

# EST-46115: Modelación Bayesiana

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

**Objetivo.** Estudiaremos diagnósticos típicos de métodos de simulación Markoviana (por Metropolis-Hastings o HMC) con el objetivo de poder identificar problemas y posibles soluciones para simulaciones deficientes.

**Lectura recomendada:** Para una revisión de los diagnósticos de MCMC busca la sección 11.4 de [3]. El Capítulo 3 de [1] tiene una revisión muy bonita de series de tiempo que es útil para algunos términos que estudiaremos. El artículo de Roy [5] tiene un breve resumen sobre diagnósticos. Además, se mencionan algunas técnicas modernas de diagnóstico como los métodos basados en los principios de Stein o distancias adecuadas para espacios de probabilidad.

## 1. INTRODUCCIÓN

El avance en poder computacional ha permitido la proliferación de métodos Bayesianos. El poder generar cadenas de Markov es múltiples procesadores nos ayuda a relajar los requisitos de convergencia y ayuda a explotar los recursos computacionales disponibles.

Cuando generamos una muestra de la distribución posterior usando MCMC, sin importar el método (Metropolis, Gibbs, HMC), buscamos que:

- Los valores simulados **no estén influenciados** por el valor inicial (arbitrario) y deben explorar todo el rango de la posterior.
- Debemos tener **suficientes simulaciones** de tal manera que las estimaciones sean precisas y estables.
- Queremos tener **métodos y resúmenes informativos** que nos ayuden diagnosticar correctamente el desempeño de nuestras simulaciones.

En la **práctica** intentamos cumplir lo más posible estos objetivos. Debemos de tener un criterio para considerar cadenas de **longitud finita** y **evaluar la calidad** de las simulaciones.

Primero estudiaremos **diagnósticos generales** para métodos que utilicen MCMC y después *mencionaremos* particularidades del método de simulación HMC.

En general, el problema es doble:

1. Determinar si la cadena de Markov ha alcanzado el estado estacionario;
2. Determinar si los estimadores Monte Carlo convergen a los valores esperados.

## 2. DIAGNÓSTICOS GENERALES

Una forma que tenemos de evaluar la (o identificar la falta de) convergencia es considerar distintas secuencias independientes.

### 2.1. Monitoreo de convergencia

**Burn-in e iteraciones iniciales.** En primer lugar, en muchas ocasiones las condiciones iniciales de las cadenas las escogemos de tal forma que que son *“típicos”* en relación a la posterior.

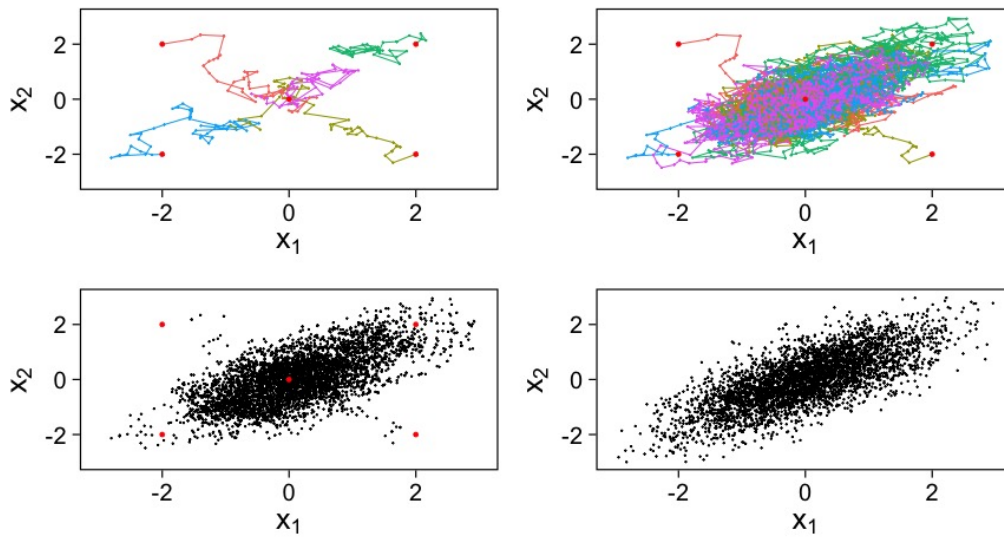
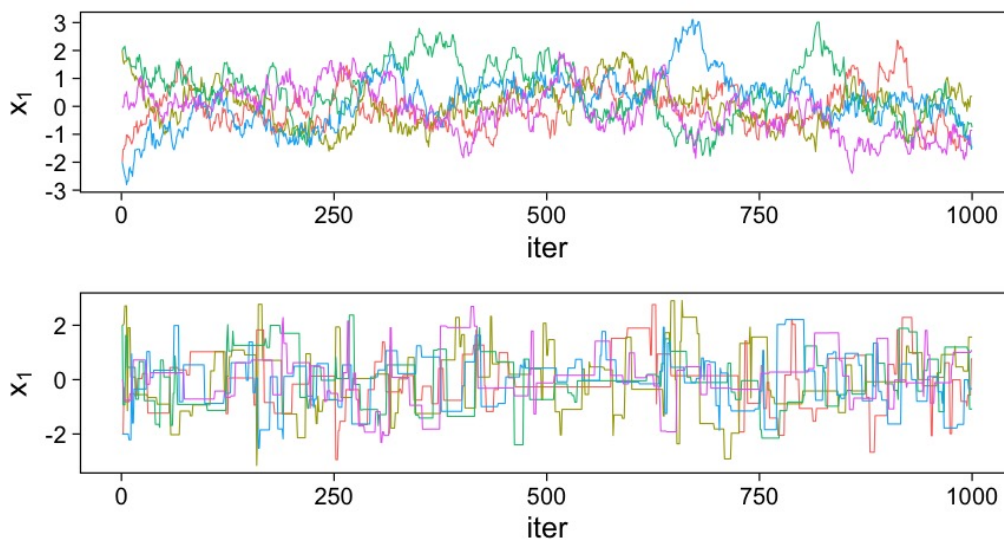


FIGURA 1. Distintas cadenas de Markov.

FIGURA 2. Trayectorias de simulación para  $X_1$ .

Estrategias de selección de puntos iniciales pueden ser valores aleatorios de la previa o perturbaciones aleatorias a estimadores MLE.

Correr varias cadenas en puntos dispersos tienen la ventaja de explorar desde distintas regiones de la posterior. Eventualmente, esperamos que todas las *cadenas mezclen bien* y representen realizaciones independientes del mismo proceso estocástico (Markoviano).

Para contrarrestar la dependencia en los distintos puntos iniciales se descarta parte de la cadena en un periodo inicial (periodo de calentamiento).

### 2.1.1. Datos: Cantantes de ópera

```

1 ## Datos: cantantes de opera -----
2 set.seed(3413)
3 cantantes <- lattice::singer >
4   mutate(estatura_cm = round(2.54 * height)) >
5   filter(str_detect(voice.part, "Tenor")) >
6   select(voice.part, estatura_cm) >
7   sample_n(20) >
8   as_tibble()

1 # A tibble: 20 × 2
2   voice.part estatura_cm
3   <fct>      <dbl>
4 1 Tenor 1      178
5 2 Tenor 2      173
6 3 Tenor 1      165
7 # ... with 17 more rows
8 # Use 'print(n = ...)' to see more rows

```

Denotamos por  $x_i$  la estatura de tenores (cantantes de ópera). Asumimos un modelo Normal con parámetros poblacionales no observados:  $\mu$  y  $\sigma$ . El modelo previo lo asumimos como

$$\mu|\sigma \sim \text{Normal}\left(\mu_0, \frac{\sigma}{n_0}\right), \quad (1)$$

$$\sigma^{-1} \sim \text{Gamma}(a_0, b_0). \quad (2)$$

En esta simulación es evidente<sup>†</sup> que necesitamos descartar una parte inicial de la simulación.

Gelman et al. [3] recomiendan descartar la mitad de las iteraciones de cada una de las cadenas que se simularon. Para problemas en dimensiones altas, incluso se podría esperar descartar hasta un 80 % de simulaciones (en especial para métodos basados en Metropolis-Hastings).

## 2.2. Monitoreo de mezcla dentro y entre cadenas

Podemos utilizar todas las simulaciones como si vinieran de una sola cadena (argumentando por estacionariedad)

```

1 cadenas.cantantes >
2   unnest(cadenas) >
3   filter(iter > 100) >
4   summarise(.estimate = mean(sigma), .variance = var(sigma), .error_mc = sqrt
5             (.variance/n()))

```

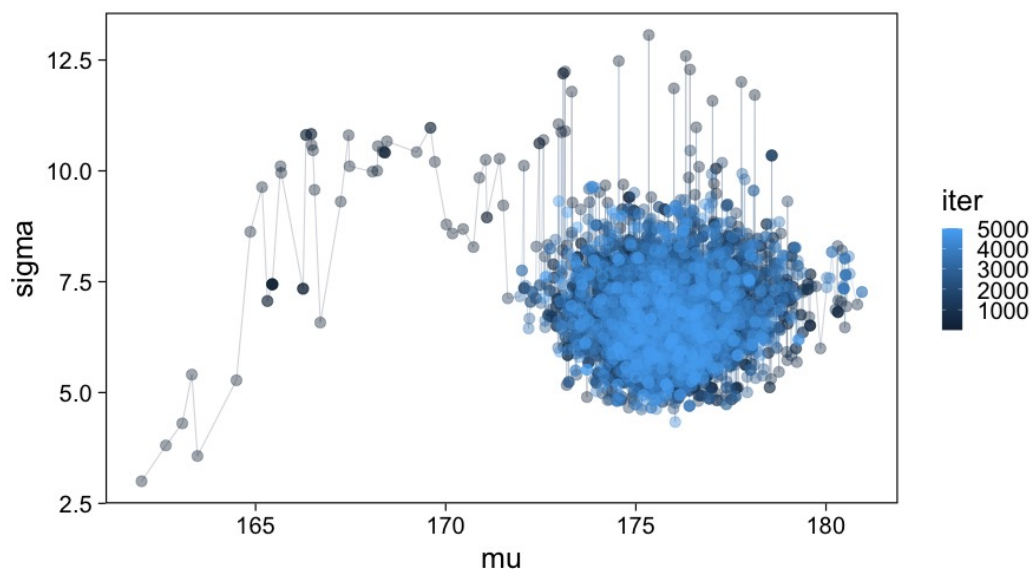


FIGURA 3. Cadena de Markov simulando de la posterior como distribución objetivo.

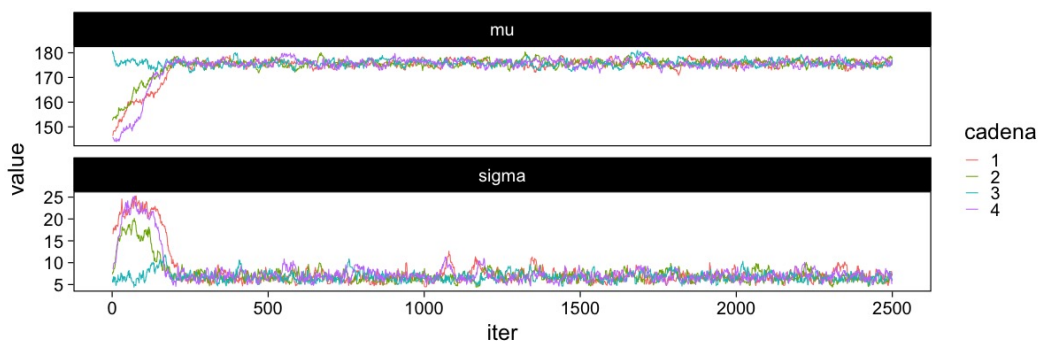


FIGURA 4. Trayectorias con dependencias iniciales.

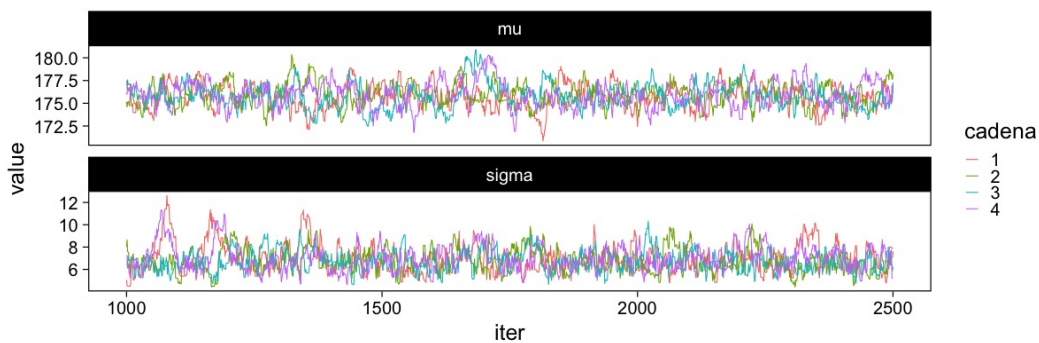


FIGURA 5. Trayectorias estacionarias.

```

1 # A tibble: 1 × 3
2   .estimate .variance .error_mc
3   <dbl>      <dbl>      <dbl>
4 1      7.11      4.14      0.0208

```

Sin embargo, al calcular la varianza como si fueran 4 cadenas independientes vemos que nuestro estimador del error Monte Carlo es mucho mas elevado de lo que esperamos ¿por qué?

```

1 cadenas.cantantes >
2   unnest(cadenas) >
3   filter(iter > 100) >
4   group_by(cadena) >
5   summarise(media = mean(sigma), varianza = var(sigma)) >
6   summarise(.estimate = mean(media), .error_mc = sd(media))

```

```

1 # A tibble: 1 × 2
2   .estimate .error_mc
3   <dbl>      <dbl>
4 1      7.11      0.272

```

Al inspeccionar cada cadena tenemos los siguientes resúmenes

```

1 cadenas.cantantes >
2   unnest(cadenas) >
3   filter(iter > 100) >
4   group_by(cadena) >
5   summarise(media = mean(sigma), varianza = var(sigma))

```

```

1 # A tibble: 4 × 3
2   cadena media varianza
3   <fct> <dbl>      <dbl>
4 1 1      7.34      7.74
5 2 2      6.93      2.35
6 3 3      6.82      1.13
7 4 4      7.35      5.12

```

Podemos partir cada cadena a la mitad y calcular nuestra estimación del error Monte Carlo. Ahora tenemos 8 cadenas que *esperamos* sean **estacionarias** (*idénticamente distribuidas*).

```

1 # A tibble: 1 × 2
2   .estimate .error_mc
3   <dbl>      <dbl>
4 1      7.11      0.418

```

Nota cómo está sucediendo algo contraintuitivo. Tenemos mas observaciones (pasamos de 1 cadena a 8) y el error Monte Carlo no decrece. Lo cual indica que nuestras cadenas realmente no han terminado de converger y tienen comportamiento distinto aunque en promedio parecen estar cercanas.

Gelman y diversos de sus coautores han desarrollado un diagnóstico numérico para evaluar implementaciones de MCMC al considerar múltiples cadenas. Aunque éste estadístico se ha ido refinando con los años, su desarrollo muestra un entendimiento gradual de éstos métodos en la práctica. La medida  $\hat{R}$  se conoce como el **factor de reducción potencial de escala**.

El estadístico  $\hat{R}$  pretende ser una estimación de la posible **reducción en la longitud** de un intervalo de confianza si las simulaciones continuaran infinitamente. Recuerda que la varianza de un estimador nos ayuda a construir intervalos en el sentido frecuentista.

La  $\hat{R}$  estudia de manera simultánea la **mezcla** de todas las cadenas (cada cadena, y fracciones de ella, deberían de haber transitado el soporte de la distribución objetivo) y **estacionariedad** (de haberse logrado cada mitad de una cadena deberían de poseer las mismas característica estadísticas).

La estrategia es descartar la **primera mitad** de cada cadena. El resto lo volvemos a dividir en dos y utilizamos cada fracción como si fuera una cadena independiente<sup>†</sup>.

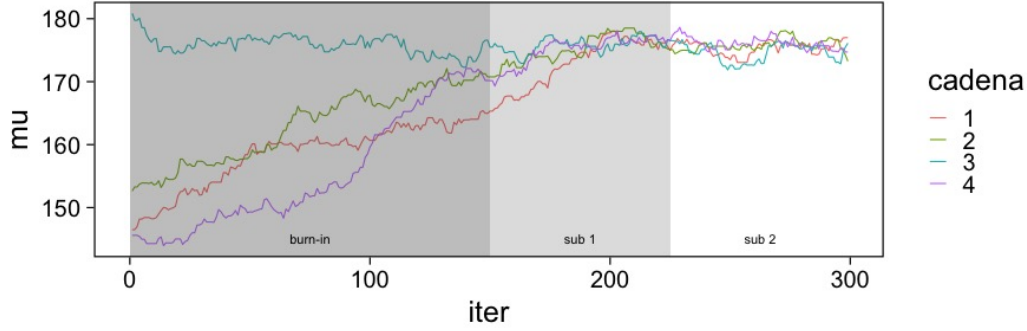


FIGURA 6. Separación de simulaciones para cálculo de  $\hat{R}$ .

Denotemos por  $m$  el número de cadenas simuladas y por  $n$  el número de simulaciones dentro de cada cadena. Cada una de las **cantidades escalares de interés** las denotamos por  $\phi$ . Éstas pueden ser los parámetros originales  $\theta$  o alguna otra cantidad derivada  $\phi = f(\theta)$ .

Ejemplos de esto puede ser en un modelo Beta-Binomial donde nos interesa la tasa de éxitos  $\theta$  pero necesitamos monitorear  $\phi = \log(\theta/(1-\theta))$ . Otra situación puede ser el caso de un modelo normal con varianza desconocida  $\sigma^2$  y necesitamos monitorear  $\phi = \log \sigma^2$ .

Ahora denotemos por  $\phi_{ij}$  las simulaciones que tenemos disponibles con  $i = 1, \dots, n$ , y  $j = 1, \dots, m$ . Calculamos  $B$  y  $W$ , la variabilidad **entre** (*between*) y **dentro** (*within*) cadenas, respectivamente, por medio de

$$W = \frac{1}{m} \sum_{j=1}^m s_j^2, \quad \text{con} \quad s_j^2 = \frac{1}{n-1} \sum_{i=1}^n (\phi_{ij} - \bar{\phi}_{\cdot j})^2, \quad \text{donde} \quad \bar{\phi}_{\cdot j} = \frac{1}{n} \sum_{i=1}^n \phi_{ij}, \quad (3)$$

$$B = \frac{n}{m-1} \sum_{j=1}^m (\bar{\phi}_{\cdot j} - \bar{\phi}_{\cdot \cdot})^2, \quad \text{donde} \quad \bar{\phi}_{\cdot \cdot} = \frac{1}{m} \sum_{j=1}^m \bar{\phi}_{\cdot j}. \quad (4)$$

La varianza entre cadenas,  $B$ , se multiplica por  $n$  dado que ésta se calcula por medio de promedios y sin este factor de corrección no reflejaría la variabilidad de las cantidades de interés  $\phi$ .

La varianza de  $\phi$  se puede estimar por medio del **estimador agregado de varianza**

$$\hat{V}(\phi)^+ = \frac{n-1}{n} W + \frac{1}{n} B. \quad (5)$$

Este estimador **sobre-estima** la varianza pues los puntos iniciales pueden estar sobre-dispersos, mientras que es un **estimador insesgado** una vez que se haya alcanzado el estado estacionario (realizaciones de la distribución objetivo)

Por otro lado, la varianza estimada por  $W$  será un sub-estimador pues podría ser el caso de que cada cadena no ha tenido la oportunidad de recorrer todo el soporte de la distribución. En el límite  $n \rightarrow \infty$ , el valor esperado de  $W$  aproxima  $\mathbb{V}(\phi)$ .

Se utiliza como diagnostico el factor por el cual la escala de la distribución actual de  $\phi$  se puede reducir si se continua con el procedimiento en el límite  $n \rightarrow \infty$ . Esto es,

$$\hat{R} = \sqrt{\frac{\hat{\mathbb{V}}(\phi)^+}{W}},$$

por construcción converge a 1 cuando  $n \rightarrow \infty$ .

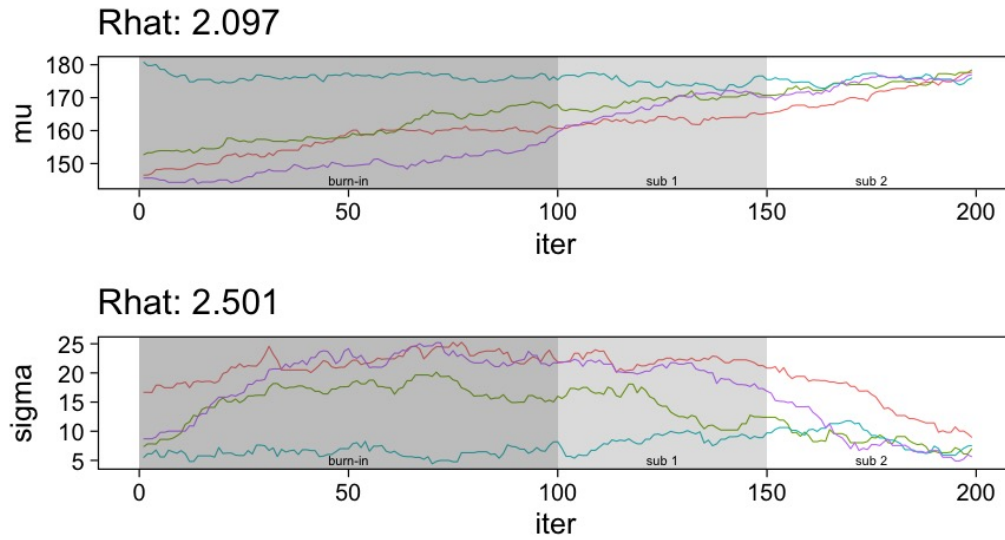


FIGURA 7. Diagnóstico de reducción de escala. Sugerencia: generar mas simulaciones.

Problemas con  $\hat{R}$ . El estimador de reducción de escala funciona bien para monitorear estimadores y cantidades de interés basados en medias y varianzas, o bien, cuando la distribución es simétrica y cercana a una Gaussiana. Es decir, colas ligeras. Sin embargo, para percentiles, o distribuciones lejos del supuesto de normalidad no es un buen indicador. Es por esto que también se recomienda incorporar transformaciones que nos permitan generar un buen estimador. Puedes leer mas de esto en el artículo de Vehtari et al. [6].



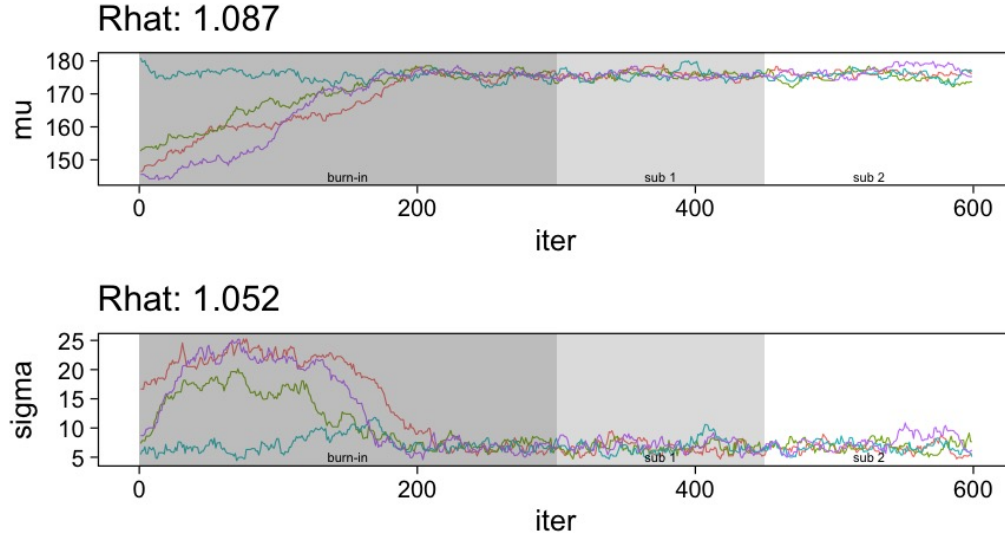


FIGURA 8. Diagnóstico de reducción de escala. Observaciones: parece estar bien.

Existe una versión multivariada de este estadístico que busca resumir la misma información a lo largo de todos los componentes de  $\theta \in \mathbb{R}^p$ . Es prácticamente el mismo estimador con las versiones multivariadas de varianza dentro de cada cadena  $\Sigma_W$  y la varianza agregada  $\Sigma_{V+}$ . Es decir, el estadístico se calcula como

$$\hat{R}_p = \max_{a \in \mathbb{R}^p} \frac{a^\top \Sigma_{V+} a}{a^\top \Sigma_W a}, \quad (6)$$

el cual, se puede probar, está relacionado con el valor propio mas grande de la matriz  $\Sigma_W^{-1} \Sigma_{V+} / n$  [5].

### 2.3. Número efectivo de simulaciones

Queremos que los recursos que hemos asignado a generar simulaciones sean representativos de la distribución objetivo. Si las  $n$  simulaciones dentro de cada cadena en verdad son realizaciones independientes entonces la estimación de  $B$  sería un estimador insesgado de  $\mathbb{V}(\phi)$ .

En esta situación tendríamos  $n \times m$  realizaciones de la distribución que queremos simular. Sin embargo, la correlación entre las muestras hacen que  $B$  sea mayor que  $\mathbb{V}(\phi)$  en promedio.

Una manera para definir el tamaño efectivo de simulaciones es por medio del estudio del estimador

$$\bar{\phi}_{..} \approx \mathbb{E}(\phi). \quad (7)$$

Del cual podemos derivar que

$$\mathbb{V}(\bar{\phi}_{..}) = \frac{\mathbb{V}(\phi)}{m \cdot n}.$$

El problema es que la correlación en las cadenas implica el denominador  $(m \cdot n)$  realmente sea una fracción del total de muestras, digamos  $\lambda$ . De tal forma que el número efectivo de



simulaciones es

$$N_{\text{eff}} = \lambda \cdot (m n),$$

donde

$$\lambda = \frac{1}{\sum_{t=-\infty}^{\infty} \rho_t} = \frac{1}{1 + 2 \sum_{t=1}^{\infty} \rho_t}.$$

El término  $\rho_t$  denota la **auto-correlación** con diferencia en  $t$  unidades de tiempo.

**Definición (autocorrelación):** La autocovarianza y autocorrelación de una serie temporal **estacionaria**  $\{Y_t : t = 0, \dots\}$  están definidas (respectivamente) como

$$C_\tau = \mathbb{E}[(Y_{t+\tau} - \mu)(Y_t - \mu)], \quad \rho_\tau = \frac{\mathbb{E}[(Y_{t+\tau} - \mu)(Y_t - \mu)]}{\sigma^2}. \quad (8)$$

**Definición (estimador de autocorrelación):** La función de autocorrelación se estima utilizando

$$\hat{C}_\tau = \frac{1}{T - \tau} \sum_{t=1}^{T-\tau} (Y_{t+\tau} - \hat{\mu})(Y_t - \hat{\mu}), \quad \hat{\rho}_\tau = \frac{\hat{C}_\tau}{\hat{C}_0}. \quad (9)$$

**Definición (variograma):** El variograma de una serie temporal **estacionaria**  $\{Y_t : t = 0, \dots\}$  está definido como

$$V_\tau = \mathbb{E}[(Y_{t+\tau} - Y_t)^2]. \quad (10)$$

**Nota** que  $V_\tau = C_0 - C_\tau$ .

Regresando a nuestro contexto... para estimar  $\rho_t$  partimos de nuestro estimador  $\hat{V}(\phi)^+$ ; y utilizamos el **variograma**  $V_t$  para cada **retraso**  $t$

$$V_t = \frac{1}{m(n-t)} \sum_{j=1}^m \sum_{i=t+1}^n (\phi_{i,j} - \phi_{i-t,j})^2.$$

Utilizando la igualdad  $\mathbb{E}(\phi_i - \phi_{i-t})^2 = 2(1 - \rho_t)\mathbb{V}(\phi)$ , podemos estimar

$$\hat{\rho}_t = 1 - \frac{V_t}{2 \hat{V}(\phi)^+}.$$

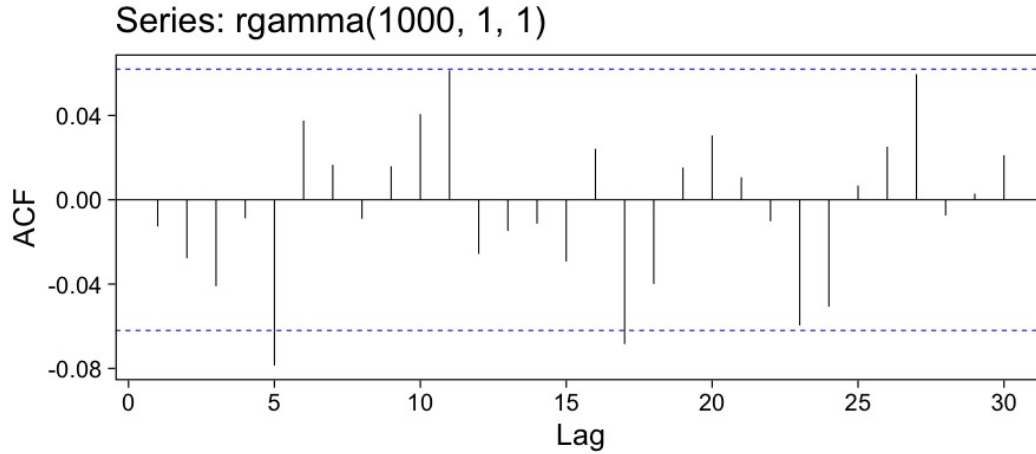
La mayor dificultad que presenta el estimador es considerar **todos** los retrasos posibles. Eventualmente agotaremos la longitud de las cadenas para ello. Por otro lado, para  $t$  eventualmente grande nuestros estimadores del variograma  $V_t$  serán muy ruidosos (¿por qué?). En la práctica truncamos la serie de acuerdo a las observaciones [4]. La serie tiene la propiedad de que para cada par  $\rho_{2t} + \rho_{2t+1} > 0$ . Por lo tanto, la serie se trunca cuando observamos  $\hat{\rho}_{2t} + \hat{\rho}_{2t+1} < 0$  para dos retrasos sucesivos.

Si denotamos por  $T$  el **tiempo de paro** (el máximo número de rezagos que podemos considerar), el estimador para el número efectivo de simulaciones es

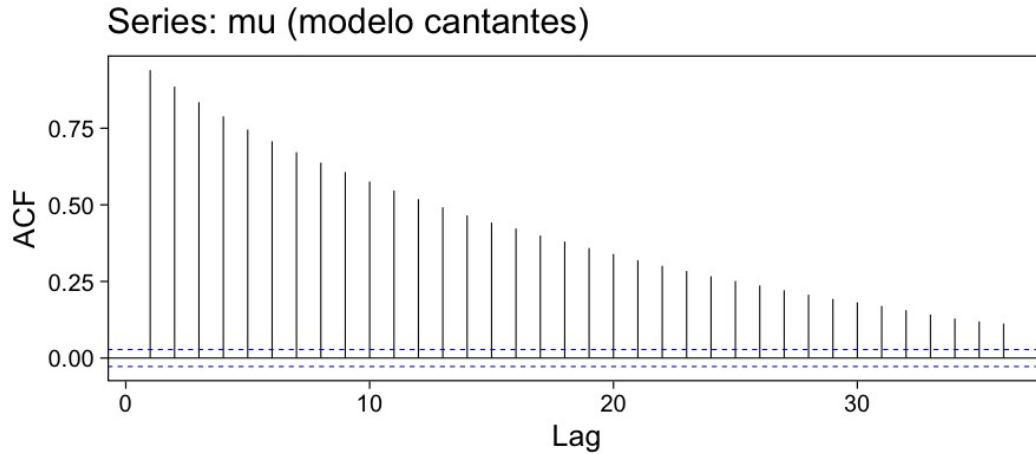
$$\widehat{\text{ESS}} = \frac{m n}{1 + 2 \sum_{t=1}^T \hat{\rho}_t}.$$

El **tamaño efectivo de simulaciones** nos ayuda a monitorear lo siguiente. Si las simulaciones fueran independientes  $N_{\text{eff}}$  sería el número total de simulaciones; sin embargo, las simulaciones de MCMC suelen estar correlacionadas, de modo que cada iteración de MCMC es menos informativa que si fueran independientes.

Por ejemplo si graficáramos simulaciones independientes, esperaríamos valores de autocorrelación chicos:



Sin embargo, los valores que simulamos tienen el siguiente perfil de autocorrelación:



```

1      mu      sigma  accept
2 0.02921 0.04822 0.69000

```

LISTING 1. Fracción ESS/nm y tasas de aceptación para la simulación de la posterior los cantantes de ópera.

```

1      mu      sigma  accept
2 0.08824 0.14556 0.38020

```

LISTING 2. Fracción ESS/nm y tasas de aceptación para la simulación (calibrada) de la posterior los cantantes de ópera.

## 2.4. Relación con error Monte Carlo

Conocer el número efectivo de simulaciones nos permite calcular el error estándar de una aproximación Monte Carlo por medio de expresiones como

$$\text{SE}(\pi(f)) \approx \left( \frac{\hat{V}_{\pi}(f)}{\widehat{\text{ESS}}} \right)^{\frac{1}{2}}. \quad (11)$$

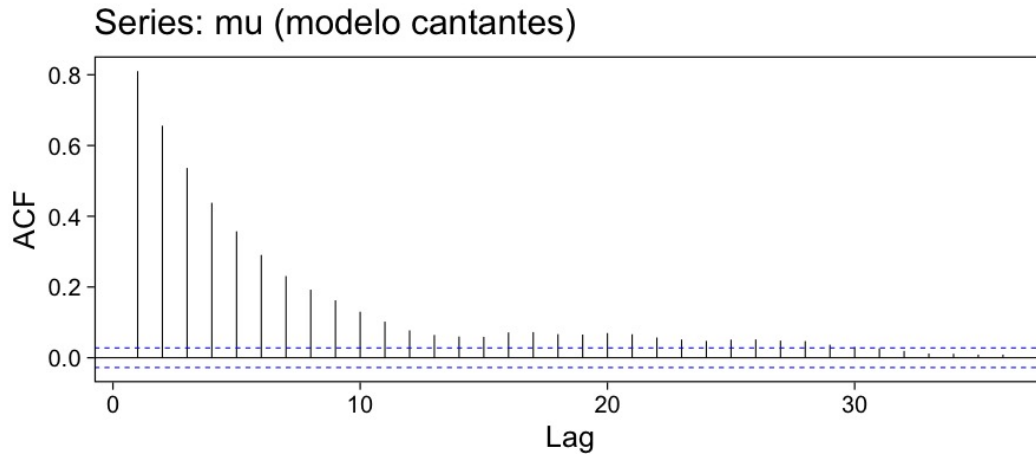


FIGURA 9. Perfil de correlación para la simulación calibrada.

Geyer [4] menciona que realizar estimaciones puntuales sin medidas de incertidumbre es el *deber* ser de un analista/estadístico. Esto es, independientemente de si la técnica es Bayesiana o frecuentista. El artículo de Flegal et al. [2] también discute la importancia de reportar errores estándar.

## 2.5. Adelgazamiento de cadenas

El método de simulación por medio de cadenas de Markov es computacionalmente intensivo. Sin embargo, los estimadores de varianza pueden ser muy volátiles (debido a altas correlaciones entre muestras). Adelgazar la cadena tiene como objetivo buscar quedarse con muestras de **alta calidad** para realizar estimaciones Monte Carlo. Se puede utilizar el ESS como una noción de cuántas muestras preservar.

## 3. CONCLUSIONES

- Ambos estadísticos asumen la existencia de un teorema de límite central Markoviano.
- En la práctica, las condiciones teóricas para garantizar su existencia (CLT) se necesitan probar caso a caso. Sin embargo, no es común que estos no existan bajo un mecanismo de muestreo tipo Metropolis-Hastings.
- Existen otras alternativas para diagnosticar un buen comportamiento de la cadena de Markov. Por ejemplo, se pueden utilizar pruebas de hipótesis para diferencias en medias con dos pedazos de cadenas.
- Para casos multivariados se pueden ajustar los análisis univariados por medio de pruebas de hipótesis múltiples con sus ajustes correspondientes (tipo Bonferroni, por nombrar uno).

## REFERENCIAS

- [1] N. Cressie and C. K. Wikle. *Statistics for Spatio-Temporal Data*. John Wiley & Sons, 2015. [1](#)
- [2] J. M. Flegal, M. Haran, and G. L. Jones. Markov chain monte carlo: Can we trust the third significant figure? *Statistical Science*, pages 250–260, 2008. [11](#)
- [3] A. Gelman, J. B. Carlin, H. S. Stern, D. B. Dunson, A. Vehtari, and D. B. Rubin. *Bayesian Data Analysis*, volume 2. CRC press Boca Raton, FL, 2014. [1](#), [3](#)
- [4] C. J. Geyer. Introduction to Markov Chain Monte Carlo. *Handbook of Markov Chain Monte Carlo*, pages 3–48, 2002. [9](#), [11](#)

- 
- [5] V. Roy. Convergence diagnostics for Markov chain Monte Carlo, oct 2019. [1](#), [8](#)
- [6] A. Vehtari, A. Gelman, D. Simpson, B. Carpenter, and P.-C. Bürkner. Rank-Normalization, Folding, and Localization: An Improved  $\hat{R}$  for Assessing Convergence of MCMC (with Discussion). *Bayesian Analysis*, 16(2), jun 2021. ISSN 1936-0975. . [7](#)