

Antithetic Variables

Two identically r.v. γ and γ'

$$\frac{\gamma + \gamma'}{2} \quad E\left[\frac{\gamma + \gamma'}{2}\right] = E[\gamma] = \mu$$

$$\text{Var}\left(\frac{\gamma + \gamma'}{2}\right) = \frac{1}{4} [\text{Var}(\gamma) + \text{Var}(\gamma') + 2 \text{Cov}(\gamma, \gamma')]$$

to make $\text{Cov}(\gamma, \gamma') < 0$

$$\text{Suppose } X_j = F_x^{-1}(u), j=1, \dots, m \quad u \sim \text{Unif}(0,1)$$

$$Y_j = g[F_x^{-1}(u)], j=1, \dots, m$$

u and $1-u \sim \text{Unif}(0,1)$

$$\text{Corr}(u, 1-u) = -1$$

$$Y'_j = g[F_x^{-1}(1-u)], j=1, \dots, m$$

Q: Under what conditions are Y_j and Y'_j negatively correlated?

Ans: If $g(x)$ is monotone, Y_j and Y'_j are negatively correlated.

$$Y_j = g[F_x^{-1}(u)] = f(u) \quad \text{assume } g(u) \text{ is increasing func.}$$

$$-Y'_j = -g[F_x^{-1}(1-u)] = h(u) \quad \text{Cov}(Y_j, Y'_j) < 0$$

$$E[f(u)h(u)] \geq E[f(u)] E[h(u)] = E(Y_j Y'_j) - E(Y_j)E(Y'_j)$$

$$E[-Y_j Y'_j] \geq E(Y_j) E(-Y'_j)$$

$$E(Y_j Y'_j) \leq E(Y_j) E(Y'_j)$$

Algorithm

$$u, u_2, \dots, u_m \sim \text{Unif}(0,1)$$

$$Y_j = f(u_j)$$

$$Y'_j = h(u_j) = f(1-u_j)$$

$$\hat{\theta} = \sum_{j=1}^m \frac{Y_j + Y'_j}{m}$$

E.g. $x \geq 0, f_x(x) + 0.5, g(u)$

$$\hat{\theta} = \int_0^1 \left(\frac{x}{m} e^{-\frac{1}{2}(xu)} \right) du$$

$u \sim \text{Unif}(0,1)$

$$Y_j = g(u_j) \Rightarrow f(u)$$

$$Y'_j = g(1-u_j) \Rightarrow h(u)$$

$$\hat{\theta} = \sum_{j=1}^m \frac{Y_j + Y'_j}{m}$$

MC_phi <- function(x, m = 10000, anti = TRUE){

```
  u <- runif(m / 2)
  if(!anti)
    v <- runif(m / 2)
  else
    v <- 1 - u

  u <- c(u, v)
  mc <- numeric(length(x))
  for(i in 1:length(x)){
    g <- x[i] * exp(-(u * x[i]) ^ 2 / 2)
    mc[i] <- mean(g) / sqrt(2 * pi) + 0.5
  }
  mc
```

l2022

The R code to compute the approximate reduction in variance at a

```
x = 1.3
n <- 50
MC1 <- MC2 <- numeric(n)
x0 <- 1.3
for(i in 1:n){
  MC1[i] <- MC_phi(x0, m = 1000, anti = FALSE)
  MC2[i] <- MC_phi(x0, m = 1000)
}
(var(MC1) - var(MC2)) / var(MC1)

## [1] 0.9569375 ~96% reduction in variance
```

Control Variates

► $\text{Var}(\hat{\theta}_c)$ is minimized at $c = c^*$, where

$$c^* = -\frac{\text{Cov}[g(X), f(X)]}{\text{Var}[f(X)]}$$

► and minimum variance is

$$\text{Var}(\hat{\theta}_{c^*}) = \text{Var}[g(X)] - \frac{\text{Cov}[g(X), f(X)]^2}{\text{Var}[f(X)]}$$

$$\hat{\theta}_c = g(x) + c[f(x) - \mu]$$

$E[f(x)] = \mu$ is known

$$E[\hat{\theta}_c] = \theta$$

$$\text{Var}(\hat{\theta}_c) \leq \text{Var}(\hat{\theta})$$

$$\text{Var}(\hat{\theta}_c) = \text{Var}(g(x) + c[f(x) - \mu])$$

$$= \text{Var}(g(x)) + 2c \text{Cov}[g(x), f(x)] + c^2 \text{Var}[f(x)]$$

$$\min_c \text{Var}(\hat{\theta}_c)$$

$\frac{\partial}{\partial c} \text{Var}(\hat{\theta}_c) = 0$ and solve it for c^*

$$c^* = -\frac{\text{Cov}[g(x), f(x)]}{\text{Var}[f(x)]}$$

$$\text{Var}(\hat{\theta}_c) = \text{Var}[g(x)] - \frac{\text{Cov}[g(x), f(x)]^2}{\text{Var}[f(x)]} + \frac{\text{Cov}[g(x), f(x)]}{\text{Var}[f(x)]}$$

Control Variance c

$$= \text{Var}[g(x)] - \frac{\text{Cov}[g(x), f(x)]^2}{\text{Var}[f(x)]} \leq \text{Var}[g(x)]$$

$$\text{The percent reduction} = 100 \cdot \frac{\text{Cov}[g(x), f(x)]^2}{\text{Var}[f(x)]}$$

$$= 100 \cdot \frac{\text{Cov}[g(x), f(x)]^2}{\text{Var}[g(x)] \cdot \text{Var}[f(x)]}$$

$$= 100 \cdot \text{Cov}[g(x), f(x)]^2$$

$$g(u) = \frac{e^{-u}}{1+u^2} \quad u \sim \text{Unif}(0,1)$$

$\text{corr}[g(u), f(u)] \uparrow$

$$f(u) = e^{-u}$$

$$E[f(u)] = \int_0^1 e^{-u} du = 1 - e^{-1}$$

Algorithm

① find $f(x)$ and $E[f(x)] = \mu$

② $u \sim \text{Unif}(0,1)$ u_1, \dots, u_m

③ Compute $c^* = -\frac{\text{Cov}[g(u), f(u)]}{\text{Var}[f(u)]}$

④ $\hat{\theta}_c = \sum_{i=1}^m g(u_i) + c^* [f(u_i) - \mu]$

f <- function(u){

exp(u)

}

g <- function(u){

u * exp(-3 * u)

}

u <- runif(10000, 0, 2)

B <- f(u)

A <- g(u)

a <- -cov(A, B) / var(B)

u <- runif(10000, 0, 2)

T1 <- g(u) * 2

T2 <- T1 + a * (f(u) - mean(B))

cbind(mean(T1), mean(T2))

[1] [,1] [,2]

[1] 0.1082181 0.108405

con_var <- (var(T1) - var(T2)) / var(T1)

con_var

[1] 0.529722

Simple MC

- Define the Function: You've already defined the function $fx(x) = x * \exp(-3 * x)$.
- Generate Uniform Random Numbers: Generate a large number of uniform random numbers over the interval $[0, 2]$.
- Apply the Function to Random Numbers: Apply fx to each of these random numbers.
- Calculate the Average: Compute the average of these function values.
- Multiply by the Length of the Interval: Multiply this average by the length of the interval (which is 2 in this case) to get the estimate of the integral.

fx <- function(x) { x * exp(-3 * x) }

m <- 10000

u <- runif(m, 0, 2)

mci_estimates <- fx(u)

theta_mci <- mean(mci_estimates)*2

theta_mci

- Initialize Parameters: Set the total number of samples m and the number of strata n_{strata} .
- Initialize Storage for Stratified Estimates: Create a numeric vector `stratified_estimates` to store the mean estimates from each stratum.
- Stratified Sampling Loop: For each stratum, calculate the lower and upper bounds, generate uniform random samples within these bounds, apply the function fx to these samples, and store the mean of these function values.
- Calculate the Overall Estimate: Average the stratified means and multiply by the length of the interval (which is 2 in this case).

Var($\hat{\theta}_{mc}$) = Var($\bar{g}(x)$) = $\frac{1}{m} \text{Var}(g(x))$

$$\text{Var}(\hat{\theta}_{mc}) = \frac{1}{m} [\text{Var}[E[g(x)|J]] + E[\text{Var}(g(x)|J)]]$$

between-group variance

within-group variance

$$\text{Var}(\hat{\theta}_{mc}) = \frac{1}{m} [\text{Var}(\hat{\theta}_1) + \frac{1}{K} \sum_{j=1}^K \text{Var}(\hat{\theta}_j)]$$

$$= \frac{1}{m} [\text{Var}(\hat{\theta}_1) + \frac{1}{K} \sum_{j=1}^K \frac{1}{m_j} \text{Var}(g(x|J=j))]$$

$$= \frac{1}{m} [\text{Var}(\hat{\theta}_1) + \frac{1}{K} \sum_{j=1}^K \frac{1}{m_j} \text{Var}(g(x|J=j))] = \frac{1}{m} [\text{Var}(\hat{\theta}_1) + \frac{1}{K} \sum_{j=1}^K \frac{1}{m_j} \text{Var}(g(x|J=j))]$$

$$= \frac{1}{m} [\text{Var}(\hat{\theta}_1) + \frac{1}{K} \sum_{j=1}^K \frac{1}{m_j} \text{Var}(g(x|J=j))] = \frac{1}{m} [\text{Var}(\hat{\theta}_1) + \frac{1}{K} \sum_{j=1}^K \frac{1}{m_j} \text{Var}(g(x|J=j))]$$

$$= \frac{1}{m} [\text{Var}(\hat{\theta}_1) + \frac{1}{K} \sum_{j=1}^K \frac{1}{m_j} \text{Var}(g(x|J=j))] = \frac{1}{m} [\text{Var}(\hat{\theta}_1) + \frac{1}{K} \sum_{j=1}^K \frac{1}{m_j} \text{Var}(g(x|J=j))]$$

$$= \frac{1}{m} [\text{Var}(\hat{\theta}_1) + \frac{1}{K} \sum_{j=1}^K \frac{1}{m_j} \text{Var}(g(x|J=j))] = \frac{1}{m} [\text{Var}(\hat{\theta}_1) + \frac{1}{K} \sum_{j=1}^K \frac{1}{m_j} \text{Var}(g(x|J=j))]$$

$$= \frac{1}{m} [\text{Var}(\hat{\theta}_1) + \frac{1}{K} \sum_{j=1}^K \frac{1}{m_j} \text{Var}(g(x|J=j))] = \frac{1}{m} [\text{Var}(\hat{\theta}_1) + \frac{1}{K} \sum_{j=1}^K \frac{1}{m_j} \text{Var}(g(x|J=j))]$$

$$= \frac{1}{m} [\text{Var}(\hat{\theta}_1) + \frac{1}{K} \sum_{j=1}^K \frac{1}{m_j} \text{Var}(g(x|J=j))] = \frac{1}{m} [\text{Var}(\hat{\theta}_1) + \frac{1}{K} \sum_{j=1}^K \frac{1}{m_j} \text{Var}(g(x|J=j))]$$

$$= \frac{1}{m} [\text{Var}(\hat{\theta}_1) + \frac{1}{K} \sum_{j=1}^K \frac{1}{m_j} \text{Var}(g(x|J=j))] = \frac{1}{m} [\text{Var}(\hat{\theta}_1) + \frac{1}{K} \sum_{j=1}^K \frac{1}{m_j} \text{Var}(g(x|J=j))]$$

$$= \frac{1}{m} [\text{Var}(\hat{\theta}_1) + \frac{1}{K} \sum_{j=1}^K \frac{1}{m_j} \text{Var}(g(x|J=j))] = \frac{1}{m} [\text{Var}(\hat{\theta}_1) + \frac{1}{K} \sum_{j=1}^K \frac{1}{m_j} \text{Var}(g(x|J=j))]$$

$$= \frac{1}{m} [\text{Var}(\hat{\theta}_1) + \frac{1}{K} \sum_{j=1}^K \frac{1}{m_j} \text{Var}(g(x|J=j))] = \frac{1}{m} [\text{Var}(\hat{\theta}_1) + \frac{1}{K} \sum_{j=1}^K \frac{1}{m_j} \text{Var}(g(x|J=j))]$$

$$= \frac{1}{m} [\text{Var}(\hat{\theta}_1) + \frac{1}{K} \sum_{j=1}^K \frac{1}{m_j} \text{Var}(g(x|J=j))] = \frac{1}{m} [\text{Var}(\hat{\theta}_1) + \frac{1}{K} \sum_{j=1}^K \frac{1}{m_j} \text{Var}(g(x|J=j))]$$

$$= \frac{1}{m} [\text{Var}(\hat{\theta}_1) + \frac{1}{K} \sum_{j=1}^K \frac{1}{m_j} \text{Var}(g(x|J=j))] = \frac{1}{m} [\text{Var}(\hat{\theta}_1) + \frac{1}{K} \sum_{j=1}^K \frac{1}{m_j} \text{Var}(g(x|J=j))]$$

$$= \frac{1}{m} [\text{Var}(\hat{\theta}_1) + \frac{1}{K} \sum_{j=1}^K \frac{1}{m_j} \text{Var}(g(x|J=j))] = \frac{1}{m} [\text{Var}(\hat{\theta}_1) + \frac{1}{K} \sum_{j=1}^K \frac{1}{m_j} \text{Var}(g(x|J=j))]$$

$$\begin{aligned}\Theta &= E_f[g(x)] = \int_D g(x) f(x) dx \\ &= \int_D g(x) \cdot \frac{f(x)}{\phi(x)} \phi(x) dx \\ &= \int_D g(x) \cdot \phi(x) dx \\ &= E_{\phi}[\int g(x) dx]\end{aligned}$$

$$\begin{aligned}① \quad x_1, \dots, x_m &\sim \phi(x) \\ ② \quad \frac{\sum_{i=1}^m g(x_i) f(x_i)}{\sum_{i=1}^m \phi(x_i)} \cdot \frac{1}{m} &= \hat{\theta} \\ \text{MSE}(\hat{\theta}) &= \text{Var}_{\phi}(\hat{\theta}) + [\mathbb{E}[\hat{\theta}] - \theta]^2\end{aligned}$$

$$\begin{aligned}\text{Var}_{\phi}(\hat{\theta}) &= \text{Var}_{\phi}\left[\frac{\sum g(x) f(x)}{\sum \phi(x)}\right] \cdot \frac{1}{m} \\ &= \left\{ \mathbb{E}_{\phi}\left[\frac{g(x) f(x)}{\phi(x)}\right] \right\}^2 - \mathbb{E}_{\phi}\left[\frac{g(x) f(x)}{\phi(x)}\right]^2 \cdot \frac{1}{m}\end{aligned}$$

Find the optimal ϕ that leads $\int_D \frac{g(x)^2 f(x)}{\phi(x)^2} \phi(x) dx = \theta^2$
 $\Rightarrow \text{Var}_{\phi}(\hat{\theta}) = 0$

important sampling mean estimate

$$\hat{\theta} = \frac{1}{m} \sum g(x) w(x) \rightarrow \mathbb{E}_f[g(x)]$$

If $\phi(x) \approx f(x)$, $w(x) \approx 1$ $\mathbb{E}[g(x) | w(x)]$

If $\phi(x) \gg f(x)$, $w(x) \ll 1$ $\mathbb{E}[g(x) \cdot w(x)]$

$\phi(x) \approx f(x) \Rightarrow w(x) \approx 1$

$$\text{Define efficiency} = \frac{1}{\text{Var}_{\phi}(w(x))} \quad \mathbb{E}_f[w(x)] = \int_D w(x) \cdot \phi(x) dx$$

Stratified Importance Sampling

$$\mathbb{E}[g(x)] = \int_D g(x) dx = \int_D \frac{g(x)}{f(x)} f(x) dx$$

$$\text{Var}_{\phi}(\hat{\theta}_1) = \frac{1}{M} \sum_{i=1}^M \frac{g(x_i)}{f(x_i)}$$

$$\text{Var}_{\phi}(\hat{\theta}_2) = \frac{1}{M} \sum_{i=1}^M \frac{g(x_i)}{f(x_i)} \cdot \frac{f(x_i)}{f(x_i)} = \frac{1}{M} \sum_{i=1}^M g(x_i)$$

$$\Theta = \Theta_1 + \Theta_2 + \dots + \Theta_k$$

$$\frac{\sum \Theta_i}{M} - \frac{K}{M} \sum_{j=1}^k \frac{\Theta_j}{f(x_j)} \geq 0$$

$$\frac{\sum \Theta_i}{M} - K \frac{\sum \Theta_j}{\sum f(x_j)} \geq 0$$

$$\mathbb{E}(\hat{\theta}^{str}) = \Theta$$

$$\text{Var}_{\phi}(\hat{\theta}^{str}) \leq \text{Var}(\hat{\theta}^i)$$

$$\Theta = \sum_{j=1}^k \mathbb{E}_{f|I_j} \left[\frac{g_j(x)}{f_j(x)} \right]$$

$$\text{Var}_{\phi}(\frac{1}{M} \sum_{j=1}^k \hat{\theta}_j) \leq \text{Var}(\hat{\theta}^i)$$

$$\sum_{j=1}^k \text{Var}(\hat{\theta}_j) \leq \text{Var}(\hat{\theta}^i)$$



Example: Folded Normal Distribution

Let $q(x)$ be an unnormalized density for $f(x)$, given by $q(x) = e^{-x^2/2}$, for $x \geq 0$.

We want to use self-normalized importance sampling to estimate $\hat{\theta} = \mathbb{E}_f(X) = \int_0^\infty x f(x) dx$

$$\hat{\theta} = \mathbb{E}_f(X) = \int_0^\infty x f(x) dx = \int_0^\infty x \frac{q(x)}{Z_q} dx = \frac{1}{Z_q} \int_0^\infty x e^{-x^2/2} dx$$

We previously found the theoretical $E_f(X) = \sqrt{\frac{2}{\pi}}$.

trial density: $\exp(2)$ $h(x) = \frac{x}{\sqrt{2\pi}} e^{-x^2/2}$

$\mathbb{E}_\phi(\hat{\theta}) = \mathbb{E}_\phi(\hat{\theta}_1) = \mathbb{E}_\phi(\hat{\theta}_2) = \dots = \mathbb{E}_\phi(\hat{\theta}_k) = \dots$

$$w(x) = \frac{g(x)}{h(x)} = \frac{e^{-x^2/2}}{\frac{x}{\sqrt{2\pi}} e^{-x^2/2}} = e^{-\frac{x^2}{2} + \frac{1}{2}}$$

$$\hat{\theta} = \frac{\sum x_i w(x_i)}{\sum w(x_i)}$$

R Code to estimate $E_f(X)$ (self-normalized importance sampling):

```
set.seed(9999) # for reproducibility
n <- 10000 # Specify the number of points to generate
# Generate n points from Exp(lambda = 2)
x <- rexp(n, rate = 2)

# Compute importance weights
w <- exp(-x^2 / 2 + 2 * x) / sqrt(2 * pi)

# Compute sum(w(x) * x) / sum(w(x))
sum(w * x) / sum(w)
```

[1] 0.800945

Stratified sampling

$$\Theta = \int_{a_0}^{a_k} g(x) dx \quad \text{assume } m_1 = m_2 = \dots = m_k = m$$

$$\Theta = (\Theta_1 + \Theta_2 + \dots + \Theta_k) \frac{1}{m}$$

$$= \Theta_1 \frac{1}{m} + \Theta_2 \frac{1}{m} + \dots + \Theta_k \frac{1}{m}$$

$$= \frac{\sum_{i=1}^m g(x_i)}{m} \cdot \frac{1}{m} + \dots + \frac{\sum_{i=1}^m g(x_i)}{m} \cdot \frac{1}{m} = \mathbb{E}[g(x) | J=j] = \mathbb{E}[g(x)]$$

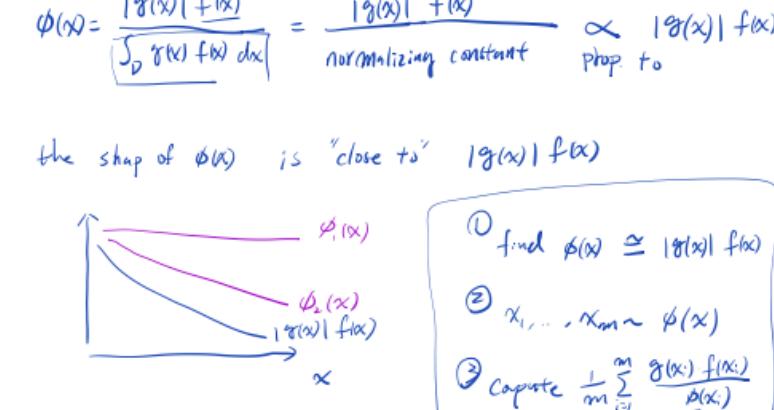
$$\mathbb{E}[g(x) | J=j] = \sum_{i=1}^k \mathbb{E}[g(x) | J=j] \cdot P(J=j)$$

$$= \sum_{j=1}^k \sum_{i=1}^m g(x_i) \cdot P(x_i | J=j) \cdot P(J=j)$$

$$= \sum_{j=1}^k \sum_{i=1}^m g_j(x_i) \cdot p_j(x_i)$$

$$= \sum_{j=1}^k \mathbb{E}_{f|I_j} [g_j(x)] = \mathbb{E}[g(x)]$$

$$\mathbb{E}\left[\frac{g(x)}{f(x)}\right] = \mathbb{E}\left[\mathbb{E}_f\left[\frac{g(x) | J=j}{f(x) | J=j}\right]\right] = \sum_{j=1}^k \mathbb{E}_{f|I_j} \left[\frac{g(x) | J=j}{f(x) | J=j} \right] = \sum_{j=1}^k \mathbb{E}_{f|I_j} \left[\frac{g(x) | J=j}{f_j(x)} \right]$$



the shape of $\phi(x)$ is "close to" $g(x)/f(x)$

- ① find $\phi(x) \approx 1/g(x) f(x)$
- ② $x_1, \dots, x_m \sim \phi(x)$
- ③ Compute $\frac{1}{m} \sum_{i=1}^m \frac{g(x_i) f(x_i)}{\phi(x_i)}$
 $\hat{\theta} = \frac{1}{m} \sum_{i=1}^m \frac{g(x_i) f(x_i)}{\phi(x_i)}$
 Let $w(x) = \frac{f(x)}{\phi(x)}$

$$\hat{\theta} = \frac{1}{m} \sum_{i=1}^m g(x_i) w(x_i)$$

```
g <- function(x){exp(-x - log(1 + x^2)) * (x > 0) * (x < 1)}
m <- 10000
# Use Uniform(0, 1) as the candidate function
is1 <- replicate(1000, expr = {
  x <- runif(m)
  phi <- 1
  mean(g(x) / phi)
})
# Use Exponential(1) as the candidate function
is2 <- replicate(1000, expr = {
  u <- runif(m)
  x <- -log(u)
  x <- x[x <= 1]
  phi <- exp(-x)
  sum(g(x) / phi) / m
})
# Use exp(-x) / (1 - exp(-1)) as the candidate function
is3 <- replicate(1000, expr = {
  u <- runif(m)
  x <- -log(1 - u * (1 - exp(-1)))
  phi <- exp(-x) / (1 - exp(-1))
  mean(g(x) / phi)
})
# Use Standard Cauchy as the candidate function
is4 <- replicate(1000, expr = {
  x <- rcauchy(m)
  x <- x[x > 0 & x <= 1]
  phi <- dcauchy(x)
  outcome <- g(x) / phi
  sum(outcome) / m
})
c(mean(is1), mean(is2), mean(is3), mean(is4))
```

```
# Use exp(-x) / (1 - exp(-1)) as the candidate function
is3 <- replicate(1000, expr = {
  u <- runif(m)
  x <- -log(1 - u * (1 - exp(-1)))
  phi <- exp(-x) / (1 - exp(-1))
  mean(g(x) / phi)
})
# Use Standard Cauchy as the candidate function
is4 <- replicate(1000, expr = {
  x <- rcauchy(m)
  x <- x[x > 0 & x <= 1]
  phi <- dcauchy(x)
  outcome <- g(x) / phi
  sum(outcome) / m
})
c(mean(is1), mean(is2), mean(is3), mean(is4))
```

[1] 0.5246914 0.5245130 0.5247545 0.5248252
 c(var(is1), var(is2), var(is3), var(is4))
 ## [1] 6.011627e-06 1.681950e-05 9.240642e-07 9.484520e-05

Example: Folded Normal Distribution

Suppose we want to estimate $E_f(X)$, where $f(x)$ is the PDF of the folded normal distribution, $f(x) = 2 \cdot \frac{1}{\sqrt{2\pi}} e^{-x^2/2} = \frac{\sqrt{2}}{\pi} e^{-x^2/2}$, for $x \geq 0$.

The support of $f(x)$ is $D = [0, \infty)$:

$f(x) > 0$ for all $x \in D$

$\int_0^\infty f(x) dx = 1$

We want to use importance sampling to estimate $E_f(X)$.

① choose $\phi(x) = 2 e^{-x^2}$, $x \geq 0 \Rightarrow \exp(x)$
 ② $w(x) = \frac{f(x)}{\phi(x)} = \frac{\sqrt{2}}{\pi} e^{-x^2/2} / 2 e^{-x^2} = \frac{1}{\pi} e^{-x^2/2 + x^2}$
 ③ $x_1, \dots, x_m \sim \exp(2)$ $x = -\log(u) / 2$ $\hat{\theta} = \frac{\sum x_i \cdot w(x_i)}{m}$

uan Wu, 2022

set.seed(9999)

n <- 10000 # Specify the number of points to generate

Generate n points from Exp(lambda = 2)

x <- rexp(n, rate = 2)

Compute importance weights

w <- exp(-x^2 / 2 + 2 * x) / sqrt(2 * pi)

Compute mean(w(x) * g(x)) (g(x) = x here)

mean(w * x)

Suppose we sample X_1, \dots, X_m from an importance sampling function $\phi(x)$. Please show that $E(\sum_i w_i) = E(f(x_i)/\phi(x_i)) = m$, where f is target density.

$$\mathbb{E}[\sum_{i=1}^m w_i] = \mathbb{E}\left[\sum_{i=1}^m \frac{f(x_i)}{\phi(x_i)}\right]$$

$$\mathbb{E}\left[\frac{f(x_i)}{\phi(x_i)}\right] = \int \frac{f(x)}{\phi(x)} \phi(x) dx$$

$$\mathbb{E}\left[\frac{f(x_i)}{\phi(x_i)}\right] = \int f(x) dx$$

$$\mathbb{E}\left[\frac{f(x_i)}{\phi(x_i)}\right] = 1$$

$$\sum_{i=1}^m \mathbb{E}\left[\frac{f(x_i)}{\phi(x_i)}\right] = m \times 1 = m$$

We showed that the variance of importance sampling estimator $\text{Var}_f(\hat{\theta}) = \frac{1}{m} \text{Var}_{\phi}\left(\frac{g(X)f(X)}{\phi(X)}\right)$, where $\hat{\theta} = \frac{1}{m} \sum_{i=1}^m \frac{g(x_i)f(x_i)}{\phi(x_i)}$ and $x_i \sim \phi(x)$. Please show that $\text{Var}_{\phi}\left(\frac{g(X)f(X)}{\phi(X)}\right)$ is equal to $\mathbb{E}_{\phi}\left[\left(\frac{g(X)f(X) - \theta\phi(X)}{\phi^2(X)}\right)^2\right]$. Here $f(x)$ is importance function, and $\phi(x)$ is importance sampling function.

Your answer:

$$\hat{\theta} = \frac{1}{m} \sum_{i=1}^m \frac{g(x_i)f(x_i)}{\phi(x_i)}$$

$$\text{Var}_{\phi}(\hat{\theta}) = \frac{1}{m} \text{Var}_{\phi}\left(\frac{g(X)f(X)}{\phi(X)}\right)$$

$$\mathbb{E}_{\phi}\left[\left(\frac{g(X)f(X)}{\phi(X)} - \hat{\theta}\right)^2\right]$$

$$\text{Var}_{\phi}(\hat{\theta}) = \mathbb{E}_{\phi}\left[\left(\frac{g(X)f(X)}{\phi(X)} - \mathbb{E}_{\phi}\left[\frac{g(X)f(X)}{\phi(X)}\right]\right)^2\right]$$

$$\text{Var}_{\phi}(\hat{\theta}) = \mathbb{E}_{\phi}\left[\left(\frac{g(X)f(X)}{\phi(X)} - \hat{\theta}\right)^2\right]$$

$$\Theta = \int_{a_0}^{a_k} g(x) dx$$

$$= \sum_{j=1}^k \mathbb{E}_{f|I_j} \left[\frac{g(x)}{f(x)} \right]$$

$$= \sum_{j=1}^k \mathbb{E}_{f|I_j} \left[\frac{g(x)}{f_j(x)} \right]$$

$$= \sum_{j=1}^k \frac{1}{m_j} \sum_{i=1}^{m_j} \frac{g(x_i)}{f_j(x_i)}$$

$$= \sum_{j=1}^k \frac{1}{m_j} \sum_{i=1}^{m_j} \frac{g(x_i)}{f_j(x_i)} \cdot \frac{m_j}{m_j} = \frac{1}{m} \sum_{j=1}^k \frac{1}{m_j} \sum_{i=1}^{m_j} g(x_i)$$

$$= \sum_{j=1}^k \frac{1}{m_j} \sum_{i=1}^{m_j} g(x_i) \cdot \frac{1}{m_j} = \frac{1}{m} \sum_{j=1}^k \frac{1}{m_j} \sum_{i=1}^{m_j} g(x_i)$$

$$= \frac{1}{m} \sum_{j=1}^k \mathbb{E}_{f|I_j} [g(x)] = \mathbb{E}[g(x)]$$

Transition Probabilities (2)

If the state space is finite, the transition probabilities P_{ij} can be represented by a transition matrix

$$\mathbb{P} = \begin{bmatrix} P_{00} & P_{01} & P_{02} & \dots & P_{0N} \\ P_{10} & P_{11} & P_{12} & \dots & P_{1N} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ P_{i0} & P_{i1} & P_{i2} & \dots & P_{iN} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ P_{N0} & P_{N1} & P_{N2} & \dots & P_{NN} \end{bmatrix}, \quad \sum_{j=0}^N P_{ij} = 1$$

where

- All the entries are non-negative:

$$P_{ij} \geq 0, \text{ for all } i, j. \quad (3)$$

- The sum of each row is 1:

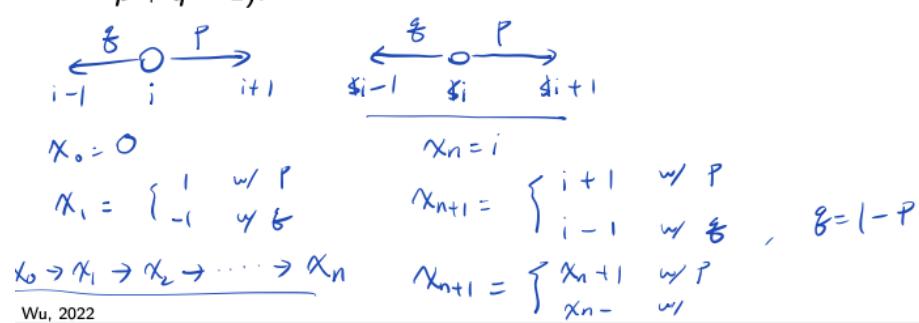
$$\sum_{j=0}^N P_{ij} = \sum_{j=0}^N P(X_{n+1} = j | X_n = i) = 1, \text{ for all } i. \quad (4)$$

A random walk is a stochastic process:

$$\{X_0, X_1, X_2, \dots, X_n, \dots\},$$

defined on the integers \mathbb{Z} , such that:

- The walk starts at 0: $X_0 = 0$.
- At each step, the random walk moves to the right 1 unit with probability p and moves to the left 1 unit with probability q (so $p + q = 1$).



Consider a two-state Markov chain, with state space $\{0, 1\}$ and the transition matrix

$$\mathbb{P} = \begin{bmatrix} 0.7 & 0.3 \\ 0.4 & 0.6 \end{bmatrix}. \quad P(X_1 = 0 | X_0 = 0) = P_{00}^{(1)}$$

Suppose we are given that the chain starts at 0: $X_0 = 0$. Calculate the probability that the state 0 in 3 steps, $P(X_3 = 0) = P_{00}^{(3)}$.

If we generate X_1, X_2, \dots, X_n , for some large n (for example, $n = 10000$). How often is the Markov chain in state 0? How often is the Markov chain in state 1?

$$\begin{aligned} X_0 &= 0 \\ X_0 \rightarrow X_1 \rightarrow X_2 \rightarrow X_3 & \\ P(X_3 = 0 | X_0 = 0) & \\ \begin{array}{c} \text{state 0} \\ \text{0.7} \xrightarrow{\text{0.1}} \text{0.3} \xrightarrow{\text{0.7}} \text{0.1} \xrightarrow{\text{0.3}} \text{0.7} \end{array} & \\ P_{00}^{(1)} &= 0.7 + 0.7 + 0.3 + 0.4 \\ P_{01}^{(1)} &= 0.7 + 0.3 + 0.3 + 0.6 \\ \left[\begin{array}{cc} 0.7 & 0.3 \\ 0.4 & 0.6 \end{array} \right] \left[\begin{array}{cc} 0.7 & 0.3 \\ 0.4 & 0.6 \end{array} \right] & \\ &= \left[\begin{array}{cc} 0.7 \cdot 0.7 + 0.3 \cdot 0.4 & 0.7 \cdot 0.3 + 0.3 \cdot 0.6 \\ 0 & 0 \end{array} \right] \\ P^{(n)} &= P^n \end{aligned}$$

```

n <- 1000
P <- matrix(c(0.7, 0.4, 0.3, 0.6), nrow = 2)
P2 <- P
for(i in 1:(n - 1)){
  P2 <- P2 %*% P
}
P2
  
```

Problem 2: Two urns A and B contain a total of N balls. Assume that at time t there were exactly k balls in A. At time $t+1$, an urn is selected at random in proportion to its contents (i.e., A is chosen with probability k/N and B is chosen with probability $(N-k)/N$). Then one of the N balls is randomly selected and placed in the chosen urn.

Let X_t denote the number of balls in A at time t , so $X_t, t \geq 0$ defines a Markov chain. Determine the transition matrix for this Markov chain.

Since we have that

$$\begin{aligned} P(X_{t+1} = k-1 | X_t = k) &= \frac{k}{N} \\ P(X_{t+1} = k+1 | X_t = k) &= \frac{N-k}{N} \end{aligned}$$

Then when $k=0$, $P(X_{t+1} = k+1) = 1$; $P(X_{t+1} = k-1) = 0$, and when $k=1$, $P(X_{t+1} = k+1) = 0$; $P(X_{t+1} = k-1) = 1$; \dots So on.

Therefore, we will have the following transition matrix:

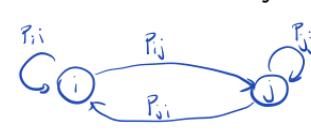
$$\begin{pmatrix} 0 & 1 & 0 & \dots & 0 \\ 1/N & 0 & (N-1)/N & \dots & 0 \\ 0 & 2/N & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{pmatrix}$$

That is, $P(X_{t+1} = k+1 | X_t = k) = \frac{N-k}{N}$ and $P(X_{t+1} = k-1 | X_t = k) = \frac{k}{N}$.

Transition Probabilities (3)

Another way to visualize the transition probabilities of a Markov chain is with a **transition state diagram**:

- Each state in the state space is represented by a node/vertex.
- Each nonzero transition probability P_{ij} is represented by an arrow from vertex i to vertex j .



where

- All the entries are non-negative:

$$P_{ij} \geq 0, \text{ for all } i, j. \quad (3)$$

- The sum of each row is 1:

$$\sum_{j=0}^N P_{ij} = \sum_{j=0}^N P(X_{n+1} = j | X_n = i) = 1, \text{ for all } i. \quad (4)$$

A random walk is a stochastic process:

$$\{X_0, X_1, X_2, \dots, X_n, \dots\},$$

(5)

$$X_1 = X_0 + Y_1 \quad Y_1 \sim \text{Ber}(P(Y_1=1) = p, P(Y_1=-1) = q)$$

$$X_2 = X_1 + Y_2 \quad \vdots$$

$$X_n = X_{n-1} + Y_n$$

$$Y = (-1, 1, \dots)$$

$$\Rightarrow X_0 + \sum_{i=1}^n Y_i$$

$$P(X_{n+1} = j | X_n = i)$$

$$P(X_{n+k} = j | X_n = i)$$

set.seed(999) # for reproducibility

```

n <- 1000 # specify length of random walk
p <- 0.5 # specify P(Y = 1)
Y <- sample(c(1, -1), size = n, replace = TRUE,
            prob = c(p, 1 - p)) # Generate n iid samples from Y
X <- c(0, cumsum(Y)) # Compute the random walk X
plot(X, type = "l", las = 1) # Plot the random walk over time
  
```

The Limiting Distribution

Let $\pi = (\pi_0, \pi_1, \pi_2, \dots, \pi_N)$ be a probability distribution on the state space $\{0, 1, 2, \dots, N\}$. We say π is the **limiting distribution** of a Markov chain $\{X_0, X_1, X_2, \dots\}$ if

$$\lim_{n \rightarrow \infty} P(X_n = j | X_0 = i) = \pi_j, \quad (7)$$

for all $i, j \in \{0, 1, 2, \dots, N\}$.

Assume the limiting dist. exists. Then for n large enough

$$P(X_n = k | X_0 = i) = \pi_k \quad \text{and} \quad P(X_{n+1} = j | X_0 = i) = \pi_j$$

find $\pi = (\pi_0, \pi_1, \dots, \pi_N)$

$$\begin{aligned} \pi_j &= P(X_{n+1} = j | X_0 = i) \\ &= \sum_{k=0}^N P(X_{n+1} = j, X_n = k | X_0 = i) \\ &= \sum_{k=0}^N P(X_{n+1} = j | X_n = k, X_0 = i) \cdot P(X_n = k | X_0 = i) \\ &= \sum_{k=0}^N P(X_{n+1} = j | X_n = k) \pi_k \\ &= \sum_{k=0}^N P_{kj} \pi_k = \sum_{k=0}^N \pi_k P_{kj} \\ \left\{ \begin{array}{l} \pi_j = \sum_{k=0}^N \pi_k P_{kj} \\ \pi_0 + \pi_1 + \dots + \pi_N = 1 \end{array} \right. \quad \text{①} \quad \text{②} \quad \pi = \pi_0 \pi_1 \dots \pi_N \end{aligned}$$

$$\begin{aligned} \text{Ex: } P &= \begin{bmatrix} 0.7 & 0.3 \\ 0.4 & 0.6 \end{bmatrix} \\ &\begin{array}{c} \text{0.7} \xrightarrow{\text{0.3}} \text{0.3} \xrightarrow{\text{0.6}} \text{0.6} \\ \text{0.4} \end{array} \\ \pi_0 &= \pi_0 P_{00} + \pi_1 P_{01} \\ \pi_1 &= \pi_0 P_{10} + \pi_1 P_{11} \\ &= \pi_0 0.7 + \pi_1 0.4 \\ \pi_0 + \pi_1 &= 1 \\ 0.3 \pi_0 &= 0.4 \pi_1 \Rightarrow \pi_1 = \frac{3}{4} \pi_0 \\ \pi_0 + \frac{3}{4} \pi_0 &= 1 \Rightarrow \pi_0 = \frac{4}{7} \\ \pi_1 &= 1 - \frac{4}{7} = \frac{3}{7} \end{aligned}$$

$$\pi = \left\{ \frac{4}{7}, \frac{3}{7} \right\}$$

Suppose $X_0 \sim \pi$ we can show $X_i \sim \pi$

$$P(X_1 = 0) = \pi_0$$

$$= \sum_{k=0}^N P(X_1 = 0 | X_0 = k) \cdot P(X_0 = k)$$

$$= P_{00} \pi_0 + P_{10} \pi_1$$

$$= 0.7 \frac{4}{7} + 0.4 \frac{3}{7} = \frac{4}{7} = P(X_0 = 0) = \pi_0$$

If a Markov chain $\{X_0, X_1, X_2, \dots\}$ is irreducible, aperiodic, and has a stationary distribution $\pi = (\pi_0, \pi_1, \pi_2, \dots, \pi_N)$, then

$$\lim_{n \rightarrow \infty} P(X_n = j | X_0 = i) = \pi_j.$$

That is, π is the **limiting distribution** of the Markov chain, and π is uniquely determined by the system of equations

$$\begin{cases} \pi = \pi \mathbb{P} \\ \sum_{i=0}^N \pi_i = 1 \end{cases} \quad \text{where } \mathbb{P} \text{ is the transition matrix of the Markov chain.}$$

Chapman-Kolmogorov Equations

We have defined the one-step transition probabilities P_{ij} . We now define the probability that the chain moves from state i to state j in n steps is $P_{ij}^{(n)}$.

$$P_{ij}^{(n)} = P(X_{n+k} = j | X_k = i), \quad n \neq 0, i, j \neq 0$$

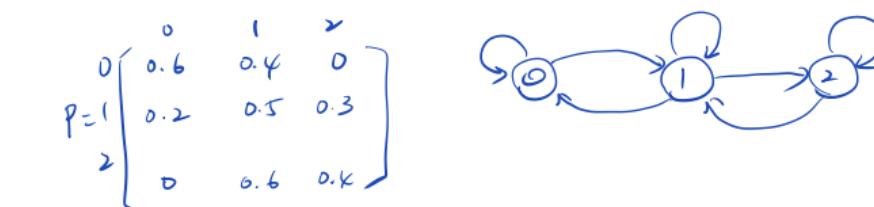
The Chapman-Kolmogorov equations provide a method for computing these n -step transition probabilities:

$$P_{ij}^{(n+m)} = \sum_{k=0}^{\infty} P_{ik}^{(n)} P_{kj}^{(m)} \text{ for all } n, m \neq 0, \text{ all } i, j. \quad (6)$$

Irreducible Markov Chains (1)

Definition:

- State j is **accessible** from state i , denoted by $i \rightarrow j$, if there is a path from state i to state j .
- Two states i and j **communicate**, denoted by $i \leftrightarrow j$, if each state is accessible from the other. In other words, if $i \rightarrow j$ and $j \rightarrow i$, then $i \leftrightarrow j$.
- A Markov chain is **irreducible** if all states communicate with each other. A Markov chain is **reducible** if it is not irreducible.



Irreducible Markov Chains (2)

Accessibility can be defined more rigorously using n -step transition probabilities:

Let $\{X_0, X_1, X_2, \dots\}$ be a Markov chain with state space $\{0, 1, 2, \dots, N\}$. The **n -step transition probability** from state i to j is defined by

$$P_{ij}^{(n)} := P(X_n = j | X_0 = i). \quad (8)$$

Note that when $n = 1$, $P_{ij}^{(1)} = P_{ij}$.

State j is **accessible** from state i , denoted by $i \rightarrow j$, if there is positive probability that state j can be reached from state i in a finite number of transitions. In other words, $P_{ij}^{(n)} > 0$ for some n .

Reducible

$$P = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0.6 & 0.4 & 0 & 0 \\ 0.2 & 0.5 & 0.3 & 0 \\ 0 & 0 & 0.6 & 0.4 \end{bmatrix}$$



$$d(0) = \gcd(1, 2, 3, 4, \dots) = 1$$

$$d(3) = 1$$

$$d(1) = ?$$

$$1 \rightarrow$$

Bayesian Perspective

In the Bayesian perspective, we are able to take our prior beliefs into account. We represent our beliefs about θ prior to observing data by a **prior distribution** $P(\theta)$.

- The prior distribution can represent past information, such as past experiments or literature, or subjective beliefs from a knowledgeable person.
- If no prior information is available (or we do not want to take it into account), we can use an **uninformative (or flat) prior**, which assigns equal density to all possibilities of the parameter.
- When using an uninformative prior, Bayesian estimators tends to be similar (sometimes identical) to frequentist estimators: The data easily outweighs a prior with no information.
- The **posterior distribution** $f(\theta|y)$ represents our updated beliefs about θ after observing the data. If the posterior is in the same parametric family as the prior, the prior and posterior are called **conjugate distributions**, and the prior is called a **conjugate prior** for the likelihood.

$$\begin{aligned} \theta | Y &\sim f(\cdot) & P(A|B) &= \frac{P(A,B)}{P(B)} \\ f(\theta|Y) &= \frac{f(Y|\theta) \cdot f(\theta)}{\int f(Y|\theta) \cdot f(\theta) d\theta} & = \frac{P(B|A) \cdot P(A)}{P(B)} \\ &\propto \frac{f(Y|\theta) \cdot f(\theta | a_1, a_2, \dots, a_n)}{\text{likelihood} \quad \text{prior}} \\ &f(Y|\theta_1, \theta_2, \dots, \theta_k) \cdot P(\theta_1) P(\theta_2) \cdots P(\theta_k) \end{aligned}$$

Example: Beta-Binomial Model

For our coin flipping example, suppose our prior is $P(\theta) \sim \text{Beta}(\alpha, \beta)$. Find the posterior mean.

$$\begin{aligned} Y &\text{iid } \text{Bin}(n, \theta) & P(\theta) &= \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} \\ f(\theta|Y) &\propto \binom{n}{Y} \theta^Y (1-\theta)^{n-Y} \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} \\ &= \binom{n}{Y} \frac{1}{B(\alpha, \beta)} \theta^{Y+\alpha-1} (1-\theta)^{n-Y+\beta-1} \\ &\propto \theta^{Y+\alpha-1} (1-\theta)^{n-Y+\beta-1} & \text{let } \lambda' = Y + \alpha \\ &= \theta^{\lambda'-1} (1-\theta)^{\beta'-1} & \beta' = n - Y + \beta \\ f(\theta|Y) &= \text{Beta}(\lambda', \beta') \end{aligned}$$

$$\begin{aligned} E[\theta|Y] &= \frac{\lambda'}{\lambda' + \beta'} \\ &= \frac{Y + \alpha}{Y + \alpha + n - Y + \beta} \\ &= \frac{\alpha}{n + \beta} + \frac{1 - \alpha}{n + \beta} \frac{Y + \alpha}{n + \beta} \\ &= w \frac{Y}{n} + (1-w) \frac{\alpha}{n + \beta} & \text{Recall: } \theta \sim \text{Beta}(\lambda, \beta) \\ & \quad \boxed{w} \quad \boxed{1-w} & E(\theta) = \frac{\lambda}{\lambda + \beta} \\ & \quad \boxed{\frac{\alpha}{n + \beta}} \quad \boxed{\frac{1 - \alpha}{n + \beta}} & \text{Var}(\theta) = \frac{\lambda \beta}{(\lambda + \beta)^2 (\lambda + \beta + 1)} \\ & \quad \boxed{\frac{\alpha}{n + \beta}} & \text{Mode} = \frac{\lambda - 1}{\lambda + \beta - 2} \\ & \quad \boxed{\frac{\alpha}{n + \beta}} & \end{aligned}$$

if Y, n are large, $E(\theta) \approx E(\theta|Y)$
if Y, n are small, $E(\theta) \approx E(\theta|Y)$

```
alpha <- 6
beta <- 2
sum_y1 <- 20
n1 <- 5
sum_y2 <- 80
n2 <- 20

posterior_alpha1 <- alpha + sum_y1
posterior_beta1 <- beta + n1
posterior_alpha2 <- alpha + sum_y2
posterior_beta2 <- beta + n2

x_seq <- seq(0, 15, length.out = 1000)
plot(x_seq, dgamma(x_seq, shape = alpha, rate = beta), type = 'l',
      col = 'blue', lwd = 2, ylim = c(0, 0.2), ylab = 'Density',
      xlab = expression(lambda), main = 'Gamma Distributions')

# Posterior plot for first scenario
lines(x_seq, dgamma(x_seq, shape = posterior_alpha1, rate = posterior_beta1),
      col = 'red', lwd = 2)

# Posterior plot for second scenario
lines(x_seq, dgamma(x_seq, shape = posterior_alpha2, rate = posterior_beta2),
      col = 'green', lwd = 2)

legend("topright",
      legend = c("Prior", "Posterior with n=5", "Posterior with n=20"),
      col = c("blue", "red", "green"), lwd = 2)
```

Problem 4: We want to model the number of siblings people have in a certain population. We can model the number of siblings a person has as a Poisson random variable Y for some unknown mean parameter λ . The probability mass function of $Y \sim \text{Pois}(\lambda)$ is

$$f(y|\lambda) = \frac{e^{-\lambda} \lambda^y}{y!}, \text{ for } y = 0, 1, 2, \dots$$

Suppose, before observing any data, we model our prior beliefs about λ by a gamma distribution $\text{Gamma}(\alpha, \beta)$ with hyperparameters $\alpha, \beta > 0$, i.e,

$$\pi(\lambda) = \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\beta\lambda}, \text{ for } \lambda > 0$$

where $\Gamma(z) = \int x^{z-1} e^{-x} dx$. Note that the prior mean is $E(\lambda) = \frac{\alpha}{\beta}$ and the prior variance is $\text{Var}(\lambda) = \frac{\alpha}{\beta^2}$. The prior mode is $\text{mode}(\lambda) = \frac{\alpha-1}{\beta}$. Suppose we observe data $y_1, y_2, \dots, y_n \sim \text{iid Pois}(\lambda)$. Let $y = (y_1, y_2, \dots, y_n)$.

- (a) Show that the gamma distribution is a conjugate prior for the Poisson likelihood. More specifically, show that the posterior distribution $\pi(\lambda|y)$ is of the form

$$\pi(\lambda|y) \sim \text{Gamma}(\alpha + \sum_{i=1}^n y_i, \beta + n)$$

Proof:

$$\begin{aligned} L(\lambda|y) &= \frac{e^{-n\lambda} \lambda^{\sum y_i}}{\prod y!} \\ \pi(\lambda|y) &\propto \frac{e^{-n\lambda} \lambda^{\sum y_i}}{\prod y!} \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\beta\lambda} \\ \pi(\lambda|y) &\propto \lambda^{\sum y_i + \alpha - 1} e^{-(n+\beta)\lambda} \\ \pi(\lambda|y) &\sim \text{Gamma}(\alpha + \sum_{i=1}^n y_i, \beta + n) \end{aligned}$$

- (b) Show that the posterior mean $E(\lambda|y)$ is a weighted average of the prior mean and the sample mean $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$. That is, find the weight w so that

$$E(\lambda|y) = w \frac{\alpha}{\beta} + (1-w) \bar{y}$$

Proof:

$$\begin{aligned} E(\lambda|y) &= \frac{\alpha + \sum_{i=1}^n y_i}{\beta + n} \\ w \frac{\alpha}{\beta} + (1-w) \bar{y} &= \frac{\alpha + n \bar{y}}{\beta + n} \\ w &= \frac{\beta}{\beta + n} \end{aligned}$$

- (c) What is $\lim_{n \rightarrow \infty} E(\lambda|y)$? What does this limit represent?

$$\lim_{n \rightarrow \infty} E(\lambda|y) = \bar{y}$$

This limit represents that as we get more data, the posterior mean converges to the sample mean, and the influence of the prior diminishes.

An interval $[\ell(y), u(y)]$, based on the observed data $Y = y$, is a $100(1 - \alpha)\%$ credible interval for θ if

$$P[\ell(y) < \theta < u(y) | Y = y] = 1 - \alpha.$$

The probability $1 - \alpha$ is called the **(Bayesian) coverage probability**. The interpretation of a credible interval is that it describes the information about the location of the true value of θ after you have observed $Y = y$.

This is different from the frequentist interpretation of coverage probability, which describes the probability that the interval will cover the true value *before* the data is observed.

A Bayesian analogue to a frequentist confidence interval is to use posterior quantiles. If $\theta_{\alpha/2}$ and $\theta_{1-\alpha/2}$ are the $\alpha/2$ and $1 - \alpha/2$ posterior quantiles of θ , then

$$P(\theta_{\alpha/2} < \theta < \theta_{1-\alpha/2} | Y = y) = 1 - \alpha,$$

so $[\theta_{\alpha/2}, \theta_{1-\alpha/2}]$ is a $100(1 - \alpha)\%$ quantile-based credible interval for θ .

Here the quantile function for the Beta distribution in R is `qbeta()`.

The Metropolis-Hastings algorithm induces a Markov chain that converges to the stationary (limiting) distribution $\pi(x)$ if the Markov chain satisfies certain regularity conditions:

- Irreducibility
(All states communicate with each other.)
 - Aperiodicity
(Every state has a period of 1.)
 - Positive recurrence
(Expected time until the chain returns to any state is finite.)
- Positive recurrence is trivially satisfied for irreducible and aperiodic Markov chains with finite state spaces.
- ### Good Mixing Behavior
- If the spread of the proposal is too large:
 - The proposal moves through the state space quickly (good mixing behavior).
 - The proposed state may be far from the current state, and the probability of acceptance will be low.
 - If the spread of the proposal is too small:
 - The probability of acceptance will be high.
 - The chain will take a long time to move through the state space, and low density regions will be undersampled.
 - Both situations above will exhibit high **autocorrelation**: Correlation between subsequent values in the chain.
 - An ideal proposal distribution will balance between good mixing behavior and a high acceptance rate.

The idea of Markov Chain Monte Carlo (MCMC) is to generate a Markov Chain such that its stationary distribution is the target distribution.

Main Goals:

- To generate a sequence of correlated samples from $\pi(x)$.
- To estimate $E_\pi[g(X)] = \int g(x)\pi(x)dx$.

Let G be any specified irreducible Markov transition matrix (kernel)

$$G = i \left(\begin{array}{c} g(j|i) \end{array} \right) \text{ and } g(\cdot|i) \text{ is a proposal dist. funct. given } i$$

when $X_t = i$, generate $Y \sim g(\cdot|i)$

If $Y = j$, then set $X_{t+1} = j$ w/ prob. $\alpha(i,j)$

If $Y \neq j$, then set $X_{t+1} = i$ w/ prob. $1 - \alpha(i,j)$

$$0 \leq \alpha(i,j) \leq 1$$

$$P(X_{t+1} = j | X_t = i) = g(j|i)$$

$$\text{Recall: } \pi(x) Q(x,y) = \pi(y) Q(y,x)$$

$$\pi(i) g(j|i) \alpha(i,j) = \pi(j) g(i|j) \alpha(j,i)$$

$$\text{If set } \alpha(i,j) = \frac{\pi(j) g(i|j)}{\pi(i) g(j|i)}, \text{ then } \alpha(j,i) = 1$$

$$\text{If set } \alpha(j,i) = \frac{\pi(i) g(j|i)}{\pi(j) g(i|j)}, \text{ then } \alpha(i,j) = 1$$

$$\text{Set } \alpha(i,j) = \min \left\{ \frac{\pi(j) g(i|j)}{\pi(i) g(j|i)}, 1 \right\}$$

$$Y = m_h = \frac{\pi(j) g(i|j)}{\pi(i) g(j|i)}$$

Metropolis-Hastings (M-H) sampling algorithm

There is a candidate point Y generated from a proposal distribution $g(\cdot|X_t)$. If this candidate point is accepted, the chain moves to state Y at time $t + 1$ and $X_{t+1} = Y$; otherwise the chain stays in state X_t and $X_{t+1} = X_t$. (Rizzo's book)

$$\begin{aligned} \alpha(i,j) &= \min \left\{ \frac{\pi(j) g(i|j)}{\pi(i) g(j|i)}, 1 \right\} \\ \text{if } \frac{\pi(j) g(i|j)}{\pi(i) g(j|i)} &\uparrow, \text{ most likely accept } X_{t+1} = j \\ &\downarrow, \text{ most likely accept } X_{t+1} = i \\ U &\sim \text{Unif}(0,1) \\ \text{if } (U < m_h) \text{ accept } X_{t+1} = j \\ &\text{ow let } X_{t+1} = i \end{aligned}$$

M-H Algorithm $\pi(x)$

- Choose a proposal distribution $g(\cdot|X_t)$.
- Generate X_0 from the distribution g .
- Repeat until the chain has converged to stationary distribution according to some criterion:
 - Generate Y from $g(\cdot|X_t)$
 - Generate U from Uniform(0, 1).
 - Compute $r = \frac{\pi(Y) g(X_t|Y)}{\pi(X_t) g(Y|X_t)}$
 - If $(U < r)$, accept Y and set $X_{t+1} = Y$; otherwise set $X_{t+1} = X_t$
 - Increment t .

Use the Metropolis-Hastings sampler to generate a sample from Standard Normal Distribution $N(0, 1)$.

```
set.seed(9999) # for reproducibility
n <- 10000 # specifying length of chain
X <- 0 # initialize chain
C <- 1
for (t in 2:n) {
    # Generate Y from proposal
    Y <- runif(1, X[t-1] - C, X[t-1] + C)
    # Compute MH ratio
    r <- min(1, exp(-0.5 * (Y^2 - X[t-1]^2)))
    U <- runif(1, 0, 1) # Generate U from Unif(0,1)
    if (U <= r) {
        X[t] <- Y # Move to Y if U <= r
    } else{
        X[t] <- X[t-1] # Stay at X[t-1] if U > r
    }
}
```

Consider two urns A and B, where we place n black balls in urn A and n white balls in urn B. At each step, we randomly choose a ball from each urn and interchange the two balls. Let X_t , $t \geq 0$, be the number of black balls in urn A at step t . Please find the transition probabilities.

```
##      1   2   3
## 1  0.6  0.3  0.1
## 2  0.6  0.3  0.1
## 3  0.2  0.2  0.6
```

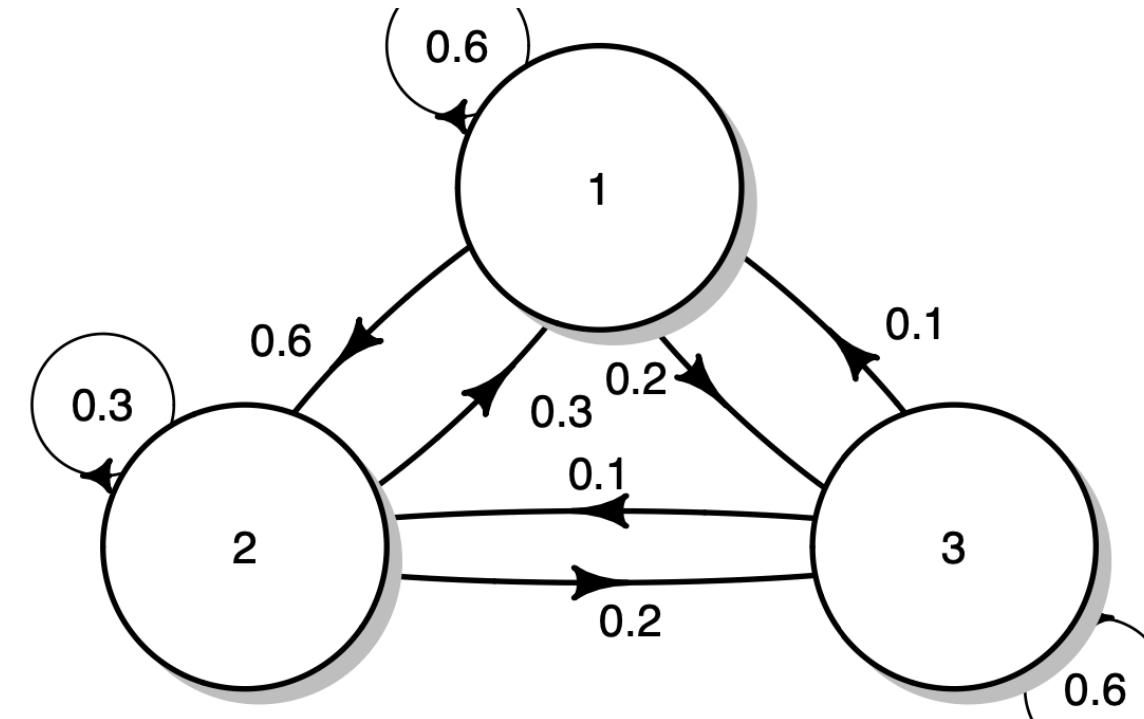
$$P(X_{t+1} = k+1 | X_t = k) = \frac{m-(n-k)}{m} \times \frac{k}{n}$$

$$P(X_{t+1} = k-1 | X_t = k) = \frac{(n-k)}{m} \times \frac{m-(n-k)}{n}$$

$$P(X_{t+1} = k | X_t = k) = 1 - P(X_{t+1} = k+1 | X_t = k) - P(X_{t+1} = k-1 | X_t = k)$$

```
plotmat(M,
  pos = c(1,2),
  self.cex = 0.5,
  self.shiftx = c(-0.1, -0.1, 0.1),
  self.shifty = c(0.1, 0.1, -0.1),
  main = "Transition Diagram")
```

Transition Diagram



Consider the transition matrix for the state space $\{0, 1\}$ below.

$\begin{pmatrix} 1-a & a \\ b & 1-b \end{pmatrix}$. Can we find a stationary distribution $\pi = [\pi_0, \pi_1]$ with this transition matrix? If so, please find the equations for computing π_0 and π_1 .

Edit View Insert Format Tools Table

12pt ▾ Paragraph ▾ B I ∪ A ▾ T² ▾ :

1. $\pi_0(1-a) + \pi_1 b = \pi_0$
2. $\pi_0 a + \pi_1 (1-b) = \pi_1$
3. $\pi_0 + \pi_1 = 1$

transition matrix:

$$\pi_0 = \frac{b}{a+b}$$

$$\pi_1 = \frac{a}{a+b}$$

University of California, Los Angeles
Department of Statistics

Statistics 100B

Instructor: Nicolas Christou

Probability Distributions - Summary

Discrete Distributions				
Distribution	Probability Mass Function	Mean	Variance	Moment-generating Function
Binomial	$P(X = x) = \binom{n}{x} p^x (1-p)^{n-x}$ $x = 0, 1, \dots, n$	np	$np(1-p)$	$[pe^t + (1-p)]^n$
Geometric	$P(X = x) = (1-p)^{x-1} p$ $x = 1, 2, \dots$	$\frac{1}{p}$	$\frac{1-p}{p^2}$	$\frac{pe^t}{1-(1-p)e^t}$
Negative Binomial	$P(X = x) = \binom{x-1}{r-1} p^r (1-p)^{x-r}$ $x = r, r+1, \dots$	$\frac{r}{p}$	$\frac{r(1-p)}{p^2}$	$[\frac{pe^t}{1-(1-p)e^t}]^r$
Hypergeometric	$P(X = x) = \frac{\binom{r}{x} \binom{N-r}{n-x}}{\binom{N}{n}}$ $x = 0, 1, \dots, n$ if $n \leq r$, $x = 0, 1, \dots, r$ if $n > r$	$\frac{nr}{N}$	$n \frac{r}{N} \frac{N-r}{N} \frac{N-n}{N-1}$	Fairly complicated!
Poisson	$P(X = x) = \frac{\lambda^x e^{-\lambda}}{x!}$ $x = 0, 1, \dots$	λ	λ	$\exp[\lambda(e^t - 1)]$
Continuous Distributions				
Distribution	Probability Density Function	Mean	Variance	Moment-generating Function
Uniform	$f(x) = \frac{1}{b-a}$ $a \leq x \leq b$	$\frac{a+b}{2}$	$\frac{(b-a)^2}{12}$	$\frac{e^{tb} - e^{ta}}{t(b-a)}$
Gamma	$f(x) = \frac{x^{\alpha-1} e^{-\frac{x}{\beta}}}{\beta^\alpha \Gamma(\alpha)}$, $\alpha, \beta > 0$, $x \geq 0$	$\alpha\beta$	$\alpha\beta^2$	$(1 - \beta t)^{-\alpha}$
Exponential	$f(x) = \lambda e^{-\lambda x}$, $\lambda > 0$, $x \geq 0$	$\frac{1}{\lambda}$	$\frac{1}{\lambda^2}$	$(1 - \frac{1}{\lambda} t)^{-1}$
Beta	$f(x) = \frac{x^{\alpha-1} (1-x)^{\beta-1}}{B(\alpha, \beta)}$ $\alpha > 0$, $\beta > 0$, $0 \leq x \leq 1$	$\frac{\alpha}{\alpha+\beta}$	$\frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$	
Normal	$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2}$ $-\infty < x < +\infty$	μ	σ^2	$e^{\mu t + \frac{t^2 \sigma^2}{2}}$

Remarks:

- Binomial:
- Geometric:
- Negative Binomial:
- Hypergeometric:
- Poisson:

- X represents the number of successes among n trials.
- X represents the number of trials needed until the first success.
- X represents the number of trials needed until r successes occur.
- X represents the number of items among the n selected that comes from the r group.
- X represents the number of events that occur in time, area, etc.