

# Bayesian statistics, simulation and software

## Module 6: The Gibbs sampler

Jesper Møller and Ege Rubak

Department of Mathematical Sciences  
Aalborg University

# The Gibbs sampler — the general algorithm

**Aim:** We want to sample  $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_k)$  from a density  $\pi(\boldsymbol{\theta})$ , e.g. the prior or the posterior density (in the latter case, suppressing in the notation the dependence of the data  $x$ :  $\pi(\boldsymbol{\theta}) = \pi(\boldsymbol{\theta}|x)$ ).

Assume  $\theta_i \in \Omega_i \subseteq \mathbf{R}^{d_i}$  and  $\boldsymbol{\theta} \in \Omega_1 \times \Omega_2 \times \dots \times \Omega_k \subseteq \mathbf{R}^{d_1+d_2+\dots+d_k}$

We can then generate an *approximate* sample from  $\pi(\boldsymbol{\theta})$  (provided some technical conditions are satisfied) as follows:

## Gibbs Sampler

- Choose initial value  $\boldsymbol{\theta}^{(0)} = (\theta_1^{(0)}, \theta_2^{(0)}, \dots, \theta_k^{(0)})$ .
- For  $i = 1, 2, \dots, t$ 
  1. Generate  $\theta_1^{(i)} \sim \pi(\theta_1 | \theta_2^{(i-1)}, \theta_3^{(i-1)}, \dots, \theta_k^{(i-1)})$
  2. Generate  $\theta_2^{(i)} \sim \pi(\theta_2 | \theta_1^{(i)}, \theta_3^{(i-1)}, \dots, \theta_k^{(i-1)})$
  - ⋮
  - k. Generate  $\theta_k^{(i)} \sim \pi(\theta_k | \theta_1^{(i)}, \theta_2^{(i)}, \dots, \theta_{k-1}^{(i)})$

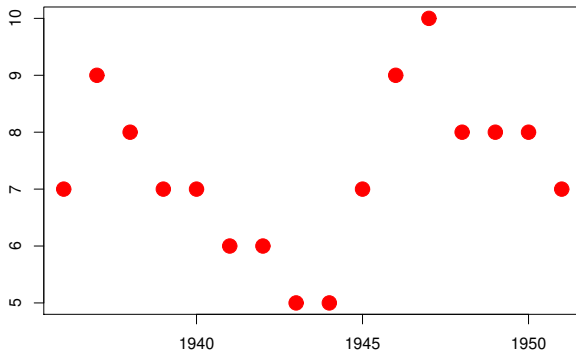
The higher  $i$  is the closer  $\boldsymbol{\theta}^{(i)} = (\theta_1^{(i)}, \theta_2^{(i)}, \dots, \theta_k^{(i)})$  is to being a sample from  $\pi(\boldsymbol{\theta})$ .

When  $d_1, \dots, d_k$  are small, Gibbs sampling may be easy to use.

## Example: Marriage rates in Italy

For the years 1936 to 1951 (16 years) the marriage rates per 1000 of the population in Italy have been observed. Is it practical to model marriage rates that occurred during WW2 to rates just before and after?

**Data:**  $\mathbf{y} = (y_1, y_2, \dots, y_{16})$ .



# Italian marriages: Model

**Model:** Conditional on (true) rates  $\lambda_1, \lambda_2, \dots, \lambda_{16}$  the observed rates  $y_1, y_2, \dots, y_{16}$  are independent and  $y_i \sim Pois(\lambda_i)$ :

- Joint density of data  $\mathbf{y}$ :

$$\pi(\mathbf{y}|\boldsymbol{\lambda}) = \prod_{i=1}^{16} \pi(y_i|\lambda_i) = \prod_{i=1}^{16} \frac{e^{-\lambda_i} \lambda_i^{y_i}}{y_i!}.$$

## Italian marriages: Prior and hyper prior

**Prior:** Conditional on a *hyper parameter*  $\beta > 0$  the rates  $\lambda_1, \lambda_2, \dots, \lambda_{16}$  are i.i.d. with  $\lambda_i | \beta \sim \text{Exp}(\beta)$ :

- The prior density of  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_{16})$  conditional on  $\beta$  is

$$\pi(\boldsymbol{\lambda} | \beta) = \prod_{i=1}^{16} \pi(\lambda_i | \beta) = \prod_{i=1}^{16} \beta \exp(-\beta \lambda_i).$$

As we are not sure which value the common parameter  $\beta$  should take, we assume a so-called *hyper prior* on  $\beta$ :

- $\beta \sim \text{Exp}(1)$ , i.e.  $\pi(\beta) = e^{-\beta}$  for  $\beta > 0$ .

Thus the prior density for  $(\boldsymbol{\lambda}, \beta)$  is

$$\pi(\boldsymbol{\lambda}, \beta) = \pi(\beta)\pi(\boldsymbol{\lambda} | \beta) = e^{-\beta} \prod_{i=1}^{16} \beta \exp(-\beta \lambda_i).$$

**Posterior** density:

$$\begin{aligned}\pi(\boldsymbol{\lambda}, \beta | \mathbf{y}) &\propto \pi(\mathbf{y} | \boldsymbol{\lambda}, \beta) \pi(\boldsymbol{\lambda}, \beta) \\ &= \prod_{i=1}^{16} \pi(y_i | \lambda_i) \prod_{i=1}^{16} \pi(\lambda_i | \beta) \pi(\beta) \\ &= \left( \prod_{i=1}^{16} \frac{e^{-\lambda_i} \lambda_i^{y_i}}{y_i!} \right) \left( \prod_{i=1}^{16} \beta e^{-\beta \lambda_i} \right) e^{-\beta}, \quad \lambda_1, \dots, \lambda_{16}, \beta > 0.\end{aligned}$$

This looks complicated. Therefore to explore the posterior we make use of a Gibbs sampler with low dimensional distributions – these are called full conditionals and specified as follows.

## Full conditionals — $\lambda_i$

- Let  $\boldsymbol{\lambda}_{-i} = (\lambda_1, \dots, \lambda_{i-1}, \lambda_{i+1}, \dots, \lambda_{16})$ ,  $i = 1, \dots, 16$ .
- The full conditional for  $\lambda_i$  has density

$$\begin{aligned}\pi(\lambda_i | \boldsymbol{\lambda}_{-i}, \mathbf{y}, \beta) &= \frac{\pi(\lambda_i, \boldsymbol{\lambda}_{-i}, \mathbf{y}, \beta)}{\pi(\boldsymbol{\lambda}_{-i}, \mathbf{y}, \beta)} \\ &\propto \left( \prod_{j=1}^{16} \pi(y_j | \lambda_j) \right) \left( \prod_{j=1}^{16} \pi(\lambda_j | \beta) \right) \pi(\beta) \\ &\propto \pi(y_i | \lambda_i) \pi(\lambda_i | \beta) \\ &= \frac{e^{-\lambda_i} \lambda_i^{y_i}}{y_i!} \cdot \beta e^{-\beta \lambda_i} \\ &\propto e^{-\lambda_i(1+\beta)} \lambda_i^{y_i+1-1} \\ &\sim \text{Gamma}(y_i + 1, (1 + \beta)^{-1}),\end{aligned}$$

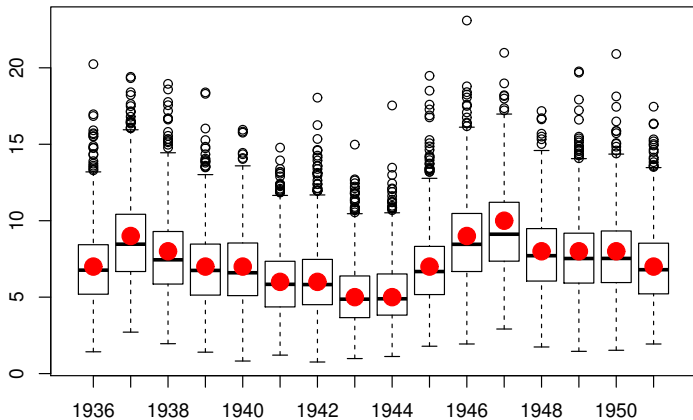
- The full conditional for  $\beta$  has density

$$\begin{aligned}\pi(\beta|\boldsymbol{\lambda}, \mathbf{y}) &\propto \left( \prod_{i=1}^{16} \pi(y_i|\lambda_i) \right) \left( \prod_{i=1}^{16} \pi(\lambda_i|\beta) \right) \pi(\beta) \\ &\propto \left( \prod_{i=1}^{16} \pi(\lambda_i|\beta) \right) \pi(\beta) \\ &= \left( \prod_{i=1}^{16} \beta e^{-\beta\lambda_i} \right) e^{-\beta} \\ &\propto \beta^{16+1-1} e^{-\beta(1+\sum_{i=1}^{16} \lambda_i)} \\ &\sim \text{Gamma}\left(17, \left(1 + \sum_{i=1}^n \lambda_i\right)^{-1}\right).\end{aligned}$$



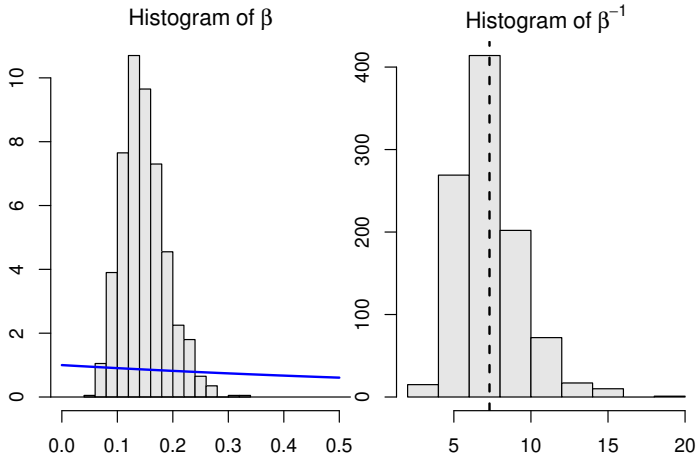
# Posterior marriage rates: Boxplots

Although there is a clear trend of a drop during WW2 it is not extreme:



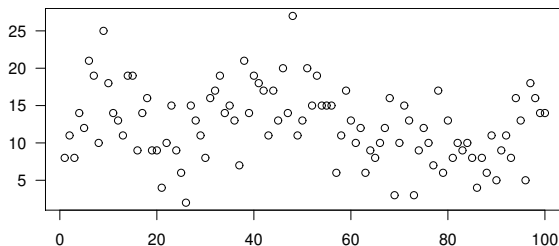
# Posterior distribution of $\beta$

Note that  $\beta^{-1}$  is the prior mean of a marriage rate.



## Example: Airport mishandling of luggage

Every hour the number of mishandled bags have been recorded:



### Notation:

- Let  $y_t \in \mathbb{N}_0$  denote the number of mishandled bags at time (hour)  $t$ .
- The airport is in (so to say) one of two states: **Normal** or **broken**.  
Let  $x_t \in \{1, 2\}$  denote the state of the airport at time  $t$  (1=normal, 2=broken).

### Objective:

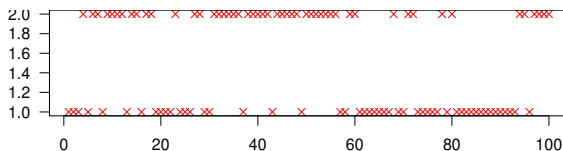
- Estimate the state of the airport at each hour.

## Mishandling: Data model

- Conditional on  $\mathbf{x} = (x_1, \dots, x_{100})$  the number of mishandlings are independent, and the conditional distribution of  $y_t | \mathbf{x}$  depends only on  $x_t$ .
- The number of mishandlings is assumed to follow a Poisson distribution:
  - ▶  $y_t | x_t = 1 \sim \text{Pois}(10)$  Normal state
  - ▶  $y_t | x_t = 2 \sim \text{Pois}(15)$  Broken state

Maximum likelihood estimate (most likely state according to data model):  $x_t = 1$  is most likely

$$\Leftrightarrow \frac{e^{-10} 10^{y_t}}{y_t!} > \frac{e^{-15} 15^{y_t}}{y_t!} \Leftrightarrow y_t > \frac{5}{\ln 15 - \ln 10}.$$

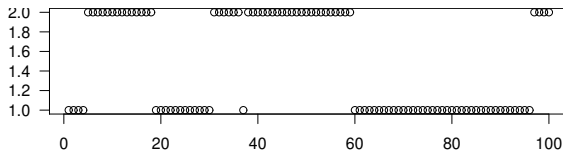


# Mishandling: Prior

It is known that the airport tends to “stick” in the same state. Thus the prior for  $\mathbf{x}$  is assumed to be a Markov chain:

- $P(x_1 = 1) = P(x_1 = 2) = \frac{1}{2}$  (probabilities for initial state)
- $P(x_{t+1} = x_t | x_t) = 0.9$  (probability of staying)
- $P(x_{t+1} \neq x_t | x_t) = 0.1$  (probability of switching)

Example of a realisation from the prior:



# Mishandling: Posterior

The posterior density is

$$\begin{aligned}\pi(\mathbf{x}|\mathbf{y}) &\propto \pi(\mathbf{y}|\mathbf{x})\pi(\mathbf{x}) \\ &= \prod_{t=1}^{100} \pi(y_t|x_t)\pi(x_1) \prod_{t=1}^{99} \pi(x_{t+1}|x_t)\end{aligned}$$

Thus we obtain a full conditional for each  $x_t$ :

$$\pi(x_t|y_t, \mathbf{x}_{-t}) \propto \pi(y_t|x_t)\pi(x_{t+1}|x_t)\pi(x_t|x_{t-1}).$$

for  $1 < t < 99$  with obvious modifications for  $t = 1$  and  $t = 100$ .

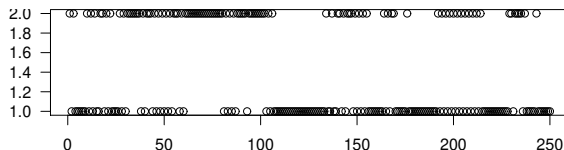
So  $x_t|y_t, \mathbf{x}_{-t}$  is a 1-2 random variable with probabilities

$$\pi(x_t = i|y_t, \mathbf{x}_{-t}) = \frac{\pi(y_t|x_t = i)\pi(x_{t+1}|x_t = i)\pi(x_t = i|x_{t-1})}{\sum_{j=1}^2 \pi(y_t|x_t = j)\pi(x_{t+1}|x_t = j)\pi(x_t = j|x_{t-1})}$$

for  $i = 1, 2$ . It is of course easy to simulate from this distribution.

# Posterior results

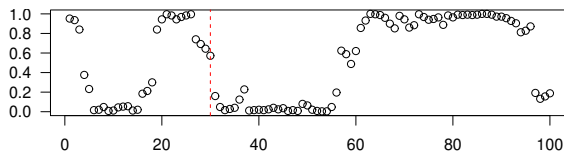
*Example:* Plot of  $x_{30}$  during  $I = 250$  “sweeps” of the Gibbs sampler



Estimate of the posterior probability that  $x_{30} = 1$ :

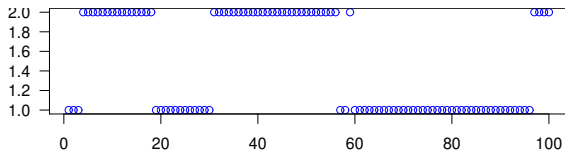
$$P(x_{30} = 1|\mathbf{y}) \approx \frac{1}{I} \sum_{i=1}^I 1[x_{30,i} = 1] = 57.2\%.$$

**For all hours:** Plot of posterior probabilities  $P(x_t = 1|\mathbf{y})$ ,  $t = 1, \dots, 100$ .



# Comparison

Most likely state according to the posterior distribution



Compare this to the MLE (the most likely state using only the data model):

