

EST-46115: Modelación Bayesiana

Profesor: Alfredo Garbuno Iñigo — Primavera, 2022 — Integración Monte Carlo.

Objetivo: Estudiar integración numérica en el contexto probabilístico. Estudiar, en particular, el método Monte Carlo y entender sus bondades y limitaciones en el contexto de inferencia Bayesiana.

Lectura recomendada: Sección 6.1 de Johnson et al. [1]. Una lectura mas técnica sobre reglas de cuadratura se puede encontrar en la sección 3.1 de Reich and Cotter [2]. Y una buena referencia (técnica) sobre el método Monte Carlo lo encuentran en Sanz-Alonso et al. [3].

1. INTRODUCCIÓN

En inferencia bayesiana lo que queremos es poder resolver

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

Lo que necesitamos es resolver integrales con respecto a la distribución de interés.

- La pregunta clave (I) es: ¿qué distribución?
- La pregunta clave (II) es: ¿con qué método numérico resuelvo la integral?
- La pregunta clave (III) es: ¿y si no hay método numérico?

2. ¿POR QUÉ INTEGRAR?

Consideremos de interés estimar la proporción de volumen de la hiper-esfera contenida en un hiper-cubo unitario conforme aumenta la dimensión del problema.

```
1 distancia_euclideana <- function(u) sqrt(sum(u * u));
2
3 experimento <- function(ndim){
4   nsamples <- 1e5;
5   y <- matrix(runif(nsamples * ndim, -0.5, 0.5), nsamples, ndim);
6   mean(apply(y, 1, distancia_euclideana) < 0.5)
7 }
```

```
1 puntos.grafica <- tibble(dims = c(1, 5, 10, 25, 50, 100)) ▷
2   mutate(points = map(dims, function(dim){
3     tibble(x = seq(0, 15, length.out = 1000)) ▷
4       mutate(y = chi.pdf(x, dim))
5   }), dimensions = factor(dims))
```

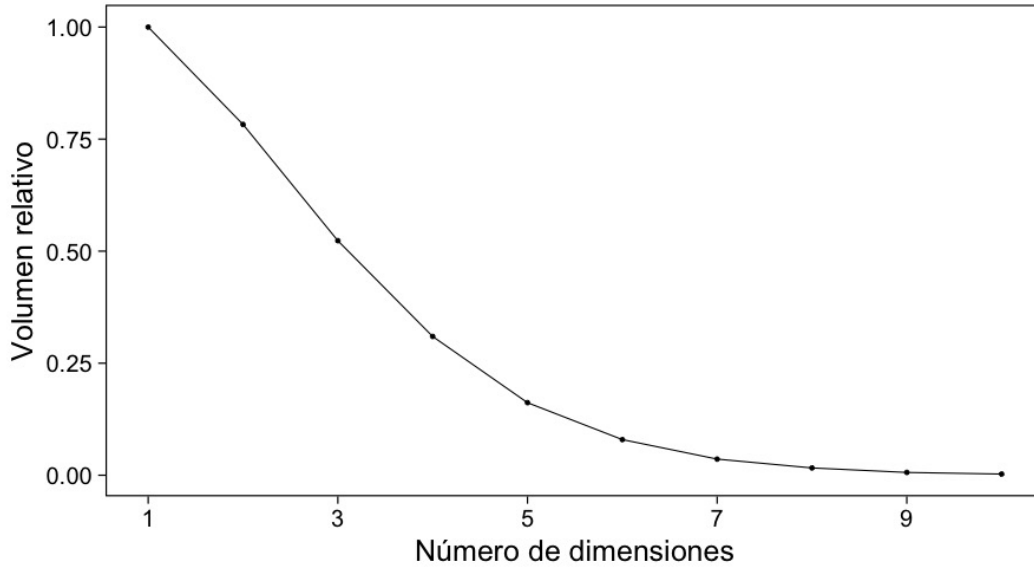


FIGURA 1. Evolución del volumen relativo de la hiper-esfera circunscrita dentro del hiper-cubo unitario.

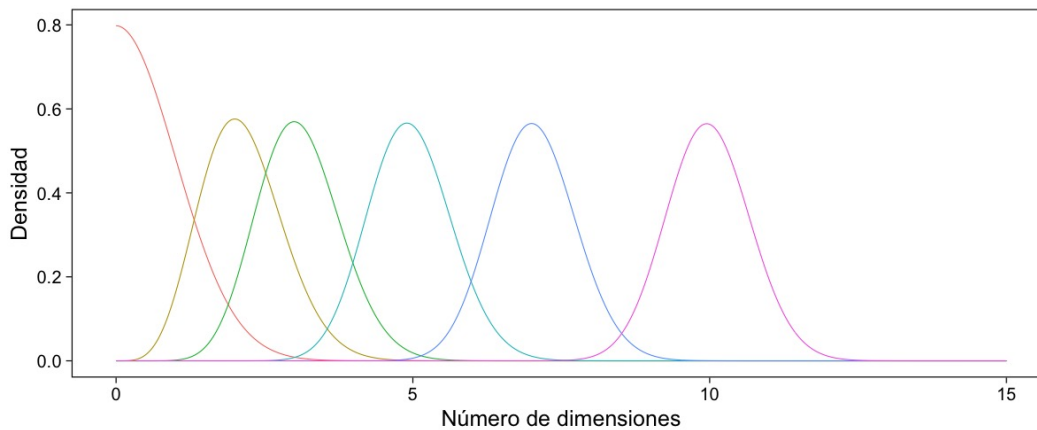


FIGURA 2. Densidad de los vecinos de la moda para una Normal multivariada estándar $N(0, I_p)$.

Otro detalle interesante de altas dimensiones es la poca intuición probabilística que tenemos de estos espacios y de lo que es una muestra típica de una distribución.

Por ejemplo, para $X \sim N(0, 1)$ estamos acostumbrados a asociar la moda como el valor de mas alta densidad. Lo cual es un error terrible en varias dimensiones.

Consideremos un análisis analítico. Por ejemplo, sabemos que si $X \sim N(0, I_p)$, entonces tenemos que

$$\sum_{i=1}^p X_i^2 \sim \chi_p^2. \quad (2)$$

Gráfiqemos el valor central de estas variables aleatorias y sus percentiles del 2.5 % y 97.5 %. Lo que observamos es que una **muestra típica** no se comporta como el promedio de nuestra distirbución, *aka* el individuo promedio no es tan común.

Lo que sucede se conoce como el fenómeno de **concentración de medida** donde los puntos de más alta densidad no corresponden a los puntos de mayor volumen (probabilidad).

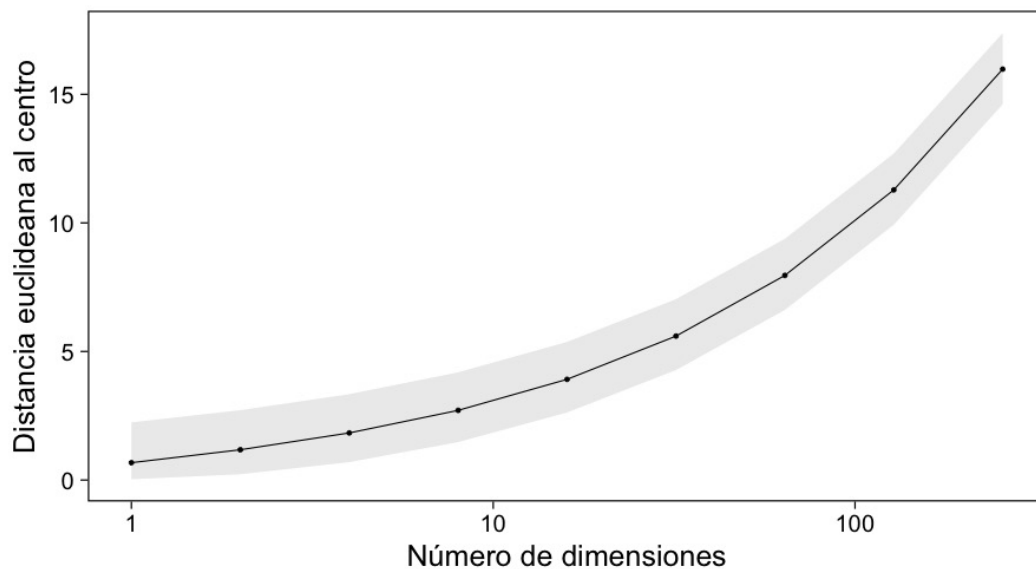


FIGURA 3. Distancia euclídea de puntos aleatorios de una Gaussiana multivariada al centro de la distribución. Esto ilustra que aunque el centro es el comportamiento promedio, los puntos típicos de una Gaussiana se encuentran lejos.

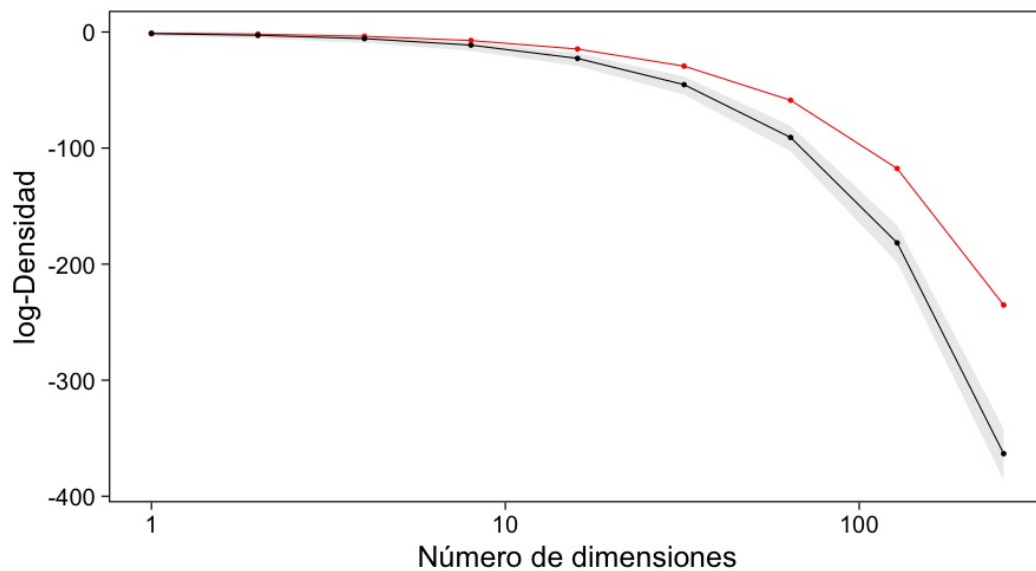


FIGURA 4. En rojo la log-densidad de la moda de una Gaussiana multivariada. En negro la log-densidad de muestras aleatorias de una Gaussiana multivariada. Esto muestra que los elementos con mayor densidad no corresponden a vecindades de mayor volumen.

Esto explica por qué no queremos realizar la aproximación

$$\pi(f) \approx f(\theta^*), \quad \text{donde} \quad \theta^* = \arg \max_{\theta \in \Theta} \pi(\theta). \quad (3)$$

3. INTEGRACIÓN NUMÉRICA

Recordemos la definición de integrales Riemann:

$$\int_a^b f(x) dx.$$

La aproximación utilizando una malla de N puntos sería:

$$\sum_{n=1}^N f(u_n) \Delta u_n.$$

El método útil cuando las integrales se realizan cuando tenemos pocos parámetros. Es decir, $\theta \in \mathbb{R}^p$ con p pequeña.

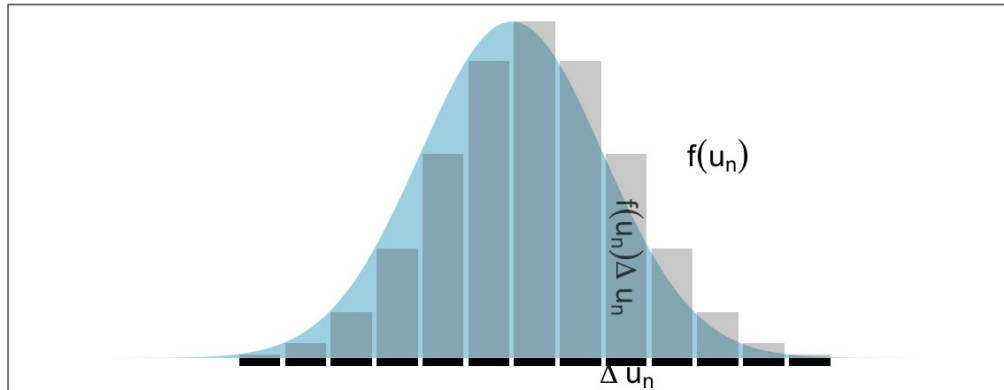


FIGURA 5. Integral por medio de discretización con $N = 11$.

3.1. Ejemplo: Proporción

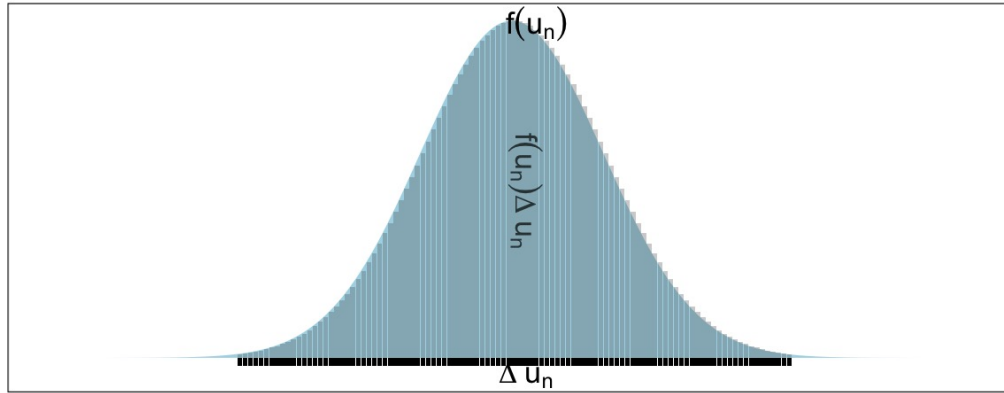
Supongamos que $p(S_n = k|\theta) \propto \theta^k(1 - \theta)^{n-k}$ cuando observamos k éxitos en n pruebas independientes. Supongamos que nuestra inicial es $p(\theta) = 2\theta$ (checha que es una densidad).

```

1 crear_log_post ← function(n, k){
2   function(theta){
3     verosim ← k * log(theta) + (n - k) * log(1 - theta)
4     inicial ← log(theta)
5     verosim + inicial
6   }
7 }
```

```

1 # observamos 3 exitos en 4 pruebas:
2 log_post ← crear_log_post(4, 3)
3 prob_post ← function(x) { exp(log_post(x)) }
```

FIGURA 6. Integral por medio de una malla fina, $N = 101$.

```

4 # integramos numericamente
5 p_x <- integrate(prob_post, lower = 0, upper = 1, subdivisions = 100L)
6 p_x

```

```

1 0.03333 with absolute error < 3.7e-16

```

Y ahora podemos calcular la media posterior:

$$\mathbb{E}[\theta|S_n] = \int_{\Theta} \theta \pi(\theta|S_n) d\theta. \quad (4)$$

```

1 media_funcion <- function(theta){
2   theta * prob_post(theta) / p_x$value
3 }
4 integral_media <- integrate(media_funcion,
5                               lower = 0, upper = 1,
6                               subdivisions = 100L)
7 media_post <- integral_media$value
8 c(Numerico = media_post, Analitico = 5/(2+5))

```

```

1 Numerico Analitico
2 0.7143 0.7143

```

3.2. Más de un parámetro

Consideramos ahora un espacio con $\theta \in \mathbb{R}^p$. Si conservamos N puntos por cada dimensión, ¿cuántos puntos en la malla necesitaríamos? Lo que tenemos son recursos computacionales limitados y hay que buscar hacer el mejor uso de ellos. En el ejemplo, hay zonas donde no habrá contribución en la integral.

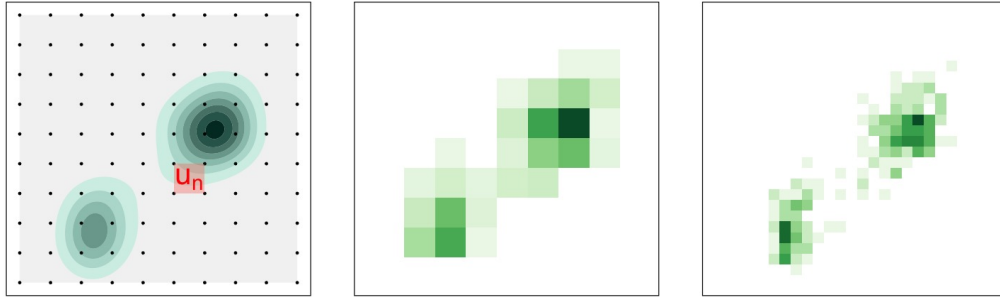


FIGURA 7. Integral por método de malla.

3.3. Reglas de cuadratura

Por el momento hemos escogido aproximar las integrales por medio de una aproximación con una **mall**a **uniforme**. Sin embargo, se pueden utilizar aproximaciones

$$\int_a^b f(x)dx \approx \sum_{n=1}^N f(\xi_n) \omega_n .$$

Estas aproximaciones usualmente se realizan para integrales en intervalos cerrados $[a, b]$. La regla de cuadratura determina los pesos ω_n y los centros ξ_n pues se escogen de acuerdo a **ciertos criterios de convergencia**.

Por ejemplo, se consideran polinomios que aproximen con cierto grado de precisión el integrando. Los pesos y los centros se escogen de acuerdo a la familia de polinomios. Pues para cada familia se tienen identificadas las mallas que optimizan la aproximación. Ver sección 3.1 de Reich and Cotter [2].

4. INTEGRACIÓN MONTE CARLO

$$\pi(f) = \mathbb{E}_\pi[f] = \int f(x)\pi(x)dx ,$$

$$\hat{\pi}_N^{\text{MC}}(f) = \frac{1}{N} \sum_{n=1}^N f(x^{(n)}) , \quad \text{donde } x^{(n)} \stackrel{\text{iid}}{\sim} \pi, \quad \text{con } n = 1, \dots, N ,$$

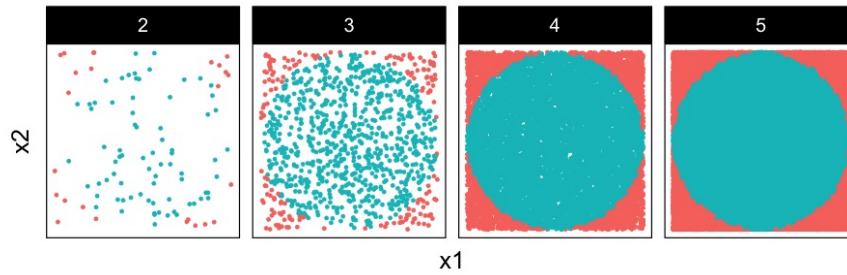
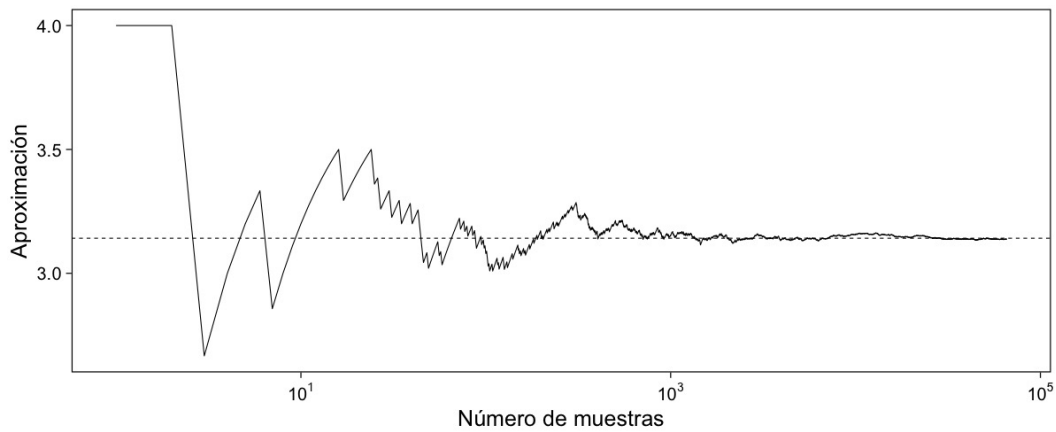
$$\pi(f) \approx \hat{\pi}_N^{\text{MC}}(f) .$$

4.1. Ejemplo: Dardos

Consideremos el experimento de lanzar dardos uniformemente en un cuadrado de tamaño 2, el cual contiene un círculo de radio 1.

Si escogemos N suficientemente grande entonces nuestro promedio converge a la integral. En Fig. 9 se muestra para cada n en el eje horizontal cómo cambia nuestra estimación $\hat{\pi}_n^{\text{MC}}(f)$.

También podemos replicar el experimento unas M veces y observar cómo cambiaría nuestra estimación con distintas semillas. Por ejemplo, podemos replicar el experimento 10 veces. En R y python lo usual es utilizar **arreglos multidimensionales** para poder guardar muestras bajo distintas replicas.

FIGURA 8. Integración Monte Carlo para aproximar π .FIGURA 9. Estimación $\hat{\pi}_N^{MC}(f)$ con $N \rightarrow \infty$.

```

1 set.seed(108)
2 nsamples <- 10**4; nexp <- 100
3 U <- runif(nexp * 2 * nsamples)
4 U <- array(U, dim = c(nexp, 2, nsamples))
5 apply(U[1:5,,], 1, str)

```

```

1 num [1:2, 1:10000] 0.455 0.164 0.529 0.415 0.474 ...
2 num [1:2, 1:10000] 0.404 0.9282 0.0883 0.3126 0.4386 ...
3 num [1:2, 1:10000] 0.351 0.449 0.369 0.814 0.695 ...
4 num [1:2, 1:10000] 0.664 0.627 0.185 0.882 0.113 ...
5 num [1:2, 1:10000] 0.4635 0.0115 0.0117 0.5086 0.9384 ...
6 NULL

```

```

1 resultados <- apply(U, 1, function(x){
2   dardos <- apply(x**2, 2, sum)
3   exitos <- ifelse(dardos <= 1, 1, 0)
4   prop <- cummean(exitos)
5   4 * prop
6 })

```

Lo cual nos permite realizar distintos escenarios posibles.

Bajo ciertas consideraciones teóricas podemos esperar un buen comportamiento de nuestro estimador de la integral. E incluso podríamos (si el número de simulaciones lo permite)

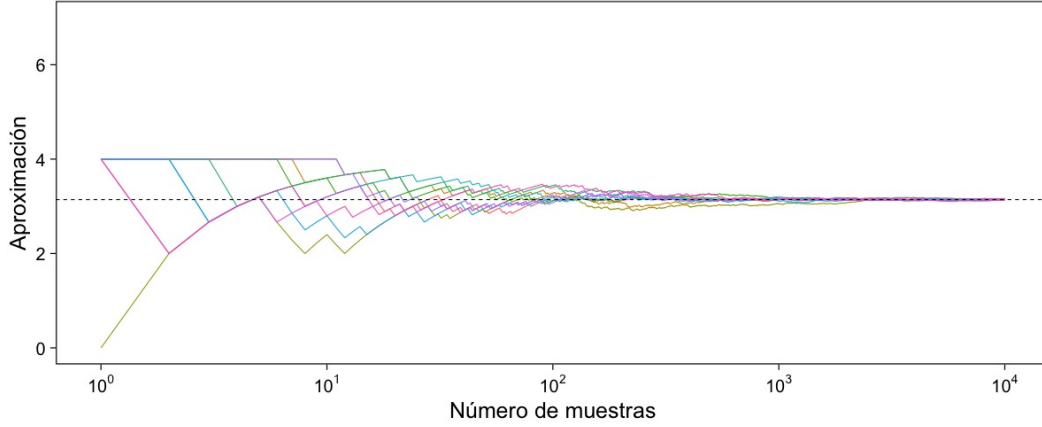


FIGURA 10. Réplica de las trayectorias de diversas realizaciones de la aproximación de la integral.

aproximar dicho comportamiento utilizando distribuciones asintóticas, (TLC).

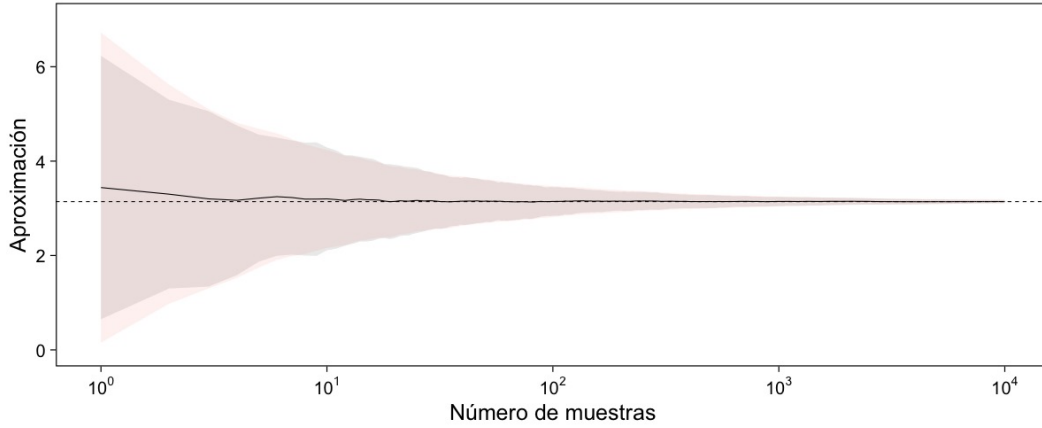


FIGURA 11. Comportamiento promedio e intervalos de confianza.

Podemos explicar la reducción de los intervalos de confianza por medio de la varianza de la estimación de la integral en las distintas réplicas que tenemos. Mas adelante explicaremos de dónde viene la línea punteada. Vemos cómo, aunque captura bien la reducción en varianza, puede sub- o sobre-estimarla.

4.2. Propiedades

A continuación enunciaremos algunas propiedades clave del método Monte Carlo. Poco a poco las iremos explicando y en particular discutiremos algunas de ellas.

4.2.1. Teorema [Error Monte Carlo] Sea $f : \mathbb{R}^p \rightarrow \mathbb{R}$ cualquier función bien comportada[†]. Entonces, el estimador Monte Carlo es **insesgado**. Es decir, se satisface

$$\mathbb{E} \left[\hat{\pi}_N^{\text{MC}}(f) - \pi(f) \right] = 0, \quad (5)$$

para cualquier N . Usualmente estudiamos el error en un escenario pesimista donde medimos el **error cuadrático medio** en el peor escenario

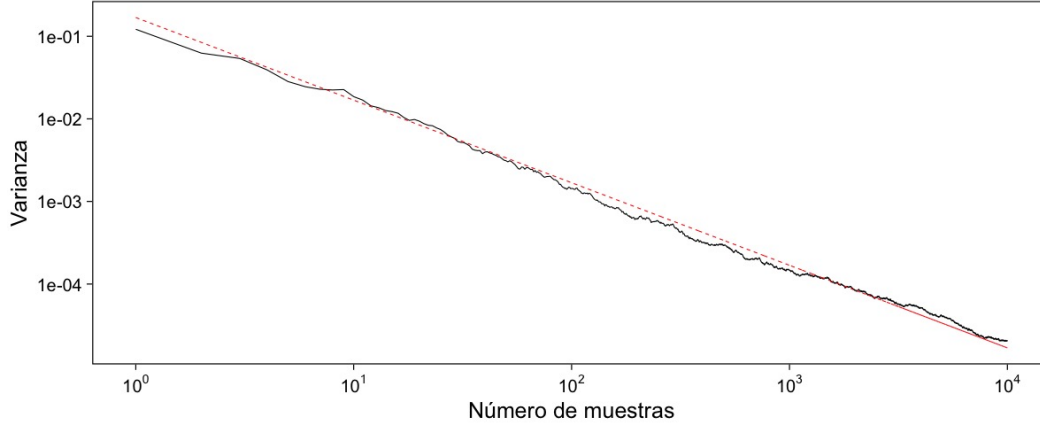


FIGURA 12. Comportamiento promedio e intervalos de confianza.

$$\sup_{f \in \mathcal{F}} \mathbb{E} \left[\left(\hat{\pi}_N^{\text{MC}}(f) - \pi(f) \right)^2 \right] \leq \frac{1}{N}.$$

Esta desigualdad nos muestra una de las propiedades que usualmente se celebran de los métodos Monte Carlo. La integral y nuestra aproximación de ella por medio de simulaciones tiene un error acotado proporcionalmente por el número de simulaciones.

En particular, la varianza del estimador (**error estándar**) satisface la igualdad

$$\text{ee}^2 \left(\hat{\pi}_N^{\text{MC}}(f) \right) = \frac{\mathbb{V}_\pi(f)}{N}.$$

Esta igualdad, aunque consistente con nuestra desigualdad anterior, nos dice algo más. El error de nuestra aproximación **depende** de la varianza de f bajo la distribución π .

4.2.2. Teorema [TLC para estimadores Monte Carlo] Sea f una función **bien comportada** ^{††}, entonces bajo una N suficientemente grande tenemos

$$\sqrt{N} \left(\hat{\pi}_N^{\text{MC}}(f) - \pi(f) \right) \sim \mathcal{N}(0, \mathbb{V}_\pi(f)). \quad (6)$$

4.2.3. Nota: El estimador Monte Carlo del que hablamos, $\hat{\pi}_N^{\text{MC}}(f)$, es una estimación con una **muestra finita de simulaciones**. En ese sentido podemos pensar que tenemos un *mapeo* de muestras a estimador

$$(x^{(1)}, \dots, x^{(N)}) \mapsto \hat{\pi}_N^{\text{MC}}(f), \quad (7)$$

con $x^{(i)} \stackrel{\text{iid}}{\sim} \pi$.

De lo cual es natural pensar: ¿y si hubiéramos observado otro conjunto de simulaciones? Nuestro proceso de estimación es el mismo pero la muestra puede cambiar.

En este sentido nos preguntamos por el **comportamiento promedio** bajo distintas muestras observadas

$$\mathbb{E}[\hat{\pi}_N^{\text{MC}}(f)] = \mathbb{E}_{x_1, \dots, x_N}[\hat{\pi}_N^{\text{MC}}(f)]. \quad (8)$$

De la misma manera nos podemos preguntar sobre la **dispersión alrededor de dicho promedio** (varianza)

$$\mathbb{V}[\hat{\pi}_N^{\text{MC}}(f)] = \mathbb{V}_{x_1, \dots, x_N}[\hat{\pi}_N^{\text{MC}}(f)]. \quad (9)$$

Al ser un ejercicio de **estimación** la desviación estándar del estimador recibe el nombre de **error estándar**. Lo cual denotamos por

$$\text{ee}[\hat{\pi}_N^{\text{MC}}(f)] = \left(\mathbb{V}[\hat{\pi}_N^{\text{MC}}(f)] \right)^{1/2} = \left(\frac{\mathbb{V}_{\pi}(f)}{N} \right)^{1/2}. \quad (10)$$

4.2.4. Nota: Para algunos estimadores la fórmula del error estándar se puede obtener de manera analítica. Para otro tipo, tenemos que utilizar propiedades asintóticas (p.e. cota de Cramer-Rao).

Hay casos en los que no existe una fórmula asintótica o resultado analítico, pero podemos usar simulación [8)] para cuantificar dicha dispersión.

4.2.5. Nota: Hay situaciones en las que la distribución normal asintótica no tiene sentido. Para este tipo de situaciones también veremos cómo podemos utilizar simulación para cuantificar dicha dispersión.

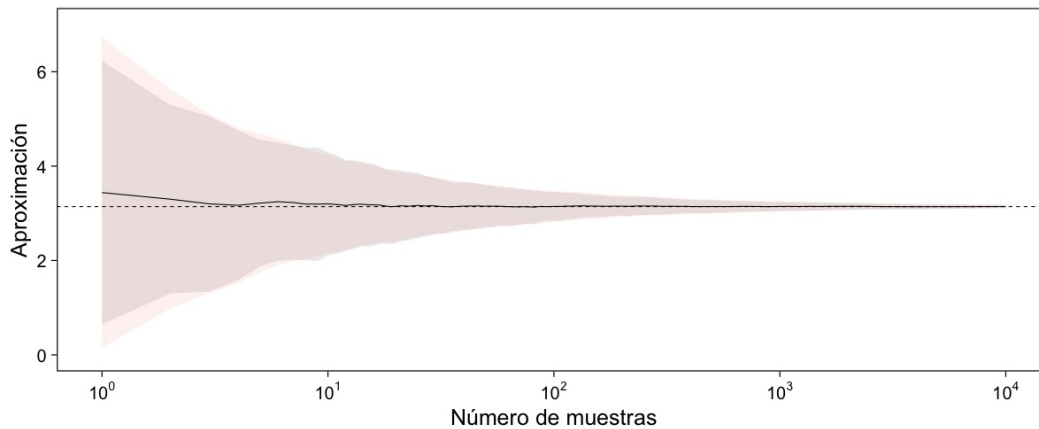


FIGURA 13. Comportamiento promedio e intervalos de confianza con aproximación asintótica.

4.3. Ejemplo: Proporciones

Consideramos la estimación de una proporción θ , tenemos como inicial $p(\theta) \propto \theta$, que es una $\text{Beta}(2, 1)$. Si observamos 3 éxitos en 4 pruebas, entonces sabemos que la posterior es $p(\theta|x) \propto \theta^4(1-\theta)$, que es una $\text{Beta}(5, 2)$. Si queremos calcular la media y el segundo momento posterior para θ , en teoría necesitamos calcular

$$\mu_1 = \int_0^1 \theta p(\theta|X=3) d\theta, \quad \mu_2 = \int_0^1 \theta^2 p(\theta|X=3) d\theta. \quad (11)$$

Utilizando el método Monte Carlo:

```
1 theta <- rbeta(10000, 5, 2)
2 media_post <- mean(theta)
3 momento_2_post <- mean(theta^2)
4 c(mu_1 = media_post, mu_2 = momento_2_post)
```

```
1 mu_1 mu_2
2 0.7131 0.5343
```

Incluso, podemos calcular cosas mas *exóticas* como

$$P(e^\theta > 2|x). \quad (12)$$

```
1 mean(exp(theta) > 2)
```

```
1 [1] 0.5918
```

4.4. Ejemplo: Sabores de helados

Supongamos que probamos el nivel de gusto para 4 sabores distintos de una paleta. Usamos 4 muestras de aproximadamente 50 personas diferentes para cada sabor, y cada uno evalúa si le gustó mucho o no. Obtenemos los siguientes resultados:

```
1 # A tibble: 4 × 4
2   sabor      n gusto prop_gust
3   <chr>    <dbl> <dbl>    <dbl>
4 1 fresa      50    36    0.72
5 2 limon      45    35    0.778
6 3 mango      51    42    0.824
7 4 guanabana  50    29    0.58
```

LISTING 1. Resultados de las encuestas.

Usaremos como inicial $\text{Beta}(2, 1)$ (pues hemos observado cierto sesgo de cortesía en la calificación de sabores, y no es tan probable tener valores muy bajos) para todos los sabores, es decir $p(\theta_i)$ es la funcion de densidad de una $\text{Beta}(2, 1)$. La inicial conjunta la definimos entonces, usando `independencia inicial`, como

$$p(\theta_1, \theta_2, \theta_3, \theta_4) = p(\theta_1)p(\theta_2)p(\theta_3)p(\theta_4).$$

Pues inicialmente establecemos que ningún parámetro da información sobre otro: saber que mango es muy gustado no nos dice nada acerca del gusto por fresa. Bajo este supuesto, y el supuesto adicional de que las muestras de cada sabor son independientes, podemos mostrar que las **posteriores son independientes**:

$$p(\theta_1, \theta_2, \theta_3, \theta_4 | k_1, k_2, k_3, k_4) = p(\theta_1 | k_1)p(\theta_2 | k_2)p(\theta_3 | k_3)p(\theta_4 | k_4)$$

```
1 # A tibble: 4 × 7
2   sabor      n gusto prop_gust a_post b_post media_post
3   <chr>    <dbl> <dbl>    <dbl> <dbl> <dbl>    <dbl>
4 1 fresa      50    36    0.72    38    15    0.717
5 2 limon      45    35    0.778   37    11    0.771
6 3 mango      51    42    0.824   44    10    0.815
7 4 guanabana  50    29    0.58    31    22    0.585
```

LISTING 2. Resultado de inferencia Bayesiana.

Podemos hacer preguntas interesantes como: ¿cuál es la probabilidad de que mango sea el sabor preferido? Para contestar esta pregunta podemos utilizar simulación y responder por medio de un procedimiento Monte Carlo.

```

1 ## Generamos muestras de la posterior
2 paletas <- datos >
3 mutate(alpha = a_post, beta = b_post) >
4 nest(params.posterior = c(alpha, beta)) >
5 mutate(muestras.posterior = map(params.posterior, modelo_beta)) >
6 select(sabor, muestras.posterior)

```

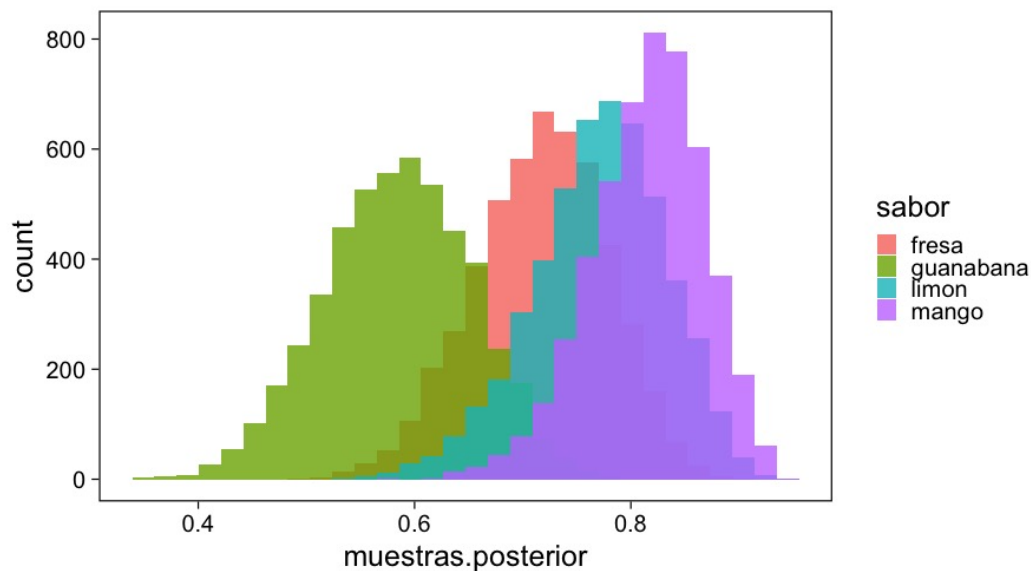


FIGURA 14. Histogramas de la distribución predictiva marginal para cada θ_j .

```

1 ## Utilizamos el metodo Monte Carlo para aproximar la integral.
2 paletas >
3 unnest(muestras.posterior) >
4 mutate(id = rep(seq(1, 5000), 4)) > group_by(id) >
5 summarise(favorito = sabor[which.max(muestras.posterior)]) >
6 group_by(favorito) > tally() >
7 mutate(prop = n/sum(n))

```

```

1 # A tibble: 4 × 3
2   favorito     n   prop
3   <chr>   <int> <dbl>
4 1 fresa     291 0.0582
5 2 guanabana    2 0.0004
6 3 limon    1331 0.266
7 4 mango    3376 0.675

```

LISTING 3. Aproximación Monte Carlo.

Escencialmente estamos preguntándonos sobre calcular la integral:

$$\mathbb{P}(\text{mango sea preferido}) = \int_{\Theta} f(\theta_1, \dots, \theta_4) p(\theta_1, \dots, \theta_4 | X_1, \dots, X_n) d\theta, \quad (13)$$

donde $f(\theta_1, \dots, \theta_4) = \mathbb{I}_{[\theta_4 \geq \theta_j, j \neq 4]}(\theta_1, \dots, \theta_4)$.

4.5. Tarea: Sabores de helados

- ¿Cuál es la probabilidad a priori de que cada sabor sea el preferido?
- Con los datos de arriba, calcula la probabilidad de que la gente prefiera el sabor de mango sobre limón.

5. EXTENSIONES: MUESTREO POR IMPORTANCIA

Incluso cuando tenemos una integral **complicada** podemos **relajar** el problema de integración. De tal forma que podemos **sustituir**

$$\int f(x) \pi(x) dx = \int f(x) \frac{\pi(x)}{\rho(x)} \rho(x) dx = \int f(x) w(x) \rho(x) dx,$$

donde ρ es una densidad de una variable aleatoria **adecuada**.

Esto nos permite utilizar lo que sabemos de las propiedades del método Monte Carlo para resolver la integral

$$\pi(f) = \int f(x) \pi(x) dx = \int f(x) w(x) \rho(x) dx =: \rho(fw),$$

por medio de una aproximación

$$\pi(f) \approx \sum_{n=1}^N \bar{w}^{(n)} f(x^{(n)}), \quad x^{(n)} \stackrel{\text{iid}}{\sim} \rho. \quad (14)$$

Al estimador le llamamos el estimador por importancia y lo denotamos por

$$\pi_N^{\text{IS}}(f) = \sum_{n=1}^N \bar{w}^{(n)} f(x^{(n)}), \quad \bar{w}^{(n)} = \frac{w(x^{(n)})}{\sum_{m=1}^N w(x^{(m)})}. \quad (15)$$

5.1. Propiedades: muestreo por importancia

Lamentablemente, utilizar muestreo por importancia **impacta la calidad de la estimación** (medida, por ejemplo, en términos del **peor error cuadrático medio cometido**). El impacto es un factor que incorpora la **diferencia** entre la distribución **objetivo** –para integrales de la forma $\int f(x) dx$, implica la distribución uniforme– y la distribución **sustituto**. Puedes leer más de esto (aunque a un nivel mas técnico) en la sección 5 de las notas de Sanz-Alonso et al. [3].

5.2. Ejemplo

El análisis del error en la sección anterior habla en del error cuadrático medio en el peor escenario posible bajo una familia de funciones de prueba (resumen). El ejemplo anterior muestra el error Monte Carlo cometido con respecto a una función resumen $f(\theta) = \theta$ con la cual, vemos, se reduce la varianza. Esto no contradice lo anterior pues para esta función resumen nuestra distribución instrumental satisface el criterio de reducción de varianza. En general, lo complicado es encontrar dicha distribución que podamos usar en la estimación Monte Carlo.

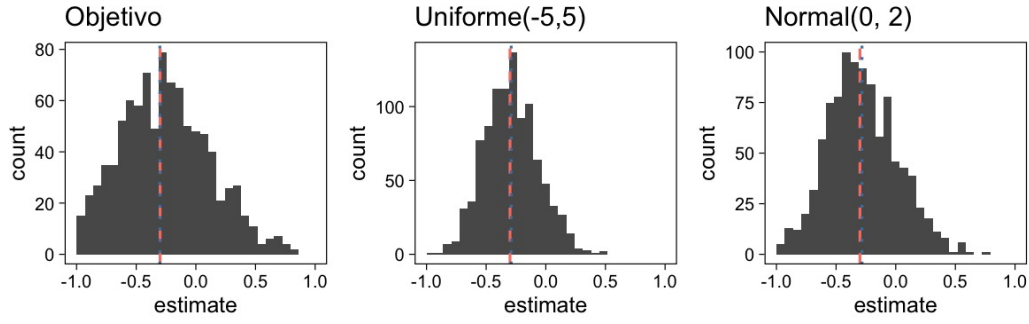


FIGURA 15. Muestreo por importancia utilizando distintas distribuciones instrumentales. Distribución bootstrap de π_N^{IS} con $B = 10,000$ y $n = 100$.

REFERENCIAS

- [1] A. Johnson, M. Ott, and M. Dogucu. *Bayes Rules! An Introduction to Applied Bayesian Modeling*. 2021.
- [2] 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.
- [3] D. Sanz-Alonso, A. M. Stuart, and A. Taeb. Inverse Problems and Data Assimilation. *arXiv:1810.06191 [stat]*, jul 2019.