

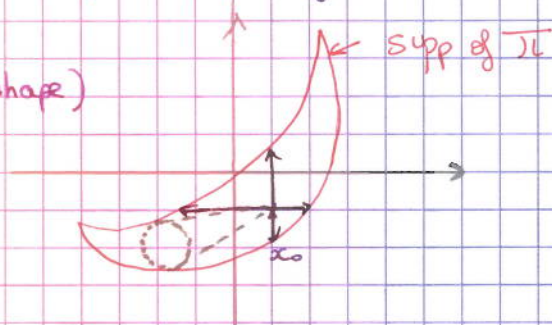
## The Gibbs Sampler:

The idea here is to create a sampler which is not a rejection type one and most particularly which will substitute the simulation of our r.v.  $x \in \mathbb{R}^d$  to sampling in smaller dimension.

The intuition of this is easy to understand in  $\mathbb{R}^2$ . Consider the following

target:  $\pi$

(the banana shape)



→ small angle in  $\mathbb{R}^2$  which will not make the sample be rejected with a MH sampler for eg.

→ larger possibilities in 1D and iterate.

### 1) Algorithm:

Assumptions: •  $X$  can be decomposed into  $X_1 \dots X_k$  so that  $\forall x \in X, \exists (x^1, \dots, x^k)$  in  $X_1 \times \dots \times X_k$  such that  $x = (x^1, \dots, x^k)$ . When  $k=d$ , each  $x^i \in \mathbb{R}$ .

We denote  $x^i$  the element of the  $i^{\text{th}}$  bloc and  $x^{(-i)}$  a vector  $x$  where we have removed the  $i^{\text{th}}$  bloc.

• We know how to simulate from the conditional distributions of the bloc  $i$  given the other when targeting  $\pi$ . That is to say that  $x \sim \pi$  is not tractable directly BUT:  $\forall i, x^i \sim \pi_i(x^i | x^{(-i)}) = \frac{\pi(x)}{\pi_{(-i)}(x)}$

is manageable directly, where  $\pi_{(-i)}(x)$  is the marginal distribution

We will see examples in the following where this clearly appears.

Pseudo code:  $x_0$  given;  $n=0$   
then iterate:

for  $j=1:k$

$x_{n+1}^j \sim \pi_j^i(x^j | x_{n+1}^1, \dots, x_{n+1}^{j-1}, x_n^{j+1}, \dots, x_n^k)$

end;  $n = n+1$ .

end

### 2) Transition kernel and invariant measure:

We will focus on measure which are absolutely continuous w.r.t. the Lebesgue measure and will still denote  $\pi$  the associated p.d.f.



the associated kernel is the composition of  $k$  kernels, each for one bloc. This writes <sup>(2)</sup>

$$P(x, y) = \prod_{i=1}^k \pi_i(y_i | y_1, \dots, y_{i-1}, x_{i+1}, \dots, x_k)$$

Now of course we expect that  $\pi$  will be the invariant measure. And this holds.

Prop:  $\pi P = \pi$

Proof: For sake of simplicity and clarity, we only consider the case where  $k=2$ , although the proof generalises to higher  $k$ .

Let  $B \subset \mathcal{X}$

$$\begin{aligned} \int_{\mathcal{X} \times B} \pi(x) P(x, y) dx dy &= \int_{\mathcal{X} \times B} \pi(x_1, x_2) \pi_1(y_1 | x_2) \pi_2(y_2 | y_1) dx dy \\ &= \int_{\mathcal{X} \times B} \pi(x_1, x_2) \frac{\pi(y_1, x_2)}{\int \pi(z, x_2) dz} \frac{\pi(y_2, y_1)}{\int \pi(y_1, z) dz} dx dy \\ &= \int_{\mathcal{X} \times B} \frac{\pi(x_1, x_2)}{\int \pi(z, x_2) dz} \frac{\pi(y_1, x_2)}{\int \pi(y_1, z) dz} \pi(y_1, y_2) dx dy \\ &= \int_{\mathcal{X} \times B} \pi_1(x_1 | x_2) \pi_2(x_2 | y_1) \pi(y) dx dy \\ &= \int_B \pi(y) \left[ \int_{\mathcal{X}} \pi_1(x_1 | x_2) \pi_2(x_2 | y_1) dx \right] dy \\ &= \int_B \pi(y) \left[ \int_{\mathcal{X}_1} \underbrace{\left\{ \int_{\mathcal{X}_2} \pi_1(x_1 | x_2) dx_1 \right\}}_{=1} \pi(x_2 | y_1) dx_2 \right] dy \\ &= \int_B \pi(y) dy = \pi(B) \end{aligned}$$

⚠ This kernel is not reversible! By construction, it is easy to see that the order the blocs are updated is important!

### 3) Examples:

#### a) Hierarchical models:

Let  $(Y_j)_{1 \leq j \leq n}$  iid  $\sim \mathcal{P}(\mu, \sigma^2)$ . Let  $\tau = 1/\sigma^2$ .  $(Y_j)_j$  are  $n$  observations and the goal is to estimate  $\mu$  and  $\tau$  given these observations.



In order to do this we will introduce a prior on the parameters  $(\mu, \tau)$  so that we work in the Bayesian framework.

For example:  $p(\mu, \tau) = \tau^{-1/2} \mathbb{1}_{\{\tau > 0\}}$

this prior is non informative on  $\mu$  [it is a degenerated uniform distribution on  $\mathbb{R}$ ] and forces  $\tau$  to be positive and behave like  $\frac{1}{\tau}$  in probability

Rq: these priors are not probability measures as they do not sum to 1  
this seems a problem at first sight, however, as long as the posterior distribution  $q(\tau, \mu | y)$  is well defined you are allowed to introduce degenerated priors.

Our target distribution is here  $q(\tau, \mu | (y_i^n)) = \pi(\mu, \tau)$ . This writes:

$$\pi(\mu, \tau) \propto \tau^{n/2} \exp\left(-\frac{1}{2} \tau \sum_{i=1}^n (y_i - \mu)^2\right) \tau^{-1/2} \mathbb{1}_{\{\tau > 0\}}$$

this probability measure is quite hard to sample from!

But looking at  $\pi(\mu | \tau)$  we have:

$$\begin{aligned} \pi(\mu | \tau) &\propto \exp\left(-\frac{1}{2} \tau n \mu^2 + \mu \tau n \left(\frac{1}{n} \sum_{j=1}^n y_j\right)\right) \quad \text{where we will denote} \\ &\quad \frac{1}{n} \sum_{j=1}^n y_j = \bar{y} \\ &\propto \exp\left(-\frac{1}{2 \cdot \frac{1}{n\tau}} (\mu - \bar{y})^2\right) \end{aligned}$$

We get that  $\pi(\mu | \tau) \sim \mathcal{D}\left(\bar{y}; \frac{1}{n\tau}\right)$  easy to sample from.

In the same way:  $\pi(\tau | \mu) \propto \tau^{\frac{n+1}{2}} \exp\left(-\tau \times \frac{1}{2} \sum_{j=1}^n (y_j - \mu)^2\right)$

Denoting  $B = \frac{1}{2} \sum_{j=1}^n (y_j - \mu)^2$  we recognise that  $\pi(\tau | \mu) \sim \Gamma\left(\frac{n+1}{2}, B\right)$

This produces a very simple Gibbs sampler: Given  $(\mu_j, \tau_j)$

- $\mu_{j+1} \sim \mathcal{D}\left(\bar{y}; (n\tau_j)^{-1}\right)$
- $\tau_{j+1} = \tau_{j+1} \sim \Gamma\left(\frac{n+1}{2}, B_{j+1}\right)$  where  $B_{j+1}$  is  $B$  where  $\mu$  has been substituted by  $\mu_{j+1}$ .

→ Rq sampling from a  $\Gamma$  distribution is easy in the case where  $\Gamma(a, b)$

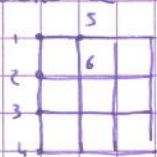
$a \in \mathbb{N}^*$  as it can be done by summing iid  $\text{Exp}(b)$

Otherwise use a rejection method with  $\Gamma([a], b)$  as proposal.



## b) the Ising model for image processing:

Let  $V$  be a set of connected vertices. the connections between vertices are called edges. this forms a graph. think as an example of the pixels in a grid of an image:



The Ising model is a distribution which assigns to each vertex a r.v. in  $\{0, 1\}$  or  $\{-1, 1\}$  with additional constraints: the value at a vertex is influenced by its immediate neighbours.

To introduce things formally, we will need some definitions:

def: let  $S$  be the set of vertices (of finite number)  
 $\mathcal{X}_s$  is the state space at vertex  $s \in S$  (typically for images  $\mathcal{X}_s = \{0, 255\}$ )  
 we assume the random variables take only a finite number of values.  
 i.e.  $|\mathcal{X}_s| < \infty$ .

let also  $\mathcal{X} = \prod_{s \in S} \mathcal{X}_s$  the space of all possible configurations

We denote  $\mathcal{X}_A = \prod_{s \in A} \mathcal{X}_s$  and  $X_A : x \mapsto x_s, s \in A$  the canonical projection

def: We call random field any measure  $\pi$  on  $\mathcal{X}$  s.t.  $\forall x \in \mathcal{X} \pi(x) > 0$

This means that all configurations in  $\mathcal{X}$  are possible with a non zero probab.

def: let  $A \subset S$ ,  $x_A \in \mathcal{X}_A$ ,  $x_{S \setminus A} \in \mathcal{X}_{S \setminus A}$  and  $\pi$  be a random field. We call local characteristic the conditional probability  $\pi(X_A = x_A | X_{S \setminus A} = x_{S \setminus A})$

Re: As  $\forall x \in \mathcal{X} \pi(x) > 0$  this is always well-defined.

def: A neighbour system  $\mathcal{V} = \{V_s, s \in S\}$  is a collection of sub sets of  $S$  s.t.

(i)  $\forall s \in S, s \notin V_s$

(ii)  $s \in V_t \iff t \in V_s$

def: A clique is a subset  $C$  of  $S$  s.t. all its elements are neighbours from each other

def: A Markov random field  $\pi$  with respect to  $\mathcal{V}$  is a random field such that

$$\forall s \in S: \pi(x_s = x_s | x_t = x_t, \forall t \neq s) = \pi(x_s = x_s | x_t = x_t, \forall t \in V_s)$$



Proof:  $\pi = \pi(x_A = x_A \mid x_{S \setminus A} = x_{S \setminus A}) = \frac{\pi(x_A = x_A \mid x_{S \setminus A})}{\pi(x_{S \setminus A} = x_{S \setminus A})}$  (Bayes' rule) (6)

$\leftarrow$  marginal on  $x_A$

$$= \frac{\exp \left[ - \sum_{c \in E} U_c(x_A, x_{S \setminus A}) \right]}{\sum_{y_A \in \mathcal{X}_A} \exp \left[ - \sum_{c \in E} U_c(y_A, x_{S \setminus A}) \right]}$$

let  $\mathcal{C} = \mathcal{C}_1 \sqcup \mathcal{C}_2$  (disjoint union) where:

$$\mathcal{C}_1 = \{c \in \mathcal{C} \text{ s.t. } c \cap A = \emptyset\} \text{ and } \mathcal{C}_2 = \{c \in \mathcal{C} \text{ s.t. } c \cap A \neq \emptyset\}$$

let  $R = S \setminus (A \cup \mathcal{V}(A))$  where  $\mathcal{V}(A) = \{\text{neighbours of elements in } A\} \setminus A$

We introduce a neutral point  $o$  (a landmark)

If  $c \in \mathcal{C}_2$ :  $U_c(x_A, x_{\mathcal{V}(A)}, x_R) = U_c(o_A, x_{\mathcal{V}(A)}, x_R)$  because as  $c \cap A = \emptyset$  the value of  $U_c$  will only depend on the projection of  $z$  on  $C$  therefore we don't care about what happens on  $A$ .

If  $c \in \mathcal{C}_1$ :  $U_c(x_A, x_{\mathcal{V}(A)}, x_R) = U_c(x_A, x_{\mathcal{V}(A)}, o_R)$ . Again what is important is the projection of the cliques which contains elements of  $A$ . This includes  $A$  and  $\mathcal{V}(A)$ . Then,

$$\pi = \frac{\exp \left[ - \sum_{c \in \mathcal{C}_1} U_c(x_A, x_{\mathcal{V}(A)}, x_R) \right] \exp \left[ - \sum_{c \in \mathcal{C}_2} U_c(x_A, x_{\mathcal{V}(A)}, x_R) \right]}{\left( \sum_{y_A \in \mathcal{X}_A} \exp \left[ - \sum_{c \in \mathcal{C}_1} U_c(y_A, x_{\mathcal{V}(A)}, x_R) \right] \right) \exp \left[ - \sum_{c \in \mathcal{C}_2} U_c(y_A, x_{\mathcal{V}(A)}, x_R) \right]}$$

$o_A$  on  $\mathcal{C}_2$

$o_A$  on  $\mathcal{C}_2$

$$= \frac{\exp \left[ - \sum_{c \in \mathcal{C}_1} U_c(x_A, x_{\mathcal{V}(A)}, x_R) \right]}{\sum_{y_A \in \mathcal{X}_A} \left[ \exp \left( - \sum_{c \in \mathcal{C}_1} U_c(y_A, x_{\mathcal{V}(A)}, x_R) \right) \right]}$$

which proves (\*)

Moreover,

$$\pi = \frac{\exp \left( - \sum_{c \in \mathcal{C}_1} U_c(x_A, x_{\mathcal{V}(A)}, x_R) \right) \times \left[ \sum_{y_R \in \mathcal{X}_R} \exp \left( - \sum_{c \in \mathcal{C}_2} U_c(o_A, x_{\mathcal{V}(A)}, y_R) \right) \right]}{\left[ \sum_{y_A \in \mathcal{X}_A} \exp \left( - \sum_{c \in \mathcal{C}_1} U_c(y_A, x_{\mathcal{V}(A)}, x_R) \right) \right] \times \left[ \sum_{y_R \in \mathcal{X}_R} \exp \left( - \sum_{c \in \mathcal{C}_2} U_c(o_A, x_{\mathcal{V}(A)}, y_R) \right) \right]}$$

$o_R = y_R$  on  $\mathcal{C}_1$

$x_A$  on  $\mathcal{C}_2$

$y_A$  on  $\mathcal{C}_2$

$$= \frac{\sum_{y_R \in \mathcal{X}_R} \exp \left[ - \sum_{c \in \mathcal{C}_1} U_c(x_A, x_{\mathcal{V}(A)}, y_R) - \sum_{c \in \mathcal{C}_2} U_c(x_A, x_{\mathcal{V}(A)}, y_R) \right]}{\sum_{y_A, y_R \in \mathcal{X}_A, \mathcal{X}_R} \exp \left[ - \sum_{c \in \mathcal{C}_1} U_c(y_A, x_{\mathcal{V}(A)}, y_R) - \sum_{c \in \mathcal{C}_2} U_c(y_A, x_{\mathcal{V}(A)}, y_A) \right]}$$

$= \pi(x_A = x_A \mid x_{\mathcal{V}(A)} = x_{\mathcal{V}(A)})$  and for  $A = \{s\}$  we get the Markov property.



def: A Gibbs random field induced by an energy function  $H$  is defined as

$$\pi(x) = \frac{\exp(-H(x))}{\sum_{y \in \mathcal{X}} \exp(-H(y))}$$

def: A potential is a family of functions  $U = \{U_A, A \in \mathcal{S}\}$  on  $\mathcal{X}$  s.t.

- (i)  $U_\emptyset = 0$
- (ii)  $U_A(x) = U_A(y)$  if  $x_A = y_A$

We can define an energy related to a potential  $U$  as  $H_U = \sum_{A \in \mathcal{S}} U_A$

(We sometimes denote  $U = H_U$  to make things simpler)

A potential  $U$  is a neighbour potential w.r.t.  $\mathcal{P}$  if  $U_A \equiv 0$  for all  $A \in \mathcal{S}$  which are not cliques for  $\mathcal{P}$ .

ex: Ising:  $\mathcal{X} = \{-1, 1\}$   $\mathcal{S} = \mathbb{Z}^2$  or a sub grid of  $\mathbb{Z}^2$ . Let  $c = (s, t)$   
 $\mathcal{P} = 4\text{-connectivity}$

$$U_c(x_s, x_t) = \begin{cases} -\beta & \text{if } x_s = x_t \\ +\beta & \text{if } x_s \neq x_t \end{cases} = \beta x_s x_t$$

And  $U_{c \rightarrow s}(x_s) = -\beta x_s$ . This defines an energy as:  $\forall x \in \mathcal{X}$

$$U(x) = - \sum_{c=(s,t)} \beta x_s x_t = \sum_{s \in \mathcal{S}} \beta x_s$$

We will prove the Hammersley-Clifford theorem which says that if we have a Gibbs random field w.r.t. a potential coming from a neighbour system  $\mathcal{P}$ , then it is a Markov random field and we can easily calculate the local characteristics.

Theorem: Let  $\pi$  a random field given by a neighbour potential  $U$  w.r.t. a neighbour system  $\mathcal{P}$  i.e:

$$\pi(x) = \frac{\exp(-\sum_{c \in \mathcal{C}} U_c(x))}{\sum_{y \in \mathcal{X}} \exp(-\sum_{c \in \mathcal{C}} U_c(y))} \quad \text{where } \mathcal{C} \text{ is the set of cliques for } \mathcal{P}.$$

Then the local characteristics are:

$$\pi(x_s = x_s, \forall s \in A \mid x_s = x_s, s \in \mathcal{S} \setminus A) = \frac{\exp[-\sum_{c \in \mathcal{C}, c \cap A \neq \emptyset} U_c(x)]}{\sum_{y \in \mathcal{X}} \exp[-\sum_{c \in \mathcal{C}, c \cap A \neq \emptyset} U_c(y_A, x_{\mathcal{S} \setminus A})]} \quad (*)$$

where  $x_y = x - x_A, x_{y \setminus A}$

$$\text{Moreover } \pi(x_s = x_s, s \in A \mid x_s = x_s, s \in \mathcal{S} \setminus A) = \pi(x_s = x_s, s \in A \mid x_s = x_s, s \in \mathcal{P}(A))$$



the Ising model is very often used for many applications as it is a simple Markov random field model. It is a very interesting choice as prior law on images as they are given a natural neighbour system thanks to pixels. This is particularly appealing for segmentation and it is the choice that may be used in our 1<sup>st</sup> example of estimation of  $x^*$  as our "cost" on images. It favours the homogeneous areas as we tried to.

The Ising model which has only 2 classes can be extended to more and is then called the Potts model.

As the local characteristics are very easy to compute, it is a very interesting choice for a Gibbs sampler.

#### 4) Around the traditional Gibbs sampler:

The previous algorithm is sometimes called the deterministic updating GS (DUGS). And we have seen that it is not reversible. Some other versions of the sampler have been proposed to face this problem.

##### a) The reversible Gibbs sampler

Pseudo code: Given  $x_n \in \mathbb{R}^k$

$$2k-1 \text{ simulations } \left\{ \begin{array}{l} \tilde{x}^1 \sim \pi(\cdot | x_n^2, \dots, x_n^k) \\ \vdots \\ \tilde{x}^k \sim \pi(\cdot | \tilde{x}^1, \dots, \tilde{x}^{k-1}) \rightarrow \tilde{x}_{n+1}^k = x_n^k \\ x_{n+1}^{k+1} \sim \pi(\cdot | \tilde{x}^1, \dots, \tilde{x}^{k-1}, x_{n+1}^k) \\ \vdots \\ x_{n+1}^1 \sim \pi(\cdot | x_{n+1}^2, \dots, x_{n+1}^k) \end{array} \right.$$

As expected it is reversible:

**Proof:** Again for sake of simplicity, we prove it for  $k=2$  but it easily generalises to any  $k \geq 2$

Let  $(A, B) \in \mathbb{R}^2$ :



$$\begin{aligned}
 & \int_{A \times B} \pi(y_1, y_2) \left[ \int_{\mathcal{X}} \pi(w | y_2) \pi(y_2' | w) dw \right] \pi(y_1' | y_2') dy dy' \\
 &= \int_{A \times B} \pi(y_1 | y_2) \pi(y_2) \left[ \int_{\mathcal{X}} \frac{\pi(w, y_2)}{\pi(y_2)} \pi(y_2' | w) dw \right] \frac{\pi(y_1', y_2')}{\pi(y_2')} dy dy' \\
 &= \int_{A \times B} \pi(y_1', y_2') \left[ \int_{\mathcal{X}} \pi(w, y_2) \frac{\pi(y_2' | w)}{\pi(y_2')} dw \right] \pi(y_1 | y_2) dy dy' \\
 &= \int_{A \times B} \pi(y_1', y_2') \left[ \int_{\mathcal{X}} \pi(w, y_2) \frac{\pi(w | y_2')}{\pi(w)} dw \right] \pi(y_1 | y_2) dy dy' \\
 &= \int_{A \times B} \pi(y_1', y_2') \left[ \int_{\mathcal{X}} \pi(y_2 | w) \pi(w | y_2') dw \right] \pi(y_1 | y_2) dy dy'
 \end{aligned}$$

### b) Random Scan Gibbs Sampler:

the simulations of each component is done in a random order. But not one at a time. We need to scan all coordinate each time otherwise it fails to be reversible.

Pseudo code: 1) Sample a permutation  $\sigma \in \mathcal{S}_k$

2) Sample  $y_{\sigma_1}^{n+1} \sim \pi_{\sigma_1}(y_{\sigma_1} | y^{(-\sigma_1)^n})$

$\vdots$

$y_{\sigma_k}^{n+1} \sim \pi_{\sigma_k}(y_{\sigma_k} | y^{(-\sigma_k)^{n+1}})$

Proof of reversibility:  $\Sigma$  over all permutations in  $\mathcal{S}_k$

Rg: It is better than the previous one as one does not sample coordinates twice. All samples are used!

### 5) The Gibbs Sampler as an MCMC sampler:

We said that the GS is not a rejection based sampler. However it can be seen that it is a particular case of the MCMC algorithm.

Let consider for  $1 \leq i \leq k$  the following transition kernel: It only updates the  $i^{th}$  coordinate and  $P_i(x, y) = \pi_i(y^i | x^{(-i)}) \mathbb{1}_{y^i = x^{(-i)}}$

If we use this proposal in a MCMC it appears that  $\alpha(x, y) = 1$  and it is therefore always accepted.



By iterating over the coordinates we get that the GS is an iteration of HN with proposal evolving with  $1 \leq i \leq k$

### 8) Metropolis - within - Gibbs sampler:

The GS assumes we are able to sample from  $\pi_i(x_i | x^{(-i)})$ . There are unfortunately many models where this is not as trivial as the previous examples. To face this, we can use a HN algorithm at each step of the GS to sample each coordinate. Let  $q(x_i; x^{(-i)}, y_i; x^{(-i)})$  be a proposal; a candidate  $y_i$  is proposed. One computes the acceptance ratio  $\alpha_i$  and accept or not with probability  $\alpha_i$ . Then, we change the coordinate.

The transition kernel is a little more complicated:

$$P_j(z, dz) = \left( \prod_{m \neq j} \delta_{x_m}(dz^m) \times \left[ q_j(dz^j | z^{(-j)}) \alpha_j(x^j, dz^j) + \delta_{x^j}(dz^j) \right] (1 - \alpha_j(x^j, b)) q_j(x^j, b) \right)$$

### 9) Which algo for high dimensional random variable?

Gibbs Sampler designed for this But loops over  $k$

Itala is a good alternative

tradeoff between  $k$  and the calculation on  $\nabla \log \pi$

↓

$k$  & to compute

↓

1 x only.