

Estadística Bayesiana

# **Métodos de Simulación Estocástica**

Lizbeth Naranjo Albarrán

Facultad de Ciencias, UNAM

March 25, 2022

# Índice

<b>1</b>	<b>Métodos de Simulación Estocástica en Estadística Bayesiana</b>	<b>1</b>
1.1	Integración Monte Carlo . . . . .	1
1.1.1	Muestreo por Importancia . . . . .	4
1.1.2	Método de Aceptación y Rechazo . . . . .	6
1.2	Monte Carlo vía Cadenas de Markov . . . . .	7
1.2.1	Algoritmo Metropolis-Hastings . . . . .	8
1.2.2	Muestreo de Gibbs . . . . .	11
1.2.3	Convergencia . . . . .	14
1.3	Software . . . . .	15
1.3.1	WinBUGS . . . . .	15
1.3.2	JAGS . . . . .	18
1.3.3	STAN . . . . .	18
1.4	INLA . . . . .	18

# Capítulo 1

## Métodos de Simulación Estocástica en Estadística Bayesiana

En este capítulo estudiaremos algunos métodos de simulación estocástica, necesarios para la estimación Bayesiana, usaremos principalmente métodos de Monte Carlo vía cadenas de Markov, y daremos una breve descripción de WinBUGS, JAGS y STAN, y finalmente trataremos de usar métodos INLA.

Hay varios libros que pueden ayudar en estos temas, algunos de los que se usarán son Chen et al. (2000), Gamerman and Lopes (2006), Gilks et al. (1996) y Robert and Casella (2000), entre otros.

### 1.1 Integración Monte Carlo

Como se mencionó anteriormente, en el enfoque Bayesiano de la Estadística, la incertidumbre presente en un modelo dado,  $p(x|\theta)$ , se representa a través de una distribución de probabilidad  $p(\theta)$  sobre el espacio parametral  $\Theta$  (generalmente multidimensional) que define al modelo. El teorema de Bayes,

$$p(\theta|x) = \frac{p(\theta)p(x|\theta)}{p(x)},$$

permite entonces incorporar la información contenida en un conjunto de datos  $x = (x_1, \dots, x_n)$ , produciendo una descripción conjunta de la incertidumbre sobre los valores de los parámetros del modelo a través de la distribución final  $p(\theta|x)$ . Desafortunadamente, la implementación de los métodos Bayesianos ocasionalmente requieren de un esfuerzo computacional muy alto. La mayor parte de este esfuerzo se concentra en el cálculo de ciertas características de la distribución final del parámetro de interés. Además, en la práctica es común que la dimensión de  $\theta$  sea muy grande. Por otro lado, excepto en aplicaciones muy sencillas tanto  $p(x|\theta)$  como

$p(\theta)$  pueden llegar a tener formas muy complicadas. En la mayoría de los problemas las integrales requeridas no pueden resolverse analíticamente, por lo que es necesario contar con métodos numéricos eficientes que permitan calcular o aproximar integrales en varias dimensiones. Algunos de estos métodos de integración son las técnicas de Monte Carlo vía cadenas de Markov. Para mayores detalles consultar Gilks, Richardson y Spiegelhalter (1996).

La integración de Monte Carlo evalúa  $E[g(\theta)]$  obteniendo muestras denotadas por  $\{\theta^{(t)}, t = 1, \dots, n\}$  de la distribución  $p(\theta|x)$  y entonces aproximando

$$E[g(\theta)] \approx \frac{1}{n} \sum_{t=1}^n g(\theta^{(t)}).$$

De esta manera, la media poblacional de  $g(\theta)$  se estima por medio de una media muestral. Cuando las muestras  $\{\theta^{(t)}\}$  son independientes, la ley de los grandes números asegura que la aproximación puede hacerse tan precisa como se desee incrementando el tamaño de la muestra,  $n$ . Note que aquí  $n$  está bajo el control del analista: no es el tamaño fijo para una muestra de datos.

En general, seleccionar muestras  $\{\theta^{(t)}\}$  independientes de  $p(\theta|x)$  no es factible, ya que  $p(\theta|x)$  puede ser no estándar. Sin embargo, las  $\{\theta^{(t)}\}$  no necesariamente necesitan ser independientes. Las  $\{\theta^{(t)}\}$  pueden ser generadas por cualquier proceso que selecciona muestras por todo el soporte de  $p(\theta|x)$  en las proporciones correctas. Una forma de hacer esto es a través de una cadena de Markov teniendo a  $p(\theta|x)$  como su distribución estacionaria. Esto es entonces un método de Monte Carlo vía cadenas de Markov.

En resumen, las técnicas de Monte Carlo vía cadenas de Markov permiten generar, de manera iterativa, observaciones de distribuciones multivariadas que difícilmente podrían simularse utilizando métodos directos. La idea básica es muy simple: construir una cadena de Markov que sea fácil de simular y cuya distribución de equilibrio corresponda a la distribución final que nos interesa.

**Proposición 1.1.** Sea  $\theta^{(1)}, \theta^{(2)}, \dots$  una cadena de Markov homogénea, irreducible y aperiódica, con espacio de estados  $\Theta$  y distribución de equilibrio  $p(\theta|x)$ . Entonces, conforme  $t \rightarrow \infty$ ,

$$\begin{aligned} (i) \quad & \theta^{(t)} \xrightarrow{\mathcal{D}} \theta, \quad \text{donde } \theta \sim p(\theta|x); \\ (ii) \quad & \frac{1}{t} \sum_{i=1}^t g(\theta^{(i)}) \rightarrow E(g(\theta)|x). \end{aligned}$$

Ejemplos: Bayes7\_1IntegracionMonteCarlo.pdf

Código R: Bayes7\_1IntegracionMonteCarlo.R

**Ejemplo 1.1** (Normal). Dada una muestra de tamaño  $n$  de una distribución  $Normal(0, 1)$ ,  $x_1, \dots, x_n$ , y considere la función de distribución acumulada  $\Phi(t)$  igual a:

$$\Phi(t) = \int_{-\infty}^t \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx$$

Por medio del método de Monte Carlo se puede obtener una aproximación como:

$$\hat{\Phi}(t) \approx \frac{1}{n} \sum_{i=1}^n \mathbb{I}[x_i \leq t]$$

con varianza exacta

$$\frac{\Phi(t)[1 - \Phi(t)]}{n}$$

¿Cómo se obtiene?

La integración Monte Carlo permite calcular la esperanza poblacional por medio de una media muestral:

$$E[g(X)] \approx \frac{1}{n} \sum_{t=1}^n g(X^{(t)}).$$

Entonces, sea  $X \sim Normal(0, 1)$ , y defina la función  $g$  como  $g(X) = \mathbb{I}[X \leq t]$ , entonces, la media poblacional de  $g(X)$  es:

$$\begin{aligned} E[g(X)] &= \int g(x)f(x)dx = \int \mathbb{I}[x \leq t]f(x)dx \\ &= \int \mathbb{I}[x \leq t] \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx = \int_{-\infty}^t \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx = \Phi(t) \end{aligned}$$

Entonces se obtiene una muestra de tamaño  $n$  de una distribución  $Normal(0, 1)$ ,  $x_1, \dots, x_n$ , y se calcula  $g(x_1) = \mathbb{I}[x_1 \leq t], \dots, g(x_n) = \mathbb{I}[x_n \leq t]$ , y se calcula la media muestral:

$$\frac{1}{n} \sum_{t=1}^n g(X^{(t)}) = \frac{1}{n} \sum_{t=1}^n \mathbb{I}[x_i \leq t]$$

Entonces, usando la integración Monte Carlo,

$$\Phi(t) \approx \frac{1}{n} \sum_{t=1}^n \mathbb{I}[x_i \leq t]$$

Además cada una de estas funciones  $g(x) = \mathbb{I}[x \leq t]$  valen 0 o 1, es decir, que son v.a. Bernoulli con probabilidad de éxito  $\Phi(t)$ , porque

$$P(\mathbb{I}[X \leq t] = 1) = P(X \leq t) = \int_{-\infty}^t \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx = \Phi(t)$$

■

### 1.1.1 Muestreo por Importancia

El muestreo por Importancia se basa en una representación alternativa de

$$E_f[h(x)] = \int_{\mathcal{X}} h(x)f_X(x)dx$$

Dada una función de densidad  $g$ , estrictamente positiva cuando  $h(x)f_X(x) \neq 0$ , se reescribe la esperanza,

$$\begin{aligned} E_f[h(x)] &= \int_{\mathcal{X}} h(x)f_X(x)\frac{g(x)}{g(x)}dx \\ &= \int_{\mathcal{X}} h(x)\frac{f_X(x)}{g(x)}g(x)dx \\ &= E_g\left[h(x)\frac{f_X(x)}{g(x)}\right] \end{aligned}$$

donde  $\mathcal{X}$  es el soporte de  $X$ , es decir, el conjunto donde  $X$  toma sus valores bajo  $f_X$  y  $h$ , puede ser más pequeño que el soporte de la densidad  $g$ , es decir,

$$\text{sup}(hf_X) \subseteq \text{sup}(g)$$

La restricción del soporte se debe cumplir, porque de lo contrario se podría truncar el soporte, y producir resultados insesgados.

El estimador de  $E_f[h(x)]$  usando la integración Monte Carlo:

$$E_f[h(x)] \approx \frac{1}{n} \sum_{i=1}^n h(x_i) \frac{f_X(x_i)}{g(x_i)}$$

donde  $x_1, \dots, x_n$  es una muestra generada a partir de  $g$ .

**Ejemplo 1.2.** Calcular  $E[h(X)]$ , donde  $X \sim \text{Uniforme}(0, 10)$  y  $h(X) = 10 \exp(-2|x - 5|)$ . Usando: (a) integración Monte Carlo y (b) muestreo por importancia. ■

**Ejemplo 1.3.** Estimar los momentos  $E[X^k]$  donde  $X$  se distribuye Doble Exponencial o Laplace:

$$\begin{aligned} f(x) &= \frac{1}{2} \exp(-|x|) \\ F(x) &= \exp(x)/2 I[x \leq 0] + (1 - \exp(-x)/2) I[x > 0] \end{aligned}$$

Opción 1: simulando de Doble Exponencial y usando Integración Monte Carlo.

Opción 2: simulando de  $N(0, 4)$  y usando muestreo por Importancia. ■

**Ejemplo 1.4.** Sea  $Z \sim Normal(0, 1)$ , calcular  $P[Z > 4.5] = 3.398 \times 10^{-6}$

Esta probabilidad es demasiado pequeña y el método de Monte Carlo estándar se rompe/descompone cuando se tienen que calcular probabilidades de las colas. Se requiere un número grande de simulaciones.

Usando muestreo por Importancia, se considera el soporte restringido a  $(4.5, \infty)$ . Tomar las funciones  $g$  como la función de densidad *Exponencial*(1), que es la selección 'natural', pero que esté truncada en 4.5,

$$g(X) = \frac{e^{-x}}{P[x \geq 4.5]} = \frac{e^{-x}}{\int_{4.5}^{\infty} e^{-x} dx} = \frac{e^{-x}}{e^{-4.5}} = e^{-(x-4.5)}$$

El estimador por el muestreo por importancia es

$$\begin{aligned} P[X > 4.5] &= E[\mathbb{I}[X > 4.5]] = \int_{4.5}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2} dx \\ &= \int_{4.5}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2} \frac{1}{e^{-(x-4.5)}} e^{-(x-4.5)} dx \\ &\approx \frac{1}{n} \sum_{i=1}^n \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x_i^2} \frac{1}{e^{-(x_i-4.5)}} \end{aligned}$$

donde  $x_1, \dots, x_n$  son v.a.i.i.d. *Exponencial*(1). ■

Código R: Bayes7\_1IntegracionMonteCarlo.R

### 1.1.2 Método de Aceptación y Rechazo

El método de aceptación y rechazo es un método indirecto para generar variables aleatorias:  $f_X(\cdot)$  función objetivo (*target*) y  $g_Y(\cdot)$  función de densidad candidato.

Restricciones:

- $f_X(\cdot)$  y  $g_Y(\cdot)$  tienen soportes compatibles,  $g_Y(x) > 0$  cuando  $f_X(x) > 0$ .
- Existe una constante  $M$  con  $\frac{f_X(x)}{g_Y(x)} \leq M \forall x$ .

Algoritmo de Aceptación y Rechazo para simular  $X$ :

- (1) Generar  $Y \sim g_Y(y)$ .
- (2) Generar  $U \sim U(0, 1)$ , independiente de  $Y$ .
- (3) Si  $u \leq \frac{1}{M} \frac{f_X(y)}{g_Y(y)}$  entonces  $X = y$ . En caso contrario, descartar  $y$ , y regresando a (1).

¿Cómo funciona el algoritmo de Aceptación-Rechazo?

La función de densidad acumulada de la v.a. aceptada es:

$$\begin{aligned} P \left[ Y \leq x | u \leq \frac{1}{M} \frac{f_X(y)}{g_Y(y)} \right] &= \frac{P \left[ Y \leq x, U \leq \frac{1}{M} \frac{f_X(y)}{g_Y(y)} \right]}{P \left[ U \leq \frac{1}{M} \frac{f_X(y)}{g_Y(y)} \right]} = \frac{\int_{-\infty}^x \int_0^{\frac{1}{M} \frac{f_X(y)}{g_Y(y)}} du g_Y(y) dy}{\int_{-\infty}^{\infty} \int_0^{\frac{1}{M} \frac{f_X(y)}{g_Y(y)}} du g_Y(y) dy} \\ &= \frac{\int_{-\infty}^x \frac{1}{M} \frac{f_X(y)}{g_Y(y)} g_Y(y) dy}{\int_{-\infty}^{\infty} \frac{1}{M} \frac{f_X(y)}{g_Y(y)} g_Y(y) dy} = \frac{\int_{-\infty}^x f_X(y) dy}{\int_{-\infty}^{\infty} f_X(y) dy} \\ &= \int_{-\infty}^x f_X(y) dy = P[X \leq x] \end{aligned}$$

**Ejemplo 1.5.** Sea  $X \sim \text{Beta}(a, b)$  el target, usamos  $U(0, 1)$  como densidad candidata, y  $M = \max\{f(x)\}$ . ■

Código R: Bayes7.2DistribucionBeta.R



## 1.2 Monte Carlo vía Cadenas de Markov

En estadística Bayesiana,  $P(X|\boldsymbol{\theta})$  y  $P(\boldsymbol{\theta})$  pueden llegar a tener formas muy complicadas. En la mayoría de los problemas las integrales requeridas para calcular  $P(\mathbf{X})$ ,  $P(X^*)$  y  $P(X^*|\mathbf{X})$  no pueden resolverse de manera analítica, por lo que es necesario usar métodos numéricos eficientes que permitan calcular o aproximar estas integrales. Algunos de estos métodos son las técnicas de Monte Carlo vía cadenas de Markov (*Markov chain Monte Carlo*, MCMC) (Gilks et al., 1996; Chen et al., 2000; Gamerman and Lopes, 2006).

La intergración de Monte Carlo evalúa la media de  $g(\boldsymbol{\theta})$ ,  $E[g(\boldsymbol{\theta})]$ , obteniendo muestras de la distribución  $P(\boldsymbol{\theta}|\mathbf{X})$ , denotadas por  $\{\boldsymbol{\theta}^{(t)}; t = 1, \dots, m\}$ , y aproximando

$$E[g(\boldsymbol{\theta})] \approx \frac{1}{m} \sum_{t=1}^m g(\boldsymbol{\theta}^{(t)}),$$

es decir, la media poblacional de  $g(\boldsymbol{\theta})$  se estima por medio de la media muestral.

Las muestras  $\{\boldsymbol{\theta}^{(t)}; t = 1, \dots, m\}$  pueden ser generadas por cualquier proceso que selecciona muestras por todo el soporte de  $P(\boldsymbol{\theta}|\mathbf{X})$  en las proporciones correctas. En particular, las muestras se pueden generar a través de una cadena de Markov teniendo a  $P(\boldsymbol{\theta}|\mathbf{X})$  como su distribución estacionaria, y a este proceso se le conoce como método de Monte Carlo vía cadenas de Markov. Los dos algoritmos más comunes en los métodos MCMC son muestreo de Gibbs y Metropolis-Hastings.

Código R: Bayes7.2DistribucionBeta.R

Código R: Bayes7.3DistribucionNormal.R

### 1.2.1 Algoritmo Metropolis-Hastings

El algoritmo de Metropolis-Hastings construye una cadena de Markov definiendo las probabilidades de transición de la siguiente manera.

Sea  $Q(\boldsymbol{\theta}^*|\boldsymbol{\theta})$  una densidad de transición (arbitraria) y se define

$$\alpha(\boldsymbol{\theta}^*, \boldsymbol{\theta}) = \min \left\{ \frac{P(\boldsymbol{\theta}^*|\mathbf{X})Q(\boldsymbol{\theta}|\boldsymbol{\theta}^*)}{P(\boldsymbol{\theta}|\mathbf{X})Q(\boldsymbol{\theta}^*|\boldsymbol{\theta})}, 1 \right\}.$$

Dado un valor inicial  $\boldsymbol{\theta}^{(0)}$ , la  $t$ -ésima iteración consta de:

- 1) Generar una observación  $\boldsymbol{\theta}^*$  de  $Q(\boldsymbol{\theta}^*|\boldsymbol{\theta}^{(t)})$ .
- 2) Generar una variable  $u \sim U(0, 1)$ .
- 3) Si  $\alpha(\boldsymbol{\theta}^*, \boldsymbol{\theta}^{(t)}) \geq u$ , hacer  $\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^*$ , en caso contrario  $\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)}$ .

De esta manera se genera una cadena de Markov con distribución de transición

$$\begin{aligned} P(\boldsymbol{\theta}^{(t+1)}|\boldsymbol{\theta}^{(t)}) &= \alpha(\boldsymbol{\theta}^{(t)}, \boldsymbol{\theta}^{(t+1)})Q(\boldsymbol{\theta}^{(t+1)}|\boldsymbol{\theta}^{(t)})\mathbf{I}[\boldsymbol{\theta}^{(t+1)} \neq \boldsymbol{\theta}^{(t)}] \\ &+ \left[ 1 - \int \alpha(\boldsymbol{\theta}^{(t)}, \boldsymbol{\theta}^*)Q(\boldsymbol{\theta}^*|\boldsymbol{\theta}^{(t)})d\boldsymbol{\theta}^* \right] \mathbf{I}[\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)}]. \end{aligned}$$

La probabilidad de aceptación  $\alpha(\boldsymbol{\theta}^{(t)}, \boldsymbol{\theta}^*)$  sólo depende de  $P(\boldsymbol{\theta}|\mathbf{X})$  a través de un cociente, por lo que la constante de normalización no es necesaria.

Este algoritmo construye una cadena de Markov definiendo las probabilidades de transición de la siguiente manera.

Sea  $Q(\theta^*|\theta)$  una densidad de transición (arbitraria) y definamos

$$\alpha(\theta^*, \theta) = \min \left\{ \frac{p(\theta^*|x)Q(\theta|\theta^*)}{p(\theta|x)Q(\theta^*|\theta)}, 1 \right\}.$$

*Algoritmo:*

Dado un valor inicial  $\theta^{(0)}$ , la  $t$ -ésima iteración consiste en:

1. generar una observación  $\theta^*$  de  $Q(\theta^*|\theta^{(t)})$ ;
2. generar una variable  $u \sim U(0, 1)$ ;

3. si  $u \leq \alpha(\theta^*, \theta^{(t)})$ , hacer  $\theta^{(t+1)} = \theta^*$ ; en caso contrario, hacer  $\theta^{(t+1)} = \theta^{(t)}$ .

Este procedimiento genera una cadena de Markov con distribución de transición

$$p(\theta^{(t+1)}|\theta^{(t)}) = \alpha(\theta^{(t)}, \theta^{(t+1)})Q(\theta^{(t+1)}|\theta^{(t)})I_{\{\theta^{(t+1)} \neq \theta^{(t)}\}} + \left[1 - \int \alpha(\theta^{(t)}, \theta^*)Q(\theta^*|\theta^{(t)})d\theta^*\right] I_{\{\theta^{(t+1)} = \theta^{(t)}\}}.$$

La probabilidad de aceptación  $\alpha(\theta^*, \theta)$  sólo depende de  $p(\theta|x)$  a través de un cociente, por lo que la constante de normalización no es necesaria.

*Comentario.* La versión original del algoritmo de Metropolis-Hastings toma  $Q(\theta^*|\theta) = Q(\theta|\theta^*)$ , en cuyo caso

$$\alpha(\theta^*, \theta) = \min \left\{ \frac{p(\theta^*|x)}{p(\theta|x)}, 1 \right\}.$$

Dos casos particulares en la práctica son:

- *Caminata aleatoria.* Sea  $Q(\theta^*|\theta) = Q_1(\theta^* - \theta)$ , donde  $Q_1(\cdot)$  es una densidad de probabilidad simétrica centrada en el origen. Entonces

$$\alpha(\theta^*, \theta) = \min \left\{ \frac{p(\theta^*|x)}{p(\theta|x)}, 1 \right\}.$$

- *Independencia.* Sea  $Q(\theta^*|\theta) = Q_0(\theta^*)$ , donde  $Q_0(\cdot)$  es una densidad de probabilidad sobre  $\Theta$ . Entonces

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

con  $\omega(\theta) = p(\theta|x)/Q_0(\theta)$ .

En la práctica es común utilizar, después de una reparametrización apropiada, distribuciones de transición normales ó  $t$  de Student ligeramente sobredispersas, por ejemplo

$$Q(\theta^*|\theta) = N_d(\theta^*|\theta, \kappa V(\hat{\theta})) \quad (\text{caminata aleatoria})$$

ó

$$Q_0(\theta^*) = N_d(\theta^*|\hat{\theta}, \kappa V(\hat{\theta})) \quad (\text{independencia}),$$

donde  $\hat{\theta}$  y  $V(\hat{\theta})$  denotan a la media y a la matriz de varianzas-covarianzas de la aproximación normal asintótica para  $p(\theta|x)$ , respectivamente, y  $\kappa \geq 1$  es una factor de sobredispersión.

**Ejemplo 1.6** (Coeficiente de correlación). Suponga que  $D = \{\mathbf{y}_i = (y_{1i}, y_{2i})^t, i = 1, \dots, n\}$  es una muestra aleatoria de una distribución normal bivariada  $N_2(0, \Sigma)$ , donde

$$\Sigma = \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}.$$

Suponga una distribución inicial uniforme  $U(-1, 1)$  para  $\rho$ ; la densidad final para  $\rho$  está dada por

$$p(\rho|D) \propto (1 - \rho^2)^{-n/2} \exp \left\{ -\frac{1}{2(1 - \rho^2)} (S_{11} - 2\rho S_{12} + S_{22}) \right\},$$

donde  $-1 < \rho < 1$ , y  $S_{rs} = \sum_{i=1}^n y_{ri}y_{si}$  para  $r, s = 1, 2$ . Generar valores de  $\rho$  de la densidad  $p(\rho|D)$  no es algo trivial, ya que no es una función log-cóncava. Por tanto, consideramos el algoritmo de Metropolis-Hastings para generar  $\rho$ . Como  $-1 < \rho < 1$ , consideremos la transformación

$$\rho = \frac{-1 + e^\xi}{1 + e^\xi}, \quad -\infty < \xi < \infty.$$

Entonces

$$p(\xi|D) = p(\rho(\xi)|D) \frac{2e^\xi}{(1 + e^\xi)^2}.$$

En vez de obtener un muestreo directo de  $\rho$ , generamos  $\xi$  eligiendo una distribución proporcional a  $N(\hat{\xi}, \hat{\sigma}_\xi^2)$ , donde  $\hat{\xi}$  es el máximo del logaritmo de  $p(\xi|D)$ , que puede obtenerse mediante el algoritmo de Newton-Raphson o el algoritmo de Nelder-Mead, y  $\hat{\sigma}_\xi^2$  es menos el inverso de la segunda derivada del  $\log p(\xi|D)$  evaluada en  $\xi = \hat{\xi}$ , es decir,

$$\hat{\sigma}_\xi^{-2} = - \left. \frac{d^2 \log p(\xi|D)}{d\xi^2} \right|_{\xi=\hat{\xi}}.$$

El algoritmo para generar  $\xi$  opera de la manera siguiente:

Dado un valor inicial  $\xi^{(0)}$ ,

1. generar  $\xi^*$  de  $N(\hat{\xi}, \hat{\sigma}_\xi^2)$ ;
2. generar una variable  $u \sim U(0, 1)$ ;
3. si  $u \leq \alpha(\xi^*, \xi^{(t)})$ , entonces  $\xi^{(t+1)} = \xi^*$ ; en caso contrario,  $\xi^{(t+1)} = \xi^{(t)}$  donde

$$\alpha(\xi^*, \xi^{(t)}) = \min \left\{ \frac{p(\xi^*|D) \phi \left( \frac{\xi^{(t)} - \hat{\xi}}{\hat{\sigma}_\xi} \right)}{p(\xi^{(t)}|D) \phi \left( \frac{\xi^* - \hat{\xi}}{\hat{\sigma}_\xi} \right)}, 1 \right\}$$

y  $\phi$  es la función de densidad de probabilidad normal estándar.

Después de obtener  $\xi$ , podemos obtener  $\rho$  aplicando la transformación correspondiente. ■

### 1.2.2 Muestreo de Gibbs

El muestreo de Gibbs permite simular una cadena de Markov  $\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)}, \dots, \boldsymbol{\theta}^{(m)}$  con distribución de equilibrio  $P(\boldsymbol{\theta}|\mathbf{X})$ . Cada nuevo valor de la cadena se puede obtener a través de un proceso iterativo que sólo requiere generar muestras de distribuciones cuya dimensión es menor que  $d$  (la dimensión de  $\boldsymbol{\theta}$ ) y que en la mayoría de los casos tienen una forma más sencilla que la de  $P(\boldsymbol{\theta}|\mathbf{X})$ .

Sea  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_k)$  una partición del vector  $\boldsymbol{\theta}$ , donde  $\theta_i \in \mathbb{R}^{d_i}$  y  $\sum_{i=1}^k d_i = d$ . Se elige un valor inicial  $\boldsymbol{\theta}^{(0)} = (\theta_1^{(0)}, \dots, \theta_k^{(0)})$ , se simula una cadena de Markov en la que  $\boldsymbol{\theta}^{(t+1)}$  se obtiene a partir de  $\boldsymbol{\theta}^{(t)}$  de las distribuciones condicionales completas de la siguiente manera:

- 1) Generar  $\theta_1^{(t+1)}$  de  $P(\theta_1|\theta_2^{(t)}, \theta_3^{(t)}, \dots, \theta_k^{(t)}, \mathbf{X})$ .
- 2) Generar  $\theta_2^{(t+1)}$  de  $P(\theta_2|\theta_1^{(t+1)}, \theta_3^{(t)}, \dots, \theta_k^{(t)}, \mathbf{X})$ .
- $\vdots$
- $k$ ) Generar  $\theta_k^{(t+1)}$  de  $P(\theta_k|\theta_1^{(t+1)}, \theta_2^{(t+1)}, \dots, \theta_{k-1}^{(t+1)}, \mathbf{X})$ .

De esta manera se genera una cadena de Markov con distribución de transición

$$P(\boldsymbol{\theta}^{(t+1)}|\boldsymbol{\theta}^{(t)}) = \prod_{i=1}^k P(\theta_i^{(t+1)}|\theta_1^{(t+1)}, \dots, \theta_{i-1}^{(t+1)}, \theta_{i+1}^{(t)}, \dots, \theta_k^{(t)}, \mathbf{X}).$$

En algunos casos la distribución final implica cierta estructura de independencia condicional entre algunos de los elementos de  $\boldsymbol{\theta}$ , y en ocasiones las distribuciones condicionales completas se simplifican.

Al igual que el algoritmo de Metropolis-Hastings, el algoritmo de Gibbs permite simular una cadena de Markov  $\theta^{(1)}, \theta^{(2)}, \dots$  con distribución de equilibrio  $p(\theta|x)$ . En este caso, sin embargo, cada nuevo valor de la cadena se puede obtener a través de un proceso iterativo que sólo requiere generar muestras de distribuciones cuya dimensión es menor que  $d$  y que en la mayoría de los casos tienen una forma más sencilla que la de  $p(\theta|x)$ .

Sea  $\theta = (\theta_1, \dots, \theta_k)$  una partición del vector  $\theta$ , donde  $\theta_i \in \mathbb{R}^{d_i}$  y  $\sum_{i=1}^k d_i = d$ . Las densidades

$$\begin{aligned} & p(\theta_1|\theta_2, \dots, \theta_k, x) \\ & \vdots \\ & p(\theta_i|\theta_1, \dots, \theta_{i-1}, \theta_{i+1}, \dots, \theta_k, x) \quad (i = 2, \dots, k-1) \\ & \vdots \\ & p(\theta_k|\theta_1, \dots, \theta_{k-1}, x) \end{aligned}$$

se conocen como *densidades condicionales completas* y en general pueden identificarse fácilmente al inspeccionar la forma de la distribución final  $p(\theta|x)$ . De hecho, para cada  $i = 1, \dots, k$ ,

$$p(\theta_i|\theta_1, \dots, \theta_{i-1}, \theta_{i+1}, \dots, \theta_k, x) \propto p(\theta|x),$$

donde  $p(\theta|x) = p(\theta_1, \dots, \theta_k|x)$  es vista sólo como función de  $\theta_i$ .

*Algoritmo:*

Dado un valor inicial  $\theta^{(0)} = (\theta_1^{(0)}, \dots, \theta_k^{(0)})$ , el algoritmo de Gibbs simula una cadena de Markov en la que  $\theta^{(t+1)}$  se obtiene a partir de  $\theta^{(t)}$  de la siguiente manera:

generar una observación  $\theta_1^{(t+1)}$  de  $p(\theta_1|\theta_2^{(t)}, \theta_3^{(t)}, \dots, \theta_k^{(t)}, x)$ ;

generar una observación  $\theta_2^{(t+1)}$  de  $p(\theta_2|\theta_1^{(t+1)}, \theta_3^{(t)}, \dots, \theta_k^{(t)}, x)$ ;

$\vdots$

generar una observación  $\theta_k^{(t+1)}$  de  $p(\theta_k|\theta_1^{(t+1)}, \theta_2^{(t+1)}, \dots, \theta_{k-1}^{(t+1)}, x)$ .

La sucesión  $\theta^{(1)}, \theta^{(2)}, \dots$  así obtenida es entonces una realización de una cadena de Markov cuya distribución de transición está dada por

$$p(\theta^{(t+1)}|\theta^{(t)}) = \prod_{i=1}^k p(\theta_i^{(t+1)}|\theta_1^{(t+1)}, \dots, \theta_{i-1}^{(t+1)}, \theta_{i+1}^{(t)}, \dots, \theta_k^{(t)}, x).$$

*Comentario.* En ocasiones la distribución final implica cierta estructura de independencia condicional entre algunos de los elementos del vector  $\theta$ . En estos casos es común que muchas de las densidades condicionales completas se simplifiquen.

**Ejemplo 1.7** (Gibbs: modelo jerárquico). Consideremos el modelo jerárquico definido por

$$\begin{aligned} \text{I.} \quad p(x|\omega) &= \prod_{i=1}^m p(x_i|\omega_i); \\ \text{II.} \quad p(\omega|\phi) &= \prod_{i=1}^m p(\omega_i|\phi); \\ \text{III.} \quad p_0(\phi). \end{aligned}$$

Esta estructura define un modelo para  $x$  parametrizado por  $\theta = (\omega, \phi) = (\omega_1, \dots, \omega_m, \phi)$  y con distribución inicial  $p(\theta) = p_0(\phi)p(\omega|\phi)$ , de manera que la distribución final está dada por

$$p(\theta|x) \propto p_0(\phi) \prod_{i=1}^m \{p(x_i|\omega_i)p(\omega_i|\phi)\}.$$

Entonces  $k = m + 1$  y las densidades condicionales completas toman la forma

$$\begin{aligned} p(\theta_1|\theta_2, \dots, \theta_{k-1}, \theta_k, x) &= p(\omega_1|\phi, x_1) \\ &\vdots \\ p(\theta_{k-1}|\theta_1, \dots, \theta_{k-2}, \theta_k, x) &= p(\omega_m|\phi, x_m) \\ p(\theta_k|\theta_1, \dots, \theta_{k-2}, \theta_{k-1}, x) &= p_0(\phi|\omega), \end{aligned}$$

donde  $p_0(\phi|\omega) \propto p_0(\phi)p(\omega|\phi)$ . ■

**Ejemplo 1.8** (Gibbs: modelo normal bivariado). El propósito de este ejemplo es examinar la estructura de la correlación exacta de una cadena de Markov inducida por el muestreo de Gibbs. Suponga que la distribución posterior  $p(\theta|D)$  es una distribución normal bivariada  $N_2(\mu, \Sigma)$  con

$$\mu = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} \quad y \quad \sigma = \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix},$$

donde  $\mu_j, \sigma_j, j = 1, 2$ , y  $\rho$  son conocidos. Entonces el muestreo de Gibbs requiere muestrear de

$$\theta_1 \sim N\left(\mu_1 + \rho\frac{\sigma_1}{\sigma_2}(\theta_2 - \mu_2), \sigma_1^2(1 - \rho^2)\right)$$

y

$$\theta_2 \sim N\left(\mu_2 + \rho\frac{\sigma_2}{\sigma_1}(\theta_1 - \mu_1), \sigma_2^2(1 - \rho^2)\right).$$

Sea  $\{\theta_i = (\theta_{1,i}, \theta_{2,i})^t, i \geq 0\}$ , que denota una cadena de Markov inducida por el muestreo de Gibbs para la distribución normal bivariada anterior. Si iniciamos a partir de la distribución estacionaria, es decir,  $\theta_0 \sim N_2(\mu, \Sigma)$ , entonces cada una de las sucesiones  $\{\theta_{1,i}, i \geq 0\}$  y  $\{\theta_{2,i}, i \geq 0\}$  es un proceso AR(1).

Veamos este resultado. Sea  $\{z_{1,i}, z_{2,i}, i \geq 0\}$  una sucesión de variables aleatorias i.i.d.  $N(0, 1)$ . Entonces la estructura del muestreo de Gibbs implica

$$\begin{aligned} \theta_{1,0} &= \mu_1 + \sigma_1 z_{1,0}, \\ \theta_{2,0} &= \mu_2 + \rho\frac{\sigma_2}{\sigma_1}(\theta_{1,0} - \mu_1) + \sigma_2\sqrt{1 - \rho^2}z_{2,0}, \end{aligned}$$

y

$$\begin{aligned} \theta_{1,i+1} &= \mu_1 + \rho\frac{\sigma_1}{\sigma_2}(\theta_{2,i} - \mu_2) + \sigma_1\sqrt{1 - \rho^2}z_{1,i+1}, \\ \theta_{2,i+1} &= \mu_2 + \rho\frac{\sigma_2}{\sigma_1}(\theta_{1,i+1} - \mu_1) + \sigma_2\sqrt{1 - \rho^2}z_{2,i+1}, \end{aligned} \tag{1.1}$$

para  $i \geq 0$ . Ahora consideremos el primer componente  $\theta_{1,i+1}$ . De la ecuación (1.1), para  $i \geq 0$ ,

$$\begin{aligned}\theta_{1,i+1} &= \mu_1 + \rho \frac{\sigma_1}{\sigma_2} \left[ \rho \frac{\sigma_2}{\sigma_1} (\theta_{1,i} - \mu_1) + \sigma_2 \sqrt{1 - \rho^2} z_{2,i} \right] + \sigma_1 \sqrt{1 - \rho^2} z_{1,i+1} \\ &= \mu_1 + \rho^2 (\theta_{1,i} - \mu_1) + \rho \sigma_1 \sqrt{1 - \rho^2} z_{2,i} + \sigma_1 \sqrt{1 - \rho^2} z_{1,i+1}.\end{aligned}\quad (1.2)$$

Sea  $\psi = \rho^2$  y  $\sigma_1^{*2} = \sigma_1^2(1 - \rho^4)$ . Sea  $\{z_i^*, i \geq 0\}$  que denota una sucesión de variables aleatorias i.i.d.  $N(0, 1)$ . Como  $z_{1,i}$  y  $z_{2,i+1}$  son independientes e idénticamente distribuidas  $N(0, 1)$ , entonces podemos reescribir (1.2) como

$$\theta_{1,0} = \mu_1 + \sigma_1 z_0^*, \quad (1.3)$$

$$\theta_{1,i+1} = \mu_1 + \psi(\theta_{1,i} - \mu_1) + \sigma_1^* z_{i+1}^* \quad \forall i \geq 0. \quad (1.4)$$

Así,  $\{\theta_{1,i}, i \geq 0\}$  es un proceso AR(1) con rezago 1 y autocorrelación  $\psi = \rho^2$ . Similarmente,  $\{\theta_{2,i}, i \geq 0\}$  es también un proceso AR(1) con rezago 1 y autocorrelación  $\psi = \rho^2$ . La única diferencia es que usamos  $\sigma_2^* = \sigma_2 \sqrt{1 - \rho^4}$  en vez de  $\sigma_1^*$  en (1.4), y usamos  $\mu_2$  y  $\sigma_2$  en vez de  $\mu_1$  y  $\sigma_1$  en (1.3). ■

### 1.2.3 Convergencia

Supongamos que se desea generar una muestra de tamaño  $N$  de la distribución  $p(\theta|x)$ . Si para cada uno de  $N$  valores iniciales  $\theta_1^{(0)}, \dots, \theta_N^{(0)}$  corremos alguno de los algoritmos discutidos en esta sección, entonces, de acuerdo con la proposición 1.1(i), después de un cierto número de iteraciones  $T$  suficientemente grande los valores  $\theta_1^{(T)}, \dots, \theta_N^{(T)}$  pueden considerarse como una muestra de tamaño  $N$  de la distribución final de  $\theta$ . Alternativamente podemos generar una sola cadena y tomar los valores  $\theta^{(T+K)}, \theta^{(T+2K)}, \dots, \theta^{(T+NK)}$  como una muestra de  $p(\theta|x)$ , donde  $K$  se elige de manera que la correlación entre las observaciones sea pequeña.

En general no es fácil determinar en qué momento la(s) cadena(s) ha(n) convergido. Un método empírico comúnmente utilizado, basado en la proposición 1.1(ii), consiste en graficar los promedios ergódicos de algunas funciones de  $\theta$  contra el número de iteraciones y elegir el valor  $T$  a partir del cual las gráficas se estabilizan. En este caso es frecuente omitir los primeros valores de la(s) cadena(s) al calcular los promedios ergódicos. La idea de este *periodo de calentamiento* es permitir que la(s) cadena(s) salga(n) de una primera fase de inestabilidad. En el caso particular del muestreo de Gibbs la velocidad de convergencia depende fuertemente de la correlación entre los componentes del vector  $\theta$  bajo la distribución final  $p(\theta|x)$ : entre más alta sea la correlación más lenta será la convergencia.



## 1.3 Software

Para el análisis de datos de estadística Bayesiana se requieren métodos numéricos, y muchos de ellos ya se encuentran en librerías de R o pueden programarse en R (Albert (2009)). A continuación presentaremos otros software utilizados en estadística Bayesiana.

### 1.3.1 WinBUGS

WinBUGS es un *software* de libre acceso en internet que se puede obtener de la página <http://www.mrc-bsu.cam.ac.uk/bugs/>. Fue diseñado por Spiegelhalter, Thomas y Best y es parte del proyecto BUGS (*Bayesian Inference Using Gibbs Sampling*) que desarrollaron estos investigadores para el análisis Bayesiano de modelos estadísticos a través de métodos de Monte Carlo vía cadenas de Markov.

WinBUGS permite que métodos complejos de simulación sean sencillos para los usuarios de la estadística Bayesiana aplicada en diversas disciplinas (Lunn et al. (2012), Lunn et al. (2000), Ntzoufras (2009)).

El *software* ofrece una interfaz con el usuario basada en cuadros de diálogo y comandos a través de los cuales se analiza el modelo, por lo que el ambiente de WinBUGS se vuelve más amigable. Además, también es posible realizar una interfaz con R<sup>1</sup> o S-Plus, los cuales manejan una sintaxis muy similar a la que usa WinBUGS, aunque WinBUGS no cuenta con tantos comandos como R o S-Plus.

La manera de operar de WinBUGS está basada en el muestreo de Gibbs (ver sección 1.2.2); es decir, dada una función de verosimilitud y una distribución inicial, el propósito es muestrear valores de los parámetros del modelo a partir de la distribución final. De estos valores muestreados, es posible obtener estimadores de los parámetros y hacer todo tipo de inferencias sobre éstos.

El desarrollo y mejoramiento de las técnicas de Monte Carlo basadas en cadenas de Markov ha hecho posible construir la distribución final de modelos Bayesianos muy complejos. Sin embargo, generar una muestra a través del algoritmo de Gibbs para obtener la distribución final de un modelo Bayesiano puede ser una tarea muy complicada, sobre todo para el caso en el que se introducen muchos parámetros en el modelo.

WinBUGS proporciona herramientas simples para llevar a cabo la simulación Monte Carlo; a través de WinBUGS es posible implementar el muestreo de Gibbs para una gran variedad de modelos.

---

<sup>1</sup>*Software* de libre acceso en la página de internet <http://www.r-project.org>.

El *software* tiene ciertas limitaciones para detectar la convergencia, resumir las muestras y realizar diagnósticos sobre el ajuste de los modelos. Sin embargo, es posible usar los paquetes *boa* o *coda* de R junto con WinBUGS para realizar un mejor análisis de las muestras simuladas.

WinBUGS intenta usar el método de muestreo más apropiado para cada parámetro del modelo estocástico de acuerdo con la tabla 1.1.

Distribución Objetivo	Método de muestreo
Discreta	Inversión de una función de distribución acumulada
Forma cerrada	Muestreo directo usando algoritmos estándar
Log-cóncava	Muestreo de rechazo adaptativo
Rango restringido	Muestreo <i>slice</i>
Rango sin restricción	Metropolis-Hastings (sección 1.2.1)

Table 1.1: Métodos de muestreo de WinBUGS.

## Ejemplo

El siguiente ejemplo se obtuvo del libro de Peter Congdon (2001), *Bayesian Statistical Modelling*. Hay preocupación por el incremento de la incidencia de cáncer de seno en las mujeres, sin embargo, en el Reino Unido el número de muertes por esta enfermedad es más o menos constante. Ambas tendencias implican una disminución de “casos fatales”, es decir, una disminución de la tasa de muerte entre los nuevos casos de la enfermedad.

Suponga que estamos interesados en el intervalo de mayor densidad de esta tasa, en una situación donde tenemos  $k = 11$  tipos de cáncer.

Sea  $\mathbf{n} = (n_1, n_2, \dots, n_k)$  el número de muertes en cada uno de los 11 tipos de cáncer (1 es cáncer de seno). El vector aleatorio  $\mathbf{n}$  tiene una distribución  $Multinomial(n_1, \dots, n_k | \pi, N)$ , donde  $N = \sum_{i=1}^k n_i$  es el número total de muertes por cáncer.

Para elegir la distribución inicial de  $\pi$  se podría tomar el patrón de muertes por cáncer de mujeres en el Reino Unido en años anteriores, o de algún otro país semejante a él.

En este ejemplo se considera una distribución inicial conjugada no informativa (ver sección ??) para  $\pi$ , la distribución  $Dirichlet(\pi | 1, 1, \dots, 1)$ .

El programa realizado en WinBUGS es el siguiente:

Modelo

```
model { # Modelo Multinomial
for(i in 1:k) {
```

```

c[i] <- 1
}
pi[1:k] ~ ddirch(c[]) # Inicial
n[1:k] ~ dmulti(pi[],N) # Verosimilitud
}

Datos
list(n = c(14080, 12990, 6440, 4350, 3420, 3190, 2600, 2420,
1820, 1760, 23610),
k = 11, # Tipos de Cáncer
N = 76680) # Total de Muertes por Cáncer

```

La distribución multinomial puede expresarse, equivalentemente, como la distribución conjunta de  $k$  variables aleatorias independientes con distribución Poisson, de tal manera que  $n_i \sim \text{Poisson}(n_i|\mu_i)$  con  $i = 1, \dots, k$ , sujetas a la condición  $\sum_{j=1}^k n_j = N$  (ver sección ??). Las probabilidades de la distribución multinomial se obtienen con:

$$\pi_j = \frac{\mu_j}{\sum_{i=1}^k \mu_i}.$$

Para este caso las distribuciones iniciales conjugadas para los parámetros  $\mu_i$  son gamma, y utilizaremos distribuciones iniciales que son no informativas,  $\text{Gamma}(\mu_i|1, 1)$ .

El programa utilizado en WinBUGS es el siguiente:

```

Modelo
model{ # Modelo Poisson
for(i in 1:k) {
mu[i] ~ dgamma(1,1) # Inicial
n[i] ~ dpois(mu[i]) # Verosimilitud
pi[i] <- mu[i]/mu.sum # Parámetro Multinomial
}
mu.sum <- sum(mu[ ])
}

```

```

Datos
list(n = c(14080, 12990, 6440, 4350, 3420, 3190, 2600, 2420,
1820, 1760, 23610),
k = 11) # Tipos de Cáncer

```

Las estadísticas básicas que se obtienen para la distribución final del vector de parámetros  $\pi$  son las mismas en ambos modelos y se presentan en la tabla 1.2.

	mean	sd	MC error	2.5%	median	97.5%
$\pi_1$	0.1836	0.00139	$1.542E-5$	0.1809	0.1836	0.1863
$\pi_2$	0.1694	0.00135	$1.285E-5$	0.1667	0.1694	0.1721
$\pi_3$	0.0840	$9.939E-4$	$8.783E-6$	0.0820	0.0839	0.0859
$\pi_4$	0.0567	$8.385E-4$	$8.441E-6$	0.0551	0.0567	0.0584
$\pi_5$	0.0446	$7.422E-4$	$6.667E-6$	0.0431	0.0445	0.0460
$\pi_6$	0.0416	$7.225E-4$	$7.654E-6$	0.0402	0.0416	0.0430
$\pi_7$	0.0339	$6.555E-4$	$7.594E-6$	0.0326	0.0339	0.0351
$\pi_8$	0.0315	$6.333E-4$	$6.818E-6$	0.0303	0.0315	0.0328
$\pi_9$	0.0237	$5.534E-4$	$4.943E-6$	0.0226	0.0237	0.0248
$\pi_{10}$	0.0229	$5.457E-4$	$5.224E-6$	0.0219	0.0229	0.0240
$\pi_{11}$	0.3079	0.00168	$1.604E-5$	0.3047	0.3079	0.3113

Table 1.2: Estadísticas básicas de la distribución final de  $\pi$ .

■

### 1.3.2 JAGS

Código R: Bayes7\_4DistNegBinom.R

### 1.3.3 STAN

## 1.4 INLA

Gomez-Rubio (2020)

# Bibliografía

- Albert, J. (2009). *Bayesian Computation with R* (second ed.). New York: Springer.
- Chen, M.-H., Q.-M. Shao, and J. G. Ibrahim (2000). *Monte Carlo Methods in Bayesian Computation*. Series in Statistics. New York: Springer.
- Congdon, P. (2006). *Bayesian Statistical Modelling* (Second ed.). Chichester: John Wiley & Sons.
- Gamerman, D. and H. F. Lopes (2006). *Markov Chain Monte Carlo: Stochastic Simulation for Bayesian Inference* (Second ed.). Boca Raton, Florida: Chapman & Hall/CRC.
- Gilks, W. R., S. Richardson, and D. J. Spiegelhalter (1996). *Markov Chain Monte Carlo in Practice*. London: Chapman and Hall.
- Gomez-Rubio, V. (2020). *Bayesian Inference with INLA* (1st ed.). Boca Raton, Florida: Chapman and Hall/CRC.
- Gutiérrez-Peña, E. (1997). *Métodos Computacionales en la Inferencia Bayesiana*. Monografías. México: IIMAS, UNAM.
- Lunn, D., C. Jackson, N. Best, A. Thomas, and D. Spiegelhalter (2012). *The BUGS Book - A Practical Introduction to Bayesian Analysis*. Boca Raton, Florida: CRC Press / Chapman & Hall.
- Lunn, D. J., A. Thomas, N. Best, and D. Spiegelhalter (2000). WinBUGS - A Bayesian modelling framework: Concepts, structure, and extensibility. *Statistics and Computing* 10(4), 325–337.
- Ntzoufras, I. (2009). *Bayesian Modeling Using WinBUGS*. New Jersey: John Wiley & Sons.
- Ríos-Insúa, D., S. Ríos-Insúa, J. Martín-Jiménez, and A. Jiménez-Martín (2009). *Simulación: métodos y aplicaciones*. Alfaomega.
- Robert, C. P. and G. Casella (2000). *Introducing Monte Carlo Methods with R*. Use R! Springer.