

MAST30027: Modern Applied Statistics

Week 9 Lab

1. Consider the following program, which performs a simulation experiment. The function `X.sim()` simulates some random variable X , and we wish to estimate $\mathbb{E}X$.

```
# set.seed(7)
# seed position 1
mu <- rep(0, 6)
for (i in 1:6) {
  # set.seed(7)
  # seed position 2
  X <- rep(0, 1000)
  for (j in 1:1000) {
    # set.seed(7)
    # seed position 3
    X[j] <- X.sim()
  }
  mu[i] <- mean(X)
}
spread <- max(mu) - min(mu)
mu.estimate <- mean(mu)
```

- (a) What is the value of `spread` used for?

Solution: To measure the variability of the estimates `mi[i]`.

- (b) If we uncomment the command `set.seed(7)` at seed position 3, then what is `spread`?

Solution: 0

- (c) If we uncomment the command `set.seed(7)` at seed position 2 (only), then what is `spread`?

Solution: 0

- (d) If we uncomment the command `set.seed(7)` at seed position 1 (only), then what is `spread`?

Solution: > 0

- (e) At which position should we set the seed?

Solution: 1

2. For $X \sim \text{Poisson}(\lambda)$ let $F(x) = \mathbb{P}(X \leq x)$ and $p(x) = \mathbb{P}(X = x)$. Show that the probability function satisfies

$$p(x+1) = \frac{\lambda}{x+1}p(x).$$

Using this write a function to calculate $p(0), p(1), \dots, p(x)$ and $F(x) = p(0) + p(1) + \dots + p(x)$.

If $X \in \mathbb{Z}_+$ is a random variable and `F(x)` is a function that returns the cdf F of X , then you can simulate X using the following program:

```
F.rand <- function () {
  u <- runif(1)
  x <- 0
  while (F(x) < u) {
    x <- x + 1
  }
  return(x)
}
```

In the case of the Poisson distribution, this program can be made more efficient by calculating F just once, instead of recalculating it every time you call the function `F(x)`. By using two new

variables, `p.x` and `F.x` for $p(x)$ and $F(x)$ respectively, modify this program so that instead of using the function `F(x)` it updates `p.x` and `F.x` within the `while` loop. Your program should have the form

```
F.rand <- function(lambda) {
  u <- runif(1)
  x <- 0
  p.x <- ?
  F.x <- ?
  while (F.x < u) {
    x <- x + 1
    p.x <- ?
    F.x <- ?
  }
  return(x)
}
```

You should ensure that at the start of the `while` loop you always have `p.x` equal to $p(x)$ and `F.x` equal to $F(x)$.

Check that your simulation works by choosing a parameter, generating a large number of random variables, using them to estimate the probability mass function, and comparing your estimates to the true values (which you can get using `dpois`).

Solution: We have

$$p(x+1) = \frac{e^{-\lambda} \lambda^{x+1}}{(x+1)!} = \frac{e^{-\lambda} \lambda^x}{x!} \frac{\lambda}{x+1} = p(x) \frac{\lambda}{x+1}$$

Thus we can calculate $F(x)$ as follows

```
> Fpois <- function(x, la) {
+   # poisson(la) cdf at x (assumed integer)
+   y <- 0
+   py <- exp(-la)
+   Fy <- py
+   while (y < x) {
+     y <- y + 1
+     py <- la*py/y
+     Fy <- Fy + py
+   }
+   return(Fy)
+ }
> sapply(0:4, Fpois, la=2) # check it works

[1] 0.1353353 0.4060058 0.6766764 0.8571235 0.9473470

> ppois(0:4, 2) # using R's built in function

[1] 0.1353353 0.4060058 0.6766764 0.8571235 0.9473470
```

To simulate from this cdf we have the following

```
> F.rand <- function(lambda) {
+   u <- runif(1)
+   x <- 0
+   p.x <- exp(-lambda)
+   F.x <- p.x
+   while (F.x < u) {
+     x <- x + 1
+     p.x <- lambda*p.x/x
+     F.x <- F.x + p.x
+   }
}
```

```

+   return(x)
+ }
> n <- 100000
> lambda <- 2
> F.sample <- rep(0, n)
> for (i in 1:n) F.sample[i] <- F.rand(lambda)
> table(F.sample)/n

F.sample
  0      1      2      3      4      5      6      7      8      9
0.13521 0.27026 0.27069 0.18002 0.08968 0.03748 0.01215 0.00337 0.00093 0.00015
10
0.00006

> round(dpois(0:9, 2), 5)

[1] 0.13534 0.27067 0.27067 0.18045 0.09022 0.03609 0.01203 0.00344 0.00086
[10] 0.00019

```

3. (a) Here is some code for simulating a discrete random variable Y . What is the probability mass function (pmf) of Y ?

```

Y.sim <- function() {
  U <- runif(1)
  Y <- 1
  while (U > 1 - 1/(1+Y)) {
    Y <- Y + 1
  }
  return(Y)
}

```

Let N be the number of times you go around the while loop when `Y.sim()` is called. What is $\mathbb{E}N$ and thus what is the expected time taken for this function to run?

- (b) Here is some code for simulating a discrete random variable Z . Show that Z has the same pmf as Y

```

Z.sim <- function() {
  Z <- ceiling(1/runif(1)) - 1
  return(Z)
}

```

Will this function be faster or slower than `Y.sim()`?

Solution: We have, for $y \geq 1$,

$$\mathbb{P}(Y = y) = 1 - \frac{1}{1+y} - \left(1 - \frac{1}{y}\right) = \frac{1}{y} - \frac{1}{1+y} = \frac{1}{y(1+y)}.$$

Thus

$$\begin{aligned}
 \mathbb{E}N &= \sum_{n=1}^{\infty} (n-1) \mathbb{P}(Y = n) \\
 &= \sum_{n=1}^{\infty} \frac{n-1}{n(1+n)} \\
 &= \left(\sum_{n=1}^{\infty} \frac{1}{1+n} \right) - 1 \\
 &= \infty.
 \end{aligned}$$

So the expected amount of time required for this algorithm to complete is infinite!

Put $Z = \lceil 1/U \rceil - 1$ where $U \sim U(0, 1)$, then

$$Z = z \iff z < 1/U \leq z + 1 \iff 1/(1+z) \leq U < 1/z.$$

Thus $\mathbb{P}(Z = z) = 1/z - 1/(1+z)$ as required. Clearly `Z.sim` will be much faster than `Y.sim` on average.

4. Consider the continuous random variable X with pdf given by:

$$f_X(x) = \frac{\exp(-x)}{(1 + \exp(-x))^2} \quad -\infty < x < \infty.$$

X is said to have a standard logistic distribution. Find the cdf for this random variable. Show how to simulate a rv with this cdf using the inversion method.

Solution: X has cdf

$$\begin{aligned} F(x) &= \int_{-\infty}^x f_X(y) dy \\ &= \left. \frac{1}{1 + e^{-y}} \right|_{-\infty}^x \\ &= \frac{1}{1 + e^{-x}} \end{aligned}$$

Thus $F^{-1}(y) = \log(y/(1-y)) = -\log((1/y) - 1)$ and we can simulate a sample of size n from the distribution of X as follows

```
> X.sim <- function(n=1) -log(1/runif(n) - 1)
```

5. The double exponential or Laplace distribution has the following density, for some $\lambda > 0$,

$$f(x) = \frac{\lambda}{2} e^{-\lambda|x|}, \text{ for } -\infty < x < \infty.$$

Plot the density and the cumulative distribution function, F .

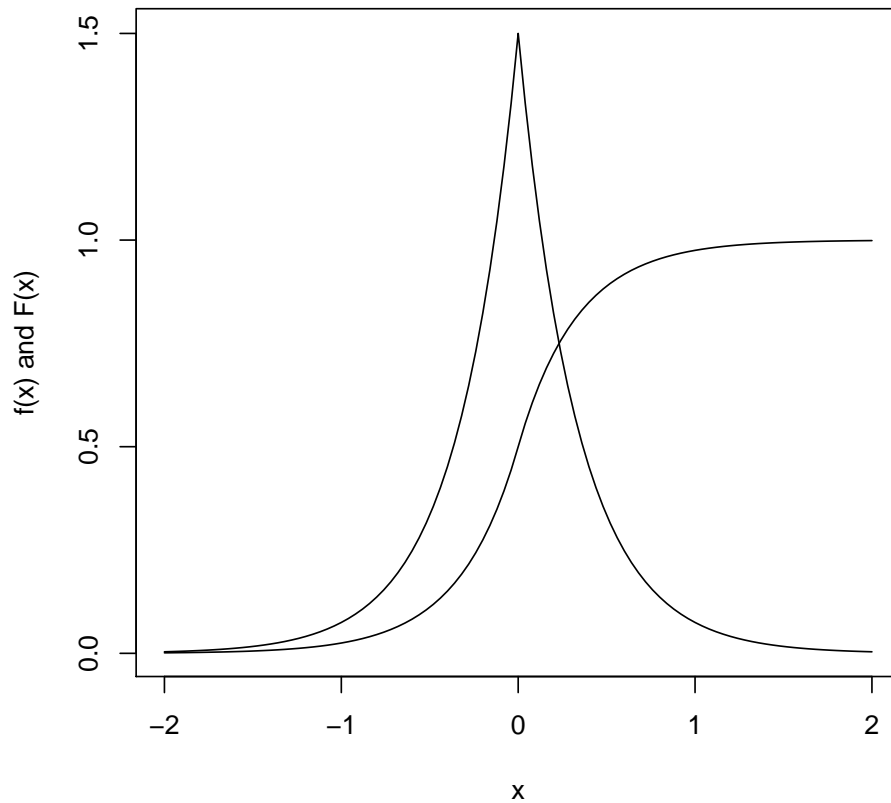
Suppose that A has an exponential distribution with rate λ , and that B is independent of A and takes on values $+1$ and -1 with equal probability. Show that AB has a Laplace distribution with parameter λ , then write a function to simulate Laplace rv's.

Solution: Integrating we get the cdf

$$F(x) = \begin{cases} \frac{1}{2} e^{\lambda x} & x < 0 \\ 1 - \frac{1}{2} e^{-\lambda x} & x \geq 0 \end{cases}$$

We plot the density and curvature for $\lambda = 3$.

```
> fx <- function(x, la) la*exp(-la*abs(x))/2
> Fx <- function(x, la) ifelse(x < 0, exp(la*x)/2, 1 - exp(-la*x)/2)
> la <- 3
> curve(fx(x, la), -2, 2, ylab="f(x) and F(x)")
> curve(Fx(x, la), -2, 2, add=TRUE)
```



Now

$$\begin{aligned}
 \mathbb{P}(AB \in (x, x + dx)) &= \frac{1}{2} \mathbb{P}(A \in (x, x + dx) | B = 1) + \frac{1}{2} \mathbb{P}(A \in (-x - dx, -x) | B = -1) \\
 &= \frac{1}{2} \lambda e^{-\lambda x} dx I_{x \geq 0} + 0 I_{x < 0} + 0 I_{x \geq 0} + \frac{1}{2} \lambda e^{-\lambda x} dx I_{x < 0} \\
 &= \frac{1}{2} \lambda e^{-\lambda |x|} dx
 \end{aligned}$$

Thus, rather than apply the inverse method to F , we can simulate an $\exp(\lambda)$ then toss a coin and with probability $1/2$ multiply by -1 . For example

```

> rLaplace <- function(la) {
+   A <- -log(runif(1))/la
+   B <- sign(runif(1) - 0.5)
+   return(A*B)
+ }

```

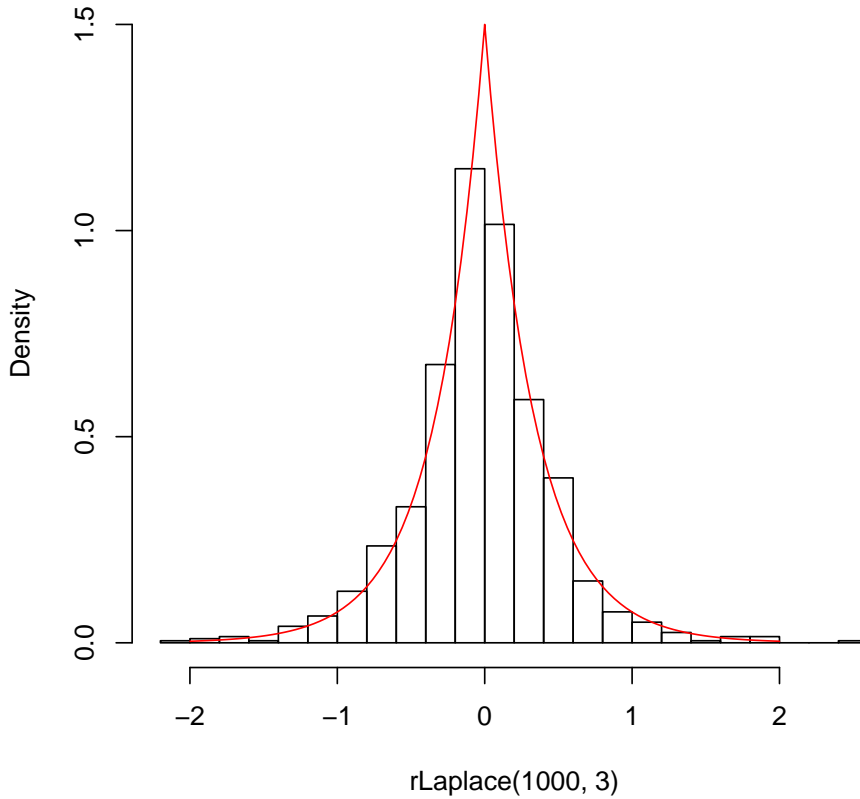
Or, in vectorised form,

```

> rLaplace <- function(n, la) {
+   A <- -log(runif(n))/la
+   B <- sign(runif(n) - 0.5)
+   return(A*B)
+ }
> hist(rLaplace(1000, 3), breaks=20, freq=FALSE, ylim=c(0,1.5))
> curve(fx(x,3), -2, 2, add=TRUE, col="red")

```

Histogram of rLaplace(1000, 3)



6. (a) Construct a rejection sampling algorithm to generate a truncated exponential distribution, which has the pdf $p(z) = \frac{e^{-z}}{1-e^{-1}}$, $0 < z < 1$.
- (b) Calculate the mean and variance for the pdf $p(z)$ in (a).
- (c) Write an R program to implement the algorithm in (a) and use it to generate a sample of 1000 observations. Plot a histogram of the sample. Calculate the sample mean and variance, and compare them with the results in (b).
- (d) Show that the following algorithm also simulates from the distribution in (a).
 - 1° Generate U from $\text{Unif}(0,1)$;
 - 2° If $U > e^{-1}$ then deliver $Z = -\ln(U)$; otherwise go to 1°.

Solution: When $0 < z < 1$, $p(z) < (1 - e^{-1})^{-1} < \infty$, so we can use a rectangular envelope to contain the density. We get the following rejection algorithm:

- 1° Generate $X \sim (0,1)$.
- 2° Generate $Y \sim (0, 1/(1 - e^{-1}))$ independently.
- 3° If $Y \leq p(X)$ return X ; otherwise go to 1°.

If Z has a truncated exponential distribution as above, then

$$\begin{aligned}
 E(Z) &= \int_0^1 z \frac{e^{-z}}{1 - e^{-1}} dz = \frac{1 - 2e^{-1}}{1 - e^{-1}} = 0.418023. \\
 E(Z^2) &= \int_0^1 z^2 \frac{e^{-z}}{1 - e^{-1}} dz = \frac{2 - 5e^{-1}}{1 - e^{-1}}. \\
 \text{Var}(Z) &= E(Z^2) - (EZ)^2 \\
 &= \frac{2 - 5e^{-1}}{1 - e^{-1}} - \left(\frac{1 - 2e^{-1}}{1 - e^{-1}} \right)^2 = 0.079326.
 \end{aligned}$$

We check our simulation algorithm by generating a sample and comparing the sample mean and variance to these theoretical values.

```
> trunc_exp <- function(){
+   # generate a truncated exponential random variable using the rejection algorithm
+   x <- runif(1)
+   y <- runif(1, 0, 1/(1 - exp(-1)))
+   while (y > exp(-x)/(1 - exp(-1))) {
+     x <- runif(1)
+     y <- runif(1, 0, 1/(1 - exp(-1)))
+   }
+   return(x)
+ }
> # generate a sample of size n
> set.seed(30027)
> n <- 10000
> z <- rep(0, n)
> for (i in 1:n) z[i] <- trunc_exp()
> # sample mean and variance
> mean(z)

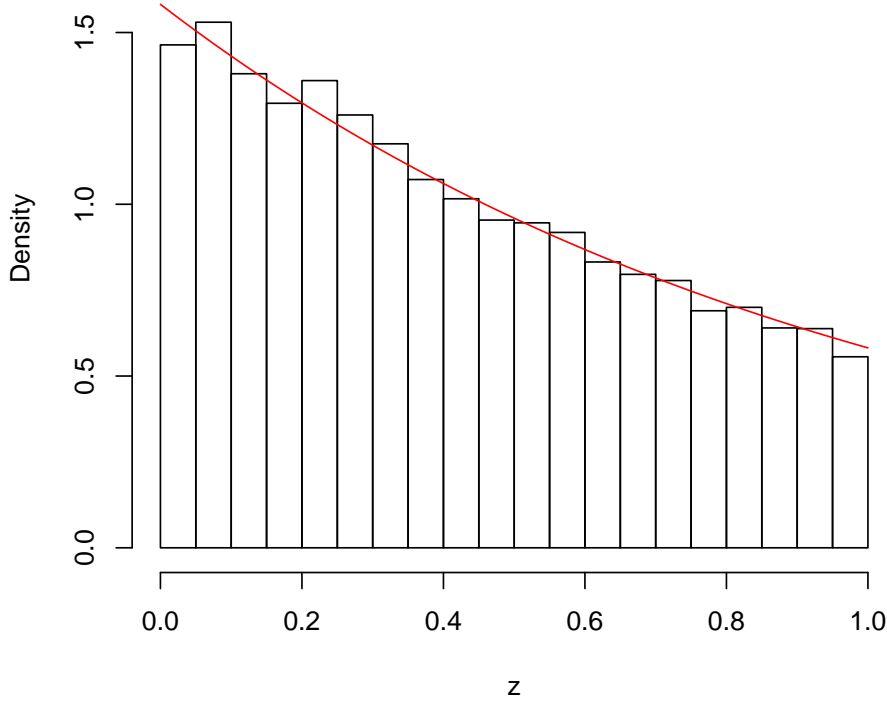
[1] 0.4159879

> var(z)

[1] 0.07836399

> # histogram with density on top
> hist(z, freq=FALSE, ylim=c(0,1.8))
> curve(exp(-x)/(1 - exp(-1)), 0, 1, add=TRUE, col="red")
```

Histogram of z



Now consider the alternative algorithm. According to the algorithm, the cdf of Z is

$$\begin{aligned}
 F_Z(z) &= P(Z \leq z) \\
 &= P(-\ln(U) \leq z | U > e^{-1}) \\
 &= P(U \geq e^{-z} | U > e^{-1}) \\
 &= \frac{P(U \geq \max\{e^{-z}, e^{-1}\})}{P(U > e^{-1})} \\
 &= \begin{cases} \frac{P(U \geq 1)}{1 - e^{-1}} = 0 & \text{if } z \leq 0, \\ \frac{P(U \geq e^{-z})}{1 - e^{-1}} = \frac{1 - e^{-z}}{1 - e^{-1}} & \text{if } 0 < z < 1, \\ \frac{P(U \geq e^{-1})}{1 - e^{-1}} = 1 & \text{if } z \geq 1 \end{cases}
 \end{aligned}$$

Thus the pdf of Z is $F'_Z(z) = \frac{e^{-z}}{1 - e^{-1}}$, $0 < z < 1$, the same as $p(z)$.

7. Suppose $g(u)$ is a decreasing function of u . Let a random number X be generated by the following algorithm:

- 1° Generate U from $\text{Unif}(0, 1)$ and V from $\text{Unif}(0, 1)$ independently.
- 2° If $U + V < 1$, then deliver $X = g(U)$; otherwise, go to 1°.

Show that the distribution function of X is given by

$$F(x) = P(X \leq x) = (1 - g^{-1}(x))^2, \quad g(1) \leq x \leq g(0).$$

Solution: According to the algorithm, the cdf of X is

$$F(x) = P(X \leq x)$$

$$\begin{aligned}
&= P(g(U) \leq x | U + V < 1) \\
&= P(U \geq g^{-1}(x) | U + V < 1) \quad (\text{because } g(u) \text{ is decreasing}) \\
&= \frac{P(U \geq g^{-1}(x), U < 1 - V)}{P(U + V < 1)} \\
&= \frac{\int_0^{1-g^{-1}(x)} \int_{g^{-1}(x)}^{1-v} 1 \, du \, dv}{\int_0^1 \int_0^{1-v} 1 \, du \, dv} \\
&= \frac{\frac{1}{2}(1 - g^{-1}(x))^2}{\frac{1}{2}} = (1 - g^{-1}(x))^2.
\end{aligned}$$

8. A random number X is to be generated by the following rejection algorithm

1° Generate two independent $\text{Unif}(0, 1)$ random numbers U and V .

2° If $U^2 + V^2 \leq 1$, deliver $X = U$; otherwise return to 1°.

(a) Find the cdf and pdf of X . Also find the mean and variance of X .

HINT: The integral $\int_0^x \sqrt{1-u^2} \, du = \frac{1}{2}(\arcsin x + x\sqrt{1-x^2})$.

(b) What is the efficiency of the algorithm? Denote by N the number of times that step 1° is to be executed to get one X value. Name the distribution of N . Also write down the mean and standard deviation of N .

(c) Write an R program to implement the algorithm. Then generate a sample of 1000 X values, and give it a numeric and a graphic summary.

Solution: The support of X is $(0, 1)$, the same as that of U . The cdf of X is

$$\begin{aligned}
F(x) &= P(X \leq x) \\
&= P(U \leq x | U^2 + V^2 \leq 1) \\
&= \frac{\int_0^x \int_0^{\sqrt{1-u^2}} 1 \, dv \, du}{\int_0^1 \int_0^{\sqrt{1-u^2}} 1 \, dv \, du} \\
&= \frac{\int_0^x \sqrt{1-u^2} \, du}{\int_0^1 \sqrt{1-u^2} \, du} \\
&= \frac{2}{\pi} (\arcsin x + x\sqrt{1-x^2}), \quad 0 < x < 1.
\end{aligned}$$

The pdf of X is

$$f(x) = F'(x) = \frac{4}{\pi} \sqrt{1-x^2}, \quad 0 < x < 1.$$

$$\begin{aligned}
E(X) &= \int_0^1 x f(x) \, dx = \int_0^1 x \frac{4}{\pi} \sqrt{1-x^2} \, dx \\
&= -\frac{4}{3\pi} (1-x^2)^{\frac{3}{2}} \Big|_0^1 \\
&= \frac{4}{3\pi}.
\end{aligned}$$

$$\begin{aligned}
E(X^2) &= \int_0^1 x^2 \frac{4}{\pi} \sqrt{1-x^2} \, dx \\
&= \frac{1}{2\pi} (\arcsin x - \frac{1}{4} \sin(4 \arcsin x)) \Big|_0^1 \\
&= \frac{1}{4}.
\end{aligned}$$

$$\text{Var}(X) = \frac{1}{4} - \left(\frac{4}{3\pi} \right)^2.$$

The efficiency $\tau = P(U^2 + V^2 \leq 1) = \int_0^1 \int_0^{\sqrt{1-u^2}} 1 dv du = \frac{\pi}{4}$.

$$\begin{aligned} N &\stackrel{d}{=} \text{GEO}\left(\frac{\pi}{4}\right) \\ E(N) &= \tau^{-1} = \frac{4}{\pi} \\ \text{Var}(X) &= \tau^{-2}(1 - \tau) = \frac{16 - 4\pi}{\pi^2} \end{aligned}$$

```
> xsim <- function(){
+   u <- runif(1)
+   v <- runif(1)
+   while(u^2 + v^2 > 1){
+     u <- runif(1)
+     v <- runif(1)
+   }
+   return(u)
+ }
> ##### Test via simulation
> set.seed(12345)
> n <- 10000
> x <- rep(0, n)
> for (i in 1:n) x[i] <- xsim()
> summary(x)
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
0.0000519	0.2018006	0.4098974	0.4264121	0.6369295	0.9971188

```
> hist(x, freq=FALSE)
> curve((4/pi)*sqrt(1-x^2), from=0, to=1, add=TRUE, lwd=2, lty=2)
```

