

# Simulación

## Integración Montecarlo.

Jorge de la Vega Góngora

Departamento de Estadística,  
Instituto Tecnológico Autónomo de México

### Semana 10

When you finally solve the integral you previously  
could only approximate with Monte Carlo methods



ITAM

# Introducción

Dos grandes clases de problemas que surgen en inferencia estadística son:

- ➊ Optimización (ML, EM)
- ➋ Integración (estimadores, valores esperados, modelos Bayesianos)

No siempre es posible resolver analíticamente estos problemas. En ambos casos se requieren soluciones numéricas, porque algunos problemas pueden ser o muy complejos o sin solución analítica.

Para integración, el problema a resolver siempre se puede escribir de la forma:

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

donde  $f$  es la densidad de la variable aleatoria  $X$ , y  $h$  es una función arbitraria.

Es importante notar que  $X$  puede ser un vector de variables aleatorias, y entonces la integral es múltiple.

Se pueden tomar los siguientes enfoques para calcular  $\theta = E[h(X)]$ :

- (a) Encontrar el valor de manera analítica. Es el método preferido si es posible.
- (b) Si la integral no se puede resolver analíticamente, se puede intentar integrales numéricas para tener un valor aproximado de la integral. Este método es bueno en dimensiones bajas, pero en espacios de dimensión grande puede ser muy costoso e ineficiente.
- (c) Integración de Monte Carlo. Esta técnica se basa en la ley de los grandes números, como se verá más adelante.

## Comentarios sobre enfoque (b)

- Hay muchas fórmulas determinísticas de cuadratura para el cálculo de integrales cuando el integrando se comporta “bien”, como la fórmula trapezoidal:

$$\int_a^b f(x) dx \approx (b - a) \left[ \frac{f(a) + f(b)}{2} \right]$$

o la regla de Simpson:

$$\int_a^b f(x) dx \approx \frac{b - a}{6} \left[ f(a) + 4f\left(\frac{a + b}{2}\right) + f(b) \right].$$

- Sin embargo, hay situaciones en donde la función tiene soporte no acotado, no es bien comportada o la integral es multivariada y las fórmulas resultan muy complicadas de aplicar.
- En esos casos, los métodos de Monte Carlo son más simples y dan buenos resultados. Estos métodos fueron inventados en 1946 por Stanislaw Ulam, un matemático polaco que trabajó junto a John von Newman en el proyecto Manhattan durante la Segunda Guerra Mundial.

# Montecarlo

El método crudo de Monte Carlo parte de lo que hemos visto para estimar áreas a través de números uniformes.

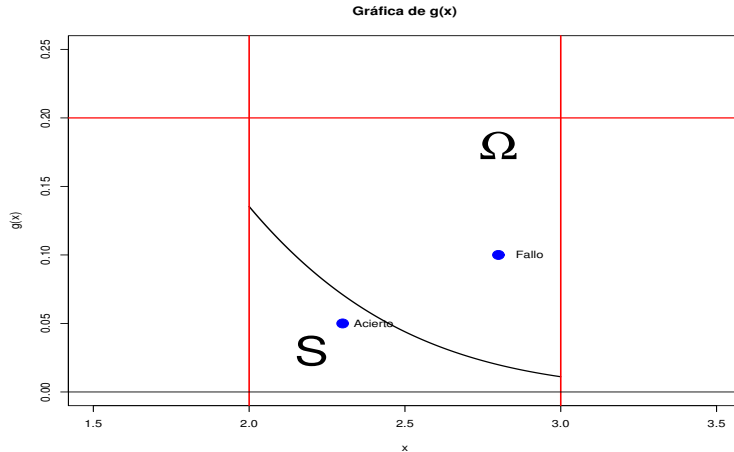
- Comencemos con un ejercicio sencillo. Supongamos que queremos calcular el valor de una integral definida, e.g.

$$\theta = \int_2^3 e^{-\frac{x^2}{2}} dx = \sqrt{2\pi}(\Phi(3) - \Phi(2))$$

El valor 'exacto' de la integral es  $\theta = 0.05364243$ .

- Habla el estadístico: cada integral puede representarse como un *valor esperado* y el problema de estimar una integral vía Monte Carlo es equivalente al problema de estimar un parámetro desconocido  $\theta$ .
- Sea  $\Omega$  el rectángulo  $\{(x, y) | 2 \leq x \leq 3, 0 \leq y \leq 0.20\}$  y sea  $S = \{(x, y) | y \leq g(x) = e^{-0.5x^2}\}$ . Ahora lanzamos dardos "al azar" sobre  $\Omega$ .

# Método crudo de Monte Carlo II



- “Al azar” significa que el vector de coordenadas  $(X, Y)$  tiene una *distribución uniforme* sobre  $\Omega$ .



- El área bajo  $g(x)$  es precisamente la integral que queremos, el área de  $S$ . La probabilidad de que un dardo pegue en  $S$  es:

$$p = \frac{\text{Área } S}{\text{Área } \Omega} = \frac{\theta}{0.2}$$

- Si lanzamos  $N$  dardos, y de éstos  $N_a$  son los que caen en el área  $S$ , podemos estimar a  $p$  con la proporción  $\hat{p} = \frac{N_a}{N}$ .
- Finalmente, una estimación de nuestra integral estaría dada por

$$\hat{\theta} = 0.2 \frac{N_a}{N}$$

- La siguiente tabla muestra valores obtenidos con diferentes lanzamientos.

$N$	$\hat{\theta}$	$ \hat{\theta} - \theta $
10	0.100	0.04635757
50	0.044	0.00964243
100	0.054	0.00035757
1,000	0.0526	0.00104243
10,000	0.05268	0.00096243
100,000	0.053684	4.157e-05

- Hay dos puntos importantes a responder aquí:
  - ¿Cuál es la variación esperada de  $\hat{\theta}$ ?, y
  - ¿Cuántos lanzamientos se requieren para lograr una precisión dada?

## ¿Cuál es la variación esperada de $\hat{\theta}$ ?

- Como cada lanzamiento de dardo acierta o falla, la distribución de cada lanzamiento es una variable de tipo Bernoulli, y suponemos que son independientes. Entonces:

$$\text{Var}(\hat{p}) = \frac{p(1-p)}{N} = \frac{\theta(0.2-\theta)}{0.04N}$$

y por lo tanto, la varianza de  $\hat{\theta}$  es

$$\text{Var}(\hat{\theta}) = 0.04\text{Var}(\hat{p}) = \frac{\theta(0.2-\theta)}{N}.$$

- El error estándar es un estimador de la desviación estándar, reemplazando  $\theta$  por  $\hat{\theta}$ :

$$se(\hat{\theta}) = \sqrt{\frac{\hat{\theta}(0.2-\hat{\theta})}{N}}.$$

## ¿Cuál es la variación esperada de $\hat{\theta}$ ?

- Ahora podemos completar la tabla anterior:

$N$	$\hat{\theta}$	$ \hat{\theta} - \theta $	error estándar = $\hat{\sigma}_{\theta}$
10	0.100	0.04635757	0.03162
50	0.044	0.00964243	0.01172
100	0.054	0.00035757	0.00888
1,000	0.0526	0.00104243	0.00278
10,000	0.05268	0.00096243	0.00088
100,000	0.053684	4.157e-05	0.00028

- La segunda pregunta es más útil en la práctica, pero requiere más teoría de probabilidad.

## ¿Cuántos lanzamientos son necesarios?

- En términos probabilísticos, queremos encontrar  $N$  tal que

$$P[|\theta - \hat{\theta}| < \epsilon] \geq \alpha,$$

donde  $\epsilon$  y  $\alpha$  son determinadas.

- Aplicando la desigualdad de Chebyshev,  $P[|\theta - \hat{\theta}| < \epsilon] \geq 1 - \frac{\text{Var}(\hat{\theta})}{\epsilon^2}$ , obtenemos que

$$\alpha \leq 1 - \frac{\text{Var}(\hat{\theta})}{\epsilon^2}$$

- En  $\text{Var}(\hat{\theta})$  aparece  $N$ . Despejando obtenemos que:

$$N \geq \frac{0.04p(1-p)}{(1-\alpha)\epsilon^2}$$

# ¿Cuántos lanzamientos son necesarios?

- La siguiente tabla muestra los resultados para diferentes combinaciones de  $\epsilon$ ,  $\alpha$  y  $p$ .

$\epsilon$	$\alpha$	$p$	$N$
0.001	0.90	0.5	100
0.00001	0.95	0.01	792
0.00001	0.95	0.6	19,200
0.0001	0.999	0.5	100,000

- La desigualdad de Chebyshev da una estimación conservadora de  $N$ . Usualmente es mejor intentar valores más grandes.
- Con el poder computacional con el que se cuenta actualmente, ya no es limitación importante encontrar valores grandes de  $N$ .
- Este método muestra una aplicación más de la técnica de *aceptación-rechazo*.

- El Método de Monte Carlo es una generalización del método crudo de Monte Carlo, para evaluar la integral sobre muestras de variables aleatorias con otras distribuciones y no sólo de la distribución uniforme.
- Introducir otras distribuciones puede ayudar a acelerar la convergencia y requerir menores tamaños de muestra.
- El soporte de MC sigue siendo la Ley de los grandes números.

El siguiente algoritmo mejora la estimación de  $\theta$ , encontrando una medida donde la integral sea más eficiente.

## Método (mejorado) de Monte Carlo

Para evaluar la integral

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

- Se obtiene una muestra  $x_1, \dots, x_n \sim f$
- Calcula  $\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n h(x_i)$



- La varianza asintótica de la aproximación  $\hat{\theta}_n$  está dada por:

$$\text{Var}(\hat{\theta}_n) = \frac{1}{n} \int_{\mathcal{X}} (h(x) - \theta)^2 f(x) dx,$$

con estimador:

$$\nu_n = \hat{\text{Var}}(\hat{\theta}_n) = \frac{1}{n^2} \sum_{i=1}^n (h(x_i) - \hat{\theta}_n)^2$$

- Además, por el TLC:

$$\frac{\hat{\theta}_n - \theta}{\sqrt{\nu_n}} \sim \mathcal{N}(0, 1)$$

El fundamento para el algoritmo anterior es la Ley fuerte de los grandes números:

## Ley (fuerte) de los grandes números

Si  $X$  es una variable aleatoria con la misma distribución que  $X_i$  y suponiendo que  $h : \mathbb{R} \rightarrow \mathbb{R}$  es una función acotada, entonces  $h(X_1), h(X_2), \dots$  es una sucesión de variables independientes e idénticamente distribuidas con media finita y

$$P \left( \lim_{n \rightarrow \infty} \frac{h(X_1) + \dots + h(X_n)}{n} = E(h(X)) \right) = 1$$

La ley débil establece que la convergencia se da en probabilidad: para cualquier  $\epsilon > 0$ ,

$$\lim_{n \rightarrow \infty} P \left( \left| \frac{h(X_1) + \dots + h(X_n)}{n} - E(h(X)) \right| < \epsilon \right) = 1$$

- Usando el mismo ejemplo que vimos en Monte Carlo crudo:

$$\theta = \int_2^3 e^{-\frac{x^2}{2}} dx = \sqrt{2\pi}(\Phi(3) - \Phi(2))$$

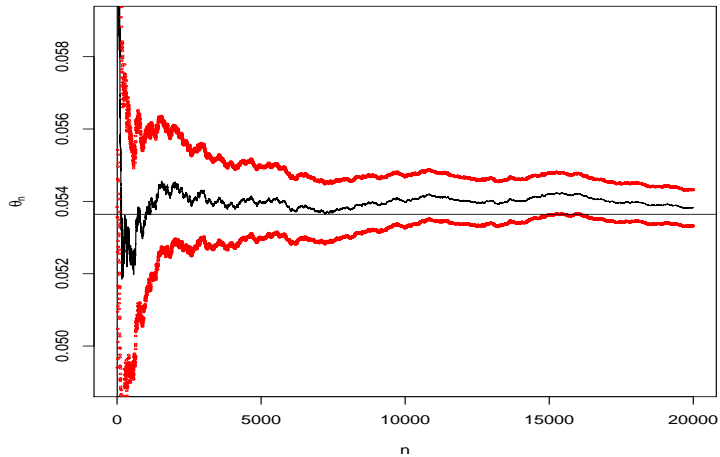
- Podemos considerar generar una muestra aleatoria de valores de una  $\mathcal{U}(2, 3)$ ,

```
n <- 20000 #máximo tamaño de muestra
h <- function(x)exp(-x^2/2)
x <- h(runif(n,2,3)) #Genera una muestra grande de observaciones
theta_n <- cumsum(x)/(1:n) #estimadores de la integral para diferentes n
vn <- sqrt(cumsum((x-theta_n)^2))/(1:n) #error estándar estimado para cada n

plot(theta_n,type="l",ylim=mean(theta_n)+c(-1,1)*20*vn[n],
      ylab=expression(theta[n]),
      xlab="n")

#Agrega líneas correspondientes a nivel de confianza del 95% para cada n
points(theta_n+c(-1,1)*2*vn,col="red",cex=0.3)
abline(h=sqrt(2*pi)*(pnorm(3)-pnorm(2)))
```

## Ejemplo MC II



- Alternativamente, podemos generar muestras de una  $\mathcal{N}(0, 1)$  (escalada por la constante) y considerar la función indicadora en el intervalo  $(2, 3)$ :

$$\theta = \int_2^3 e^{-\frac{x^2}{2}} dx = \sqrt{2\pi} \int_{-\infty}^{\infty} \phi(x) I(2 < x < 3) dx$$

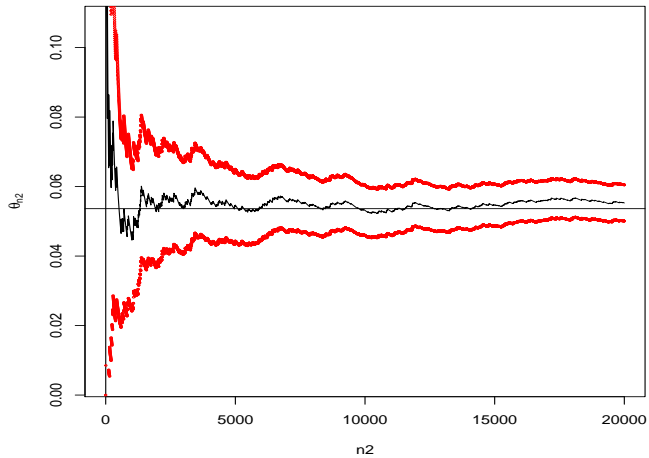
- ¿Converge a la misma velocidad?

```
#Mismo ejercicio, invirtiendo las distribuciones
g <- function(x) ifelse(x>=2 & x<=3, 1, 0) #función indicadora
x <- sqrt(2*pi)*g(rnorm(n))
theta_n2 <- cumsum(x)/(1:n)
vn2 <- sqrt(cumsum((x-theta_n2)^2))/(1:n) #error estándar estimado para cada n

plot(theta_n2, type="l", ylim=mean(theta_n2)+c(-1,1)*20*vn2[n],
      ylab=expression(theta[n2]),
      xlab="n2")

#Agrega líneas correspondientes a nivel de confianza del 95% para cada n
points(theta_n2+c(-1,1)*2*vn2, col="red", cex=0.3)
abline(h=sqrt(2*pi)*(pnorm(3)-pnorm(2)))
```

## Ejemplo MC II



## Ejemplo 2 I

Ejercicio 3.1, Casella y Robert

El siguiente ejercicio muestra que puede haber algunos problemas cuando  $\nu_n$  no es un estimador adecuado de la varianza de  $\hat{\theta}_n$  o cuando no converge o converge de manera muy lenta. En estos casos, el estimador y la región de confianza asociada no será de confianza.

Supongan que tenemos una observación  $X \sim \mathcal{N}(\theta, 1)$  y una distribución inicial para  $\theta \sim \text{Cauchy}(0, 1)$ . Queremos actualizar la información de  $\theta$  basada en la información que provee  $X$ . En este contexto, la verosimilitud es:

$$f(x|\theta) = \frac{1}{\sqrt{2\pi}} \exp \left[ -\frac{1}{2}(x - \theta)^2 \right]$$

con distribución inicial

$$\pi(\theta) = \frac{1}{\pi(1 + \theta^2)}$$

Usando el teorema de Bayes, la distribución posterior es proporcional a la verosimilitud por la inicial, esto es:

$$\pi(\theta|x) \propto \exp \left[ -\frac{1}{2}(x - \theta)^2 \right] \frac{1}{(1 + \theta^2)},$$

## Ejemplo 2 II

Ejercicio 3.1, Casella y Robert

donde la constante de proporcionalidad es el recíproco de  $C$  con

$$C = \int \exp \left[ -\frac{1}{2}(x - \theta)^2 \right] \frac{1}{(1 + \theta^2)} d\theta.$$

El estimador puntual para  $\theta$  (también conocido como estimador de Bayes para el modelo normal-Cauchy) es la media posterior, dada por

$$\delta(x) = E(\theta|x) = \frac{\int \theta \exp \left[ -\frac{1}{2}(x - \theta)^2 \right] \frac{1}{(1 + \theta^2)} d\theta}{C}$$

El ejercicio pide resolver la ecuación para  $x = 0, 2, 4$ .



# Ejemplo 2 I

Ejercicio 3.1, Casella y Robert

- Grafica los integrandos, y usa integración de Monte Carlo basada en la simulación Cauchy para calcular las integrales.

Noten que para  $C$  tenemos a la Cauchy multiplicada por *algo*, donde *algo* es

$$h(\theta) = \pi \exp \left[ -\frac{1}{2}(x - \theta)^2 \right]$$

asi que podemos aplicar MC con el siguiente algoritmo:

Dado el valor observado  $X = x$

1. simular una muestra de Cauchys:  $\theta_1, \theta_2, \dots, \theta_n$ .
2. Estimar la integral en el denominador con  $\frac{1}{n} \sum_{i=1}^n \pi \exp \left[ -\frac{(x-\theta_i)^2}{2} \right]$ .
3. Estimar la integral en el numerador con  $\frac{1}{n} \sum_{i=1}^n \pi \theta_i \exp \left[ -\frac{(x-\theta_i)^2}{2} \right]$ .
4. Define la razón  $\delta(x)$ .

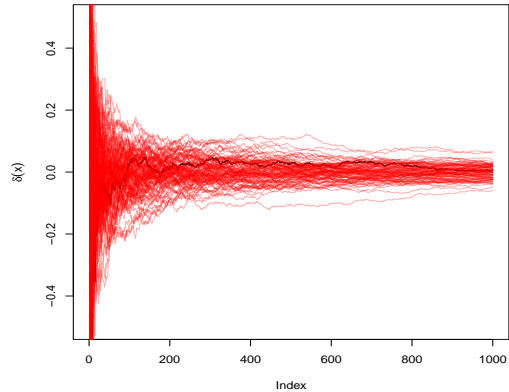
# Ejemplo 2 II

## Ejercicio 3.1, Casella y Robert

```
delta <- function(x){  
  #estimador de Bayes  
  n <- 1000 #número de valores simulados  
  x <- 0 #valor de la variable observada  
  th <- rcauchy(n)  
  h <- function(x){pi*exp(-0.5*(x-th)^2)}  
  cumsum(th*h(x))/cumsum(h(x))  
}  
rojo=rgb(1,0,0,alpha=0.3)  
plot(delta(1),type="l",ylim=c(-0.5,0.5),ylab=expression(delta(x)),lwd=2)  
for(i in 1:100)lines(delta(1),type="l",col=rojo)
```

# Ejemplo 2 III

Ejercicio 3.1, Casella y Robert



- Podemos escribir en el caso multivariado

$$\begin{aligned}\theta &= \int_{\Omega} f dV \\ &\approx \text{Volumen de } \Omega \times \text{promedio de } f \text{ en } \Omega\end{aligned}$$

- Para integración unidimensional, los métodos de MC son muy ineficientes. Pero para integración multidimensional es muy eficiente.
- Como el error es proporcional a  $1/\sqrt{n}$ , no depende de la dimensión.
- Se pueden manejar fácilmente regiones con fronteras irregulares.

## Ejemplo 3 I

### Monte Carlo multidimensional

Evaluar la integral

$$\theta = \int \int_{\Omega} \text{sen}(\sqrt{\log(x+y+1)}) dx dy$$

donde  $\Omega$  es el disco definido por la condición  $(x - 0.5)^2 + (y - 0.5)^2 \leq 0.25$ .

***Solución.***

Considerando que la región de integración está en el cuadro unitario, es posible simular de uniformes independientes y quedarnos con los valores que estén dentro de la región de integración.

# Ejemplo 3 II

## Monte Carlo multidimensional

```
N <- 1e6 #numero de simulaciones
X <- cbind(runif(N),runif(N)) #matriz de uniformes
#solo nos quedamos con las observaciones que están en la región considerada
X1 <- X[(X[,1]^2+(X[,2]-0.5)^2 <=0.25,]
N1 <- dim(X1)[1] # dimensión que queda con sólo aceptados
thetahat <- sin( sqrt( log(X1[,1] + X1[,2] + 1)))
a <- (pi/4)*cumsun(thetahat)/(1:N1) #se ajusta por el volumen del disco
a[N1]

[1] 0.5677426

sd(thetahat)/sqrt(N1) #error estándar

[1] 9.170774e-05
```



## Ejemplo 4 I

En la aproximación

$$\theta \approx \text{Volumen de } \Omega \times \text{promedio de } f \text{ en } \Omega$$

se puede estimar el volumen de la región  $\Omega$  al mismo tiempo que se estima la función  $f$ .

Se generan puntos en un volumen  $V$  que puede ser usualmente rectangular, tal que  $\Omega \subseteq V$ , y se generan  $n$  puntos en  $V$  de los cuales  $k$  están también en  $\Omega$ , entonces

$$\text{Vol}(\Omega) \approx \frac{k}{n} \text{Vol}(V)$$

Como el promedio de  $f$  en  $\Omega$  es aproximado a  $\frac{1}{k} \sum f$ , entonces las  $k$  se cancelan y se tiene:

$$\theta \approx \text{Vol}(V) \times \frac{1}{n} \sum_{p_i \in \Omega} f(p_i)$$

Por ejemplo, evaluar:  $\theta = \int \int_{\Omega} y dx dy$  con  $\Omega$  la semi-elipse  $\Omega$  dada por

$$x^2 + 4y^2 \leq 1, \quad y \geq 0$$

## Ejemplo 4 II

Se consideramos  $V$  como el rectángulo  $[-1, 1] \times [0, 0.5]$ , podemos generar  $X \sim \mathcal{U}(-1, 1)$  y  $y \sim \mathcal{U}(0, 0.5)$  El área del rectángulo es 1. Entonces:

```
N <- 1e6 #numero de simulaciones
X <- cbind(runif(N,-1,1),runif(N,0,0.5)) #matriz de uniformes
#solo nos quedamos con las observaciones que están en la región considerada
X1 <- X[X[,1]^2 + 4*X[,2]^2 <= 1,]
N1 <- dim(X1)[1] # dimensión que queda con sólo aceptados
thetahat <- X1[,2]
a <- sum(thetahat)/N
a

[1] 0.166546

sd(thetahat)/sqrt(N) #error estándar

[1] 0.0001321448
```

El valor exacto de la integral es  $1/6$ . (tarea)



# Ejemplo 1 - I

Dagpunar 1.1, Casella-Robert. 3.1

Consideremos la función gamma:

$$\Gamma(\lambda) = \int_0^{\infty} x^{\lambda-1} e^{-x} dx$$

y calculemos su valor de dos formas diferentes:

- usando la función `integrate`
- usando el método crudo de Monte Carlo.

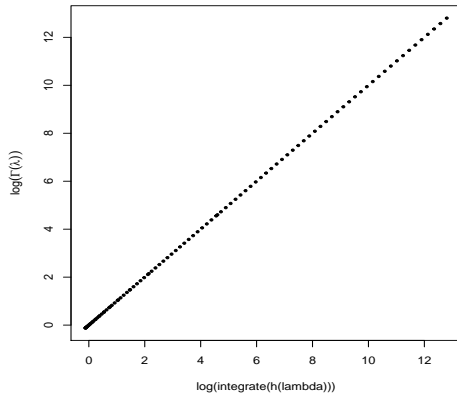
Forma 1:

En este ejemplo, la función `integrate` funciona bien.

```
# Usando la función integrate y comparando contra la función (log) gamma:
h <- function(lambda){integrate(function(x){x^{lambda-1}*exp(-x)},0,Inf)$val}
x <- seq(0.01,10,length=100) #soporte
par(pty="s") #hacemos el plot cuadrado
plot(lgamma(x),log(apply(as.matrix(x),1,h)), xlab = "log(integrate(h(lambda)))",
  ylab = expression(log(Gamma(lambda))),
  pch=19, cex=0.5)
```

# Ejemplo 1 - II

Dagpunar 1.1, Casella-Robert. 3.1



# Ejemplo I

Dagpunar 1.1, Casella-Robert. 3.1

Forma 2:

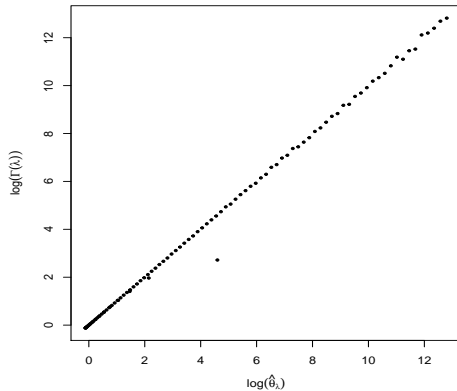
Ahora consideremos la estimación usando Monte Carlo Crudo.

- Podemos usar una variable aleatoria  $X$  exponencial:  $X \sim \exp(1)$  donde  $f(x) = e^{-x}$  en  $[0, \infty)$ .
- Entonces podemos escribir a  $\theta = E_f(X^{\lambda-1})$ .
- Extraemos una muestra de la distribución exponencial con parámetro 1 y calculamos el promedio  $\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n X_i^{\lambda-1}$ . Aquí usamos una muestra de tamaño  $n = 1,000,000$ .

```
# Usando Monte Carlo crudo y comparando contra la función (log) gamma:
h <- function(lambda){mean(rexp(1e6,1)^(lambda-1))}
x <- seq(0.01,10,length=100) #soporte
par(pty="s") #hacemos el plot cuadrado
plot(lgamma(x),log(apply(as.matrix(x),1,h)),
     xlab = expression(log(hat(theta)[lambda])),
     ylab = expression(log(Gamma(lambda))),
     pch=19, cex=0.5)
```

# Ejemplo II

Dagpunar 1.1, Casella-Robert. 3.1



## Variabilidad del estimador $\hat{\theta}_n$

- Como vemos, la estimación en general no es mala, pero tenemos variabilidad que depende de la muestra. El estimador  $\hat{\theta}_n$  es un estimador insesgado de  $\theta$ :

$$E(\hat{\theta}_n) = E\left(\frac{1}{n} \sum_{i=1}^n X_i^{\lambda-1}\right) = E(X_1^{\lambda-1}) = \theta$$

- Si la muestra es independiente, la varianza de  $\hat{\theta}_n$  está dada por:

$$\text{Var}_f(\hat{\theta}_n) = \frac{1}{n} \text{Var}_f(X^{\lambda-1})$$

y su desviación estándar es:  $\sigma_f(\hat{\theta}_n) = \sigma_f(X^{\lambda-1})/\sqrt{n}$ . El error estándar es el estimador de la desviación estándar:

$$se(\hat{\theta}_n) = \hat{\sigma}_f(X^{\lambda-1})/\sqrt{n}.$$

- Noten que para cambiar el error estándar en un factor de  $K$ , se requiere que la muestra cambie por un factor de  $1/K^2$ , lo que hace ineficiente el proceso:

$$K \cdot se(\hat{\theta}_n) = K \hat{\sigma}_f(X^{\lambda-1})/\sqrt{n} = \hat{\sigma}_f(X^{\lambda-1})/\sqrt{n/K^2}$$

## Ejemplo 2

Consideremos un caso particular de la integral anterior, con  $\lambda = 1.9$ :

$$\theta = \int_0^{\infty} x^{0.9} e^{-x}$$

Considerando una muestra de  $n = 100$  observaciones:

```
x <- rexp(100,1) #genera exponenciales con parámetro 1:
theta <- mean(x^{0.9})
theta
```

```
[1] 0.9331481
```

```
sd(x) #desviación estándar muestral
```

```
[1] 0.9937605
```

```
sd(x)/sqrt(100) #error estándar
```

```
[1] 0.09937605
```

¿De qué tamaño tiene que ser la muestra para reducir el error estándar a 0.0001? Como  $K = 0.08638555/0.0001 = 863.8556$ , entonces el tamaño de muestra tiene que ser del orden de  $100 \times 863.86^2 = 74,625,410$

# Sobre la función `integrate` y `area`

- Ambas funciones sólo son para integrales unidimensionales.
- La función `integrate` utiliza un método de cuadratura. Si la función a integrar es casi constante (o cero) en su rango, es posible que el resultado de la estimación y su error puedan ser muy equivocados. Esta función es muy frágil.
- La función `area` es parte del paquete MASS, no acepta límites infinitos, por lo que se requiere conocer de antemano el comportamiento de la función en la región de integración.
- Más adelante veremos un ejemplo de problemas que pueden surgir con ambas funciones, pero antes tenemos que revisar temas Bayesianos, para introducir el contexto de los ejemplos.

## Integración en el contexto Bayesiano



- El paradigma Bayesiano se basa en el teorema de Bayes:

$$P(A|B) = \frac{P(A, B)}{P(B)} = \frac{P(B|A)P(A)}{\int P(B|A)P(A) dA} \propto P(B|A)P(A)$$

- Este teorema se aplica para desarrollar un sistema de aprendizaje: Una persona modifica su afirmación o creencia inicial de probabilidad sobre los parámetros antes de observar datos  $y = (y_1, \dots, y_n)$  a un conocimiento posterior o actualizado que combina el conocimiento inicial y los datos que se observan.
- Sea  $\theta$  un vector de parámetros. El conocimiento inicial se  $\theta$  se resume en una distribución inicial  $\pi(\theta)$ . La verosimilitud es  $f(y|\theta)$  y el conocimiento actualizado está contenido en la distribución posterior  $\pi(\theta|y)$ . Aplicando el teorema de Bayes,

$$\pi(\theta|y) = \frac{f(y|\theta)\pi(\theta)}{m(y)} \propto f(x|\theta)\pi(\theta)$$

donde  $m(y)$  es la verosimilitud marginal. Este valor se puede expresar como una integral:  $m(y) = \int_{\Theta} f(y|\theta)\pi(\theta) d\theta$ , pero al no depender de  $\theta$  se puede considerar como una constante para normalizar  $\pi(\theta|y)$  y garantizar que sea una densidad.

- Noten que las integrales que se menciona en este modelo ya no son integrales unidimensionales, sino de la dimensión del vector de parámetros  $\theta$ .
- Con la distribución posterior surgen varias cantidades de interés a estimar. A partir de la distribución posterior, se pueden estimar funciones de  $\theta$ , usualmente de la forma  $E[h(\theta|y)]$ .
  - Probabilidad posterior de que  $h(\theta)$  esté en un cierto conjunto  $A$ :

$$P(h(\theta) \in A|y) = \frac{\int_{h(\theta) \in A} \pi(\theta)f(y|\theta) d\theta}{\int \pi(\theta)f(y|\theta) d\theta}$$

- Distribuciones marginales del vector  $\theta$ :

$$\pi(\theta_1) \propto \int \pi(\theta_{-j}, \theta_j) d\theta_{-j}$$

- Densidades predictivas:

$$m(\tilde{y}) = \int f(\tilde{y}|\boldsymbol{\theta})\pi(\boldsymbol{\theta}) d\boldsymbol{\theta}.$$

Queda claro que todos estos casos generan problemas de integración.

# Otro Ejemplo I

Casella-Robert 3.2

El siguiente ejemplo parte de un contexto Bayesiano, y podemos ver el comportamiento y problemas de las funciones `integrate` y `area`.

- Ahora se tiene una muestra de tamaño  $n = 10$  de una distribución Cauchy con parámetro de localización  $\theta = 350$ . La marginal de la muestra bajo una distribución inicial  $\pi(\theta)$  queda como:

$$m(\mathbf{x}) = \int_{-\infty}^{\infty} f(\mathbf{x}|\theta)\pi(\theta) d\theta = \int_{-\infty}^{\infty} \prod_{i=1}^{10} \frac{1}{\pi[1 + (x_i - \theta)^2]} d\theta,$$

considerando que la inicial es plana. En este caso, la función `integrate` no funciona bien, comparado con la función `area`:

# Otro Ejemplo II

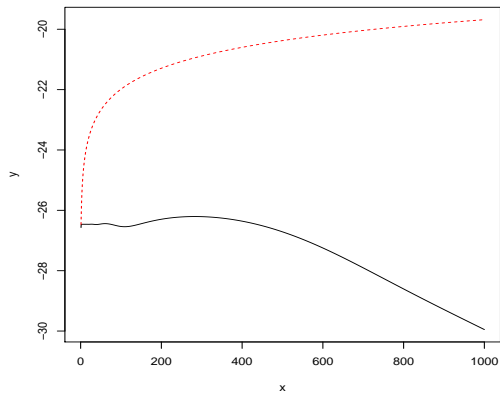
Casella-Robert 3.2

```
library(MASS) #carga para la función area
cac <- rcauchy(10)
verosimilitud <- function(th){
  #define la función de verosimilitud, a partir de la muestra dada, como función de theta
  u <- dcauchy(cac[1]-th)
  for(i in 2:10){
    u <- u*dcauchy(cac[i]-th)
  }
  return(u)
}

f1 <- function(a){integrate(verosimilitud,-a,a)$val}
f2 <- function(a){area(verosimilitud,-a,a)}
x <- seq(1,1000,length=10^4) #partición de intervalo de 1 a 1000 muy fina.
y <- log(apply(as.matrix(x),1,f1))
z <- log(apply(as.matrix(x),1,f2))
plot(x,y, type = "l", ylim = range(cbind(y,z)))
lines(x, z, lty = 2, col = "red")
```

# Otro Ejemplo III

Casella-Robert 3.2



## Técnicas de Reducción de varianza

- El objeto de las técnicas de reducción de varianza en los métodos de Montecarlo es mejorar la velocidad y eficiencia estadística de un estudio de simulación.
- Un estudio de simulación que utiliza entradas aleatorias genera salidas aleatorias, y por lo tanto tienen variabilidad que es necesario comprender.
- Por otra parte un estudio puede consumir muchos recursos como grandes cantidades de números aleatorios, y múltiples replicaciones para obtener un número puede ser demasiado costoso, por lo que es importante tratar de minimizar su variabilidad para obtener mayor precisión.



- Usualmente en simulación estocástica se desea estimar  $\theta = E(h(\mathbf{X}))$ . Como hemos visto, el algoritmo de simulación estándar dice:
  - 1 Genera  $\mathbf{X}_1, \dots, \mathbf{X}_n$
  - 2 Estima  $\theta$  con  $\hat{\theta}_n = \sum_{i=1}^n Y_i$  donde  $Y_i = h(\mathbf{X}_i)$
- Un intervalo de confianza aproximado del  $100(1 - \alpha) \%$  está dado por:

$$\left[ \hat{\theta}_n - z_{1-\alpha/2} \frac{\hat{\sigma}_n}{\sqrt{n}}, \hat{\theta}_n + z_{1-\alpha/2} \frac{\hat{\sigma}_n}{\sqrt{n}} \right]$$

donde  $\hat{\sigma}_n = \hat{\text{Var}}(\hat{\theta}_n) = \frac{\sum (Y_j - \bar{Y})^2}{n-1}$

- Una forma de medir la calidad de un estimador es con la *longitud media HW* del intervalo de confianza,

$$HW = z_{1-\alpha/2} \frac{\hat{\sigma}_n}{\sqrt{n}}$$

- Usualmente se desea que *HW* sea lo más pequeño posible, pero esto es a veces difícil de lograr.

- Una forma de reducir la longitud media puede ser a través de las **Técnicas de reducción de varianza**. Éstas incluyen:
  - Variadas antitéticas
  - Variadas de control
  - Condicionamiento
  - Números pseudo-aleatorios comunes
  - Importance sampling
  - Muestreo estratificado
- Hay muchas otras técnicas sofisticadas o generalizaciones de las que se mencionan aquí.
- En la mayoría de los casos, consideraremos la siguiente situación:  $X$  es una variable aleatoria *de salida*, es decir, es una variable relevante del estudio de simulación. En muchas situaciones lo que se desea es estimar  $E(X) = \mu$ , la media del proceso. Típicamente, entonces, este número es una integral.

- Este método trata de inducir correlación negativa entre series de números pseudo-aleatorios.
- La idea básica es generar pares de corridas de un modelo tal que una observación pequeña en la primera corrida tiende a compensarse por una observación grande en la otra corrida, para que se obtenga una correlación negativa.
- Por ejemplo, si  $u_k \sim \mathcal{U}(0, 1)$  fue generado para obtener un parámetro particular de la primera corrida (por ejemplo, un tiempo de espera) entonces  $1 - u_k$  se usa para obtener el mismo parámetro en la segunda corrida. En este caso se dice que  $u_k$  y  $1 - u_k$  están *sincronizados* en el sentido de que fueron utilizados para el mismo propósito (en este caso generar el mismo parámetro del modelo).

- Supóngase que  $(X_1^{(1)}, X_2^{(1)}, \dots, X_n^{(1)})$  y  $(X_1^{(2)}, X_2^{(2)}, \dots, X_n^{(2)})$  son 2 corridas de una simulación, en donde  $X_j^{(1)}$  se estimó con  $u$  y  $X_j^{(2)}$  se estimó con  $1 - u$ . Noten que:
  - $E(X_j^{(1)}) = E(X_j^{(2)}) = \mu$
  - Se tienen pares independientes:  $(X_{j_1}^{(1)}, X_{j_1}^{(2)}) \perp\!\!\!\perp (X_{j_2}^{(1)}, X_{j_2}^{(2)})$ .
  - Definan  $X_j = \frac{X_j^{(1)} + X_j^{(2)}}{2}$  y  $\bar{X} = \sum_{j=1}^n X_j$ .

Entonces:

- $\bar{X}$  es un estimador insesgado de  $\mu$ , y
- $\text{Var}(\bar{X}) = \text{Var}(X_1)/n = \frac{\text{Var}(X_j^{(1)}) + \text{Var}(X_j^{(2)}) + 2\text{Cov}(X_j^{(1)}, X_j^{(2)})}{4n}$
- De esta forma, si se induce correlación negativa entre  $X_j^{(1)}$  y  $X_j^{(2)}$ , entonces se puede reducir la varianza del estimador  $\bar{X}$ .
- Sin embargo, *no siempre* se puede garantizar que se logre el objetivo, depende del modelo. A veces se puede elaborar un estudio piloto para medir la magnitud de la reducción.

Si queremos estimar  $\theta = \int_a^b f(x)dx$  por el método de Montecarlo crudo, entonces

$$\hat{\theta} = \frac{(b-a)}{n} \sum_{i=1}^n f(x_i),$$

donde  $x_i \sim \mathcal{U}(a, b)$ . Si por cada  $x_i$  se usa su variable antitética  $\tilde{x}_i = a + (b - x_i)$ , entonces el estimador se convierte en

$$\hat{\theta} = \frac{(b-a)}{n} \sum_{i=1}^{n/2} (f(x_i) + f(\tilde{x}_i))$$

Como probamos antes, como la varianza de la suma es la suma de las covarianzas mas dos veces la covarianza y la covarianza es negativa para variables antitéticas, entonces se reduce la varianza de la suma.

# Variadas de control: Idea básica I

- Supóngase que queremos estimar  $\mu = E(X)$  donde  $X$  es una variable aleatoria de salida, como indicamos previamente.
- En lugar de estimar  $\mu$  directamente, se considera la diferencia entre el problema de interés y un modelo analítico:  $Y$  es otra variable relacionada con  $\mu$ , y que está correlacionada con  $X$ , pero además se conoce  $\nu = E(Y)$ . A  $Y$  se le llama **variada de control** para  $X$ .
- Sea  $X_c = X - a(Y - \nu)$  una nueva variable. Entonces
  - $E(X_c) = E(X) = \mu$ , por lo que  $X_c$  es un estimador insesgado de  $\mu$ .
  - $\text{Var}(X_c) = \text{Var}(X - a(Y - \nu)) = \text{Var}(X) + a^2 \text{Var}(Y) - 2a \text{Cov}(X, Y)$  Entonces:

$$\text{Var}(X_c) \leq \text{Var}(X) \text{ si } 2a \text{Cov}(X, Y) > a^2 \text{Var}(Y).$$

La varianza mínima se alcanza si

$$a^* = \frac{\text{Cov}(X, Y)}{\text{Var}(Y)}$$

En este caso  $\text{Var}(X_c) = (1 - \rho_{X,Y}^2) \text{Var}(Y)$ .

- En la práctica, se puede conocer o no el valor de  $\text{Var}(Y)$  y muy difícilmente conocemos  $\text{Cov}(X, Y)$ , por lo que es difícil conocer el valor de  $a$ .

- Una alternativa es estimar  $a$  a través de un estudio piloto (Lavenberg, Moeller y Welch, 1982):

$$\hat{a} = \frac{C_{\bar{X}, Y}}{S_Y^2}$$

y entonces  $\bar{X}_c^* = \bar{X} - \hat{a}(\bar{Y} - \nu)$ . Noten que  $\bar{X}_c^*$  ya no es un estimador insesgado de  $\mu$ . Para reducir el sesgo se puede utilizar, por ejemplo el *jackknife*.

- Se puede mostrar (ejercicio) que el método de variadas antitéticas es un caso particular del método de variadas de control.

## Ejemplo (a)

Supongan que  $X \sim \mathcal{N}(0, 1)$  y se requiere estimar  $E\left(\frac{X^6}{1+X^2}\right)$ . Como

$$\frac{x^6}{1+x^2} = x^4 - x^2 + 1 - \frac{1}{1+x^2}$$

entonces podemos aproximar  $\frac{x^6}{1+x^2}$  con  $Y = g(x) = x^4 - x^2 + 1$ . Para esta  $Y$ ,  
 $E(Y) = E(X^4) - E(X^2) + 1 = 3 - 1 + 1 = 3$ , ya que en una normal estándar  $E(X^4) = 3$  (curtosis).  
De este modo,

$$\theta = E\left(\frac{X^6}{1+X^2} - (X^4 - X^2 + 1)\right) + 3 = 3 - E\left(\frac{1}{1+X^2}\right)$$

Así que podemos aplicar Montecarlo crudo sólo a la función  $h(x) = \frac{1}{1+x^2}$  muestreando de una normal estándar.



## Ejemplo (b)

Se desea estimar  $\theta = E\left(e^{(U+W)^2}\right)$  donde  $U, W \sim \mathcal{U}(0, 1)$  y son independientes. Sea  $X = e^{(U+W)^2}$ . Elegimos una variable de control  $Y$ .

Una posible variable de control es  $Y_1 = U + W$ . Por la distribución de  $U$  y  $W$ , sabemos que  $\nu_1 = E(Y_1) = 1$  y se puede ver que  $\text{Cov}(X, Y_1) > 0$ . Otra posibilidad es usar  $Y_2 = (U + W)^2$ , y  $E(Y_2) = 7/6$ .

El siguiente código muestra como se puede hacer un pequeño piloto en R.

# Código R

```
#Primero hacemos un pequeño piloto
p <- 100;n <- 1000
u <- runif(p);w <- runif(p)
x <- exp((u+w)^2)
y <- (u+w)^2
covest <- cov(cbind(x,y))
a <- -covest[1,2]/covest[2,2]

# Ahora hacemos la simulación
u <- runif(n);w <- runif(n)
x <- exp((u+w)^2); y <- (u+w)^2
v <- x + a*(y-7/6)
estimadorusual <- c(mean(x),sd(x))
estimadorcont <- c(mean(v),sd(v))
CI <- c(mean(v)-1.96*sd(v)/sqrt(n),mean(v)+1.96*sd(v)/sqrt(n))
estimadorusual

[1] 4.866310 5.849764

estimadorcont

[1] 4.889896 2.853162

CI

[1] 4.713055 5.066737
```

# Condicionamiento: idea básica I

- En algunos modelos es posible reemplazar un estimado de alguna variable por su valor analítico exacto.
- Al remover esta fuente de variabilidad, se espera una salida más estable, aunque de nuevo, *no siempre hay garantía de que cumpla el objetivo*.
- Supongamos otra vez que  $X$  es una variable de salida de una simulación y  $E(X) = \mu$  es lo que se quiere estimar. Si  $Z$  es otra variable aleatoria tal que se conoce analíticamente la esperanza condicional  $E(X|Z = z)$ , entonces

$$\mu = E(X) = E(E(X|Z)).$$

Además,

$$\text{Var}(E(X|Z)) = \text{Var}(X) - E(\text{Var}(X|Z)) \leq \text{Var}(X).$$

- Entonces, conviene que  $Z$  tenga las siguientes propiedades:
  - $Z$  puede ser generado de manera eficiente.
  - $E(X|Z = z)$  se puede calcular analíticamente
  - $E(\text{Var}(X|Z))$  tiene un valor grande.

Supongan un modelo jerárquico de la siguiente forma:  $W \sim \mathcal{P}(10)$ , y  $X|W \sim \text{Beta}(w, w^2 + 1)$ . El problema es encontrar  $\theta = \mathbb{E}(X)$ .

Como sabemos que  $\mathbb{E}(X|W = w) = \frac{w}{w^2 + w + 1}$ , entonces basta con muestras  $W_1, W_2, \dots, W_n$  y construir

$$\tilde{\theta} = \frac{1}{n} \sum_{j=1}^n \mathbb{E}(X|W = w_j) = \frac{1}{n} \sum_{j=1}^n \frac{w_j}{w_j^2 + w_j + 1}.$$

Noten que en este ejemplo, no tuvimos que generar ningún valor de  $X$ , sólo valores de una distribución Poisson conocida.

En este caso es muy ineficiente utilizar el método crudo de Montecarlo, porque obliga a muestras de una distribución difícil de calcular.

- Este procedimiento se utiliza cuando se realiza un estudio y se desea comparar el desempeño de dos o más sistemas. Frecuentemente, el objetivo al comparar es estimar las diferencias en los parámetros de dos procesos estocásticos.
- Los dos parámetros pueden estar positivamente correlacionados, y en ese caso la varianza de las diferencias individuales posiblemente es menor que la varianza de la diferencia de los estimadores en conjunto. Si  $T$  y  $S$  son estimadores insesgados, queremos conocer la diferencia de las varianzas  $\text{Var}(T) - \text{Var}(S)$ , ya que el que tenga menor varianza será mejor.
- Supongan que cada estimador es una función de una muestra aleatoria  $x_1, \dots, x_n$ . Un estimador de montecarlo de la varianza de  $T$  para una muestra de tamaño  $n$  se obtiene generando  $m$  muestras de tamaño  $n$  de la distribución dada, se calcula  $T_i$  para la  $i$ -ésima muestra y luego se calcula:

$$\hat{\text{Var}}(T) = \frac{\sum_{i=1}^m (T_i - \bar{T})^2}{m - 1}$$

- Más que hacer esto separadamente para  $T$  y  $S$ , noten que como los estimadores son insesgados,  $E(S) = E(T)$  y por lo tanto

$$\begin{aligned}\text{Var}(T) - \text{Var}(S) &= E(T^2) - E(T)^2 - (E(S^2) - E(S)^2) \\ &= E(T^2) - E(S^2) \\ &= E(T^2 - S^2)\end{aligned}$$

así que podemos obtener un estimado de  $T^2 - S^2$ , utilizando números aleatorios comunes.

- Consideren un sistema de linea de espera en donde los clientes llegan de acuerdo a un proceso Poisson  $N(t)$ . El operador requiere instalar un servidor para atender las llegadas que puede ser de dos tipos,  $M$  y  $N$ . Si elige  $M$ , entonces para el cliente  $i$  sea el tiempo de servicio  $S_i^m$ ,  $X^m$  el tiempo de espera total en el sistema de todos los clientes que llegaron antes del tiempo  $T$  y  $W_i$  es el tiempo total en el sistema. Las variables se relacionan como:

$$X^m = \sum_{i=1}^{N(T)} W_i^m$$

- Noten que  $W_i^m = S_i^m + Q_i^m$  donde  $Q_i^m$  es el tiempo de espera de  $i$  antes de ser atendido.
- Para el servido  $N$ , también se tienen las variables  $S_i^n$ ,  $X^n$ ,  $W_i^n$  y  $Q_i^n$ . El operador desea estimar  $\theta = E(X^m) - E(X^n)$

# Ejemplo

Los números aleatorios comunes se pueden usar para estimar la delta de una opción cuando no puede ser calculada de modo explícito, aunque hay mejores métodos para hacerlo. Aquí se menciona este como una forma de aplicación. Si  $C(S)$  es el precio de una opción particular cuando el precio del bien subyacente es  $S$ , entonces se define a la delta como

$$\Delta = \frac{\partial C}{\partial S}$$

que mide la sensibilidad del precio de la opción a cambios en el precio  $S_t$  del bien subyacente. Entonces, si  $\epsilon > 0$ , se puede aproximar  $\Delta$  con la razón de diferencia finita

$$\Delta_\epsilon = \frac{C(S + \epsilon) - C(S - \epsilon)}{2\epsilon}$$

En este caso,  $C(S + \epsilon)$  y  $C(S - \epsilon)$  no deben ser estimados de manera independiente, sino utilizando el mismo conjunto de numeros aleatorios comunes. Para una  $\epsilon$  pequeña, la reducción de varianza puede ser dramática.



- Las regiones del espacio muestral son ponderadas de acuerdo a su contribución a la estimación de  $\theta = E(X)$ . Esta ponderación se hace a través de una densidad.
- Si  $p(x)$  es la densidad ponderadora o **función de importancia**, entonces

$$\theta = \int_D f(x) dx = \int_D \frac{f(x)}{p(x)} p(x) dx$$

Entonces, un estimador de Montecarlo es  $\hat{\theta} = \frac{1}{m} \sum_{i=1}^m \frac{f(x_i)}{p(x_i)}$  sobre  $x_i \sim p$ . La varianza de este estimador está dada por:

$$\text{Var}(\hat{\theta}) = \frac{1}{m} \text{Var} \left( \frac{f(x)}{p(x)} \right) = \frac{1}{m} \left[ E \left( \frac{f^2(x)}{p^2(x)} \right) - E^2 \left( \frac{f(x)}{p(x)} \right) \right].$$

- El objetivo de importance sampling es escoger una densidad  $p$  tal que la varianza se minimice.
- Como  $E^2\left(\frac{f(x)}{p(x)}\right) = \left(\int_D f(x)dx\right)^2$  es lo que queremos estimar y no depende de  $p$ , la elección sólo depende de  $E\left(\frac{f^2(x)}{p^2(x)}\right)$ .
- Por la desigualdad de Jensen:

$$E\left(\frac{f^2(x)}{p^2(x)}\right) \geq \left(E\left(\frac{|f(x)|}{p(x)}\right)\right)^2 = \left(\int_D |f(x)|dx\right)^2$$

- La cota se alcanza cuando  $p(x) = \frac{|f(x)|}{\int_D |f(x)|dx}$  pero otra vez, no conocemos esa integral, así que en la práctica se busca  $p(x)$  tal que  $\frac{|f(x)|}{p(x)}$  sea aproximadamente constante.
- Este método es muy parecido al método de aceptación y rechazo.

## Ejemplo I

Supongan que se tiene una estadística de prueba  $T(\mathbf{X})$  que rechaza una hipótesis  $H_0$  si  $T(\mathbf{X}) \geq c$ . Lo que podría querer estudiarse en este contexto puede ser el nivel de confianza de la prueba

$$\theta = P_{H_0}(T(\mathbf{X}) \geq c) = \int I_D(\mathbf{x})f(\mathbf{x})d\mathbf{x}$$

donde  $D = \{\mathbf{x} | T(\mathbf{x}) \geq c\}$ .

- El método de Montecarlo dice: obtén  $\mathbf{X}_1, \dots, \mathbf{X}_n$  una muestra aleatoria de la distribución conjunta  $F$  y calcula

$$\hat{\theta} = \frac{\#\{\mathbf{X}_i \in D\}}{n},$$

con  $\text{Var}(\hat{\theta}) = \frac{\theta(1-\theta)}{n}$ .

- El método de importance sampling dice: encuentra  $g(\mathbf{x})$  una densidad tal que el cociente  $\frac{f(\mathbf{x})}{g(\mathbf{x})}$  sea pequeño en  $D$ . Una vez encontrada esta  $g$  se puede mostrar de  $G$ :

$$\theta = \int I_D(\mathbf{x})f(\mathbf{x})d\mathbf{x} = \int I_D(\mathbf{x})f(\mathbf{x})/g(\mathbf{x})g(\mathbf{x})d\mathbf{x} = E\left(I_D(\mathbf{X})\frac{f(\mathbf{x})}{g(\mathbf{x})}\right)$$

- Si  $\mathbf{Y}_1, \dots, \mathbf{Y}_n$  es una muestra aleatoria de  $G$ ,

$$\tilde{\theta} = \frac{1}{n} \sum_{j=1}^n l_D(\mathbf{y}_j) \frac{f(\mathbf{y}_j)}{g(\mathbf{y}_j)}$$

tiene varianza  $\text{Var}(\tilde{\theta}) = \frac{1}{n} \int_D \frac{f^2(\mathbf{y})}{g^2(\mathbf{y})} g(\mathbf{y}) d\mathbf{y} - \theta^2/n$ , que será menor a  $\text{Var}(\hat{\theta})$ .

- Este es un caso particular de importance sampling. Si suponemos que en el problema  $\theta = \int_D f(x)dx$   $f(x)$  se puede descomponer como una mezcla de densidades, tal que  $f(x) = \sum_{j=1}^k w_j f_j(x)$  con  $w_j$  pesos convexos, entonces si  $\theta_j = \int_D f_j(x)dx$ ,

$$\theta = \sum_{j=1}^k w_j \theta_j.$$

- Se puede muestrear cada parte de manera separada, con  $n_j$  observaciones, lo que da  $\hat{\theta}_j$  con  $\text{Var}(\hat{\theta}_j) = \sigma_j^2/n_j$ , donde  $\sigma_j$  es la varianza para el estimador  $\hat{\theta}_j$ . Combinando se obtiene

$$\text{Var}(\hat{\theta}) = \sum_{j=1}^k w_j^2 \frac{\sigma_j^2}{n_j}$$

- Se puede elegir  $n_j$  de manera óptima para  $n = \sum_{j=1}^k n_j$ ,  $w_j$  y  $\sigma_j^2$  fijas, obteniendo  $n_j = \frac{w_j \sigma_j}{\sum_{j=1}^k w_j \sigma_j}$

- Consideren obtener una muestra de una *normal contaminada*:

$$X \sim \begin{cases} \mathcal{N}(0, 1) & \text{con probabilidad } 1 - \alpha \\ \mathcal{N}(0, \sigma^2) & \text{con probabilidad } \alpha \end{cases}$$

- Si  $\alpha$  es pequeña entonces en una muestra de tamaño  $n$  se tendrán una o dos observaciones del componente contaminado. Si  $M$  es el número de casos contaminados en una muestra de tamaño  $n$ , entonces  $M \sim \text{Bin}(n, \alpha)$ .
- Entonces podemos estratificar por la variable  $M$  en los casos  $M = 0, 1, \dots, n$ , con  $w_j = \binom{n}{j} \alpha^j (1 - \alpha)^{n-j}$ .
- Observación: usualmente la varianza de un estimador se incrementa conforme la contaminación aumenta, así que la  $n_j$  óptima dependerá de  $j$  así como de  $w_j$ .

# Muestreo de Importancia I

Mejorando la eficiencia de la estimación

- Hay ciertos problemas en los que el método clásico MC no nos será de utilidad, por la precisión que se necesita, o bien porque es un evento raro; por ejemplo, si deseamos calcular una probabilidad en la cola de una distribución (usualmente se requiere calcular  $p$ -values).
- El método de *muestreo de importancia* (importance sampling) sirve para ganar eficiencia y precisión en la estimación, haciendo un cambio en la medida de referencia de los datos para mejorar el estimador de la varianza. Esta es una de las técnicas de *reducción de varianza* que detallaremos más adelante.

# Muestreo de Importancia II

Mejorando la eficiencia de la estimación

## Importance sampling

$$\theta = E_f[h(X)] = \int_{\mathcal{X}} h(x)f(x) dx = \int h(x)\frac{f(x)}{g(x)}g(x) dx = E_g\left[\frac{h(X)f(X)}{g(X)}\right]$$

donde  $\text{supp}(g) \supset \text{supp}(hf)$  ( $g$  es estrictamente positiva cuando  $hf$  es diferente de cero). El estimador IS de  $\theta$  será

$$\hat{\theta}_{n,IS} = \frac{1}{n} \sum_{i=1}^n \frac{f(y_i)h(y_i)}{g(y_i)},$$

donde  $y_1, \dots, y_n \sim g$

- Entonces, el problema se cambia a estimar observaciones de una densidad  $g$  en lugar de la densidad  $f$ .



Problema : Estimar  $\theta = P(X > 20)$  donde  $X \sim \mathcal{N}(0, 1)$ .

```
pnorm(20,lower.tail=F)
```

```
[1] 2.753624e-89
```

En este caso pueden comprobar que generando observaciones de la distribución normal estándar no funciona, ya que el evento es muy raro. Pero podemos escribir la integral trasladando el problema a la cola de la distribución de interés:

$$\begin{aligned}\theta = \mathbb{E}[I(X > 20)] &= \int_{-\infty}^{\infty} I(X > 20) \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx \\&= \int_{-\infty}^{\infty} I(X > 20) \frac{\frac{1}{\sqrt{2\pi}} e^{-x^2/2}}{\frac{1}{\sqrt{2\pi}} e^{-(x-\mu)^2/2}} \frac{1}{\sqrt{2\pi}} e^{-(x-\mu)^2/2} dx \\&= \int_{-\infty}^{\infty} I(X > 20) e^{-\mu x + \mu^2/2} \frac{1}{\sqrt{2\pi}} e^{-(x-\mu)^2/2} dx \\&= \mathbb{E}_{\mu} \left[ I(X > 20) e^{-\mu x + \mu^2/2} \right]\end{aligned}$$

Ahora podemos muestrear de  $\mathcal{N}(\mu, 1)$ . Podemos escoger  $\mu = 20$  para estar cerca del punto de interés.

# Ejemplo IS III

```
#Usando importance sampling:
n <- 1e6
mu <- 20
y <- rnorm(n, mean = mu)
I <- rep(0, n)
I[which(y > 20)] <- 1
theta_is <- mean(I*exp(-mu*y+mu^2/2))
#intervalo de confianza: noten la precisión.
theta_is

[1] 2.752978e-89

theta_is + c(-1,1)*sd(I*exp(-mu*y+mu^2/2))/n

[1] 2.752964e-89 2.752991e-89
```

La varianza del estimador IS se puede obtener de la siguiente manera:

$$\begin{aligned}\text{Var}_g(\hat{\theta}_{n,IS}) &= \int (\hat{\theta}_{n,IS} - \mathbb{E}(\hat{\theta}_{n,IS}))^2 g(x) dx \\ &= \int \hat{\theta}_{n,IS}^2 g(x) dx - \theta^2 \\ &= \int \left( \frac{h(x)f(x)}{g(x)} \right)^2 g(x) dx - \theta^2 \\ &= \int \frac{h^2(x)f(x)}{g(x)} f(x) dx - \theta^2\end{aligned}$$

Recordando que la varianza del estimador usual de  $\theta$  es:

$$\text{Var}_f(\hat{\theta}_n) = \int h^2(x)f(x) dx - \theta^2,$$

la reducción de varianza que se puede dar es:

$$\text{Var}_f(\hat{\theta}_n) - \text{Var}_g(\hat{\theta}_{n,IS}) = \int h^2(x) \left(1 - \frac{f(x)}{g(x)}\right) f(x) dx$$

Necesitamos que sea positiva. Esto se puede dar cuando:

- $\frac{f(x)}{g(x)} \geq 1$  cuando  $h(x)f(x)$  es pequeño, o
- $\frac{f(x)}{g(x)} \leq 1$  cuando  $h(x)f(x)$  es grande.

Más adelante regresaremos al problema de cómo se puede tratar de elegir  $g$  de manera adecuada.

# Sampling Importance Resampling (SIR) I

(Remuestreo de muestreo de importancia)

A partir de las ideas desarrolladas para el muestreo de importancia, se puede obtener un nuevo modelo para simular observaciones de  $f$ , como se indica a continuación.

- El método IS produce una muestra  $X_1, \dots, X_n \sim g$  y un conjunto de pesos de 'importancia'  $\omega_i = f(X_i)/g(X_i)$ .
- Los pesos  $\omega_i$  no suman 1, y algunos son mayores que uno. Pero se pueden normalizar:  $w_i = \frac{\omega_i}{\sum_{i=1}^n \omega_i}$
- Si se considera obtener muestras  $X^*$  con probabilidad  $w_i$ , extraídas **con reemplazo de**  $\{X_1, \dots, X_n\}$  (esta fase se conoce como remuestreo), entonces la distribución de  $X_i^*$  es  $f$ , ya que:

$$\begin{aligned} P(X^* \in A) &= \sum_{i=1}^n P(X^* \in A \& X^* = X_i) \\ &= \int_A f(x)/g(x)g(x) dx = \int_A f(x) dx \end{aligned}$$

# Sampling Importance Resampling (SIR) II

(Remuestreo de muestreo de importancia)

## Estimador Sampling Importance Resampling

El estimador SIR está dado por:

$$\hat{\theta}_{n,SIR} = \sum_{i=1}^n w_i h(X_i)$$

donde  $w_i = \frac{f(X_i)/g(X_i)}{\sum_{j=1}^n f(X_j)/g(X_j)}$

# Ejemplo SIR I

Este método es muy general, así que podemos aplicarlo a cualquier función: Consideremos obtener la integral  $\theta = E_f(X^2)$  donde  $f(x) = e^{0.4(x-0.4)^2 - 0.08x^4}$ . Notemos que  $f$  no es una densidad, por lo que necesitamos encontrar la constante para normalizar y poder comparar valores. Por otra parte, uso  $g \sim \mathcal{U}(-4, 4)$  como cambio de medida.

```
f <- function(x, a=0.4,b=0.08){exp(a*(x-a)^2 -b*x^4)}
x <- seq(-4,4,0.1)
#Consideramos a g ~ U[-4,4]
n <- 10000
x <- runif(n,-4,4)
w <- f(x)/dunif(x,-4,4)
q <- w/sum(w)
xb <- sample(x,prob=q,replace = T) #remuestreo
thetahatSIR <- sum(q*xb^2)
cc <- integrate(function(x)f(x),-4,4)$val #constante de estandarización
integrate(function(x){(1/cc)*f(x)*x^2},-4,4)
```

2.413271 with absolute error < 1.7e-09

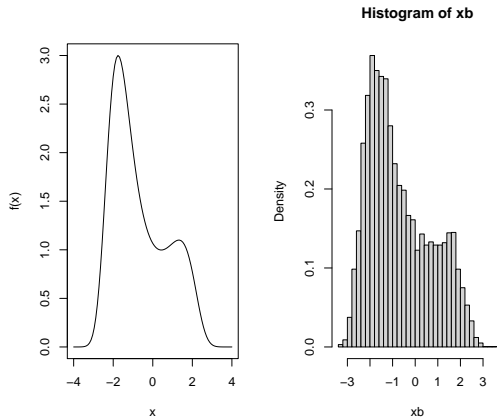
thetahatSIR

[1] 2.366468

```
par(mfrow=c(1,2))
curve(f,type="l",from = -4,to=4)
hist(xb,breaks=30,prob=T)
```



# Ejemplo SIR II



- $X \sim T(\nu, \theta, \sigma^2)$  con densidad

$$f(x) = \frac{\Gamma((\nu + 1)/2)}{\sigma \sqrt{\nu \pi} \Gamma(\nu/2)} \left( 1 + \frac{(x - \theta)^2}{\nu \sigma^2} \right)^{-(\nu+1)/2}$$

- Tomar  $\theta = 0$  y  $\sigma = 1$
- Estimar

$$\int_{2.1}^{\infty} x^5 f(x) dx$$

- Con candidatos:
  - la misma  $f$
  - Cauchy
  - Normal
  - Uniforme en  $(0,1/2.1)$