



Estadística Bayesiana

Métodos Computacionales y Algunas Aplicaciones

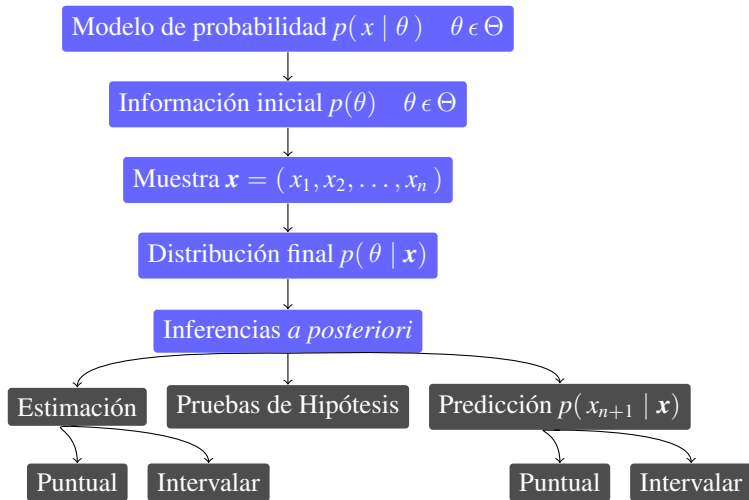
Antonio Soriano Flores

UNIVERSIDAD NACIONAL AUTÓNOMA DE MÉXICO
Instituto de Investigaciones en Matemáticas Aplicadas y en Sistemas

`asoriano@sigma.iimas.unam.mx`

XXXI Foro de Estadística - Universidad Autónoma Chapingo

Estadística Bayesiana



Introducción

En términos generales, en la estadística bayesiana surge la necesidad de llevar a cabo el cálculo de ciertas integrales que pueden llegar a ser analíticamente difíciles de resolver:

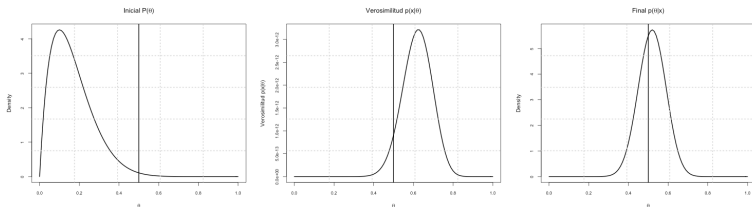
$$\begin{aligned}
 p(\theta|\underline{x}) &\propto p(\underline{x}|\theta)p(\theta) \\
 p(x^*|\underline{x}) &\propto \int_{\Theta} p(x^*|\theta)p(\theta|\underline{x})d\theta \\
 p(\text{Desconocido}|\text{Observado}) &\propto \dots
 \end{aligned}$$

Algunas veces, la solución analítica es relativamente fácil por medio del uso de las **Familias Conjugadas**

$$x \sim \text{Bernoulli}(x|\theta); \quad \theta \sim \text{Beta}(\theta|\alpha_0, \beta_0);$$

Entonces:

$$\underbrace{\theta \sim \text{Beta}(\alpha_0, \beta_0)}_{\text{Inicial}} + \underbrace{p(\underline{x}|\theta)}_{\text{Muestra}} \Rightarrow \underbrace{\theta|\underline{x} \sim \text{Beta}\left(\sum_{i=1}^n x_i + \alpha_0, n - \sum_{i=1}^n x_i + \beta_0\right)}_{\text{Final}}$$



Ver video: [Final.mp4](#)

En este ejemplo, si estamos interesados en modelar una nueva observación x^* , se puede calcular:

- Densidad predictiva inicial:

$$\begin{aligned}
 p(x^*) &= \int p(x^*, \theta) d\theta = \int p(x^*|\theta) p(\theta) d\theta \\
 &= \frac{1}{B(\alpha_0, \beta_0)} \int_0^1 \theta^{x^*} (1-\theta)^{1-x^*} \theta^{\alpha_0-1} (1-\theta)^{\beta_0-1} d\theta \\
 &= \frac{1}{B(\alpha_0, \beta_0)} \int_0^1 \theta^{x^*+\alpha_0-1} (1-\theta)^{\beta_0-x^*} d\theta \\
 &= \frac{B(x^* + \alpha_0, \beta_0 - x^* + 1)}{B(\alpha_0, \beta_0)} \\
 &= \frac{\Gamma(x^* + \alpha_0)\Gamma(\beta_0 - x^* + 1)}{\Gamma(\alpha_0 + \beta_0 + 1)} \frac{\Gamma(\alpha_0 + \beta_0)}{\Gamma(\alpha_0)\Gamma(\beta_0)}
 \end{aligned}$$

Entonces:

$$p(0) = \frac{\beta_0}{\alpha_0 + \beta_0}; \quad p(1) = \frac{\alpha_0}{\alpha_0 + \beta_0}$$

*

Haciendo $\alpha_1 = \sum_{i=1}^n x_i + \alpha_0$ y $\beta_1 = n - \sum_{i=1}^n x_i + \beta_0$,

- Densidad predictiva final:

$$\begin{aligned}
 p(x^*|\underline{x}) &= \int p(x^*|\theta) p(\theta|\underline{x}) d\theta \\
 &= \frac{1}{B(\alpha_1, \beta_1)} \int_0^1 \theta^{x^*} (1-\theta)^{1-x^*} \theta^{\alpha_1-1} (1-\theta)^{\beta_1-1} d\theta \\
 &= \frac{1}{B(\alpha_1, \beta_1)} \int_0^1 \theta^{x^*+\alpha_1-1} (1-\theta)^{\beta_1-x^*-1} d\theta \\
 &= \frac{B(x^* + \alpha_1, \beta_1 - x^* + 1)}{B(\alpha_1, \beta_1)} \\
 &= \frac{\Gamma(x^* + \alpha_1)\Gamma(\beta_0 - x^* + 1)}{\Gamma(\alpha_1 + \beta_1 + 1)} \frac{\Gamma(\alpha_1 + \beta_1)}{\Gamma(\alpha_1)\Gamma(\beta_1)}
 \end{aligned}$$

Entonces:

$$p(0) = \frac{\beta_1}{\alpha_1 + \beta_1} = \frac{n - \sum_{i=1}^n x_i + \beta_0}{n + \alpha_0 + \beta_0}; \quad p(1) = \frac{\alpha_1}{\alpha_1 + \beta_1} = \frac{\sum_{i=1}^n x_i + \alpha_0}{n + \alpha_0 + \beta_0}$$

*

Algunas familias conjugadas:

- $x \sim \text{Poisson}(x|\lambda) \rightarrow \lambda \sim \text{Gamma}(\lambda|\alpha_0, \beta_0)$
- $x \sim \text{Binomial}(x|n^*, \theta) \rightarrow \theta \sim \text{Beta}(\theta|\alpha_0, \beta_0)$
- $x \sim \text{Multinomial}(x|p_1, \dots, p_k, n^*) \rightarrow (p_1, \dots, p_k) \sim \text{Dir}(p_1, \dots, p_k|\alpha_1, \dots, \alpha_k)$
- $x \sim \text{Geometric}(x|\theta) \rightarrow \theta \sim \text{Beta}(\theta|\alpha_0, \beta_0)$
- $x \sim U(x|0, \theta) \rightarrow \theta \sim \text{Pareto}(\theta|\alpha_0, \beta_0)$
- $x \sim \text{Gamma}(x|\alpha^*, \beta) \rightarrow \beta \sim \text{Gamma}(\beta|\alpha_0, \beta_0)$
- $x \sim \text{Normal}(x|\mu, \tau^*) \rightarrow \mu \sim \text{Normal}(\mu|\mu_0, \tau_0)$, donde: $\tau = \frac{1}{\sigma^2}$ se conoce como la precisión del modelo.
- $x \sim \text{Normal}(x|\mu^{*1}, \tau) \rightarrow \tau \sim \text{Gamma}(\tau|\alpha_0, \beta_0)$
- $x \sim \text{Normal}(x|\mu, \tau) \rightarrow (\mu, \tau) \sim \text{Normal} - \text{Gamma}(\mu, \tau|\mu_0, \tau_0, \alpha_0, \beta_0)$

¹* Parámetros que se consideran conocidos

Sin embargo, en muchas ocasiones existen problemas que involucran el cálculo de integrales analíticamente complicadas.

Ejemplo 1:

$$x \sim \text{Gamma}(x|\alpha, 2); \quad \alpha \sim \text{Gamma}(\alpha|\alpha_0, \beta_0)$$

Dada una muestra $\mathbf{x} = (x_1, \dots, x_n)$, hacer inferencias sobre α así como para una nueva observación (x_F).

$$\begin{aligned} p(\alpha|\mathbf{x}) &\propto p(\mathbf{x}|\alpha) p(\alpha) \\ &\propto \left(\prod_{i=1}^n \frac{2^\alpha}{\Gamma(\alpha)} x_i^{(\alpha-1)} e^{-2x_i} \right) \alpha^{(\alpha_0-1)} e^{-\beta_0\alpha} \\ &\propto \frac{2^{n\alpha}}{\Gamma(\alpha)^n} \left(\prod_{i=1}^n x_i \right)^{(\alpha-1)} e^{-2 \sum_{i=1}^n x_i} \alpha^{(\alpha_0-1)} e^{-\beta_0\alpha} \\ &\propto \frac{2^{n\alpha}}{\Gamma(\alpha)^n} \left(\prod_{i=1}^n x_i \right)^{(\alpha-1)} \alpha^{(\alpha_0-1)} e^{-\beta_0\alpha} \\ &\int_0^\infty \frac{2^{n\alpha}}{\Gamma(\alpha)^n} \left(\prod_{i=1}^n x_i \right)^{(\alpha-1)} \alpha^{(\alpha_0-1)} e^{-\beta_0\alpha} d\alpha = ?????? \end{aligned}$$

Por otro lado, para una observación futura:

$$\begin{aligned}
 p(x_F|\mathbf{x}) &= \int_0^\infty p(x_F|\alpha) p(\alpha|\mathbf{x}) d\alpha \\
 &\propto \int_0^\infty \frac{2^\alpha}{\Gamma(\alpha)} x_F^{\alpha-1} e^{-2x_F} \frac{2^{n\alpha}}{\Gamma(\alpha)^n} \left(\prod_{i=1}^n x_i \right)^{(\alpha-1)} \alpha^{(\alpha_0-1)} e^{-\beta_0 \alpha} d\alpha \\
 &\int_0^\infty \frac{2^\alpha}{\Gamma(\alpha)} x_F^{\alpha-1} \frac{2^{n\alpha}}{\Gamma(\alpha)^n} \left(\prod_{i=1}^n x_i \right)^{(\alpha-1)} \alpha^{(\alpha_0-1)} e^{-\beta_0 \alpha} d\alpha =?????
 \end{aligned}$$

Ejemplo 2: Suponga que se tienen x_1 y x_2 observaciones independientes del modelo Cauchy:

$$p(x|\theta) \propto \frac{1}{1 + (x - \theta)^2} \quad (x \in \mathbb{R}; \theta \in \mathbb{R})$$

Se desea hacer inferencias sobre el valor de θ suponiendo una inicial no informativa de la forma:

$$p(\theta) \propto 1$$

La regla de Bayes nos dice entonces que:

$$p(\theta|x_1, x_2) \propto \left(\frac{1}{1 + (x_1 - \theta)^2} \right) \left(\frac{1}{1 + (x_2 - \theta)^2} \right)$$

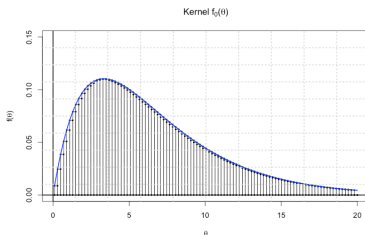
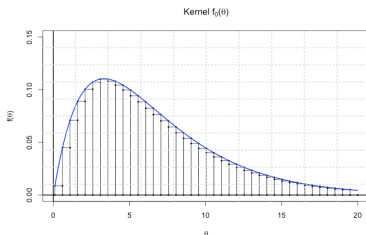
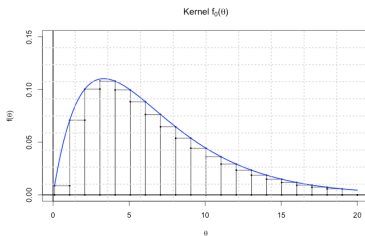
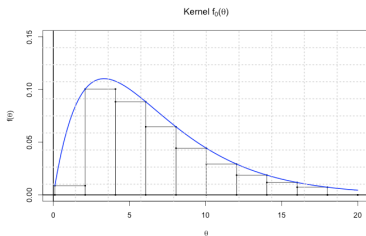
El kernel de la distribución final toma una forma complicada. Es conveniente notar, sin embargo, que $p(x|\theta)$ corresponde a la densidad marginal de x respecto a la densidad conjunta:

$$p(x, \lambda|\theta) = N\left(x \mid \theta, \frac{1}{\lambda}\right) Ga\left(\lambda \mid \frac{1}{2}, \frac{1}{2}\right)$$

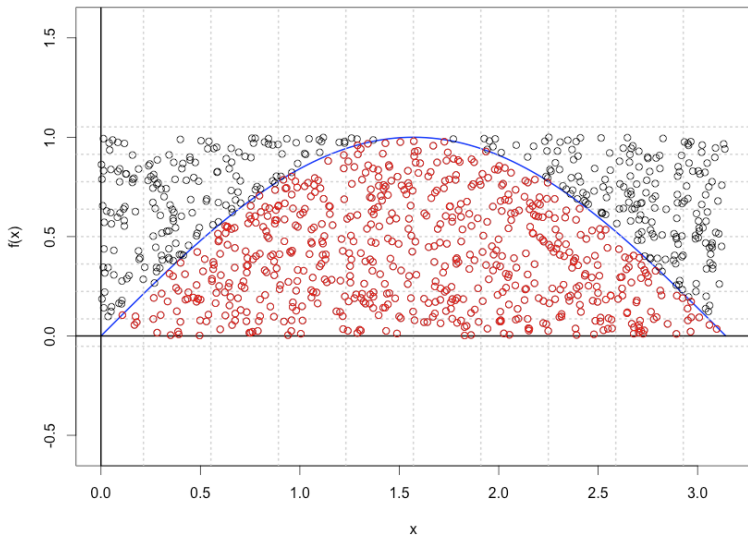
Nota: A veces incluir variables que inicialmente parece complicar el problema, en realidad lo hace más fácil de resolver!!!

Soluciones Numéricas:

- 1 Cuadratura: Regla del Punto Medio, Regla Trapezoidal
- 2 Cuadratura de Gauss-Hermite

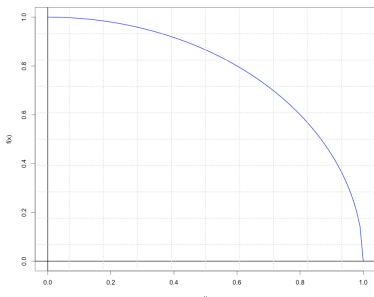


Aproximación vía simulación



Ideas básicas de la aproximación

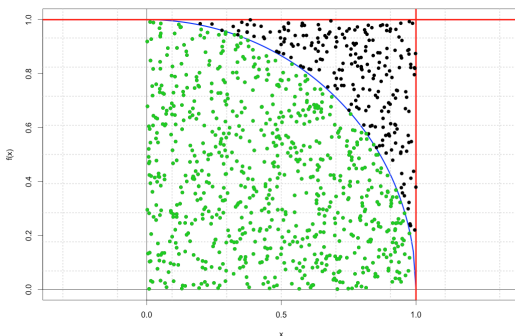
Problema 1: Se desea encontrar el área bajo la curva $f(x) = \sqrt{1-x^2}$ en el intervalo $(0, 1)$



Esta integral se puede resolver vía cambio de variable y el resultado es:

$$I = \int_0^1 \sqrt{1-x^2} dx = \frac{\pi}{4} \approx 0.7853982$$

La idea es construir un rectángulo de área A que cubra completamente a la función, luego simular N observaciones de forma uniforme y contar la proporción de puntos que caen debajo de la curva.



Definiendo N_f como el número de observaciones dentro de la curva, se tendría que el área estimada bajo la curva es:

$$\hat{I} = A \frac{N_f}{N}$$

Intuitivamente, entre más simulaciones hagamos mejor será nuestra aproximación.

N	\hat{I}
1,000	0.8050000
10,000	0.7857000
100,000	0.7848500
1,000,000	0.7856720
10,000,000	0.7851957

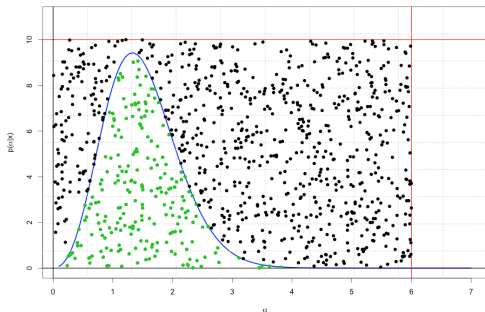
$$I = \int_0^1 \sqrt{1-x^2} dx = \frac{\pi}{4} \approx 0.7853982$$

Ver código: [Programa01.r](#)

Problema 2:

$$\int_0^{\infty} \frac{2^{n\alpha}}{\Gamma(\alpha)^n} \left(\prod_{i=1}^n x_i \right)^{(\alpha-1)} \alpha^{(\alpha_0-1)} e^{-\beta_0\alpha} d\alpha = \text{?????}$$

Asumiendo $n = 3$ y que $\mathbf{x} = (0.7615236, 0.6414124, 0.3593526)$, con hiperparámetros $(\alpha_0 = 0.001, \beta_0 = 0.001)$



En este caso la aproximación es:

$$\hat{I} = (10)(6) \frac{N_f}{N}$$

Las aproximaciones para distintos valores de N son:

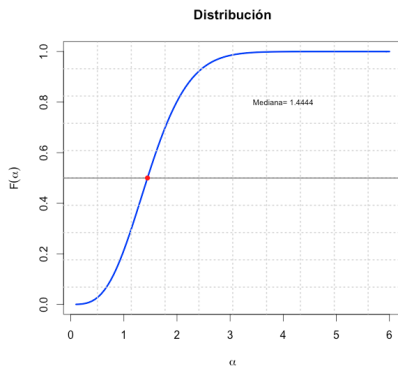
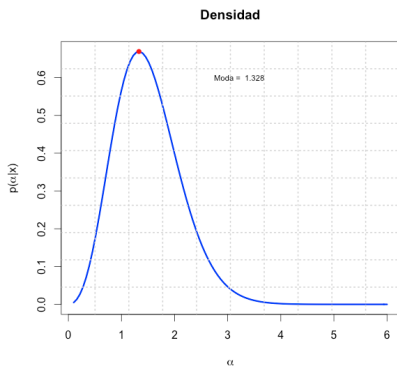
N	\hat{I}
1,000	13.38000
10,000	13.56600
100,000	14.15460
1,000,000	14.09142
10,000,000	14.08628
100,000,000	14.08861

Por lo tanto la aproximación a la constante de proporcionalidad es $1/14.08861 \approx 0.07097932$, de donde concluimos entonces que:

$$p(\alpha|\underline{x}) \approx 0.07097932 \frac{2^{n\alpha}}{\Gamma(\alpha)^n} \left(\prod_{i=1}^n x_i \right)^{(\alpha-1)} \alpha^{(\alpha_0-1)} e^{-\beta_0\alpha}$$

Ver código: [Programa02.r](#)

Calculada la constante de proporcionalidad, podemos encontrar numéricamente algunas cantidades de interés.



Función: maxLik \rightarrow Moda \rightarrow 1.327961

Función: uniroot \rightarrow Mediana \rightarrow 1.444442

Ver código: [Programa02.r](#)

Si queremos calcular la media, nuevamente necesitamos una integral!!!

$$\mathbb{E}(\alpha|\underline{x}) = \int_0^{\infty} 0.07097932 \alpha \frac{2^{n\alpha}}{\Gamma(\alpha)^n} \left(\prod_{i=1}^n x_i \right)^{(\alpha-1)} \alpha^{(\alpha_0-1)} e^{-\beta_0 \alpha} d\alpha$$

Aplicando nuevamente las simulaciones correspondientes:

N	\hat{I}
1,000	1.452000
10,000	1.492800
100,000	1.502220
1,000,000	1.500240
10,000,000	1.500289
100,000,000	1.500306

Concluimos entonces que:

$$\mathbb{E}(\alpha|\underline{x}) \approx 1.500306$$

Ver código: [Programa02.r](#)

Muestreo por importancia

Esta idea se basa en la Ley Fuerte de los Grandes Números, que nos garantiza que si se tiene x_1, \dots, x_n m.a. de un modelo $s(x)$ (puede ser multivariado) y si se asume que

$$\mathbb{E}(g(x)) = \int_{\mathbb{R}^d} g(x)s(x) dx < \infty$$

donde $g : \mathbb{R}^d \rightarrow \mathbb{R}$, entonces

$$\frac{1}{n} \sum_{i=1}^n g(x_i) \rightarrow \mathbb{E}(g(x)) = \int_{-\infty}^{\infty} g(x)s(x) dx$$

Ahora suponga que se está interesado en encontrar la integral:

$$I = \int_{\mathbb{R}^d} f(x) dx$$

Podemos reescribir dicha integral como:

$$I = \int_{\mathbb{R}^d} \left(\frac{f(x)}{s(x)} \right) s(x) dx$$

Si podemos simular x_1, \dots, x_n de la densidad $s(\cdot)$, entonces:

$$\hat{I} = \frac{1}{n} \sum_{i=1}^n \left(\frac{f(x_i)}{s(x_i)} \right) \rightarrow \mathbb{E} \left(\frac{f(x_i)}{s(x_i)} \right) = I$$

A $s(\cdot)$ se le conoce como la distribución de **muestreo por importancia** y debe tener las siguientes características:

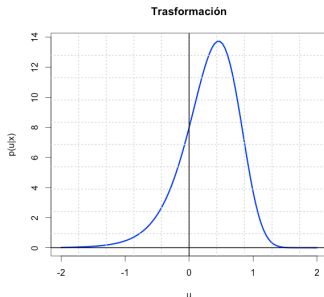
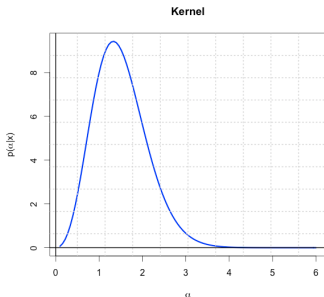
- Debe de ser fácil de simular (Normal, Gamma, Beta, Uniforme,...)
- Debe de tener una forma similar a la de $f(x)$ (la función que se desea integrar)
- Debe de tener el mismo soporte que $f(x)$
- En la práctica es común trabajar en términos de alguna reparametrización de manera que la integral esté definida en todo \mathbb{R}^d y luego utilizar como distribución de muestreo a una Normal o t de Student multivariada.

Ejemplo 1: Consideremos nuevamente el problema relacionado con la constante de proporcionalidad:

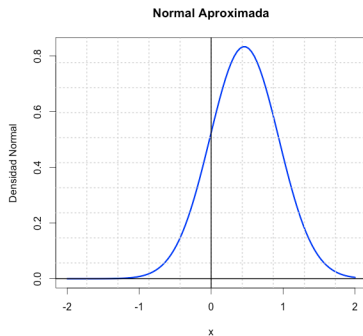
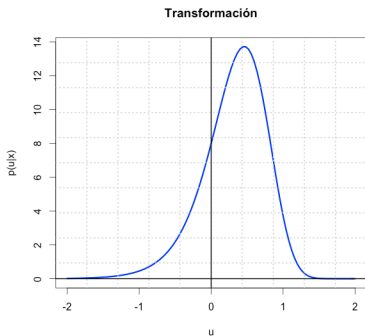
$$\int_0^{\infty} \frac{2^{n\alpha}}{\Gamma(\alpha)^n} \left(\prod_{i=1}^n x_i \right)^{(\alpha-1)} \alpha^{(\alpha_0-1)} e^{-\beta_0 \alpha} d\alpha = \text{????}$$

- Reparametrizamos el integrando para que esté definido en todo \mathbb{R} , mediante el cambio $u = \log(\alpha)$

$$\int_0^{\infty} f(\alpha) d\alpha = \int_{-\infty}^{\infty} f(e^u) e^u du$$



- Buscamos la distribución normal que mejor aproxime la “forma” de la función transformada.



Los parámetros de la normal que se obtuvieron son:

$$\mu = 0.4605435; \quad \sigma^2 = 0.228919$$

Se procede a simular observaciones de esta densidad.

Ver código: [Programa03.r](#)

Los resultados obtenidos son:

N	\hat{I}
1,000	14.23363
10,000	14.13107
100,000	14.08652
1,000,000	14.07526
10,000,000	14.07714
100,000,000	14.08597

En este ejercicio, la constante de proporcionalidad toma el valor de

$$1/14.08326 \approx 0.07099263$$

$$0.07097932 \quad \text{vs} \quad 0.07099263$$

Con esta constante se pueden realizar nuevamente las estimaciones de α .

Ver código: [Programa03.r](#)

Los resultados obtenidos son:

Estadística	Simulación basada en Uniformes	Muestreo por importancia
Constante	0.070979	0.070993
Moda	1.327961	1.327961
Mediana	1.444442	1.444299
Media	1.500306	1.500663

Ver código: [Programa03.r](#)

Simulación de una distribución

Los métodos anteriores permiten calcular resúmenes inferenciales a partir del cálculo de integrales complicadas, por ejemplo:

- La constante de normalización
- Valores esperados

Los métodos que a continuación se presentan se basan en **simular muestras de la distribución final**, las cuales permiten, en principio, aproximar cualquier característica de interés, como son cuantiles, medias, varianzas, etc.

Existe una dualidad interesante aquí: por un lado, dada la distribución podemos simular observaciones de ella; por otro lado, dada una muestra grande es posible recrear la distribución que la generó (**Teorema de Glivenko-Cantelli**).

$$\lim_{n \rightarrow \infty} \sup (|F_n(x) - F(x)|) = 0$$

La idea entonces es: Supongamos que tenemos $f_0(\theta)$ el kernel de una densidad, de tal forma que:

$$f(\theta) = \frac{f_0(\theta)}{\int f_0(\theta') d\theta'}$$

El problema es generar un algoritmo que pueda simular observaciones de la densidad $f(\theta)$ a través de simulaciones de otra densidad $s(\cdot)$. En este curso veremos dos casos:

- Muestreo por importancia
- Simulación vía cadenas de Markov (Uso de JAGS)

Muestreo por importancia

Este método asume que podemos encontrar una constante M tal que $\frac{f_0(\theta)}{s(\theta)} \leq M$ para todo θ y que $s(\cdot)$ es una función fácil de simular. El algoritmo que se propone es el siguiente:

Algoritmo

- 1 Generar una observación $\tilde{\theta}$ de $s(\theta)$
- 2 Generar una variable $u \sim U(0, 1)$
- 3 si $u \leq \frac{f_0(\tilde{\theta})}{M s(\tilde{\theta})}$, aceptar $\tilde{\theta}$ como una observación de la densidad $f(\theta)$; en caso contrario, repetir los pasos 1 a 3.

Observación: Suponiendo que se desea una muestra de tamaño N de $f(\theta)$, el valor esperado para el tamaño de la correspondiente muestra de $s(\theta)$ es:

$$N_0 = \frac{MN}{\int f_0(\theta) d\theta}$$

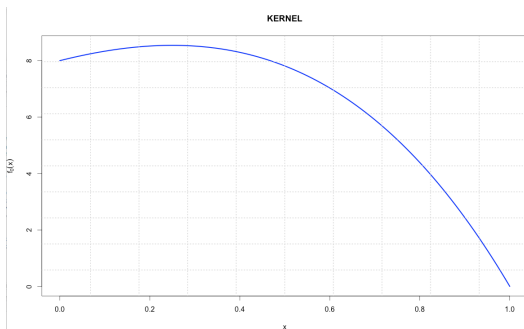
Ejemplo 3 : Suponga que se tiene el siguiente kernel:

$$f_0(x) = (x + 2)^3(1 - x) \quad x \in (0, 1)$$

Se desea simular observaciones de dicho kernel para aproximar la media $\mathbb{E}(X)$ y la mediana ($q_{0.5}$).

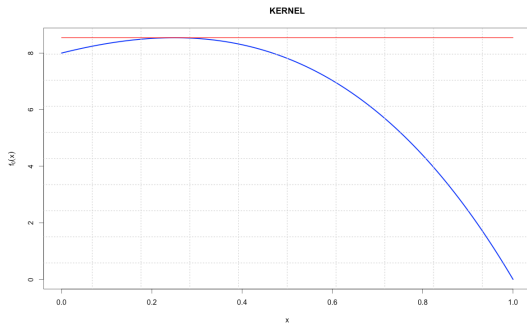
Se puede probar que, en este caso, la constante de proporcionalidad así como las cantidades de interés son:

$$k = \frac{20}{131} \approx 0.1527; \quad \mathbb{E}(x) = \frac{160}{393} \approx 0.4071; \quad q_{0.5} = 0.3900$$



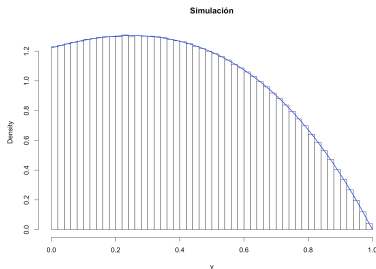
Proponemos $s \sim U(x|0, 1)$ y $M = \sup\{f_0(x)\}$ El algoritmo es:

- 1 Generar una observación \tilde{x} de $s(x) = 1$, es decir $\tilde{x} \sim U(0, 1)$
- 2 Generar una variable $u \sim U(0, 1)$
- 3 si $u \leq \frac{f_0(\tilde{x})}{M s(\tilde{x})}$, aceptar \tilde{x} como una observación de la densidad $f(x)$; en caso contrario, repetir los pasos 1 a 3. En este caso $s(\tilde{x}) = 1$ y $M = 8.542969$



Corriendo el algoritmo anterior para $N_0 = 100,000$, se obtuvo que:

$$N = 76,771; \quad \frac{N}{N_0} = 0.76776; \quad k \approx \frac{N_0}{MN} = 0.152463$$

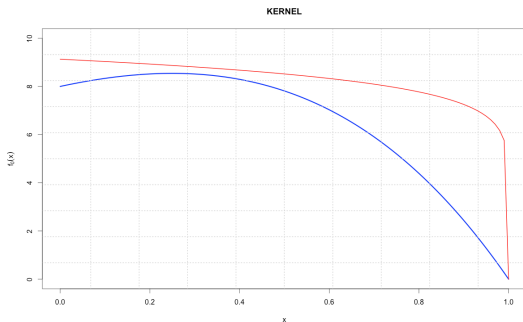


En este caso vía simulación las características distribucionales son aproximadas como:

$$\mathbb{E}(X) \approx \frac{1}{N} \sum_{i=1}^N y_i = 0.4071616 \quad q_{0.5} \approx \text{mediana}(y_1, \dots, y_N) = 0.3897104$$

Ver código: [Programa04.r](#)

Es deseable mejorar la tasa de aceptación. Para ello se puede proponer una densidad $s(\cdot)$ más parecida al kernel $f_0(x)$. En el ejemplo anterior, supongamos entonces que $s(x) = \text{Beta}(x|1, 1.1)$ y que se propone $M = 8.1$



El algoritmo es:

- 1 Generar una observación \tilde{x} de $s(x) = \text{Beta}(x|1, 1.1)$
- 2 Generar una variable $u \sim U(0, 1)$
- 3 si $u \leq \frac{f_0(\tilde{x})}{M s(\tilde{x})}$, aceptar \tilde{x} como una observación de la densidad $f(x)$; en caso contrario, repetir los pasos 1 a 3. Corriendo el algoritmo anterior para $N_0 = 100,000$ se obtuvo que:

$$N = 80653; \quad \frac{N}{N_0} = 0.80653; \quad k \approx \frac{N_0}{MN} = 0.1530715$$

Mientras que las características distribucionales son aproximadas como:

$$\mathbb{E}(X) \approx \frac{1}{N} \sum_{i=1}^N y_i = 0.4066467 \quad q_{0.5} \approx \text{mediana}(y_1, \dots, y_N) = 0.3884548$$

Ver código: [Programa05.r](#)

En resumen tenemos que las simulaciones arrojan lo siguiente:

Estadística		Exacto	Uniforme	Beta
Constante	$\frac{20}{131} \approx$	0.1526718	0.1524630	0.1530715
Mediana		0.3900254	0.3897104	0.3884548
Media	$\frac{160}{393} \approx$	0.4071247	0.4071616	0.4066467
Tasa de Aceptación			0.76776	0.80653

Ver código: [Programa05.r](#)

Ejercicio: Simular observaciones de la distribución $N(0, 1)$ utilizando $s(x) = \text{Cauchy}(x|0, 1)$

Ver código: [Programa06.r](#)

*

Ejemplo Bayesiano:

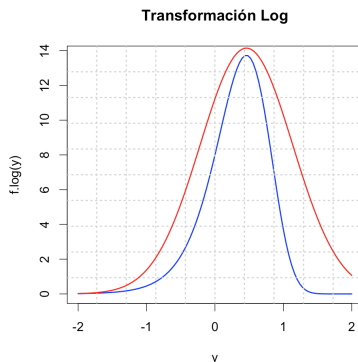
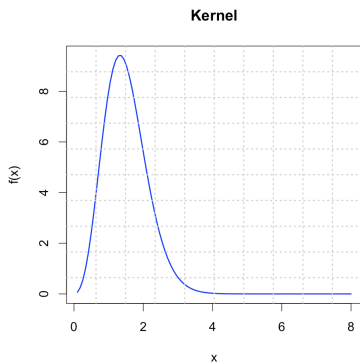
$$x \sim \text{Gamma}(x|\alpha, 2); \quad \alpha \sim \text{Gamma}(\alpha|\alpha_0, \beta_0)$$

Dada $\mathbf{x} = (x_1, \dots, x_n)$ una muestra, hacer inferencia sobre α

$$p(\alpha|\underline{x}) \propto \frac{2^{n\alpha}}{\Gamma(\alpha)^n} \left(\prod_{i=1}^n x_i \right)^{(\alpha-1)} \alpha^{(\alpha_0-1)} e^{-\beta_0\alpha}$$

Asumiendo $n = 3$ y que $\mathbf{x} = (0.7615236, 0.6414124, 0.3593526)$, con hiperparámetros ($\alpha_0 = 0.001, \beta_0 = 0.001$), simularemos observaciones de la distribución $p(\alpha|\mathbf{x})$ para posteriormente hacer la inferencia correspondiente

El proceso consiste en transformar el kernel para tomar valores en todos los reales ($Y = \log(\alpha)$) y proponer a la densidad normal que ajuste lo mejor posible a dicha transformación



Se procede entonces a simular observaciones de la v.a. Y .

Se propone $s(x) = \text{Normal}(x|0.4605435, 0.457838)$, $M = 24$

Ver código: [Programa07.r](#)

El algoritmo es:

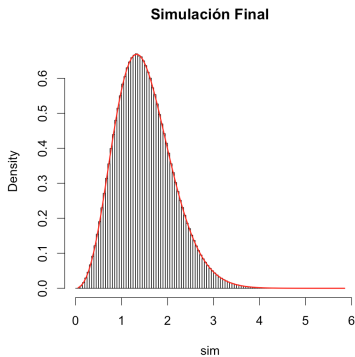
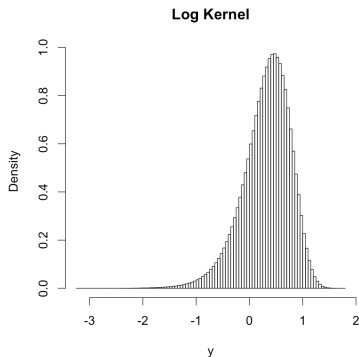
- 1 Generar una observación \tilde{x} de $s(x) = \text{Normal}(x|0.4605435, 0.457838)$
- 2 Generar una variable $u \sim U(0, 1)$
- 3 si $u \leq \frac{f_0(\tilde{x})}{M s(\tilde{x})}$, aceptar \tilde{x} como una observación de la densidad $f(x)$; en caso contrario, repetir los pasos 1 a 3.

Corriendo el algoritmo anterior para $N_0 = 10,000,000$ se obtuvo que:

$$N_0 = 5,869,603; \quad \frac{N}{N_0} = 0.5869603;$$

Lo anterior genera simulaciones del kernel transformado

Ver código: [Programa07.r](#)



La esperanza a posterior es:

$$\mathbb{E}(\alpha | \underline{x}) \approx 1.501118;$$

El intervalo de máxima densidad al 95 % es:

$$(0.3937097, 2.6891916)$$

Ver código: [Programa07.r](#)

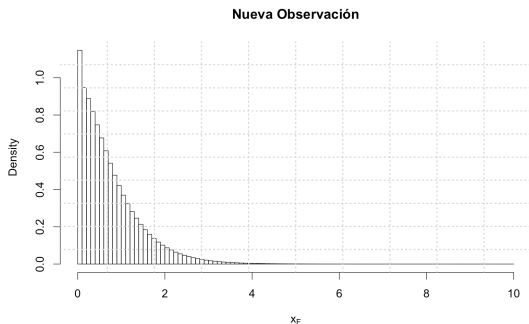
Una de las ventajas de tener simulaciones de la distribución final es que podemos simular fácilmente observaciones de la distribución de x_F (una nueva observación). Para ello recordemos que:

$$\begin{aligned} p(x_F, \alpha | \mathbf{x}) &= \frac{p(x_F, \alpha, \mathbf{x})}{p(\mathbf{x})} \\ &= \frac{p(x_F | \alpha, \mathbf{x}) p(\alpha, \mathbf{x})}{p(\mathbf{x})} \\ &= p(x_F | \alpha) p(\alpha | \mathbf{x}) \end{aligned}$$

Es decir, simular observaciones del vector (x_F, α) dada la muestra \underline{x} es simple mediante el siguiente algoritmo (muestreo condicional):

- 1 Generar una observación $\tilde{\alpha}$ de $p(\alpha | \mathbf{x})$
- 2 Generar una observación \tilde{x} de la densidad $Gamma(x | \tilde{\alpha}, 2)$
- 3 Repetir paso 1 y 2. Entonces $(\tilde{x}, \tilde{\alpha})$ son simulaciones del vector aleatorio (x_F, α) dada la muestra \mathbf{x}
Ver código: [Programa07.r](#)

Simulación para una nueva observación del modelo:



En este caso:

$$\mathbb{E}(x_F | \mathbf{x}) = 0.75045$$

y un intervalo de credibilidad al 95 % es:

$$(0, 2.104152)$$

Ver código: [Programa07.r](#)

Ejercicio: Sea $(0, 0, 0, 1, 0)$ una muestra observada del modelo $Bernoulli(x|\theta)$, asumiendo que la inicial es $p(\theta) = -\log(\theta)$, simular observaciones de la densidad final $p(\theta|\underline{x})$ y encontrar un intervalo de credibilidad al 95 % para θ .

Ver código: [Programa08.r](#)

*

Monte Carlo vía cadenas de Markov

Esta técnica permite generar, de manera iterativa, observaciones de distribuciones multivariadas que difícilmente podrían simularse utilizando métodos directos. La idea es simple: **Construir una cadena de Markov que sea fácil de simular y cuya distribución de equilibrio corresponda a la distribución que nos interesa.**

Teorema

Sea $\theta^{(1)}, \theta^{(2)}, \dots$, una cadena de Markov homogénea, irreducible y aperiódica, con espacio de estados Θ y distribución de equilibrio $p(\theta|x)$. Entonces, conforme $t \rightarrow \infty$ se tiene que:

- $\theta^{(t)} \rightarrow \theta \sim p(\theta|x)$

Es decir, si dejamos correr por mucho tiempo a la cadena, eventualmente estaremos simulando observaciones de la distribución $p(\theta|x)$

Consideremos un ejemplo para el caso discreto. Se tiene sólo dos posibles estados $\{0, 1\}$, con las siguientes probabilidades de transición:

$$\begin{array}{cc} & \begin{array}{cc} 0 & 1 \end{array} \\ \begin{array}{c} 0 \\ 1 \end{array} & \begin{pmatrix} 0.2 & 0.8 \\ 0.1 & 0.9 \end{pmatrix} \end{array}$$

Se demuestra que una forma de obtener la distribución estacionaria es resolver el sistema

$$\boldsymbol{\pi} = \boldsymbol{\pi}P,$$

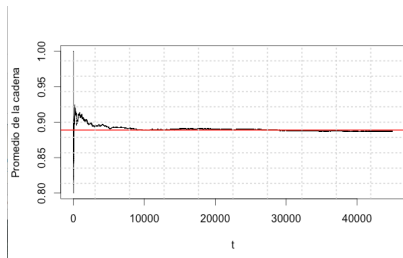
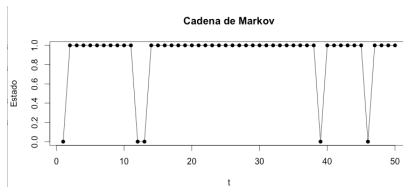
de donde se encuentra que la distribución de equilibrio es $\boldsymbol{\pi} = \left(\frac{0.1}{0.9}, \frac{0.8}{0.9}\right)$.

Teóricamente sabemos entonces que no importa dónde inicie la cadena, después de varias iteraciones la cadena estará simulando observaciones de la distribución estacionaria dada por:

$$\mathbb{P}(X = 0) = \frac{1}{9} \quad \mathbb{P}(X = 1) = \frac{8}{9}$$

Es decir, en este caso, la distribución estacionaria es una *Bernoulli* $\left(\frac{8}{9}\right)$

Simulación de la cadena:

Ver código: [Programa09.r](#)

Ejercicio: Suponga que se tiene $\{X_n : n \in \mathbb{N}\}$, un proceso estocástico con espacio de estados $\{0, 1, 2\}$ y con matriz de transición dada por:

$$\begin{array}{c} \\ 0 \\ 1 \\ 2 \end{array} \begin{pmatrix} 0 & 1 & 2 \\ 1/2 & 1/2 & 0 \\ 0 & 1/3 & 2/3 \\ 1/3 & 1/3 & 1/3 \end{pmatrix}$$

Realizar un programa en R que simule la cadena y verificar que la cadena converge a la distribución estacionaria dada por:

$$\pi_0 = 1/4; \quad \pi_1 = 3/8; \quad \pi_2 = 3/8$$

Ver código: [Programa10.r](#)

Problema: Dada una densidad $p(\theta|x)$, ¿Cómo construir una cadena que tenga como distribución de equilibrio precisamente a $p(\theta|x)$?

*

Algoritmo de Metropolis-Hasting (1953)

El algoritmo construye una cadena de Markov **apropiada** definiendo probabilidades de transición tal que se cumplan las condiciones para tener una cadena estacionaria con distribución límite $p(\theta|x)$.

El algoritmo requiere una distribución de transición (en principio arbitraria) $Q(\theta^*|\theta)$ que sea fácil de simular. Se define:

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

- 1 Inicializar la cadena $\theta^{(0)}$ (en un principio de forma arbitraria)
- 2 Para la t -ésima iteración. Simular θ^* de $Q(\theta^*|\theta^{(t-1)})$;
- 3 Genera una variable $u \sim U(0, 1)$
- 4 Si $u \leq \alpha(\theta^*, \theta^{(t-1)})$, hacer $\theta^{(t)} = \theta^*$; en caso contrario, hacer $\theta^{(t)} = \theta^{(t-1)}$.

Obs: $\alpha(\theta^*, \theta^{(t-1)})$, en este contexto, es la probabilidad de que la cadena se mueva de $\theta^{(t-1)}$ a θ^* y para su cálculo **no es necesario tener la constante de proporcionalidad** de la densidad $p(\theta|x)$.

Técnicamente, la distribución de transición $Q(\theta^*|\theta)$ puede ser arbitraria. Sin embargo, una forma práctica es suponer independencia y asumir que:

$$Q(\theta^*|\theta) = Q_0(\theta^*)$$

Donde para evitar que la cadena se quede estancada, se sugiere que Q_0 tengo una forma similar al kernel de $p(\theta|x)$ En este caso:

$$\alpha(\theta^*, \theta) = \min \left\{ \frac{p(\theta^*|x)Q_0(\theta)}{p(\theta|x)Q_0(\theta^*)}, 1 \right\}$$

En la práctica es común utilizar, después de una reparametrización apropiada, distribuciones de transición normales o t de Student sobredispersas.

Por ejemplo:

$$Q_0(\theta^*) = N_d \left(\theta^* \mid \hat{\theta}, kV(\hat{\theta}) \right)$$

donde $\hat{\theta}$ y $V(\hat{\theta})$ denotan la media y la matriz de varianzas y covarianzas de la aproximación normal asintótica para $p(\theta|x)$ y k es un factor de dispersión para lograr explorar mejor el soporte de la densidad que se desea simular.

Por construcción, después de un determinado número de iteraciones, la cadena debe empezar a estabilizarse y comenzar a simular observaciones de la distribución estacionaria, en este caso de la final $p(\theta|x)$. Una cuestión interesante es el momento en el cual la cadena se ha estabilizado. El objetivo es simular N observaciones de $p(\theta|x)$; para ello dos posibles opciones son:

- Fijar T suficiente grande, luego inicializar N cadenas $\theta_1^{(0)}, \dots, \theta_N^{(0)}$ y correrlas durante T pasos. Finalmente, considerar los valores $\theta_1^{(T)}, \theta_2^{(T)}, \dots, \theta_N^{(T)}$ como una muestra de f_Y . (Computacionalmente es demandante)
- Correr una sola cadena, luego fijar T suficientemente grande y tomar $\theta^{(T+K)}, \theta^{(T+2K)}, \dots, \theta^{(T+NK)}$ como m.a. de $p(\theta|x)$, donde K se elige de manera que la correlación entre las observaciones sea pequeña.

Nota: No es fácil determinar en qué momento la cadena converge, por lo que comúnmente se hacen pruebas empíricas, por ejemplo sobre los promedios ergódicos.

Ejemplo 1 : (Consideremos nuevamente el caso Gamma)

$$x \sim \text{Gamma}(x|\alpha, 2); \quad \alpha \sim \text{Gamma}(\alpha|\alpha_0, \beta_0)$$

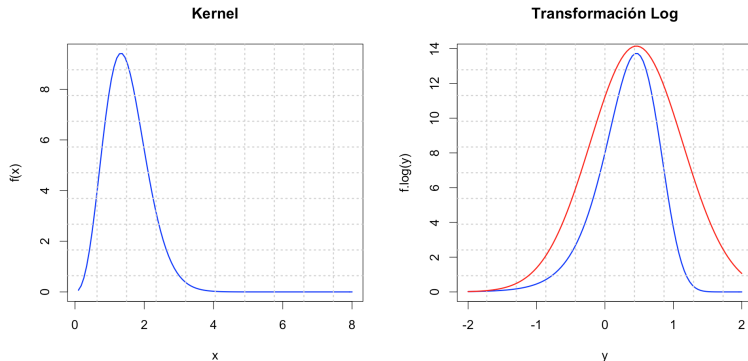
Dada $\mathbf{x} = (x_1, \dots, x_n)$ una muestra, hacer inferencia sobre α

$$p(\alpha|\underline{x}) \propto \frac{2^{n\alpha}}{\Gamma(\alpha)^n} \left(\prod_{i=1}^n x_i \right)^{(\alpha-1)} \alpha^{(\alpha_0-1)} e^{-\beta_0\alpha}$$

Asumiendo $n = 3$ y que $\mathbf{x} = (0.7615236, 0.6414124, 0.3593526)$, con hiperparámetros $(\alpha_0 = 0.001, \beta_0 = 0.001)$, simularemos observaciones de la distribución $p(\alpha|\underline{x})$ utilizando el Algoritmo de Metropolis-Hasting para posteriormente hacer la inferencia correspondiente.

Ver código: [Programa 11.r](#)

El proceso consiste en transformar el kernel para tomar valores en todos los reales ($Y = \log(\alpha)$) y proponer a la densidad normal que ajuste lo mejor posible a dicha transformación



Se procede entonces a simular observaciones de la v.a. Y .

Se propone $Q_0(\theta) = \text{Normal}(\theta|0.4605435, 0.457838)$

Ver código: [Programa11.r](#)

El algoritmo de Metropolis-Hasting es:

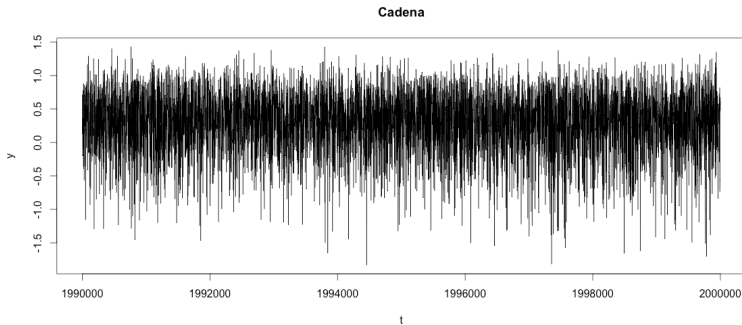
- 1 Inicializar la cadena $y^{(0)}$ (en un principio de forma arbitraria, pero cuidando evitar errores numéricos)
- 2 Para la t -ésima iteración. Simular y^* de $N(\theta|0.4605435, 0.457838)$;
- 3 Genera una variable $u \sim U(0, 1)$
- 4 Si $u \leq \alpha(y^*, y^{(t-1)})$, hacer $y^{(t)} = y^*$; en caso contrario, hacer $y^{(t)} = y^{(t-1)}$.

En este caso:

$$\alpha(y^*, y^{(t-1)}) = \min \left\{ \frac{f_Y(y^*) \text{Normal}(y^{(t-1)} | 0.4605435, 0.457838)}{f_Y(y^{(t-1)}) \text{Normal}(y^* | 0.4605435, 0.457838)}, 1 \right\}$$

Ver código: [Programa 11.r](#)

Se corrió el algoritmo durante 2,000,000 iteraciones. Calentamiento: 100,000.



Ver código: [Programa11.r](#)

Ejemplo 3: Supongamos ahora 2 parámetros desconocidos:

$$x \sim \text{Gamma}(x|\alpha, \beta); \quad p(\alpha, \beta) = \text{Gamma}(\alpha|\alpha_0, \beta_0)\text{Gamma}(\beta|\alpha_1, \beta_1)$$

Asumiendo que:

$$\alpha_0 = \beta_0 = \alpha_1 = \beta_1 = 0.001$$

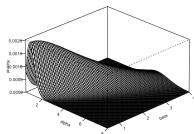
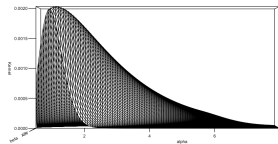
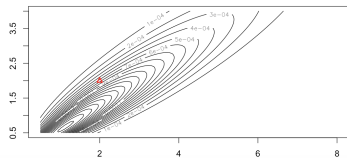
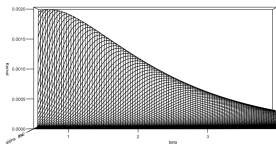
y que se observa la muestra de tamaño 5:

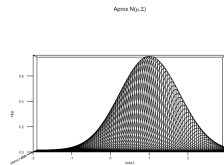
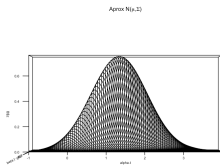
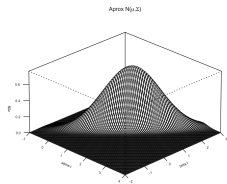
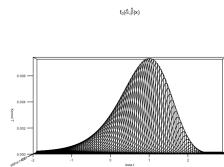
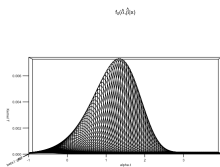
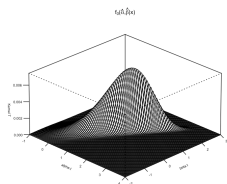
$$\mathbf{x} = (0.4154325, 1.7853782, 1.7315852, 1.0254059, 1.9427045)$$

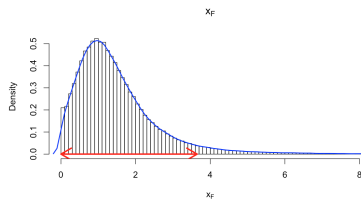
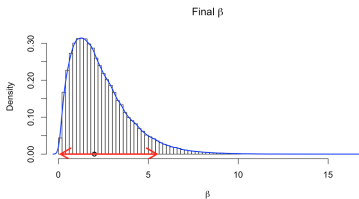
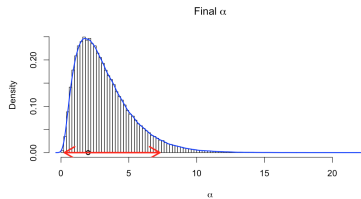
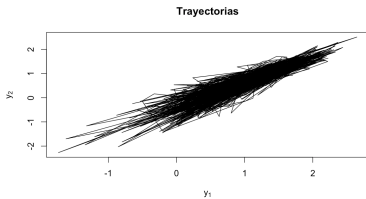
Hacer inferencia sobre α, β y x_F

Solución clásica: Mediante métodos numéricos

$$\hat{\alpha}_{MV} = 3.882556; \quad \hat{\beta}_{MV} = 2.813241; \quad x_F \sim ?$$

$f_3(\alpha, \beta)$  $f_3(\alpha, \beta)$  $f_3(\alpha, \beta)$ 





Ver video: [CADENA MARKOV.mp4](#)

Ver código: [Programa12.r](#)

Algoritmo de Gibbs

En el Algoritmo MH, el valor $\alpha(\theta^*, \theta)$ se interpreta como la probabilidad de que la cadena se mueva de θ a θ^* . Lo ideal es construir un algoritmo en el cual $\alpha(\theta^*, \theta) = 1$, es decir, que el movimiento esté siempre garantizado para favorecer posteriormente a la convergencia de la cadena. El Algoritmo de Gibbs no es más que un caso especial del Algoritmo MH que precisamente garantiza lo mencionado anteriormente. La principal característica de este algoritmo es:

- Se requiere poder simular de las condicionales completas

$$p(\theta_i | \theta_{[i]}, \mathbf{x}); \quad \theta_{[i]} = (\theta_1, \dots, \theta_{[i-1]}, \theta_{[i+1]}, \dots, \theta_k)$$

- Cada valor de la cadena se obtiene a través de un proceso iterativo que utiliza a las condicionales completas.

Algoritmo de Gibbs

El algoritmo propuesto es el siguiente:

- 0 Inicializar la cadena en un valor inicial $\theta^{(0)} = (\theta_1^{(0)}, \dots, \theta_k^{(0)})$
 Para $t \in \{0, 1, 2, \dots, \}$
- 1 Generar una muestra $\theta_1^{(t+1)}$ de $p(\theta_1^{(t+1)} \mid \theta_2^{(t)}, \theta_3^{(t)}, \dots, \theta_k^{(t)}, \mathbf{x})$
- 2 Generar una muestra $\theta_2^{(t+1)}$ de $p(\theta_2^{(t+1)} \mid \theta_1^{(t+1)}, \theta_3^{(t)}, \dots, \theta_k^{(t)}, \mathbf{x})$
- 3 Generar una muestra $\theta_3^{(t+1)}$ de $p(\theta_3^{(t+1)} \mid \theta_1^{(t+1)}, \theta_2^{(t+1)}, \dots, \theta_k^{(t)}, \mathbf{x})$
- \vdots
- k Generar una muestra $\theta_k^{(t+1)}$ de $p(\theta_k^{(t+1)} \mid \theta_1^{(t+1)}, \theta_2^{(t+1)}, \dots, \theta_{k-1}^{(t+1)}, \mathbf{x})$
- k+1 Con los pasos 1 a k construir $\theta^{(t+1)} = (\theta_1^{(t+1)}, \theta_2^{(t+1)}, \dots, \theta_k^{(t+1)})$
- k+2 Repetir los pasos 1 a k + 1, y generar la cadena $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}, \dots,$

Algoritmo de Gibbs

Ejemplo 2: Recordando el modelo Cauchy con inicial no informativa.

$$p(x|\theta) \propto \frac{1}{1 + (x - \theta)^2}; \quad p(\theta) \propto 1; \quad (x \in \mathbb{R}; \theta \in \mathbb{R})$$

Suponiendo que se observan 2 muestras independientes $\mathbf{x} = (x_1, x_2)$, la densidad final es:

$$p(\theta|\mathbf{x}) \propto \left(\frac{1}{1 + (x_1 - \theta)^2} \right) \left(\frac{1}{1 + (x_2 - \theta)^2} \right)$$

Nos interesa simular observaciones de este kernel para aproximar algunas características de interés (media, mediana, moda, IC al 95 %.)

Algoritmo de Gibbs

Para solucionar este problema se introducen variables latentes para facilitar el problema de simulación.

Al inicio de la sesión se mencionó que $p(x_1|\theta)$ se obtiene al marginalizar la siguiente densidad conjunta

$$p(x_1, \lambda_1|\theta) = N\left(x_1 \mid \theta, \frac{1}{\lambda_1}\right) Ga\left(\lambda_1 \mid \frac{1}{2}, \frac{1}{2}\right)$$

Lo anterior se puede generalizar para dos observaciones del modelo. En este caso se prueba que $p(x_1, x_2|\theta)$ se obtiene al marginalizar la siguiente densidad conjunta:

$$p(x_1, x_2, \lambda_1, \lambda_2|\theta) = p(\mathbf{x}, \lambda_1, \lambda_2|\theta) = \prod_{i=1}^2 N\left(x_i \mid \theta, \frac{1}{\lambda_i}\right) Ga\left(\lambda_i \mid \frac{1}{2}, \frac{1}{2}\right)$$

Algoritmo de Gibbs

El problema parece haberse **triplicado**, pues de tener sólo a θ como parámetro de interés, ahora han aparecido dos nuevas variables λ_1 y λ_2 .

Sin embargo notemos que gracias a las variables introducidas, es posible simular observaciones del vector $(\theta, \lambda_1, \lambda_2)$ dada la muestra x_1, x_2 mediante el Algoritmo de Gibbs. Para ello necesitaremos las distribuciones condicionales completas:

- Para λ_1 :

$$\begin{aligned}
 p(\lambda_1 | \lambda_2, \theta, \mathbf{x}) &= \frac{p(\lambda_1, \lambda_2, \theta, \mathbf{x})}{p(\lambda_2, \theta, \mathbf{x})} = \frac{p(\mathbf{x}, \lambda_1, \lambda_2 | \theta) p(\theta)}{p(\lambda_2, \theta, \mathbf{x})} \\
 &\propto p(\mathbf{x}, \lambda_1, \lambda_2 | \theta) = \prod_{i=1}^2 N\left(x_i \mid \theta, \frac{1}{\lambda_i}\right) Ga\left(\lambda_i \mid \frac{1}{2}, \frac{1}{2}\right) \\
 &\propto N\left(x_1 \mid \theta, \frac{1}{\lambda_1}\right) Ga\left(\lambda_1 \mid \frac{1}{2}, \frac{1}{2}\right) \\
 &\propto \lambda_1^{\frac{1}{2}} e^{-\frac{\lambda_1}{2}(x_1 - \theta)^2} \lambda_1^{\frac{1}{2} - 1} e^{-\frac{1}{2}\lambda_1} \\
 &\propto Ga\left(\lambda_1 \mid 1, \frac{1 + (x_1 - \theta)^2}{2}\right)
 \end{aligned}$$

Lo anterior es conocido como el **Método de Variables Latentes**

Algoritmo de Gibbs

- Para λ_2 , de forma similar se prueba que:

$$p(\lambda_2 | \lambda_1, \theta, \mathbf{x}) = Ga\left(\lambda_2 \mid 1, \frac{1 + (x_2 - \theta)^2}{2}\right)$$

- Para θ :

$$\begin{aligned} p(\theta | \lambda_1, \lambda_2, \mathbf{x}) &= \frac{p(\theta, \lambda_1, \lambda_2, \mathbf{x})}{p(\lambda_1, \lambda_2, \mathbf{x})} = \frac{p(\mathbf{x}, \lambda_1, \lambda_2 | \theta) p(\theta)}{p(\lambda_1, \lambda_2, \mathbf{x})} \\ &\propto p(\mathbf{x}, \lambda_1, \lambda_2 | \theta) = \prod_{i=1}^2 N\left(x_i \mid \theta, \frac{1}{\lambda_i}\right) Ga\left(\lambda_i \mid \frac{1}{2}, \frac{1}{2}\right) \\ &\propto N\left(x_1 \mid \theta, \frac{1}{\lambda_1}\right) N\left(x_2 \mid \theta, \frac{1}{\lambda_2}\right) \\ &\propto e^{-\frac{\lambda_1}{2}(x_1 - \theta)^2} e^{-\frac{\lambda_2}{2}(x_2 - \theta)^2} \\ &\propto N\left(\theta \mid \mu_0, \sigma_0^2\right) \end{aligned}$$

donde:

$$\mu_0 = \frac{\lambda_1 x_1 + \lambda_2 x_2}{\lambda_1 + \lambda_2}; \quad \sigma_0^2 = \frac{1}{\lambda_1 + \lambda_2}$$

Algoritmo de Gibbs

En resumen, las condicionales completas son

$$p(\lambda_1 | \lambda_2, \theta, \mathbf{x}) = \text{Ga} \left(\lambda_1 \left| 1, \frac{1 + (x_1 - \theta)^2}{2} \right. \right)$$

$$p(\lambda_2 | \lambda_1, \theta, \mathbf{x}) = \text{Ga} \left(\lambda_2 \left| 1, \frac{1 + (x_2 - \theta)^2}{2} \right. \right)$$

$$p(\theta | \lambda_1, \lambda_2, \mathbf{x}) = N \left(\theta \left| \frac{\lambda_1 x_1 + \lambda_2 x_2}{\lambda_1 + \lambda_2}, \frac{1}{\lambda_1 + \lambda_2} \right. \right)$$

¡Todas las condicionales son fáciles de simular!!! Por lo que se puede aplicar Gibbs.

Algoritmo de Gibbs

El Algoritmo de Gibbs para simular de $p(\theta, \lambda_1, \lambda_2 | \mathbf{x})$ es el siguiente:

0 Inicializar la cadena en un valor inicial $\boldsymbol{\theta}^{(0)} = (\theta^{(0)}, \lambda_1^{(0)}, \lambda_2^{(0)})$

Para $t \in \{0, 1, 2, \dots, \}$

1 Generar una muestra $\lambda_1^{(t+1)}$ de

$$p\left(\lambda_1^{(t)} \mid \lambda_2^{(t)}, \theta^{(t)}, \mathbf{x}\right) = \text{Ga}\left(\lambda_1^{(t)} \mid 1, \frac{1+(x_1-\theta^{(t)})^2}{2}\right)$$

2 Generar una muestra $\lambda_2^{(t+1)}$ de

$$p\left(\lambda_2^{(t)} \mid \lambda_1^{(t)}, \theta^{(t)}, \mathbf{x}\right) = \text{Ga}\left(\lambda_2^{(t)} \mid 1, \frac{1+(x_2-\theta^{(t)})^2}{2}\right)$$

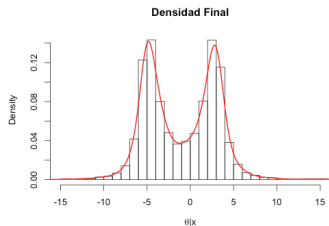
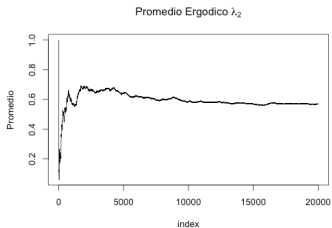
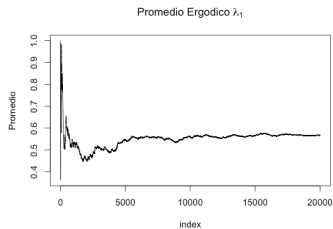
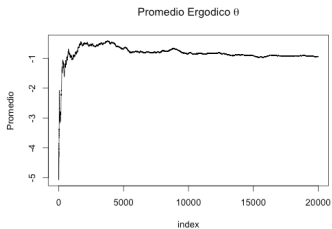
3 Generar una muestra $\theta^{(t+1)}$ de

$$p\left(\theta^{(t+1)} \mid \lambda_1^{(t+1)}, \lambda_2^{(t+1)}, \mathbf{x}\right) = N\left(\theta \mid \frac{\lambda_1^{(t+1)}x_1 + \lambda_2^{(t+1)}x_2}{\lambda_1^{(t+1)} + \lambda_2^{(t+1)}}, \frac{1}{\lambda_1^{(t+1)} + \lambda_2^{(t+1)}}\right)$$

4 Con los pasos 1 a 3 construir $\boldsymbol{\theta}^{(t+1)} = (\theta^{(t+1)}, \lambda_1^{(t+1)}, \lambda_2^{(t+1)})$

5 Repetir los pasos 1 a 4, y generar la cadena $\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)}, \boldsymbol{\theta}^{(3)}, \dots,$

Algoritmo de Gibbs



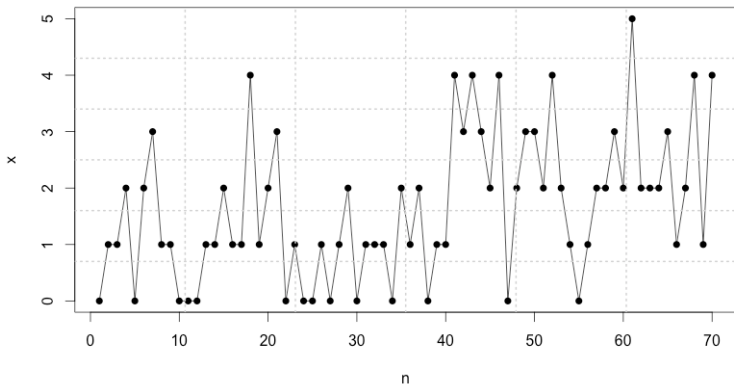
Ver código: [Programa13.r](#)

Algoritmo de Gibbs

Ejemplo 4: Se tienen $\{x_1, \dots, x_n\}$ observaciones de una distribución Poisson en la que existe un punto de cambio en el proceso de observación, digamos en algún $m \in \{1, \dots, n\}$. Es decir, condicional en el valor de m tenemos que:

$$x_i \sim \text{Poi}(x|\lambda), \quad i \in \{1, \dots, m\}; \quad x_i \sim \text{Poi}(x|\phi), \quad i \in \{m+1, \dots, n\}$$

Observaciones Poisson con punto de cambio



Algoritmo de Gibbs

El problema anterior cuenta con tres parámetros a estimar: λ , ϕ y m . Desde el punto de vista Bayesiano debemos asignar probabilidades iniciales:

- $p(\lambda|\alpha, \beta) = Ga(\lambda|\alpha_\lambda, \beta_\lambda)$
- $p(\phi|\alpha, \beta) = Ga(\phi|\alpha_\phi, \beta_\phi)$
- $p(m) = \frac{1}{n}; \quad m \in \{1, 2, \dots, n\}$

La regla de Bayes nos lleva a la siguiente distribución final:

$$p(\lambda, \phi, m|\mathbf{x}) \propto \lambda^{\alpha_\lambda t_m - 1} \phi^{\alpha_\phi + u_m - 1} \exp\{-(\beta_\lambda + m)\lambda\} \exp\{-(\beta_\phi + n - m)\phi\}$$

Donde $t_m = \sum_{i=1}^m x_i$ y $u_m = \sum_{i=m+1}^n x_i$.

En este caso $\theta = (\lambda, \phi, m)$, por lo que si queremos aplicar el algoritmo requeriremos las condicionales completas correspondientes:

$$p(\lambda|\phi, m, \mathbf{x}); \quad p(\phi|\lambda, m, \mathbf{x}); \quad p(m|\lambda, \phi, \mathbf{x})$$

*

Algoritmo de Gibbs

Las densidades condicionales completas tienen una forma fácil de simular:

- $p(\lambda|\phi, m, \mathbf{x}) = Ga(\lambda|\alpha_\lambda + t_m, \beta_\lambda + m)$
- $p(\phi|\lambda, m, \mathbf{x}) = Ga(\phi|\alpha_\phi + u_m, \beta_\phi + n - m)$
- $p(m|\lambda, \phi, \mathbf{x}) = \frac{\lambda^\alpha \lambda^{t_m - 1} \phi^\alpha \phi^{u_m - 1} \exp\{-(\beta_\lambda + m)\lambda\} \exp\{-(\beta_\phi + n - m)\phi\}}{\sum_{l=1}^n \lambda^\alpha \lambda^{t_l - 1} \phi^\alpha \phi^{u_l - 1} \exp\{-(\beta_\lambda + l)\lambda\} \exp\{-(\beta_\phi + n - l)\phi\}}$

*

Algoritmo de Gibbs

El Algoritmo de Gibbs para simular de $p(\lambda, \phi, m|\mathbf{x})$ es el siguiente:

- 0 Inicializar la cadena en un valor inicial $\theta^{(0)} = (\lambda^{(0)}, \phi^{(0)}, m^{(0)})$
 Para $t \in \{0, 1, 2, \dots, \}$
 - 1 Generar una muestra $\lambda^{(t+1)}$ de

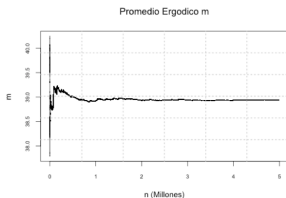
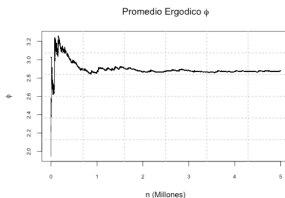
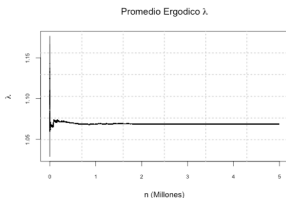
$$p(\lambda^{(t)}|\phi^{(t)}, m^{(t)}, \mathbf{x}) = Ga(\lambda^{(t)}|\alpha_\lambda + t_{m^{(t)}}, \beta_\lambda + m^{(t)})$$
 - 2 Generar una muestra $\phi^{(t+1)}$ de

$$p(\phi^{(t)}|m^{(t)}, \mathbf{x}) = Ga(\phi^{(t)}|\alpha_\phi + u_{m^{(t)}}, \beta_\phi + n - m^{(t)})$$
 - 3 Generar una muestra $m^{(t+1)}$ de $p(m^{(t+1)}|\lambda^{(t+1)}, \phi^{(t+1)}, \mathbf{x})$
 - 4 Con los pasos 1 a 3 construir $\theta^{(t+1)} = (\lambda^{(t+1)}, \phi^{(t+1)}, m^{(t+1)})$
 - 5 Repetir los pasos 1 a 4, y generar la cadena $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}, \dots,$

*

Algoritmo de Gibbs

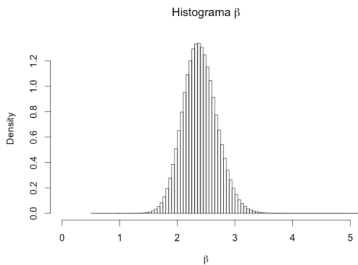
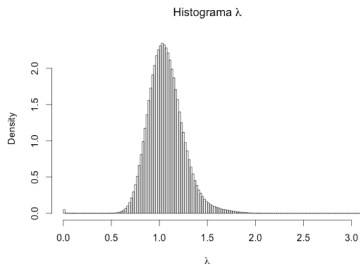
Para ejemplificar: Supongamos $n = 70$, $\alpha_\lambda = \beta_\lambda = \alpha_\phi = \beta_\phi = 0.1$. Número de iteraciones: 5 millones.



Ver código: [Programa14.r](#)

Algoritmo de Gibbs

Distribuciones finales marginales para λ y ϕ



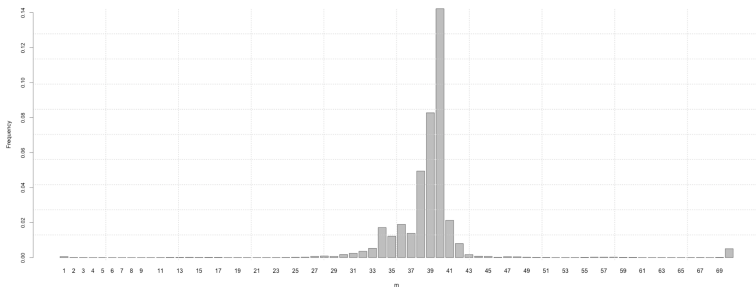
Intervalo de credibilidad de máxima densidad:

$$\mathbb{P}(0.722 < \lambda < 1.438) = 0.95 \quad \mathbb{P}(1.793 < \phi < 3.029) = 0.95$$

Ver código: [Programa14.r](#)

Algoritmo de Gibbs

Distribución final marginal para el punto de cambio:



Intervalo de credibilidad de máxima densidad:

$$\mathbb{P}(33 < m < 42) = 0.9500026$$

$$\text{Moda} = 40$$

Ver código: [Programa14.r](#)

¿Qué es JAGS?

“Just Another Gibbs Sampler”

JAGS (Plummer, 2013)

Es un programa para el análisis de modelos Bayesianos usando Monte Carlo vía Cadenas de Markov

¿Que es JAGS?

JAGS fue escrito para:

- Tener un motor para el lenguaje BUGS que corra en Unix, Mac y Windows
- Ser extendible
- Proporcionar una plataforma para experimentos con ideas de modelación Bayesiana

Corriendo un modelo en JAGS

Para obtener muestras de las distribuciones finales de los parámetros, JAGS realiza 5 pasos:

- 1 Definición del modelo
- 2 Compilación
- 3 Inicialización
- 4 Adaptación y *burn-in*
- 5 Monitoreo

Otras etapas del análisis se realizan fuera de JAGS; por ejemplo, diagnósticos de convergencia.

1. Definición del Modelo

Existen dos partes en la definición del modelo en JAGS: El modelo y los datos.

Descripción del modelo. El modelo se define en un archivo de texto usando el lenguaje BUGS.

Ej. Modelo Gamma (ambos parámetros desconocidos)

```
model.jags <- function() {  
  alpha ~ dgamma(0.01, 0.01)  
  beta ~ dgamma(0.01, 0.01)  
  for (i in 1:n){  
    x[i] ~ dgamma(alpha, beta)  
  }  
}
```

Datos

Los datos pueden ser dados en un archivo por separado o directamente en R.
Por ejemplo:

- En .txt

```
"x"<- c(1, 2, 3, 4, 5)
"Y" <- c(1, 3, 3, 3, 5)
"N" <- 5
```

- En R:

```
x<- c(1, 2, 3, 4, 5)
Y <- c(1, 3, 3, 3, 5)
N <- 5
datos<-list("Y","x","N")
```

cont..

- **2. Compilación:** Verifica si no hay errores de sintaxis

- **3. Inicialización:**
 - El usuario puede fijar los valores iniciales.
 - Si no se especifican los valores iniciales, un “valor típico” es obtenido de la distribución inicial (media, mediana o moda)

- **4. Burn-in**

- **5. Monitoreo:** Un objeto que registra los valores de los parámetros en cada iteración. (p.e. Trace monitor)

Ejemplo 3 : Gamma con 2 parámetros desconocidos.

$$x \sim \text{Gamma}(x|\alpha, \beta); \quad p(\alpha, \beta) = \text{Gamma}(\alpha|\alpha_0, \beta_0)\text{Gamma}(\beta|\alpha_1, \beta_1)$$

Asumiendo que:

$$\alpha_0 = \beta_0 = \alpha_1 = \beta_1 = 0.01$$

y que se observa la muestra de tamaño 5:

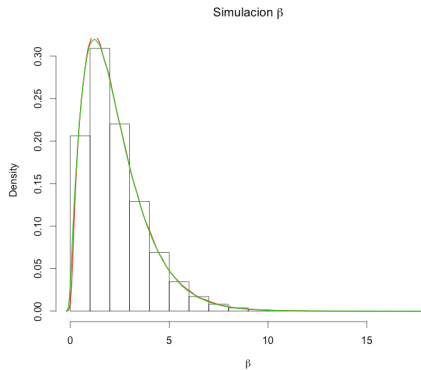
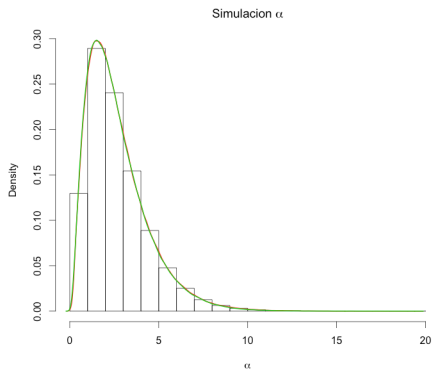
$$(0.4154325, 1.7853782, 1.7315852, 1.0254059, 1.9427045)$$

Haremos inferencias sobre α , β y x_F

```
model.jags <- function() {
  for (i in 1:n){
    sample[i] ~ dgamma(alpha.jags, beta.jags)
  }
  alpha.jags ~ dgamma(0.01, 0.01)
  beta.jags ~ dgamma(0.01, 0.01)
  x.f ~ dgamma(alpha.jags, beta.jags)
}
```

Ver código: [Programa15.r](#)

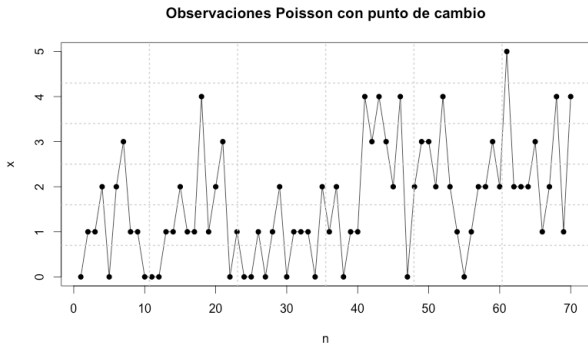
Comparación de las distribuciones finales (MH programado en R vs JAGS)



Ver código: [Programa15.r](#)

Ejemplo 4: Se tienen observaciones $\{x_1, \dots, x_n\}$ de una distribución Poisson en la que existe un punto de cambio, digamos en $m \in \{1, \dots, n\}$. Es decir, condicional en el valor de m tenemos que:

$$x_i \sim \text{Poi}(x|\lambda), \quad i \in \{1, \dots, m\}; \quad x_i \sim \text{Poi}(x|\phi), \quad i \in \{m+1, \dots, n\}$$



Ver código: [Programa16.r](#)

*

Artículos, libros, software



Plummer, Martyn. (2013).

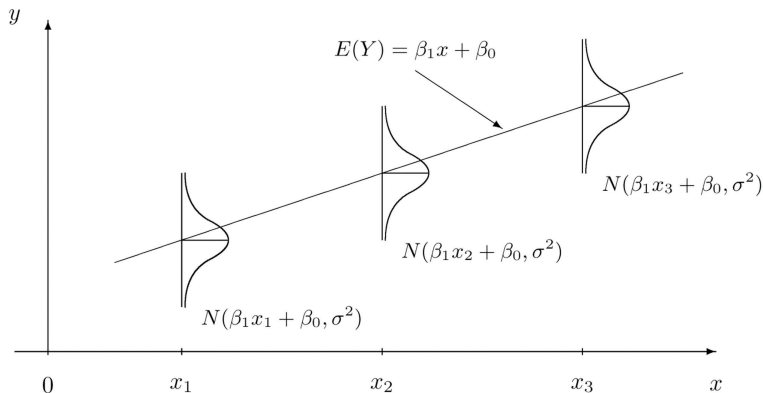
JAGS Version 3.4.0 User Manual



Plummer, Martyn y Northcott, Bill (2013).

JAGS Version 3.4.0 Installation Manual

Modelos Lineales con JAGS



Modelos Lineales Generalizados

Dada una variable respuesta y con un conjunto de covariables \mathbf{z} , surge de manera natural preguntarnos cuál podría ser la relación funcional entre ellas. Una forma de modelarla podría ser:

$$\mathbb{E}(y | \mathbf{z}) = \mu(\mathbf{z})$$

donde, en general, $\mu(\cdot)$ es una función desconocida. En la práctica es común aproximar a $\mu(\cdot)$ a través de una función más simple (paramétrica):

$$\mu(\mathbf{z}) = \psi(\mathbf{z}; \boldsymbol{\beta})$$

donde $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_k)^t$ denota a un vector de parámetros desconocidos.

La forma mas simple para modelar la relación es suponer una función lineal de β , es decir:

$$\psi(\mathbf{z}; \beta) = h(\beta_0 + \beta_1 s_1(\mathbf{z}) + \dots + \beta_k s_k(\mathbf{z}))$$

donde s_i son funciones conocidas.

Finalmente esta función $\psi(\mathbf{z}; \beta)$ es tratada como si fuera la verdadera función que modelará el valor esperado de la variable respuesta y , por lo que el problema se reduce a hacer inferencias sobre el valor del vector de parámetros β .

Es decir:

$$\mathbb{E}(y | \mathbf{z}) = h(\beta_0 + \beta_1 s_1(\mathbf{z}) + \dots + \beta_k s_k(\mathbf{z}))$$

o bien:

$$g(\mathbb{E}(y | \mathbf{z})) = \beta_0 + \beta_1 s_1(\mathbf{z}) + \dots + \beta_k s_k(\mathbf{z})$$

De forma general, si se tienen n observaciones $(y_1, \mathbf{z}_1), \dots, (y_n, \mathbf{z}_n)$, definimos el **modelo lineal generalizado** como:

$$g(\mu_i) = \eta_i$$

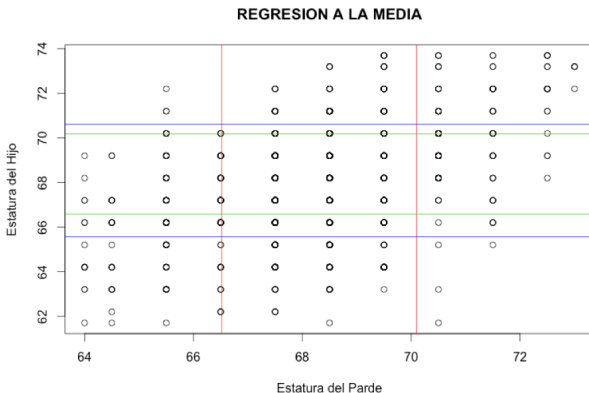
donde

- (y_1, \dots, y_n) se conoce como la **componente aleatoria** y las observaciones se asumen independientes
- $\mu_i = \mathbb{E}(y_i | \mathbf{z}_i)$ es el valor esperado de y_i condicionado en los valores de las covariables \mathbf{z}_i
- $\eta_i = \beta_0 + \beta_1 s_1(\mathbf{z}_i) + \dots + \beta_k s_k(\mathbf{z}_i)$, se conoce como la **componente sistemática** (predictor lineal)
- $g(\cdot)$ función liga (o función vínculo), la cual relaciona a las componentes aleatoria y sistemática

Un caso muy utilizado es cuando g es la función identidad y y sigue una distribución Normal, dando origen al modelo de **Regresión Lineal**

Un poco de historia....

El término **regresión** fue acuñado por Francis Galton en el Siglo XIX en su artículo **Regression towards mediocrity in hereditary stature**, en donde observó que las alturas de los descendientes de ancestros altos tienden a regresar hacia abajo, hacia un promedio normal (un fenómeno conocido como regresión a la media).



Ver: `data(Galton)`, del paquete `HistData`

Supongamos que se tiene n observaciones independientes $(y_1, \mathbf{z}_1), \dots, (y_n, \mathbf{z}_n)$ del modelo:

$$y_i | \mathbf{z}_i \sim N(y_i | \mu(\mathbf{z}_i), \sigma^2) \quad (\sigma^2 > 0, \text{desconocida})$$

donde

$$\mu(\mathbf{z}_i) = \beta_0 + \beta_1 s_1(\mathbf{z}_i) + \dots + \beta_k s_k(\mathbf{z}_i)$$

Definamos:

$$x_{ij} = s_j(\mathbf{z}_i) \quad i \in \{1, \dots, n\}; \quad j \in \{1, \dots, k\}$$

Entonces el modelo lo podemos escribir como:

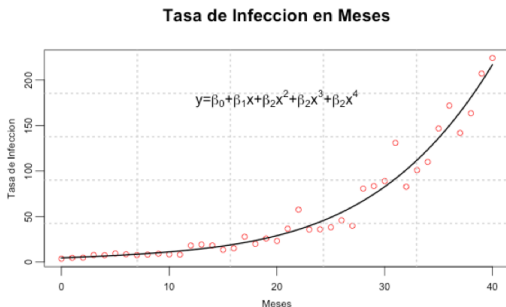
$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \varepsilon_i \quad \varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2) \quad (1)$$

Ejemplo:

Suponiendo que sólo tenemos una covariable $z \in \mathbb{R}$, y haciendo $s_j(z) = z^j$ entonces (1) toma la forma:

$$y_i = \beta_0 + \beta_1 z + \dots + \beta_k z^k + \varepsilon_i \quad \varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2) \quad (2)$$

El modelo anterior pretende modelar el valor de y a través de una función polinomial de la covariable z .



Dada la relación lineal que hemos impuesto resulta conveniente utilizar una notación matricial y escribir (1) como: (haciendo $p = k + 1$)

$$Y = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}; \quad \boldsymbol{\varepsilon} \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I}_n) \Rightarrow Y \sim N_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n)$$

donde

$$Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}_{n \times 1} \quad \mathbf{X} = \begin{pmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1k} \\ 1 & x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \ddots & \vdots & \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nk} \end{pmatrix}_{n \times p} \quad \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}_{n \times 1}$$

y

$$\boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{pmatrix}_{p \times 1}$$

Solución clásica del problema de estimación:

Dado que ahora conocemos la forma de distribución de \mathbf{Y} , podemos encontrar la función de verosimilitud

$$f(\mathbf{Y}; \boldsymbol{\beta}, \sigma^2) = \frac{1}{(2\pi)^{\frac{n}{2}} |\sigma^2 \mathbf{I}_n|^{\frac{1}{2}}} e^{-\frac{1}{2}(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})' (\sigma^2 \mathbf{I}_n)^{-1} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})}$$

Como $|\sigma^2 \mathbf{I}_n| = \sigma^{2n}$ y $(\sigma^2 \mathbf{I}_n)^{-1} = \frac{1}{\sigma^2} \mathbf{I}_n$, entonces la verosimilitud es:

$$\mathcal{L}(\boldsymbol{\beta}, \sigma^2; \mathbf{Y}) = (2\pi\sigma^2)^{-\frac{n}{2}} e^{-\frac{1}{2\sigma^2} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})' (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})}$$

Sacamos el logaritmo de la verosimilitud:

$$\log \mathcal{L}(\boldsymbol{\beta}, \sigma^2; \mathbf{Y}) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})' (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})$$

Maximizamos con respecto a $(\boldsymbol{\beta}, \sigma^2)$; para ello derivamos e igualamos a cero

$$\frac{d}{d\boldsymbol{\beta}} \log \mathcal{L}(\boldsymbol{\beta}, \sigma^2; \mathbf{Y}) = -\frac{1}{2\sigma^2} (-2\mathbf{X}'\mathbf{Y} + 2\mathbf{X}'\mathbf{X}\boldsymbol{\beta}) \quad (3)$$

$$\frac{d}{d\sigma^2} \log \mathcal{L}(\boldsymbol{\beta}, \sigma^2; \mathbf{Y}) = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})' (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \quad (4)$$

*

De la ecuación (3) obtenemos las ecuaciones normales, es decir:

$$\mathbf{X}'\mathbf{X}\boldsymbol{\beta} = \mathbf{X}'\mathbf{Y} \quad (5)$$

Notamos que (5) tiene solución única si y sólo si la matriz $\mathbf{X}'\mathbf{X}$ es invertible (\mathbf{X} de rango completo) en cuyo caso el estimador máximo verosímil para $\boldsymbol{\beta}$ es:

$$\hat{\boldsymbol{\beta}}_{MV} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y} \quad (6)$$

Para σ^2 , de la ecuación (4) obtenemos:

$$-n + \frac{1}{\sigma^2} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})' (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) = 0 \Rightarrow \sigma^2 = \frac{(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})' (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})}{n}$$

Por lo tanto, al sustituir en la última igualdad lo que obtuvimos en la ecuación (6) obtenemos que estimador máximo verosímil para σ^2 es:

$$\hat{\sigma}_{MV}^2 = \frac{(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}_{MV})' (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}_{MV})}{n}$$

Definiendo $\hat{\mathbf{Y}} := \mathbf{X}\hat{\boldsymbol{\beta}}_{MV}$, tenemos:

$$\hat{\sigma}_{MV}^2 = \frac{(\mathbf{Y} - \hat{\mathbf{Y}})' (\mathbf{Y} - \hat{\mathbf{Y}})}{n} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Los estimadores máximo verosímiles gozan de las siguientes propiedades:

- $\mathbb{E}(\hat{\beta}_{MV}) = \beta$ (insesgamiento)
- $\text{Var}(\hat{\beta}_{MV}) = \sigma^2 (\mathbf{X}'\mathbf{X})^{-1}$
- $\hat{\beta}_{MV} \sim N_p(\beta, \sigma^2 (\mathbf{X}'\mathbf{X})^{-1})$
- $\mathbb{E}(\hat{\sigma}_{MV}^2) = \sigma^2 \frac{n-p}{n}$ (sesgado)
- $\text{Var}(\hat{\sigma}_{MV}^2) = 2 \frac{n-p}{n^2} \sigma^4$
- $\hat{\sigma}_{MV}^2 \sim \text{Gamma}(\frac{n-p}{2}, \frac{n}{2\sigma^2})$
- $\hat{\sigma}^2 = \frac{1}{n-p} \sum_{i=1}^n (y_i - \hat{y}_i)^2$ (insesgamiento)
- Como $\hat{\beta}_{MV} \sim N_p(\beta, \sigma^2 (\mathbf{X}'\mathbf{X})^{-1})$, entonces haciendo $\mathbf{C} = (\mathbf{X}'\mathbf{X})^{-1}$ y definiendo C_{ij} al elemento (i, j) de la matriz \mathbf{C} se tiene que:

$$\hat{\beta}_i \sim N\left(\beta_i, \sigma^2 C_{(i+1)(i+1)}\right) \Rightarrow \frac{\hat{\beta}_i - \beta_i}{\sqrt{\hat{\sigma}^2 C_{(i+1)(i+1)}}} \sim t_{(n-p)} \quad (7)$$

De la última expresión de (7), la inferencia clásica desprende las pruebas de hipótesis e intervalos de confianza correspondientes para el parámetro β_i con $i \in \{0, \dots, k\}$

Regresión Lineal : Enfoque Bayesiano

Sea $\beta \in \mathbb{R}^p$ un vector de parámetros, y $\mathbf{X} \in \mathbb{R}^{n \times p}$ una matriz de diseño conocida. Definamos el modelo lineal:

$$\mathbf{y} = \mathbf{X}\beta + \varepsilon \quad (8)$$

donde $\varepsilon \sim N_n(\underline{\varepsilon}|\mathbf{0}, \tau\mathbf{I}_n)$ con $\tau = \sigma^{-2} > 0$. Suponemos una matriz de precisión con correlación 0, entre la variables lo que implica independencia entre las ε_i . De la ecuación (8) concluimos que:

$$\mathbf{y} \sim N_n(\mathbf{y}|\mathbf{X}\beta, \tau\mathbf{I}_n)$$

Objetivo: suponiendo que observamos \mathbf{y} , **inferir** sobre los parámetros β y $\tau = \sigma^{-2}$

Solución: (La receta) Encontrar las distribuciones finales:

$$p(\beta, \tau|\mathbf{y}) \propto p(\mathbf{y}|\beta, \tau) p(\beta, \tau)$$

$$p(\beta|\mathbf{y}) = \int_0^\infty p(\beta, \tau|\mathbf{y}) d\tau; \quad p(\tau|\mathbf{y}) = \int_{\mathbb{R}^p} p(\beta, \tau|\mathbf{y}) d\beta;$$

Bajo los supuestos que tiene este modelo, es posible construir distribuciones conjugadas que permiten encontrar distribuciones finales exactas que no requieren uso de herramientas de simulación para su estudio.

Para proponer una distribución conjugada, se estudia la verosimilitud.

$$p(\mathbf{y}|\boldsymbol{\beta}, \tau) \propto \tau^{\frac{n}{2}} e^{-\frac{\tau}{2} \left((\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) + (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) \right)} \quad (9)$$

$$\propto \tau^{\frac{n}{2}} e^{-\frac{\tau}{2} \left((\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) + \tilde{\beta}_0 \right)} \quad (10)$$

Vista como función de $(\boldsymbol{\beta}, \tau)$ aparece el denominado **kernel** de una distribución **Normal-Multivariada - Gamma !!**

*

Regresión Lineal : Enfoque Bayesiano

En la literatura se propone como inicial una distribución Normal-Multivariada - Gamma

$$p(\boldsymbol{\beta}, \tau | \boldsymbol{\mu}_0, \mathbf{P}_0, \alpha_0, \delta_0) \propto N_p(\boldsymbol{\beta} | \boldsymbol{\mu}_0, \tau \mathbf{P}_0) Ga(\tau | \alpha_0, \delta_0)$$

donde $\mathbf{P}_0 \in \mathbb{R}^{p \times p}$, $\boldsymbol{\mu}_0 \in \mathbb{R}^p$, $\alpha_0, \delta_0 \in \mathbb{R}$ son hiperparámetros.

De las propiedades de esta densidad se obtiene que:

$$\begin{aligned} \boldsymbol{\beta} &\sim T_p \left(\boldsymbol{\beta} \mid 2\alpha_0, \boldsymbol{\mu}_0, \frac{\alpha_0}{\delta_0} \mathbf{P}_0 \right) \Rightarrow \mathbb{E}(\boldsymbol{\beta}) = \boldsymbol{\mu}_0; \quad \text{Var}(\boldsymbol{\beta}) = \frac{\delta_0}{\alpha_0 - 1} \mathbf{P}_0^{-1} \\ \tau &\sim Ga(\tau | \alpha_0, \delta_0) \Rightarrow \mathbb{E}(\tau) = \frac{\alpha_0}{\delta_0}; \quad \text{Var}(\tau) = \frac{\alpha_0}{\delta_0^2} \end{aligned}$$

*

Regresión Lineal : Enfoque Bayesiano

Definida la distribución inicial tenemos todos los ingredientes para obtener la final, usando nuestra receta!!

Tras un poco de algebra se obtiene que:

$$p(\boldsymbol{\beta}, \tau | \mathbf{y}) = N_p(\boldsymbol{\beta} | \boldsymbol{\mu}_1, \tau \mathbf{P}_1) Ga(\tau | \alpha_1, \delta_1)$$

donde:

- $\mathbf{P}_1 = \mathbf{X}^T \mathbf{X} + \mathbf{P}_0$
- $\boldsymbol{\mu}_1 = (\mathbf{X}^T \mathbf{X} + \mathbf{P}_0)^{-1} (\mathbf{X}^T \underline{\mathbf{y}} + \mathbf{P}_0 \boldsymbol{\mu}_0)$
- $\alpha_1 = \frac{n}{2} + \alpha_0$
- $\delta_1 = \delta_0 + \frac{1}{2} \left((\underline{\mathbf{y}} - \mathbf{X} \boldsymbol{\mu}_1)^T (\underline{\mathbf{y}} - \mathbf{X} \boldsymbol{\mu}_1) + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)^T \mathbf{P}_0 (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) \right)$

Finalmente, sabemos que, por propiedades de esta distribución,

$$p(\boldsymbol{\beta} | \mathbf{y}) = T_p \left(\boldsymbol{\beta} \mid 2\alpha_1, \boldsymbol{\mu}_1, \frac{\alpha_1}{\delta_1} \mathbf{P}_1 \right)$$

$$p(\tau | \mathbf{y}) = Ga(\alpha_1, \delta_1) \Rightarrow p(\tau^{-1} | \underline{\mathbf{y}}) = IGa(\alpha_1, \delta_1)$$

*

Regresión Lineal: Inferencia sobre y^*

Si se desea hacer inferencia sobre nuevas observaciones del modelo, también se obtiene fórmulas cerradas:

$$\mathbf{w} = \mathbf{Z}\underline{\beta} + \mathbf{e}$$

donde

- $\mathbf{Z} \in \mathbb{R}^{k \times p}$ es una nueva matriz de covariables.
- $\mathbf{e} \sim N_k(\mathbf{e}|\mathbf{0}, \tau\mathbf{I})$
- $\mathbf{w} \sim N_k(\mathbf{w}|\mathbf{Z}\underline{\beta}, \tau\mathbf{I})$

Desde el punto de vista Bayesiano, el objetivo es determinar $p(\mathbf{w}|\mathbf{y})$

$$\begin{aligned} p(\mathbf{w}|\mathbf{y}) &= \int_0^\infty \int_{\mathbb{R}^p} p(\mathbf{w}, \underline{\beta}, \tau|\mathbf{y}) d\underline{\beta} d\tau \\ &= \dots \\ &= T_k \left(\mathbf{w} \mid 2\alpha_1, \mathbf{Z}\mu_1, \frac{\alpha_1}{\delta_1} \left(\mathbf{I} + \mathbf{ZP}_1\mathbf{Z}^T \right)^{-1} \right) \end{aligned}$$

*

Regresión Lineal: Distribuciones no Informativas

En estos casos es posible, además, construir distribuciones no informativas.

- **Jeffreys:**

$$p(\boldsymbol{\beta}, \tau) \propto \tau^{\frac{p-2}{2}}$$

La cual da origen a la final

$$p(\boldsymbol{\beta}, \tau | \underline{y}) = NG\left(\boldsymbol{\beta}, \tau \left| \hat{\boldsymbol{\beta}}, \mathbf{X}^T \mathbf{X}, \frac{n}{2}, \frac{1}{2} (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) \right.\right)$$

- Distribución de **Referencia:**

$$p(\boldsymbol{\beta}, \tau) \propto \tau^{-1}$$

La cual da origen a la final

$$p(\boldsymbol{\beta}, \tau | \underline{y}) = NG\left(\boldsymbol{\beta}, \tau \left| \hat{\boldsymbol{\beta}}, \mathbf{X}^T \mathbf{X}, \frac{n-p}{2}, \frac{1}{2} (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) \right.\right)$$

*

Regresión Lineal con JAGS

Para ilustrar el uso de la simulación en este modelo, supongamos que se propone una inicial de la siguiente forma:

$$p(\boldsymbol{\beta}, \tau | \boldsymbol{\mu}_0, \alpha_0, \delta_0) = N_p(\boldsymbol{\beta} | \boldsymbol{\mu}_0, \mathbf{I}) Ga(\tau | \alpha_0, \delta_0)$$

(Asume independencia entre los parámetros de forma inicial)

En este caso no hay conjugación con la verosimilitud, por lo que no necesariamente se puede llegar a formas cerradas para la distribución final

$$p(\boldsymbol{\beta}, \tau | \mathbf{y}) \propto p(\mathbf{y} | \boldsymbol{\beta}, \tau) N_p(\boldsymbol{\beta} | \boldsymbol{\mu}_0, \mathbf{I}) Ga(\tau | \alpha_0, \delta_0)$$

El objetivo es, entonces, construir una cadena de Markov cuya distribución final sea precisamente $p(\boldsymbol{\beta}, \tau | \mathbf{y})$.

¿Cómo describimos un modelo lineal simple en JAGS?

El modelo lineal simple es:

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i; \quad \varepsilon_i \sim N(\varepsilon_i | 0, \tau); \quad y_i \sim N(y_i | \beta_0 + \beta_1 x_i, \tau)$$

Asumiendo las iniciales

$$\tau \sim Ga(\tau | 0.01, 0.01); \quad \beta_0 \sim N(\beta_0 | 0, 0.001); \quad \beta_1 \sim N(\beta_1 | 0, 0.001).$$

(Obs: Asumimos independencia de forma inicial entre los parámetros y una precisión pequeña, lo que refleja poca información inicial)

Asumiendo que tenemos los vectores x , y de longitud n , entonces el modelo es:

```
model.jags <- function() {
  beta0 ~ dnorm(0, 0.001)
  beta1 ~ dnorm(0, 0.001)
  tau ~ dgamma(0.01, 0.01)
  for (i in 1:n){
    y[i] ~ dnorm(beta0+beta1*x[i], tau)
  }
}
```

¿Cómo describimos un modelo lineal múltiple en JAGS?

El modelo lineal múltiple es:

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \varepsilon_i \quad \varepsilon_i \sim N(\varepsilon_i | 0, \tau)$$

En este caso:

$$\tau \sim Ga(\tau | 0.01, 0.01); \quad \beta_i \sim N(\beta_i | 0, 0.001); \quad i \in \{0, 1, \dots, k\}$$

Asumiendo que tenemos la matriz de diseño X , y el vector y , entonces el modelo es:

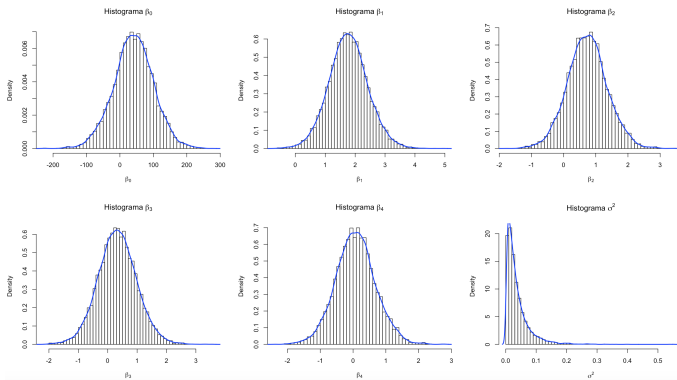
```
model.jags <- function() {
  for (i in 1:(k+1)){
    beta[i] ~ dnorm(0, 0.001)
  }
  tau ~ dgamma(0.01, 0.01)
  for (i in 1:n){
    y[i] ~ dnorm(X[i,]%*%beta, tau)
  }
}
```

Ejemplo: El siguiente conjunto de datos (Hald, 1952), también descrito en Draper and Smith (1981), consiste de 13 observaciones que relacionan el calor producido por el endurecimiento de cierto tipo de cemento con cuatro variables explicativas, cada una midiendo el contenido de un ingrediente en particular (en porcentajes).

Y	x_1	x_2	x_3	x_4
78.5	7	26	6	60
74.3	1	29	15	52
104.3	11	56	8	20
87.6	11	31	8	47
95.9	7	52	6	33
109.2	11	55	9	22
102.7	3	71	17	6
72.5	1	31	22	44
93.1	2	54	18	22
115.9	21	47	4	26
83.8	1	40	23	34
113.3	11	66	9	12
109.4	10	68	8	12

Ver código: [Programa17.r](#)

Se generó una cadena durante 20,000 iteraciones con un calentamiento de 10,000. Las densidades finales de cada uno de los parámetros son:



Ver código: [Programa17.r](#)

Modelo Logístico

El modelo **logístico** asume lo siguiente:

- (y_1, \dots, y_n) siguen una distribución Bernoulli.
- $\mu_i = \mathbb{E}(y_i | \mathbf{z}_i) = p_i \in (0, 1)$
- $\eta_i = \beta_0 + \beta_1 s_1(\mathbf{z}_i) + \dots + \beta_k s_k(\mathbf{z}_i)$ la **componente sistemática** (predictor lineal)
- $g(p_i) = \log\left(\frac{p_i}{1-p_i}\right)$, conocida como función $\text{logit}(p_i)$

Obs: En este caso se modela la probabilidad de éxito de y_i por medio de:

$$p_i = \frac{\exp(\eta_i)}{1 + \exp(\eta_i)}$$

¿Cómo describimos un modelo logístico en JAGS?

El modelo asume

$$y_i \sim \text{Bernoulli}(y_i|p_i) = \text{Bernoulli}\left(y_i \left| \frac{\exp(\eta_i)}{1 + \exp(\eta_i)} \right.\right)$$

Supongamos las iniciales:

$$\beta_i \sim N(\beta_i|0, 0.001); \quad i \in \{0, 1, \dots, k\}$$

Asumiendo que tenemos la matriz de diseño X , y el vector y , entonces el modelo es:

```
model.jags <- function() {
  for (i in 1:(k+1)){
    beta[i] ~ dnorm(0, 0.001)
  }
  for (i in 1:n){
    logit(p[i]) <- X[i,]%*% beta
    y[i] ~ dbin(p[i], 1)
  }
}
```

Cuando en este modelo tenemos repeticiones de las observaciones (*Bernoulli*) para cierto nivel de las covariables z , se utiliza la distribución Binomial en el componente aleatorio.

Ejemplo 2. (Gelman *et al.*, 1995) En el desarrollo de drogas y otros compuestos químicos, es normal realizar pruebas de toxicidad (experimentos de bioensayo) en grupos de animales para evaluar la conveniencia de utilizar las drogas de manera comercial. Los experimentos consisten en administrar diferentes dosis (niveles) de la droga o compuesto a diferentes grupos de animales. En general, la respuesta es dicotómica; por ejemplo, el animal vive o muere, tiene tumor o no lo tiene, o se contagia o no se contagia. Así, la información que se recolecta es de la forma

$$D = \{(x_i, n_i, y_i) : i \in J_k\},$$

en donde x_i representa el i -ésimo de k niveles de la droga que se administró a un grupo de n_i animales de los que y_i responden negativamente. La siguiente tabla muestra un conjunto de resultados de un experimento ficticio.

Dosis, x_i (log g/ml)	Número de animales, n_i	Número de muertes, y_i
-0.863	5	0
-0.296	5	1
-0.053	5	3
0.727	5	5

Para el problema anterior planteamos el modelo logístico de la siguiente forma:

$$\text{logit}(p_i) = \beta_0 + \beta_1 x_i$$

de donde se puede despejar p_i de tal manera que la probabilidad de que el animal responda negativamente (muera) con una dosis igual a x_i es:

$$p_i = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_i)}} = \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}}$$

Un objetivo importante en estos estudios es determinar la dosis mediana (LD50) que es la dosis en la que la probabilidad de muerte es precisamente 0.5. En este caso, se tiene que resolver la ecuación:

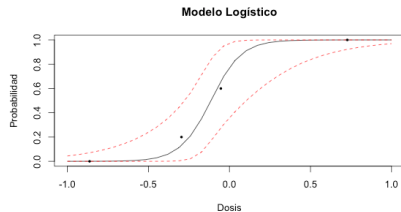
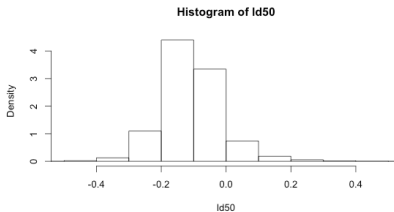
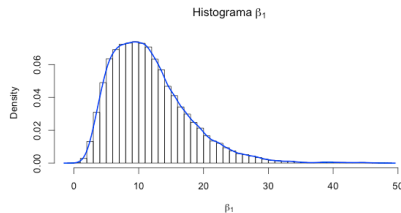
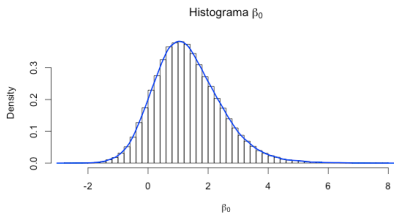
$$0.5 = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_i)}} = \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}}$$

de donde se obtiene que:

$$LD50 = -\frac{\beta_0}{\beta_1}.$$

Es entonces de interés hacer inferencias sobre esta cantidad en este tipo de problemas.

Ver código: [Programa18.r](#)



Ver código: [Programa18.r](#)

Regresión Poisson

El modelo lineal Poisson asume lo siguiente:

- (y_1, \dots, y_n) siguen una distribución Poisson.
- $\mu_i = \mathbb{E}(y_i | \mathbf{z}_i) = \lambda_i \in \mathbb{R}^+$
- $\eta_i = \beta_0 + \beta_1 s_1(\mathbf{z}_i) + \dots + \beta_k s_k(\mathbf{z}_i)$ la **componente sistemática** (predictor lineal)
- $g(\lambda_i) = \log(\lambda_i)$.

Obs: En este caso se modela la tasa λ_i de la variable y_i como

$$\lambda_i = \exp(\eta_i)$$

¿Cómo describimos un modelo Poisson en JAGS?

El modelo asume

$$y_i \sim \text{Poisson}(y_i | \lambda_i) = \text{Poisson}(y_i | \exp(\eta_i))$$

Supongamos las siguientes iniciales en los parámetros del predictor lineal

$$\beta_i \sim N(\beta_i | 0, 0.001); \quad i \in \{0, 1, \dots, k\}$$

Asumiendo que tenemos la matriz de diseño X , y el vector y , entonces el modelo es:

```
model.jags <- function() {  
  for (i in 1:(k+1)){  
    beta[i] ~ dnorm(0, 0.001)  
  }  
  for (i in 1:n){  
    log(lambda[i]) <- X[i,]%*%beta  
    y[i] ~ dpois(lambda[i])  
  }  
}
```

Ver código: [Programa19.r](#)

*

Regresión Gamma

El modelo Gamma asume lo siguiente:

- (y_1, \dots, y_n) siguen una distribución $Gamma(\alpha, \delta_i)$
- $\mu_i = \mathbb{E}(y_i | \mathbf{z}_i) = \frac{\alpha}{\delta_i} = \mu_i \in \mathbb{R}^+, \Rightarrow \delta_i = \frac{\alpha}{\mu_i}$
- $\eta_i = \beta_0 + \beta_1 s_1(\mathbf{z}_i) + \dots + \beta_k s_k(\mathbf{z}_i)$ la **componente sistemática** (predictor lineal)
- $g(\mu_i) = \log(\mu_i)$.

Obs: En este caso se modela la media μ_i de la variable y_i como

$$\mu_i = \exp(\eta_i)$$

¿Cómo describimos un modelo Gamma en JAGS?

El modelo asume

$$y_i \sim \text{Gamma}(y_i | \alpha, \delta_i) = \text{Gamma}\left(y_i \mid \alpha, \frac{\alpha}{\exp(\eta_i)}\right)$$

Asumiendo las siguientes iniciales en los parámetros del predictor lineal y para el parámetro α

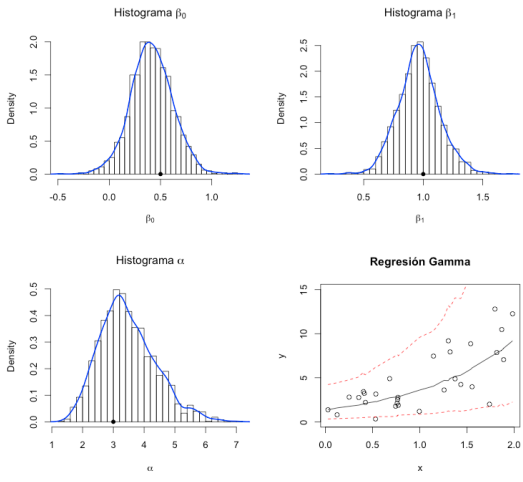
$$\beta_i \sim N(\beta_i | 0, 0.001); \quad i \in \{0, 1, \dots, k\}; \quad \alpha \sim \text{Ga}(\alpha | 0.01, 0.01)$$

```

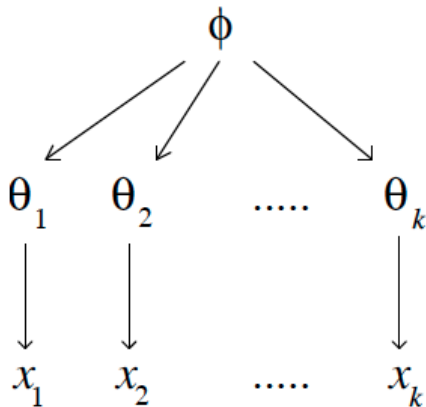
model.jags <- function() {
  for (i in 1:(k+1)){
    beta[i] ~ dnorm(0, 0.001)
  }
  alpha ~ dgamma(0.01, 0.01)
  for (i in 1:n){
    log(mu[i]) <- X[i,] %*% beta
    y[i] ~ dgamma(alpha, alpha/mu[i])
  }
}

```

Ejemplo de una regresión Gamma

Ver código: [Programa20.r](#)

Modelos Jerárquicos Lineales

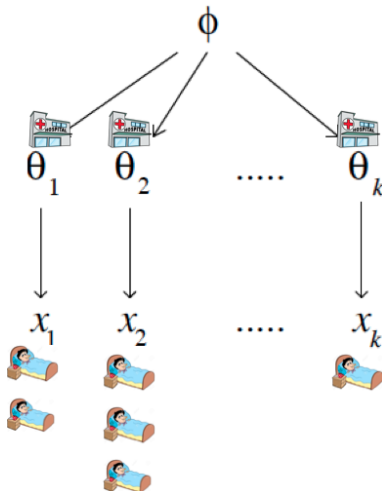


Modelos Jerárquicos Lineales: Motivación

Consideremos el siguiente problema:

- Se tiene una muestra de k hospitales
- Dentro de cada hospital, se tiene una muestra de n_i pacientes ($i = 1, \dots, k$).
- Para el Hospital i , los pacientes tiene una probabilidad de supervivencia (ante cierto padecimiento) de θ_i , de tal forma que condicionado a θ_i se tiene que $x_{ij} \sim \text{Bernoulli}(\theta_i)$.
- Si los hospitales fuera independientes, bastaría inferir el parámetro θ_i en cada hospital
- Si entre los hospitales suponemos **intercambiabilidad** (homogeneidad; por ejemplo hospitales de cierta Región del país), entonces se puede considerar que estos hospitales son una muestra de una distribución poblacional común que a su vez depende de **hiperparámetros** desconocidos ϕ .
- Con esta estructura **jerárquica**, suponiendo que se observan las muestras en cada hospital, resulta interesante entonces hacer inferencias sobre θ_i con $i = 1, \dots, k$
- Si además a esta estructura se agregan covariables z , se pueden enriquecer el modelo.

Modelos Jerárquicos Lineales



Modelos Jerárquicos Lineales

Un Modelo Jerárquico tienen la siguiente estructura:

1 Nivel I, las observaciones:

$$\begin{aligned} p(\mathbf{x}|\boldsymbol{\theta}) &= p(\mathbf{x}_1, \dots, \mathbf{x}_k | \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_k) \\ &= \prod_{i=1}^k p(\mathbf{x}_i | \boldsymbol{\theta}_i) = \prod_{i=1}^k \prod_{j=1}^{n_i} p(x_{ij} | \boldsymbol{\theta}_i) \end{aligned}$$

2 Nivel II, los parámetros:

$$\begin{aligned} p(\mathbf{x}|\boldsymbol{\theta}) &= p(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_k | \phi) \\ &= \prod_{i=1}^k p(\boldsymbol{\theta}_i | \phi) \end{aligned}$$

3 Nivel III, los hiperparámetros:

$$p(\phi)$$

Modelos Jerárquicos Lineales

La interpretación del modelo puede ser la siguiente:

- Las observaciones $\mathbf{x}_1, \dots, \mathbf{x}_k$, con $\mathbf{x}_i = (x_{i1}, \dots, x_{in_i})$, provienen de experimentos distintos pero relacionados entre sí (Ej. Experimentos realizados en k centros de investigación).
- Los parámetros $\theta_1, \dots, \theta_k$ se suponen relacionados (intercambiables, homogéneos) (Ej. θ_i puede representar la probabilidad de supervivencia en el centro de investigación i)
- Los parámetros ϕ describen alguna característica relevante de la población (Ej. $g(\phi)$ con $g: \mathbb{R}^d \rightarrow \mathbb{R}$ puede representar la probabilidad de supervivencia **global** para toda la población de cierta región del país).
- En caso de existir información adicional, por ejemplo algunas característica del paciente como edad, peso, estatura, entonces los datos vienen dados por

$$\{(\mathbf{x}_1, \mathbf{z}_1), (\mathbf{x}_2, \mathbf{z}_2), \dots, (\mathbf{x}_k, \mathbf{z}_k)\}$$

Modelos Jerárquicos Lineales

En el enfoque Bayesiano, estamos interesados en hacer inferencia sobre los parámetros,

$$(\boldsymbol{\theta}, \boldsymbol{\phi}) = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_k, \phi_1, \dots, \phi_d)$$

así como también para posibles futuras observaciones, digamos

$$x_{iF} \sim p(x|\boldsymbol{\theta}_i)$$

(una observación futura del i -ésimo centro de investigación).

Sin embargo, dada la estructura **jerárquica** también es plausible pensar en una observación futura x_F^* correspondiente a una futura θ_* que proviene de la misma población que generó a los parámetros θ_j existentes.

Modelos Jerárquicos Lineales

Dada la jerarquía que existe, es apropiado pensar en distribuciones iniciales de la siguiente forma:

$$p(\boldsymbol{\theta}, \phi) = p(\boldsymbol{\theta}|\phi)p(\phi)$$

Por otro lado, la distribución final correspondiente es

$$\begin{aligned} p(\boldsymbol{\theta}, \phi|\mathbf{x}) &= \frac{p(\boldsymbol{\theta}, \phi, \mathbf{x})}{p(\mathbf{x})} \\ &\propto p(\mathbf{x}|\boldsymbol{\theta}, \phi)p(\boldsymbol{\theta}, \phi) = p(\mathbf{x}|\boldsymbol{\theta})p(\boldsymbol{\theta}, \phi) \end{aligned}$$

de donde las marginales correspondientes para hacer inferencias son:

$$\begin{aligned} p(\boldsymbol{\theta}|\mathbf{x}) &\propto \int p(\mathbf{x}|\boldsymbol{\theta})p(\boldsymbol{\theta}, \phi) d\phi \\ &\propto \int p(\mathbf{x}|\boldsymbol{\theta})p(\boldsymbol{\theta}|\phi)p(\phi) d\phi \\ p(\phi|\mathbf{x}) &\propto p(\phi) \int p(\mathbf{x}|\boldsymbol{\theta})p(\boldsymbol{\theta}|\phi) d\boldsymbol{\theta} \end{aligned}$$

Modelos Jerárquicos Lineales

Existen casos donde se puede hacer inferencia de forma analítica, sin embargo generalmente tendremos que recurrir a aspectos computacionales para obtener aproximaciones a las densidades finales.

Ejemplo:

- Nivel I (Observaciones):

$$p(\mathbf{y}|\boldsymbol{\beta}) = N_n(\mathbf{y}|\mathbf{X}\boldsymbol{\beta}, \boldsymbol{\Sigma}_y); \quad \mathbf{y} \in \mathbb{R}^n, \quad \mathbf{X} \in \mathbb{R}^{n \times p}; \quad \boldsymbol{\beta} \in \mathbb{R}^p; \quad \boldsymbol{\Sigma}_y \in \mathbb{R}^{n \times n}$$

- Nivel II (Parámetros):

$$p(\boldsymbol{\beta}|\boldsymbol{\alpha}) = N_p(\boldsymbol{\beta}|\mathbf{H}\boldsymbol{\alpha}, \boldsymbol{\Sigma}_\beta); \quad \boldsymbol{\beta} \in \mathbb{R}^p; \quad \mathbf{H} \in \mathbb{R}^{p \times q}; \quad \boldsymbol{\alpha} \in \mathbb{R}^q; \quad \boldsymbol{\Sigma}_\beta \in \mathbb{R}^{p \times p}$$

- Nivel III (Hiperparámetros):

$$p(\boldsymbol{\alpha}) = N_q(\boldsymbol{\alpha}|\boldsymbol{\alpha}_0, \boldsymbol{\Sigma}_\alpha); \quad \boldsymbol{\alpha} \in \mathbb{R}^p; \quad \boldsymbol{\alpha}_0 \in \mathbb{R}^q; \quad \boldsymbol{\Sigma}_\alpha \in \mathbb{R}^{q \times q}$$

En este modelo, suponiendo Σ_y y Σ_β conocidas, es posible demostrar que las finales toman la siguiente forma:

$$p(\alpha|y) = N_q \left(\alpha \mid \mu_{\alpha|y}, V_{\alpha|y} \right)$$

donde

$$\mu_{\alpha|y} = \alpha_0 + \Sigma_\alpha H' X' V_y^{-1} (y - XH\alpha_0)$$

y

$$V_{\alpha|y} = \Sigma_\alpha - \Sigma_\alpha H' X' V_y^{-1} XH \Sigma_\alpha$$

con

$$V_y = X \left(\Sigma_\beta + H \Sigma_\alpha H' \right) X' + \Sigma_y$$

Mientras que

$$p(\beta|y) = N_p \left(\beta \mid \mu_{\beta|y}, V_{\beta|y} \right)$$

donde:

$$\mu_{\beta|y} = H \mu_{\alpha|y} + \Sigma_\beta X' V_y^{-1} (y - XH\alpha_0)$$

y

$$V_{\beta|y} = V_\beta - V_\beta X' V_y^{-1} X V_\beta$$

Ejemplo 5: (Simulación) Suponga el siguiente modelo jerárquico para modelar tiempos de fallas de cierto componente de aviones de una compañía determinada.

- Nivel I (Observaciones) (9 grupos, con 20 observaciones cada uno)

$$p(\mathbf{x}_i | \alpha_i, \beta_i) = Ga(\mathbf{x}_i | \alpha_i, \beta_i); \quad i = 1, 2, \dots, 9; \quad \mathbf{x}_i = (x_{i1}, \dots, x_{i20})$$

- Nivel II (Parámetros):

$$p(\alpha_i, \beta_i | \alpha_1^{(0)}, \beta_1^{(0)}, \alpha_2^{(0)}, \beta_2^{(0)}) = Ga(\alpha_i | \alpha_1^{(0)}, \beta_1^{(0)}) Ga(\beta_i | \alpha_2^{(0)}, \beta_2^{(0)});$$

- Nivel III (Hiperparámetros):

$$p(\alpha_1^{(0)}, \beta_1^{(0)}, \alpha_2^{(0)}, \beta_2^{(0)} | \boldsymbol{\theta}_0) = Ga(\alpha_1^{(0)} | 0.01, 0.01) Ga(\beta_1^{(0)} | 0.01, 0.01) \\ Ga(\alpha_2^{(0)} | 0.01, 0.01) Ga(\beta_2^{(0)} | 0.01, 0.01)$$

Obs: En este caso, colocar iniciales impropias no es adecuado pues genera distribuciones finales que no necesariamente son propias. Además, la estimación de los hiperparámetros no es muy precisa si se tienen pocos grupos.

El código de JAGS para este modelo es:

```

model.jags <- function() {
  #Nivel 3; hiperparámetros (iniciales informativas)
  alpha.3[1] ~ dgamma(100,1)
  alpha.3[2] ~ dgamma(100,1)
  beta.3[1] ~ dgamma(1,1)
  beta.3[2] ~ dgamma(1,1)

  for (i in 1:k){
    #Nivel 2; Parametros
    alpha.2[i] ~ dgamma(alpha.3[1],beta.3[1])
    beta.2[i] ~ dgamma(alpha.3[2],beta.3[2])
    for(j in 1:(n[i])){
      #Nivel 1; Las Observaciones
      X[j,i] ~ dgamma(alpha.2[i],beta.2[i])
    }
  }
  #Nueva observacion del grupo 1
  x.1 ~ dgamma(alpha.2[1],beta.2[1])
  #Nueva observacion de un nuevo grupo
  alpha.2.n ~ dgamma(alpha.3[1],beta.3[1])
  beta.2.n ~ dgamma(alpha.3[2],beta.3[2])
  x.n ~ dgamma(alpha.2.n,beta.2.n)
}

```

Ver código: [Programa21.r](#)

Resultados : (Basados en 100,000 simulaciones, con un calentamiento de 10,000)

	Reales	Media	Mediana	L0.025	U0.975
alpha.2[1]	93.35	106.07	105.35	79.88	136.07
alpha.2[2]	113.21	117.56	116.83	88.71	150.02
alpha.2[3]	112.60	116.85	116.32	88.16	149.39
alpha.2[4]	103.68	108.39	107.77	82.04	137.28
alpha.2[5]	84.73	101.02	100.51	76.50	128.46
alpha.2[6]	104.42	103.60	103.22	78.62	131.20
alpha.2[7]	107.00	106.01	105.44	80.22	134.47
alpha.2[8]	105.33	101.73	101.16	77.20	129.09
alpha.2[9]	96.48	95.80	95.41	72.85	121.26
alpha.3[1]	100.00	100.16	99.88	81.85	120.16
alpha.3[2]	100.00	101.11	100.82	82.39	121.84
beta.2[1]	91.69	106.09	105.40	79.82	136.11
beta.2[2]	93.40	96.58	96.02	72.93	123.16
beta.2[3]	96.54	97.66	97.21	73.73	124.82
beta.2[4]	99.05	103.04	102.49	77.95	130.43
beta.2[5]	90.80	110.78	110.17	83.71	141.09
beta.2[6]	107.86	106.12	105.73	80.47	134.45
beta.2[7]	105.51	103.05	102.50	77.92	130.78
beta.2[8]	108.88	106.55	105.95	80.86	135.32
beta.2[9]	107.45	110.13	109.66	83.65	139.44
beta.3[1]	1.00	0.95	0.94	0.70	1.28
beta.3[2]	1.00	0.98	0.97	0.72	1.32

Ver código: [Programa21.r](#)

Ejemplo 6: Gelfand, Hills, Racine-Poon y Smith (1990) discuten el análisis bayesiano de la siguiente tabla:

Tabla 2: Crecimiento de un grupo de ratas (controles)

Rata	x_1	x_2	x_3	x_4	x_5	Rata	x_1	x_2	x_3	x_4	x_5
1	151	199	246	283	320	16	160	207	248	288	324
2	145	199	249	293	354	17	142	187	234	280	316
3	147	214	263	312	328	18	156	203	243	283	317
4	155	200	237	272	297	19	157	212	259	307	336
5	135	188	230	280	323	20	152	203	246	286	321
6	159	210	252	298	331	21	154	205	253	298	334
7	141	189	231	275	305	22	139	190	225	267	302
8	159	201	248	297	338	23	146	191	229	272	302
9	177	236	285	340	376	24	157	211	250	285	323
10	134	182	220	260	296	25	132	185	237	286	331
11	160	208	261	313	352	26	160	207	257	303	345
12	143	188	220	273	314	27	169	216	261	295	333
13	154	200	244	289	325	28	157	205	248	289	316
14	171	221	270	326	358	29	137	180	219	258	291
15	163	216	242	281	312	30	153	200	244	286	324

Corresponde al peso (en gramos) de 30 ratas jóvenes en un grupo de controles medido en distintos días (edad: $x_1 = 8, x_2 = 15, x_3 = 22, x_4 = 29, x_5 = 36$). Se supone un crecimiento lineal del peso respecto a la edad para cada una de las ratas.

Se tienen entonces 30 grupos, todos de tamaño 5. El modelo jerárquico que se plantea es el siguiente:

- Nivel I (Observaciones) (30 grupos, con 5 observaciones en cada uno)

$$p(\mathbf{y}|\boldsymbol{\beta}_i, \tau_i) = N_5(\mathbf{X}\boldsymbol{\beta}_i, \tau_i\mathbf{I}); i = 1, \dots, 30; \mathbf{X} = \begin{pmatrix} 1 & 8 \\ 1 & 15 \\ 1 & 22 \\ 1 & 29 \\ 1 & 36 \end{pmatrix} \in \mathbb{R}^{5 \times 2}; \boldsymbol{\beta}_i \in \mathbb{R}^{2 \times 1}$$

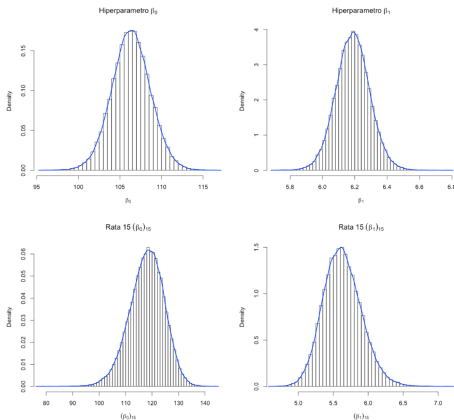
- Nivel II (Parámetros)

$$p(\boldsymbol{\beta}_i, \tau_i | \boldsymbol{\alpha}, \boldsymbol{\tau}^{(0)} \delta_1, \delta_2) = N_2(\boldsymbol{\beta}_i | \boldsymbol{\alpha}, \boldsymbol{\tau}^{(0)}) Ga(\tau_i | \delta_1, \delta_2); \boldsymbol{\tau}^{(0)} = \begin{pmatrix} \tau_1^{(0)} & 0 \\ 0 & \tau_2^{(0)} \end{pmatrix}$$

- Nivel III (Hiperparámetros)

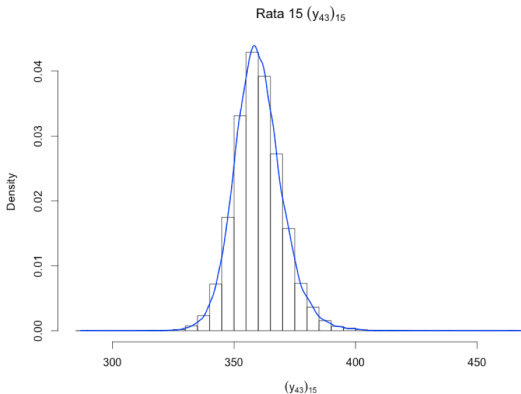
$$p(\boldsymbol{\alpha}, \boldsymbol{\tau}^{(0)}, \delta_1, \delta_2 | \boldsymbol{\theta}_0) = N_2(\boldsymbol{\alpha} | \mathbf{0}, 0.001\mathbf{I}) \\ Ga(\tau_1^{(0)} | 0.01, 0.01) Ga(\tau_2^{(0)} | 0.01, 0.01) \\ Ga(\delta_1 | 0.01, 0.01) Ga(\delta_2 | 0.01, 0.01)$$

En la siguiente figura se presentan las distribuciones finales, tanto para los coeficientes poblacionales como para los coeficientes correspondientes a la Rata 15. (Los resultados son muy parecidos a los presentados por Gelfand, Hills, Racine-Poon y Smith (1990)).



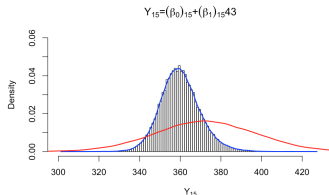
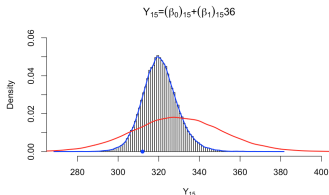
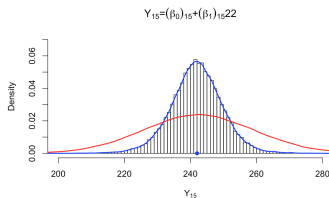
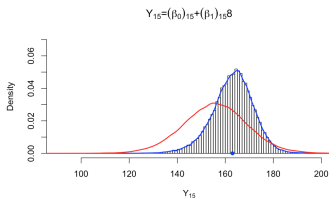
Ver código: [Programa22.r](#)

En la siguiente figura se presenta la distribución predictiva final para el peso de la Rata 15 a los 43 días de edad. El intervalo de credibilidad al 95 % obtenido vía simulación es (340.7595, 380.4825).



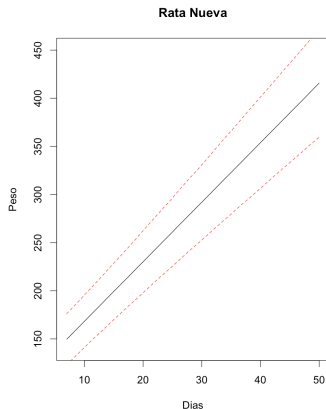
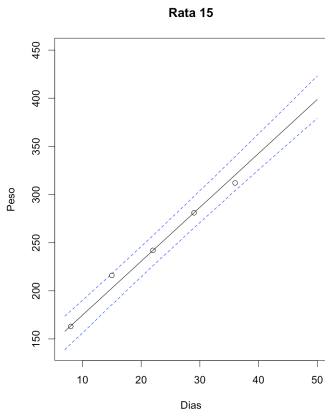
Ver código: [Programa22.r](#)

En la siguiente figura se presentan las distribuciones predictivas del peso de la Rata 15 para los días (8, 22, 36, 43). Por otro lado, la línea roja presenta la predicción en esos días para el peso de una rata nueva.



Ver código: [Programa22.r](#)

Se presenta el modelo lineal ajustado con las bandas de predicción al 95 % para la Rata 15 (líneas azules). Las líneas rojas del segundo gráfico representan la banda de predicción (95 %) para el peso de una rata nueva.



Ver código: [Programa22.r](#)

Comentarios finales

Comentarios finales

¡Muchas gracias por su atención!