

# EST-46115: Modelación Bayesiana

**Profesor:** Alfredo Garbuno Iñigo — Primavera, 2022.

**Objetivo.** Estudiar el método general de integración Monte Carlo vía cadenas de Markov (MCMC). La estrategia será construir poco a poco utilizando los principios básicos que lo componen.

**Lectura recomendada:** Capítulo 7 de [1]. Sección 6 de [4] y Capítulo 3 de [3] (avanzado).

## 1. INTRODUCCION

El interés es poder resolver

$$\mathbb{E}[f] = \int_{\Theta} f(\theta) \pi(\theta|y) d\theta. \quad (1)$$

Sin embargo, no podemos generar  $\theta^{(i)} \stackrel{\text{iid}}{\sim} \pi(\theta|y)$ .

## 2. MUESTREO POR ACEPTACIÓN RECHAZO

Podemos utilizar una versión estocástica de muestreo por importancia.

Para muestrear de  $\pi$  necesitamos utilizar una distribución sustituto (lo mismo hicimos con muestreo por importancia). Sólo que ahora permitimos rechazar muestras que no correspondan con las regiones de alta densidad de nuestra distribución objetivo. El rechazo se realiza lanzando una moneda. La tasa de éxito depende del qué tanto cubre nuestra distribución sustituto.

### 2.1. Implementación

Necesitamos algunas cosas. Ser capaces de evaluar nuestra distribución objetivo. Ser capaces de evaluar y muestrear de nuestra distribución de muestreo.

```
1 crea_mezcla <- function(weights){
2   function(x){
3     weights$w1 * dnorm(x, mean = -1.5, sd = .5) +
4     weights$w2 * dnorm(x, mean = 1.5, sd = .7)
5   }
6 }
7 objetivo <- crea_mezcla(list(w1 = .6, w2 = .4))
8 M      <- 3.3
```

LISTING 1. *Distribución objetivo.*

El objetivo del curso **no** es programación orientada a objetos. Pero todxs lxs jóvenes *cool* lo utilizan.

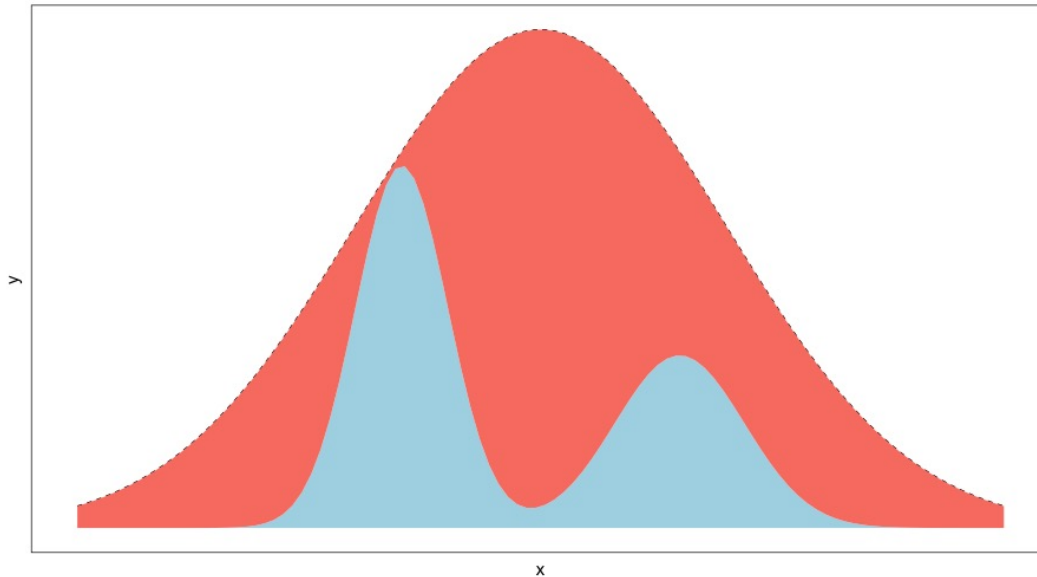


FIGURA 1. Esquema de muestreo.

```

1 library(R6)
2 ModeloNormal <-
3   R6Class("ProbabilityModel",
4     list(
5       mean = NA,
6       sd = NA,
7       ## Inicializador
8       initialize = function(mean = 0, sd = 1){
9         self$mean = mean
10        self$sd = sd
11      },
12      ## Muestreador
13      sample = function(n = 1){
14        rnorm(n, self$mean, sd = self$sd)
15      },
16      ## Evaluacion de densidad
17      density = function(x, log = TRUE){
18        dnorm(x, self$mean, sd = self$sd, log = log)
19      }
20    ))

```

LISTING 2. Distribución de muestreo.

En muestreo por aceptación rechazo necesitamos definir una distribución de la cual **si podemos** generar números aleatorios. El inconveniente es, además, **conocer** qué tanto podemos inflar la densidad de nuestra propuesta para *cubrir* la distribución objetivo.

```

1 crea_rejection_sampling <- function(objetivo, aprox, M){
2   function(niter){
3     muestras <- matrix(nrow = niter, ncol = 3)
4     for (ii in seq(1, niter)){
5       propuesta <- aprox$sample()
6       p <- objetivo(propuesta)

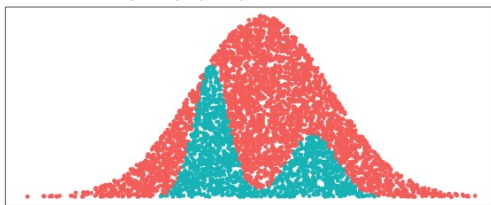
```

```

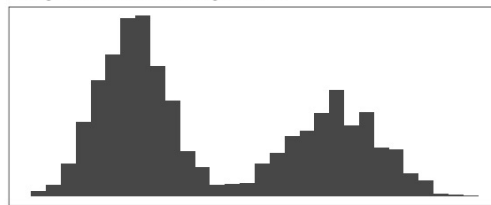
7   g <- aprox$density(propuesta, log = FALSE)
8   u <- runif(1)
9   if (u < p/(M * g)) { ## Aceptamos
10      muestras[ii, 1] <- 1
11   } else { ## Rechazamos
12      muestras[ii, 1] <- 0
13   }
14   muestras[ii, 2] <- propuesta
15   muestras[ii, 3] <- u
16 }
17 colnames(muestras) <- c("accept", "value", "auxiliar")
18 muestras
19 }
20 }

```

Muestras en el espacio (x,u), aceptación: 0.3026



Histograma de las muestras generadas



## 2.2. Tarea

- ¿Qué pasa si  $M$  es demasiado grande? Juega con el código e interpreta los resultados.
- ¿Qué pasa si  $M$  no es suficiente para cubrir la distribución objetivo? Juega con el código e interpreta los resultados.

## 2.3. Propiedades

**Lemma (Consistencia de muestreo por rechazo).** El método de muestreo de aceptación-rechazo genera muestras  $x^{(i)}$  con  $i = 1, \dots, N$  que son independientes y distribuidas acorde a la distribución objetivo  $\pi$ .

*Prueba* Usemos probabilidad condicional para medir

$$\pi(x|\text{aceptar}) = \frac{\pi(\text{aceptar}|x) \times \pi(x)}{\pi(\text{aceptar})}. \quad (2)$$

## 3. ¿QUÉ HEMOS VISTO?

- El método Monte Carlo se puede utilizar para aproximar integrales.
- Se puede utilizar una distribución sustituto para generar números aleatorios que nos interesan.
- Podemos lanzar monedas para *filtrar* sólo los aleatorios que tengan altas probabilidades.
- Hemos utilizado el supuesto de independencia.

## 4. MUESTREO POR CADENAS DE MARKOV

Vamos a **relajar** el supuesto de independencia. Es decir, vamos a generar una secuencia de números aleatorios con cierta correlación.

**Definición (Cadena de Markov).** Un **proceso estocástico** en tiempo discreto –una colección de variables aleatorias  $X_1, X_2, \dots$ – que satisface la propiedad de dependencia condicional

$$\mathbb{P}(X_{n+1} = x | X_1 = x_1, \dots, X_n = x_n) = \mathbb{P}(X_{n+1} = x | X_n = x_n), \quad (3)$$

se llama una **cadena de Markov** en tiempo discreto.

**4.1. Ejemplo:**

FIGURA 2. Problema del café.

El vendedor de galletas quiere satisfacer la demanda para acompañar un café. El vendedor:

- Viaja entre las islas.
- Decide si se queda o no se queda en la isla donde está.
- Se puede mover entre islas contiguas (a través de puentes).
- Tiene mala memoria y pregunta el número de casas en las islas aledañas (todos los días).
- Quiere visitar todas las islas y vender galletas.
- Viaja en bicicleta.

También es astuto. Sabe que en donde haya mucha gente venderá mas, pero también sabe que una isla siempre lo podría llevar a una mas grande. Asi que a veces le convendrá viajar a una isla pequeña. Asi que utilizará el **principio de aceptación rechazo** para decidir si se moverá a la siguiente isla.

1. Lanza una moneda para decidir si se mueve a la izquierda o derecha.
2. Decide si se mueve de acuerdo al cociente de poblaciones.

**4.2. Pregunta**

En el contexto de nuestro problema ¿qué cambiaría si tuviera conocimiento censal del archipiélago y pudiera viajar en avión?

**4.3. Modelación del tour de ventas**

El vendedor se encuentra en el  $t$ -ésimo día. Supongamos que va a evaluar si se cambia a la isla de la derecha. Sea  $\pi_*$  la población de la isla propuesta y  $\pi_t$  la población de la isla actual. Entonces el vendedor acepta cambiar de isla con probabilidad

$$\alpha_{\text{mover}} = \frac{\pi_*}{\pi_t}.$$

Nota que nunca dudará moverse a una isla mas grande. Por otro lado, entre mas parecidas sean las poblaciones de las islas mas **indeciso** será de moverse. Por definición  $\alpha_{\text{mover}} \in (0, 1)$ . De hecho, podemos definir la probabilidad de aceptar un viaje a otra isla por medio de

$$\alpha(t, \star) = \min \left\{ 1, \frac{\pi_{\star}}{\pi_t} \right\},$$

pues incluye los dos casos.

Una propiedad interesante es que el tiempo que el vendedor pase en cada isla corresponderá al tamaño relativo de ellas.

```

1 islas <- tibble(islas = 1:5, pob = seq(60, 100, by = 10))
2 camina_isla <- function(i){ # i: isla actual
3   u_izq <- runif(1) # Lanzamos volado para ver si nos vamos izq o der.
4   v <- ifelse(u_izq < 0.5, i - 1, i + 1) # Pedimos ndice isla vecina.
5   if (v < 1 | v > 5) { # si estas en los extremos y el volado indica salir
6     return(i)
7   }
8   u_cambio <- runif(1) # Moneda de aceptacion de cambio
9   p_cambio = min(islas$pob[v]/islas$pob[i], 1)
10  if (u_cambio < p_cambio) {
11    return(v) # isla destino
12  }
13  else {
14    return(i) # me quedo en la misma isla
15  }
16 }

```

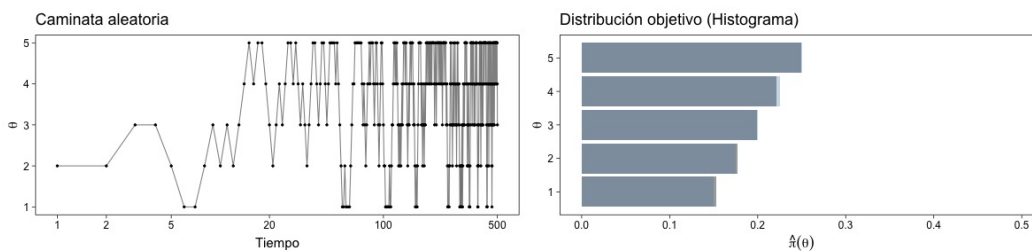


FIGURA 3. Caminata aleatoria en un archipiélago de 5 islas.

#### 4.4. Conclusiones

- La estrategia del vendedor le permitirá, en el **largo plazo**, visitar todas las islas.
- El tiempo que pasa en cada isla<sup>†</sup> corresponde a la población relativa.
- Al principio, aún no representa dicha proporción.

### 5. GENERALIZANDO...

Supongamos que tenemos un modelo

$$Y|\mu \sim N(\mu, 0.75^2), \quad (4)$$

$$\mu \sim N(0, 1^2). \quad (5)$$

Verifica que bajo la observación  $y = 6.25$  la distribución posterior que nos interesa es

$$\mu|y \sim N(4, 0.6^2). \quad (6)$$

Vamos a suponer que **no** sabemos muestrear de una Normal. Así que usaremos una estrategia parecida que con el vendedor de galletas. La estrategia será:

1. Generar una propuesta  $\mu_*$  para cambiarnos de nuestro valor actual  $\mu_t$ .
2. Decidir si nos movemos utilizando un cociente que tome en cuenta los pesos relativos.

### 5.1. Pseudo-código

- Vamos a proponer una "moneda" para lanzar la **dirección** de movimiento. Esto lo haremos con

$$\mu_*|\mu_t \sim \text{Uniforme}(\mu_t - \omega, \mu_t + \omega). \quad (7)$$

- Vamos a decidir si nos movemos de acuerdo a los pesos relativos

$$\alpha(\mu_t, \mu_*) = \min \left\{ 1, \frac{\pi(\mu_*|y)}{\pi(\mu_t|y)} \right\}. \quad (8)$$

### 5.2. Desentrañando

Escribamos el cociente en términos de la densidad de la distribución posterior y simplifiquemos. ¿Qué observas?

### 5.3. Implementación

Veamos cómo implementarlo. Vamos a suponer una distribución de muestreo con un intervalo de longitud 2. Es decir,  $\omega = 1$ .

```

1 ModeloUniforme <-
2   R6Class("ProbabilityModel",
3     list(
4       a = NA,
5       b = NA,
6       initialize = function(a = 0, b = 1){
7         self$a = a
8         self$b = b
9       },
10      sample = function(n = 1){
11        runif(n, self$a, self$b)
12      },
13      density = function(x, log = TRUE){
14        dunif(x, self$a, self$b, log = log)
15      }
16    ))

```

LISTING 3. Modelo de muestreo uniforme.

```

1 crea_cadena_markov <- function(objetivo, muestreo){
2   function(niter){
3     muestras <- matrix(nrow = niter, ncol = 2)
4     ## Empezamos en algun lugar
5     estado <- muestreo$sample()
6     muestras[1,2] <- estado
7     muestras[1,1] <- 1

```

```

8   for (ii in 2:niter){
9     propuesta  <- estado + muestreo$sample()
10    p_propuesta <- objetivo$density(propuesta, log = FALSE)
11    p_estado   <- objetivo$density(estado, log = FALSE)
12    if (runif(1) < p_propuesta/p_estado) {
13      muestras[ii, 1] <- 1 ## Aceptamos
14      muestras[ii, 2] <- propuesta
15    } else {
16      muestras[ii, 1] <- 0 ## Rechazamos
17      muestras[ii, 2] <- estado
18    }
19    estado <- muestras[ii, 2]
20  }
21  colnames(muestras) <- c("accept", "value")
22  muestras
23 }
24 }

```

```

1  objetivo <- ModeloNormal$new(mean = 4, sd = .6)
2  muestreo <- ModeloUniforme$new(a = -1, b = 1)
3
4  mcmc <- crea_cadena_markov(objetivo, muestreo)
5  muestras <- mcmc(5000)

```

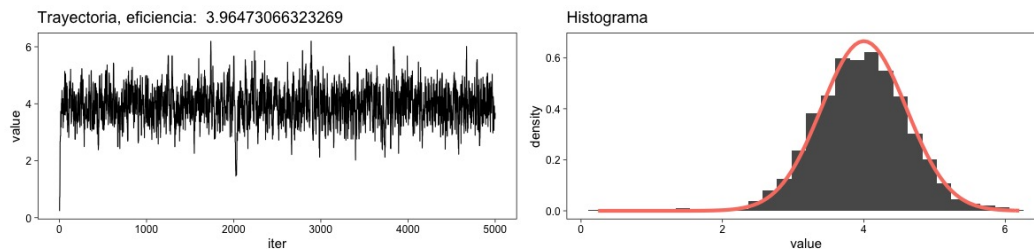


FIGURA 4. Nuestra segunda cadena de Markov.

### 5.4. Tarea (1)

Sin modificar el número de iteraciones, considera cambiar la dispersión de la distribución de muestreo.

- ¿Qué observas si  $\omega = 0.01$ ?
- ¿Qué observas si  $\omega = 100$ ?

### 5.5. Tarea (2)

Regresa a nuestro ejemplo conjugado Beta-Binomial. Considera una previa  $\theta \sim \text{Beta}(2, 3)$  y una verosimilitud  $Y|\theta \sim \text{Binomial}(2, \theta)$ . Escribe la distribución posterior asumiendo  $Y = k$ .

Para este caso tenemos un ligero inconveniente. El soporte para  $\theta$  es el intervalo cerrado  $[0, 1]$  y utilizar una propuesta como en el caso anterior nos podría colocar (casi seguramente) fuera del intervalo. Así que lo que haremos será una pequeña modificación a cómo generamos nuestra propuesta y cómo evaluamos la probabilidad de aceptar dicha propuesta.

- Vamos a generar propuestas de la siguiente manera

$$\theta_{\star}|\theta_t \sim \text{Beta}(\alpha, \beta). \quad (9)$$

□ Vamos a calcular la probabilidad de aceptar dicho movimiento a través de

$$\alpha(\theta_t, \theta_*) = \min \left\{ 1, \frac{\pi(\theta_*|y)}{\pi(\theta_t|y)} \cdot \frac{g(\theta_t)}{g(\theta_*)} \right\}, \quad (10)$$

donde  $g$  denota la densidad de la distribución de muestreo definida arriba.

Modifica el código de clase para implementar este muestreador. Compara con muestras exactas del modelo posterior bajo la observación  $Y = 1$ .

## 6. EL MÉTODO METROPOLIS-HASTINGS

La forma más general que tenemos para generar una cadena de muestras es el método de Metropolis-Hastings.

□ Generamos propuestas en cada iteración por medio de

$$\theta_*|\theta_t \sim q(\theta_*|\theta_t). \quad (11)$$

□ Calculamos la probabilidad de aceptar la propuesta como

$$\alpha(\theta_t, \theta_*) = \min \left\{ 1, \frac{\pi(\theta_*)}{\pi(\theta_t)} \cdot \frac{q(\theta_t|\theta_*)}{q(\theta_*|\theta_t)} \right\}, \quad (12)$$

donde la notación hace énfasis en que este mecanismo puede generar muestras de la distribución  $\pi$  utilizando un generador  $q$ .

### 6.1. Tarea (3)

- Repasemos los métodos anteriores.
- ¿Qué pasa si desconocemos la constante de normalización de la distribución objetivo?

### 6.2. Distribución propuesta

El *arte* está en proponer una distribución de muestreo eficiente. Como ya hemos discutido, si no está bien calibrada podríamos tener un comportamiento no deseado. Supongamos que queremos muestrear de una  $\text{Gamma}(20, 100)$ . Para esto veamos tres configuraciones de la distribución de muestreo que será  $N(\theta_t, \sigma^2)$ .

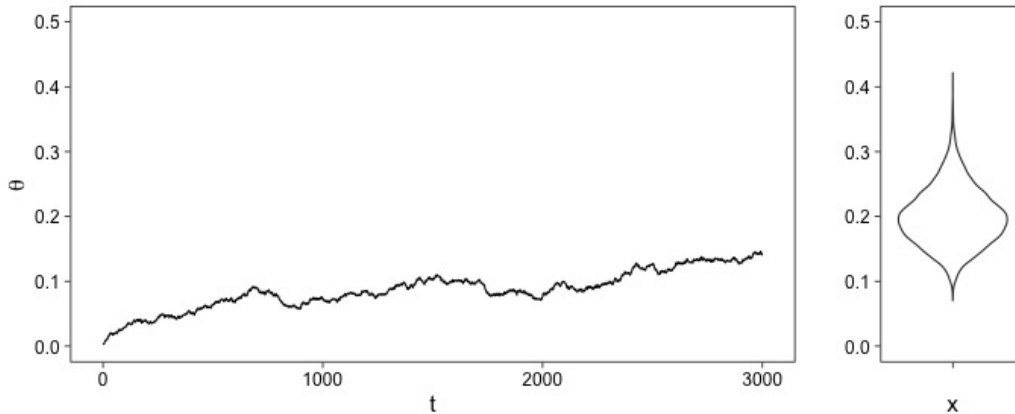


FIGURA 5. *Metropolis-Hastings en acción con un tamaño de paso muy pequeño.*



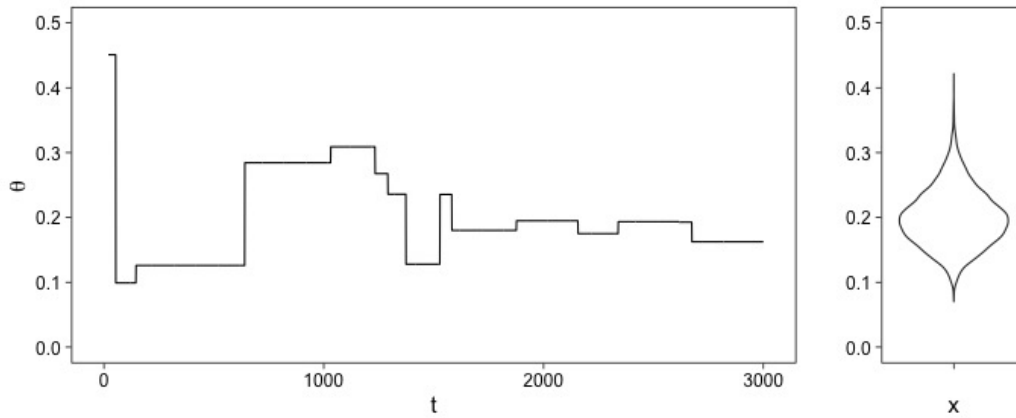


FIGURA 6. Metropolis-Hastings en acción con un tamaño de paso muy grande.

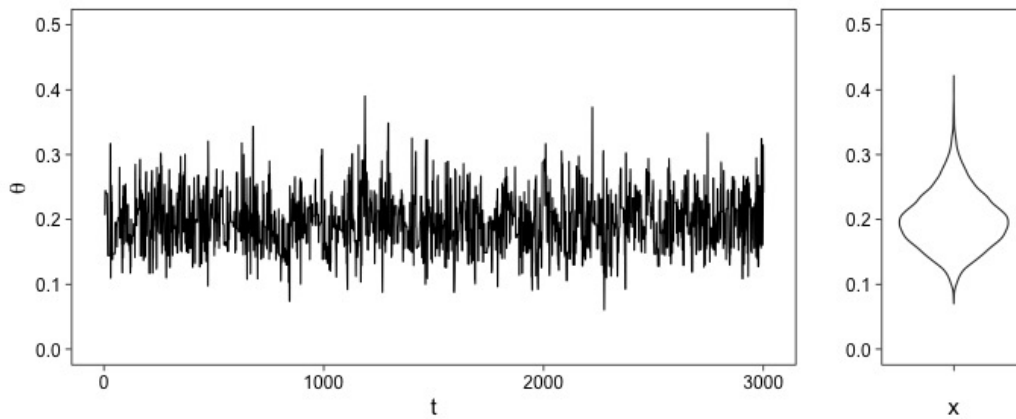


FIGURA 7. Metropolis-Hastings en acción con un tamaño de paso justo.

	configuracion	media	tasa.aceptacion
1	Paso chico	0.086	0.9440
2	Paso grande	0.309	0.0067
3	Paso justo	0.197	0.4633
4	Teorica	0.200	NA

## 7. EN MÁS DIMENSIONES

Consideremos la siguiente distribución objetivo

$$\theta \sim N(m, S), \quad m = (1, 2)^\top, \quad S = \begin{pmatrix} 1 & .75 \\ .75 & 1 \end{pmatrix}, \quad (13)$$

y utilicemos el modelo de muestreo

$$\theta \sim N(0, \Sigma), \quad 0 \in \mathbb{R}^2, \quad \Sigma = \sigma^2 \cdot \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}. \quad (14)$$

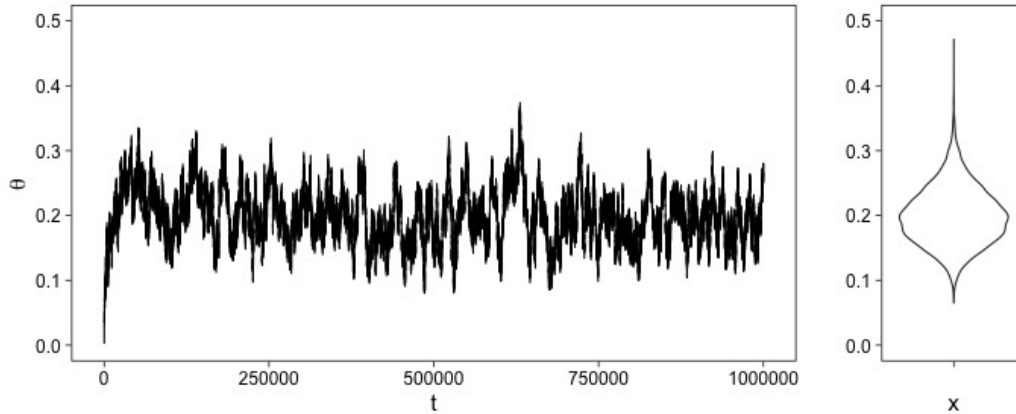


FIGURA 8. Metropolis-Hastings en acción con un tamaño de paso pequeño y un periodo suficientemente amplio.

```

1 library(mvtnorm)
2 ModeloNormalMultivariado <-
3   R6Class("ProbabilityModel",
4     list(
5       mean = NA,
6       cov = NA,
7       initialize = function(mu = 0, sigma = 1){
8         self$mean = mu
9         self$cov = sigma |> as.matrix()
10      },
11      sample = function(n = 1){
12        rmvnorm(n, mean = self$mean, sigma = self$cov)
13      },
14      density = function(x, log = TRUE){
15        dmnorm(x, self$mean, self$cov, log = log)
16      }
17    ))

```

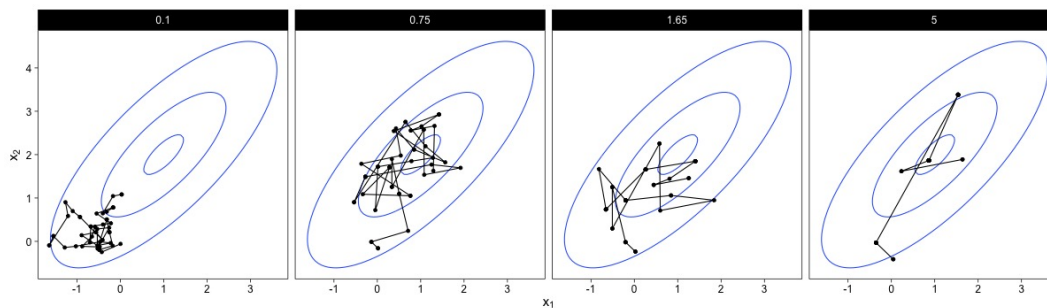


FIGURA 9. Caminata aleatoria utilizando Metropolis-Hastings para  $\theta \in \mathbb{R}^2$ .

## 8. ¿POR QUÉ FUNCIONA?

Ya vimos cómo funciona y describimos una versión suficientemente robusta. Ahora estudiaremos el por qué esa manera de operar las transiciones nos lleva a tener un mecanismo que

genera muestras de la distribución (en el largo plazo).

Para esto tenemos que preguntarnos sobre las probabilidades de transición entre dos estados. Es decir, la probabilidad de movernos al estado  $\theta_*$  condicional en estar en  $\theta$ . Lo denotamos por

$$\mathbb{P}(\theta_{t+1} = \theta_* | \theta_t = \theta) . \quad (15)$$

Si el algoritmo es capaz de mantener un balance entre las probabilidades condicionales entre dos estados de acuerdo a su frecuencia relativa, entonces el algoritmo será capaz de preservar las frecuencias.

En palabras (bueno...), buscamos que

$$\frac{\mathbb{P}(\theta_{t+1} = \theta_* | \theta_t = \theta)}{\mathbb{P}(\theta_{t+1} = \theta | \theta_t = \theta_*)} = \frac{\pi(\theta_*)}{\pi(\theta)} , \quad (16)$$

donde  $\pi(\cdot)$  denota la probabilidad objetivo.

Sólo nos falta calcular la probabilidad de transición. Esto lo logramos con dos pasos: 1) generar la propuesta y 2) aceptar o rechazar la propuesta. Por lo tanto

$$\mathbb{P}(\theta_{t+1} = \theta_* | \theta_t = \theta) = q(\theta | \theta_*) \cdot \alpha(\theta, \theta_*) = q(\theta | \theta_*) \cdot \min \left\{ 1, \frac{\pi(\theta_*)}{\pi(\theta)} \cdot \frac{q(\theta | \theta_*)}{q(\theta_* | \theta)} \right\} . \quad (17)$$

**Definición (Invarianza)** Decimos que la distribución  $\pi$  es **invariante** para un mecanismo de transición Markoviana  $(p(u, v))$  si satisface que

$$\int \pi(u) p(u, v) du = \pi(v) . \quad (18)$$

Lo que aprendemos de esto es que si tenemos un mecanismo de transición Markoviana que satisface las ecuaciones de balance entonces se mantendrá el comportamiento aleatorio de la distribución objetivo. Lo importante es que la transición preserva la distribución objetivo.

**Lemma (Comportamiento asintótico de Metropolis-Hastings).** El mecanismo de MH descrito anteriormente tiene como distribución límite  $\pi$ .

Lo que aprendemos de esto es que en particular MH preserva las ecuaciones de balance. Por lo tanto, si la cadena empieza en la distribución que nos interesa, entonces se mantendrá en ese comportamiento. Estudiar formalmente las condiciones y la tasa de convergencia para llegar a esa distribución escapa a los intereses del curso y se puede encontrar un tratamiento mas cuidadoso de esto en [2]. Sin embargo, podemos entenderlo bajo el argumento que MH busca las zonas de alta densidad. Tal como el vendedor ambulante prefería de manera consistente las islas mas grandes.

## REFERENCIAS

- [1] A. Johnson, M. Ott, and M. Dogucu. *Bayes Rules! An Introduction to Applied Bayesian Modeling*. 2021. [1](#)
- [2] S. P. Meyn and R. L. Tweedie. *Markov Chains and Stochastic Stability*. Springer London, London, 1993. ISBN 978-1-4471-3269-1 978-1-4471-3267-7. [1](#)
- [3] S. Reich and C. Cotter. *Probabilistic Forecasting and Bayesian Data Assimilation*. Cambridge University Press, Cambridge, 2015. ISBN 978-1-107-06939-8 978-1-107-66391-6. [1](#)
- [4] D. Sanz-Alonso, A. M. Stuart, and A. Taeb. Inverse Problems and Data Assimilation. *arXiv:1810.06191 [stat]*, jul 2019. [1](#)