

$(\hat{\omega}, \hat{\beta})$  : LS estimator of  $(\omega, \beta)$ .

$$\Sigma_n = \Sigma_n(\sigma) = \sigma^2 (\underline{\mathbf{X}}' \underline{\mathbf{X}})^{-1}, \quad \underline{\mathbf{X}}' = \begin{pmatrix} 1 & \cdots & 1 \\ X_0 & \cdots & X_{n-1} \end{pmatrix},$$

$$\Sigma^* = \{(\Sigma^0)^{-1} + \Sigma_n^{-1}\}^{-1},$$

$$(\omega^*, \beta^*)' = \Sigma^* \{ \Sigma_n^{-1} (\hat{\omega}, \hat{\beta})' + (\Sigma^0)^{-1} (\omega^0, \beta^0)' \}.$$

## LS estimator

We have

$$\begin{aligned}\sum_{t=1}^n (X_t - \omega - \beta X_{t-1})^2 &= \sum_{t=1}^n [X_t - \hat{\omega} - \hat{\beta} X_{t-1} - \{\omega - \hat{\omega} + (\beta - \hat{\beta}) X_{t-1}\}]^2 \\ &= \sum_{t=1}^n (X_t - \hat{\omega} - \hat{\beta} X_{t-1})^2 + \sum_{t=1}^n \{\omega - \hat{\omega} + (\beta - \hat{\beta}) X_{t-1}\}^2,\end{aligned}$$

## A posteriori distribution of $(\omega, \beta)$

$$\begin{aligned}\pi^1((\omega, \beta) | \mathbf{X}, \sigma) &\propto \exp \left\{ \frac{-1}{2\sigma^2} \sum_{t=1}^n (X_t - \omega - \beta X_{t-1})^2 \right\} \\ &\quad \times \exp \left\{ \frac{-1}{2} (\omega - \omega^0, \beta - \beta^0) (\Sigma^0)^{-1} (\omega - \omega^0, \beta - \beta^0)' \right\} \\ &\propto \exp \left\{ \frac{-1}{2\sigma^2} \sum_{t=1}^n \{(\omega - \hat{\omega}) + (\beta - \hat{\beta}) X_{t-1}\}^2 \right\} \\ &\quad \times \exp \left\{ \frac{-1}{2} (\omega - \omega^0, \beta - \beta^0) (\Sigma^0)^{-1} (\omega - \omega^0, \beta - \beta^0)' \right\} \\ &\propto \exp \left\{ \frac{-1}{2} (\omega - \hat{\omega}, \beta - \hat{\beta}) \Sigma_n^{-1} (\omega - \hat{\omega}, \beta - \hat{\beta})' \right\} \\ &\quad \times \exp \left\{ \frac{-1}{2} (\omega - \omega^0, \beta - \beta^0) (\Sigma^0)^{-1} (\omega - \omega^0, \beta - \beta^0)' \right\} \\ &\propto \exp \left\{ \frac{-1}{2} (\omega - \omega^*, \beta - \beta^*) (\Sigma^*)^{-1} (\omega - \omega^*, \beta - \beta^*)' \right\}.\end{aligned}$$

## A property of quadratic forms

For the last equality, we used a useful result on quadratic forms :

*Let  $x, a, b$  some  $k \times 1$  vectors,  $A$  and  $B$  some symmetric  $k \times k$  matrices such that  $(A+B)^{-1}$  exists. Then*

$$\begin{aligned}(x-a)'A(x-a) + (x-b)'B(x-b) &= (x-c)'(A+B)(x-c) \\ &\quad + (a-b)'A(A+B)^{-1}B(a-b)\end{aligned}$$

*where  $c = (A+B)^{-1}(Aa+Bb)$ .*

*See Box and Tiao (1973) p. 418.*

## A posteriori distributions

A *posteriori* distribution of  $(\omega, \beta) : \mathcal{N}((\omega^*, \beta^*), \Sigma^*)$ .

A *posteriori* distribution of  $\sigma^2$  :

$$\pi^2(\sigma^2 | \mathbf{X}, \omega, \beta) \propto \frac{e^{\frac{-n\bar{Y}}{2\sigma^2}}}{\sigma^n} \frac{e^{-b/\sigma^2}}{\sigma^{2(a+1)}} = \frac{e^{\frac{-1}{2\sigma^2}(n\bar{Y}+2b)}}{\sigma^{2(a+1+\frac{n}{2})}}.$$

where  $\bar{Y} = \frac{1}{n} \sum_{t=1}^n Y_t$  and  $Y_t = (X_t - \omega - \beta X_{t-1})^2$ ,  $t = 1, \dots, n$ .

Thus  $\pi^2(\sigma^2 | \mathbf{X}, \omega, \beta)$  is the density of the law  $IG(a + \frac{n}{2}, b + \frac{n\bar{Y}}{2})$ .



## Gibbs sampler for the AR(1)

- ① Specify the hyperparameters  $a, b, \omega^0, \beta^0, \Sigma^0$ .
- ② Choose an initial value  $\sigma^{(0)}$ .
- ③ For  $i=1, 2, \dots$ 
  - ① Generate  $(\omega^{(i)}, \beta^{(i)})$  in the law  $\mathcal{N}((\omega^*, \beta^*), \Sigma^*)$  obtained for  $\Sigma_n = \Sigma_n(\sigma^{(i-1)})$ .
  - ② Generate  $\sigma^{(i)}$  in the law  $IG(a + \frac{n}{2}, b + \frac{n\bar{Y}^{(i)}}{2})$  where
$$\bar{Y}^{(i)} = \frac{1}{n} \sum_{t=1}^n (X_t - \omega^{(i)} - \beta^{(i)} X_{t-1})^2.$$
- ③ Take  $\theta^{(i)} = (\omega^{(i)}, \beta^{(i)}, \sigma^{(i)})'$ .

## Comments

- ① The standard laws used in the algorithm, can be simulated efficiently using numerous statistical softwares such as Gauss, Mathematica, Matlab, R.
- ② This choice of *prior* distributions leads to standard *posterior* laws. It is possible to choose other laws, and to use the MH for the simulation of non standard *posterior* laws.
- ③ The stationarity condition  $|\beta| < 1$  can be taken into account by using, for the *prior* of  $(\omega, \beta)$ , a troncated bivariate Gaussian law on  $\mathbb{R} \times ]-1, 1[$ . The same troncature follows on the *a posteriori* law. Alternatively, one can choose for  $\beta$  a *prior* density with support  $[-1, 1]$ , for instance the law obtained for  $\beta = 2X - 1$  when  $X \sim \mathcal{B}(\beta_1, \beta_2)$ . We get the *a priori* density

$$\pi(\beta) = 0.5 \frac{\Gamma(\beta_1 + \beta_2)}{\Gamma(\beta_1)\Gamma(\beta_2)} \{0.5(1 + \beta)\}^{\beta_1-1} \{0.5(1 - \beta)\}^{\beta_2-1} \mathbf{1}_{[-1, 1]}$$

where  $\beta_1, \beta_2 > 0.5$  are hyperparameters.

## STAR(1) model

For the **STAR model** with  $c = 0$ , we have for a given  $\gamma$ ,

$$Y_t = Z_{t-1}'\beta + \epsilon_t,$$

where  $\epsilon_t$  is iid  $\mathcal{N}(0, \sigma^2)$ , and

$$Z_t = \begin{pmatrix} 1 - G_t(\gamma) \\ Y_t \{1 - G_t(\gamma)\} \\ G_t(\gamma) \\ Y_t G_t(\gamma) \end{pmatrix}, \quad \beta = \begin{pmatrix} a_0 \\ a_1 \\ b_0 \\ b_1 \end{pmatrix}, \quad G_t(\gamma) = \frac{1}{1 + \exp(-\gamma Y_t)}.$$

## STAR(1) model

Calculations similar to those done for the AR(1) show that the conditional laws of  $\beta$  and  $\sigma^2$  are explicit :

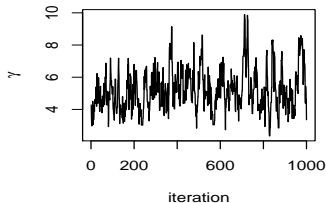
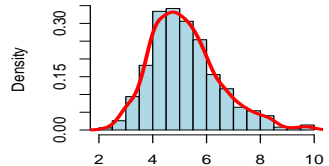
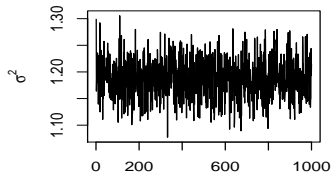
$$\begin{aligned}P_1(\beta | Y, \gamma, \sigma^2) &\sim \mathcal{N}\left\{(I_4 + \Sigma_n^{-1})^{-1} \Sigma_n^{-1} \hat{\beta}, (I_4 + \Sigma_n^{-1})^{-1}\right\}, \\P_2(\sigma^2 | Y, \beta, \sigma^2) &\sim IG\left(1 + (n-1)/2, 1 + \sum_{t=2}^n \epsilon_t^2/2\right),\end{aligned}$$

The conditional law of  $\gamma$  is not explicit but satisfies

$$\begin{aligned}P_3(\gamma | Y, \beta, \sigma^2) &\propto \exp\left\{-\sum_{t=2}^n \epsilon_t^2(\gamma)/(2\sigma^2) - \gamma\right\} \mathbf{1}_{\{\gamma > 0\}}, \\ \epsilon_t(\gamma) &= Y_t - \beta' Z_{t-1}.\end{aligned}$$

By using the MH algorithm to simulate the conditional law  $P_3$ , one gets an **hybrid method** combining the MH and Gibbs algorithms.

STAR(1) model : Sample paths and *a posteriori* laws of  $\gamma$  and  $\sigma^2$ , obtained for 1000 iterations of the hybrid algorithm

Trace of  $\gamma$ Posterior distribution of  $\gamma$ Trace of  $\sigma^2$ Posterior distribution of  $\sigma^2$ 