

# HW: Week 6

36-350 – Statistical Computing

Week 6 – Spring 2021

Name: Cherie Hua

Andrew ID: cxhua

You must submit **your own** lab as a PDF file on Gradescope.

---

```
suppressWarnings(library(tidyverse))
```

```
## -- Attaching packages ----- tidyverse 1.3.0
```

```
## v ggplot2 3.3.2    v purrr  0.3.4
## v tibble  3.0.3    v dplyr  1.0.2
## v tidyr   1.1.2    v stringr 1.4.0
## v readr   1.3.1    v forcats 0.5.0
```

```
## -- Conflicts ----- tidyverse_conflicts()
```

```
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()
```

---

## HW Length Cap Instructions

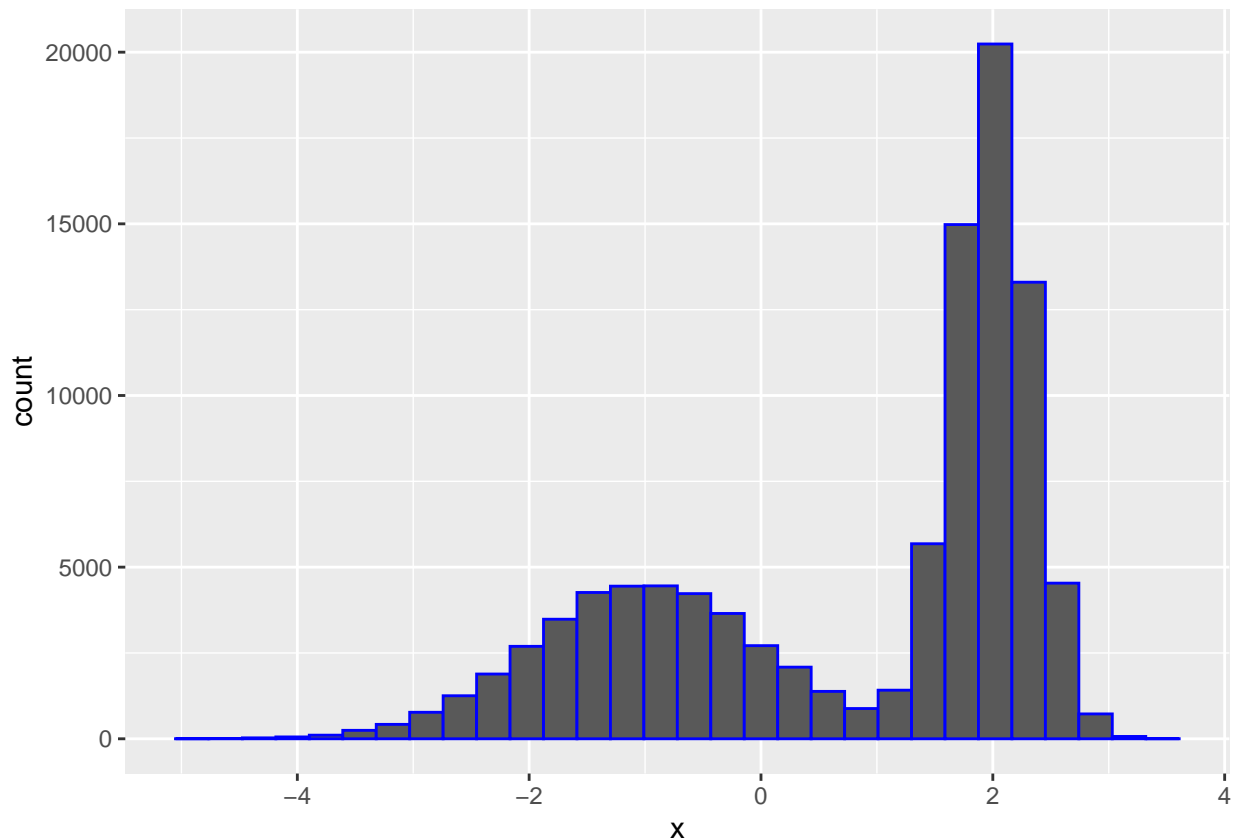
- If the question requires you to print a data frame in your solution e.g. `q1_out_df`, you must first apply `head(q1_out_df, 30)` and `dim(q1_out_df)` in the final knitted pdf output for such a data frame.
- Please note that this only applies if you are knitting the Rmd to a pdf, for Gradescope submission purposes.
- If you are using the data frame output for visualization purposes (for example), use the entire data frame in your exploration
- The **maximum allowable length** of knitted pdf HW submission is **30 pages**. Submissions exceeding this length *will not be graded* by the TAs. All pages must be tagged as usual for the required questions per the usual policy
- For any concerns about HW length for submission, please reach out on Piazza during office hours

## Question 1

(10 points)

Create a Gaussian mixture model sampler. In this sampler, a datum has a 40% chance of being sampled from a  $N(-1,1)$  distribution, and a 60% chance of being sampled from a  $N(2,1/9)$  distribution. Sample 100,000 data and create a density histogram of your result. Hint: use `sample()` with `replace` set to `TRUE` and an appropriate vector for `prob` in order to determine which of your 100,000 data should randomly be assigned to distribution 1 as opposed to distribution 2. Also note that if you create a sample of data from distribution 1 and another sample from distribution 2, you can simply combine them by doing, e.g., `x = c(sample1, sample2)`, where `x` becomes a vector of length 100,000.

```
set.seed(123)
n = 100000
ind = sample(1:2, size = n, replace = TRUE, prob = c(0.4, 0.6))
sample1 = rnorm(sum(ind == 1), mean = -1, sd = 1)
sample2 = rnorm(sum(ind == 2), mean = 2, sd = 1/3)
x = c(sample1, sample2)
ggplot(data.frame(x = x), aes(x = x)) + geom_histogram(color = "blue", bins = 30)
```



## Question 2

(10 points)

What is the mean of the mixture model in Q1? Compute this via importance sampling, with 100,000 sampled points. You should get an answer around 0.8 (which you can actually derive analytically: if  $X \sim N(-1, 1)$

and  $Y \sim N(2, 1/9)$ , then  $E[0.4X + 0.6Y] = 0.4E[X] + 0.6E[Y] = -0.4 + 1.2 = 0.8$ .

```
n = 100000
h = rnorm(100000, 1, sd = 1/2)
g = x
f = function (x) {
  0.4 * rnorm(100000, mean = -1, sd = 1) + 0.6 * rnorm(100000, mean = 2, sd = 1/3)
}
mean(g*f(x)/h)
```

```
## [1] 1.240288
```

### Question 3

(10 points)

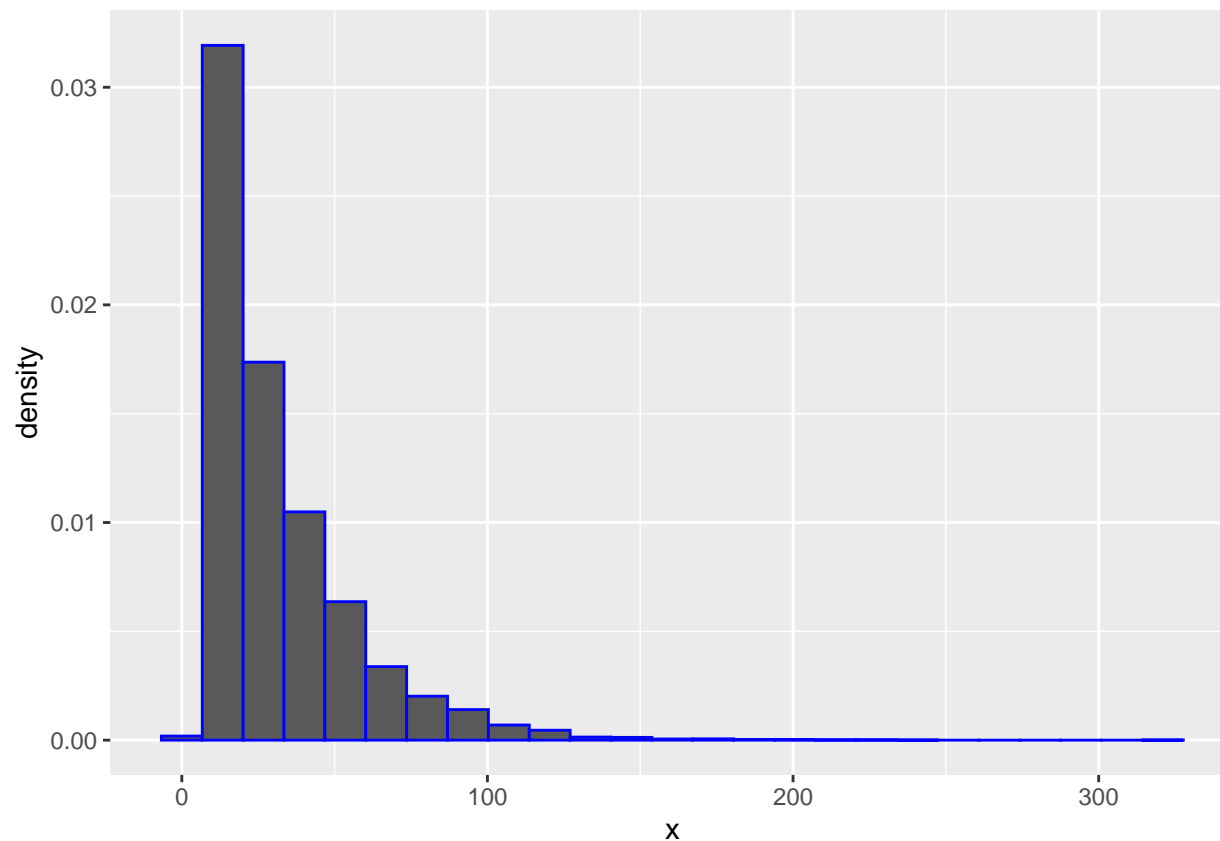
Remember the Chutes and Ladders question? (Q4 of HW 2.) Display a probability histogram that shows the empirical PDF for the number of spins, computed over 10,000 Chutes and Ladders games. Also display the average number of spins that it takes to win the game (approximately 39, give or take) and the minimum number of spins. (Display these two numbers using `cat()`, being sure to indicate which is the mean and which is the minimum number of spins.) (Free feel to use your code from HW 2 as a base for what you do here.)

```
set.seed(345)
#a game of chutes and ladders (single-player)
ladder.bottom = c(4, 10, 13, 15, 20, 24)
ladder.top = c(40, 80, 70, 69, 31, 46)
chute.bottom = c(11, 8, 22, 31, 28, 7)
chute.top = c(99, 32, 45, 73, 58, 51)

chutes.ladders <- function () {
  location = 0 #the player's position on the board
  i = 0 #keep track of iterations
  spins = 0
  while (location < 100) {
    spin = sample(1:6, 1)
    spins = spins + 1
    location = location + spin
    w = which(ladder.bottom == location)
    if (length(w) > 0) location = ladder.top[w]
    w = which(chute.top == location)
    if (length(w) > 0) location = chute.bottom[w]
    i = i + 1
  }
  return(spins)
}

all_spins <- vector(mode = "numeric", length = 10000)
for (i in 1:10000) {
  all_spins[i] = chutes.ladders()
}

ggplot(data.frame(x = all_spins), aes(x = x)) + geom_histogram(aes(y = ..density..), color = "blue", binwidth = 1)
```



```
cat("mean", mean(all_spins), "\n", "min", min(all_spins))
```

```
## mean 31.3744
## min 6
```

## Question 4

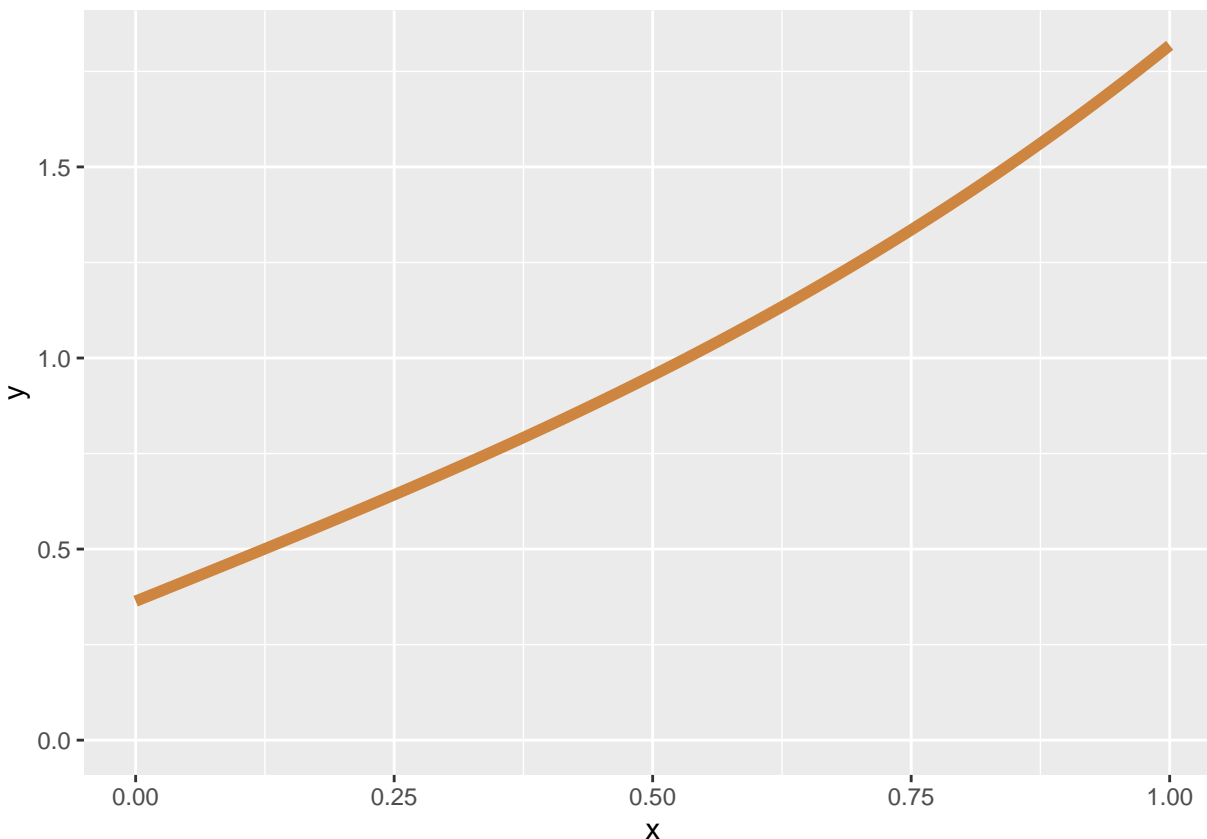
(10 points)

You are given the following distribution:

$$f(x) = \frac{4}{11}(x^3 + 3x + 1) \quad x \in [0, 1]$$

It looks like this:

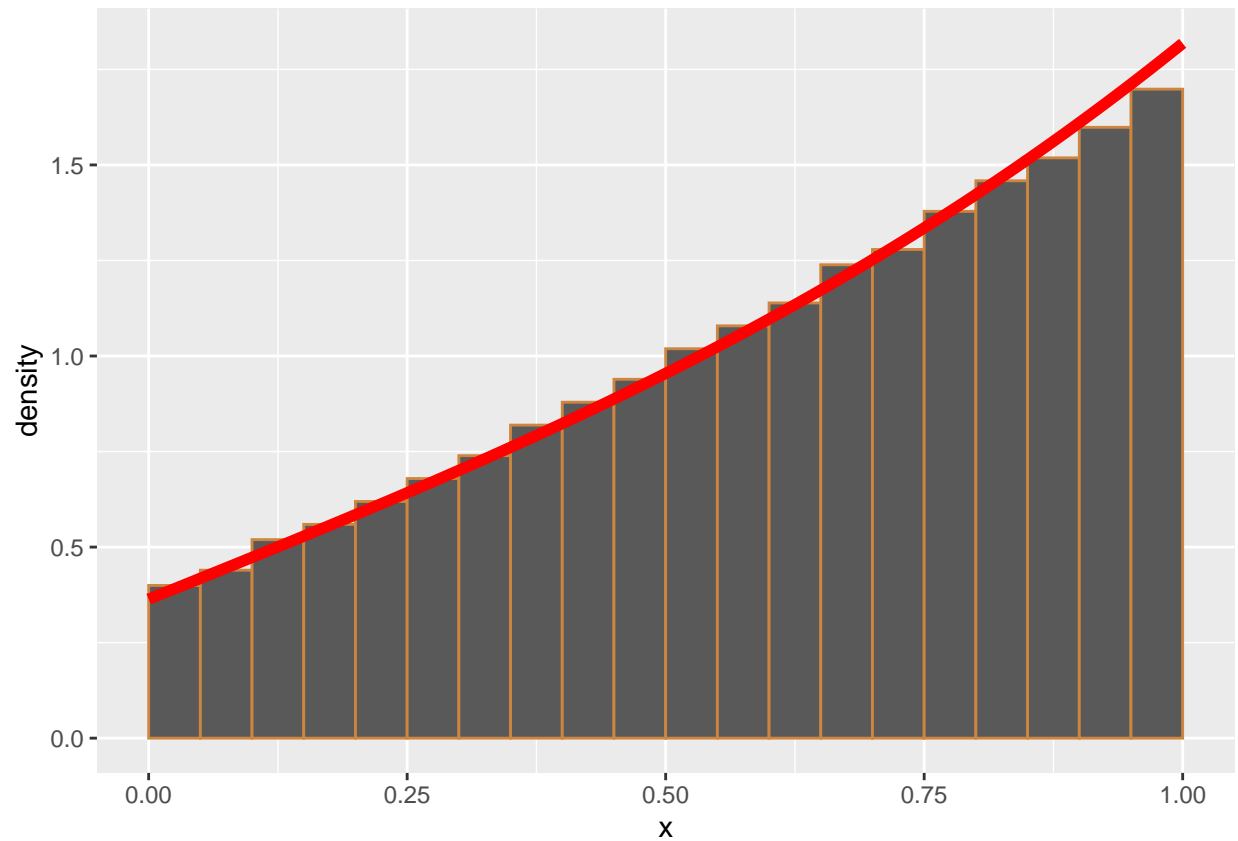
```
x = seq(0,1,by=0.001)
fx = 4*(x^3+3*x+1)/11
ggplot(data=data.frame(x=x,y=fx),mapping=aes(x=x,y=y)) +
  geom_line(col="peru",size=2) + ylim(0,max(fx))
```



Code up an inverse transform sampler that allows you to efficiently sample 1000 data from this distribution. Initially you will work with this and say “but...I cannot easily invert the CDF, because it’s a quartic and all.” To which I say “`polyroot()`, which will give you one real root between 0 and 1.” To which you will say, “how do I extract that real root?” To which I will say “If you save the output of `polyroot()` as `p`, then the real roots are given by `w = which(abs(Im(p))<1.e-6)`” (The 1.e-6 is a check against round-off error.) You then determine which value of `Re(p)[w]` is within the pdf bounds. Histogram your sample with the function line overlaid, and save your sample as `sample.it`. Note: pass a new argument to your histogram, `breaks=seq(0,1,by=0.05)`.

```
set.seed(444)
x = seq(0,1,by=0.001)
sample.it = rep(NA, 1000)
counter = 0
for (i in x) {
  p = polyroot(c(-i, 4/11, 6/11, 1/11))
  w = which(abs(Im(p)) < 1.e-6)
  val = p[which((Re(p)[w] >= 0) & (Re(p)[w] <= 1))]
  counter = counter + 1
  sample.it[counter] = val
}
ggplot(data=data.frame(x = as.numeric(sample.it)), mapping=aes(x=x)) + geom_histogram(col="peru", aes(y

## Warning in data.frame(x = as.numeric(sample.it)): imaginary parts discarded in
## coercion
```

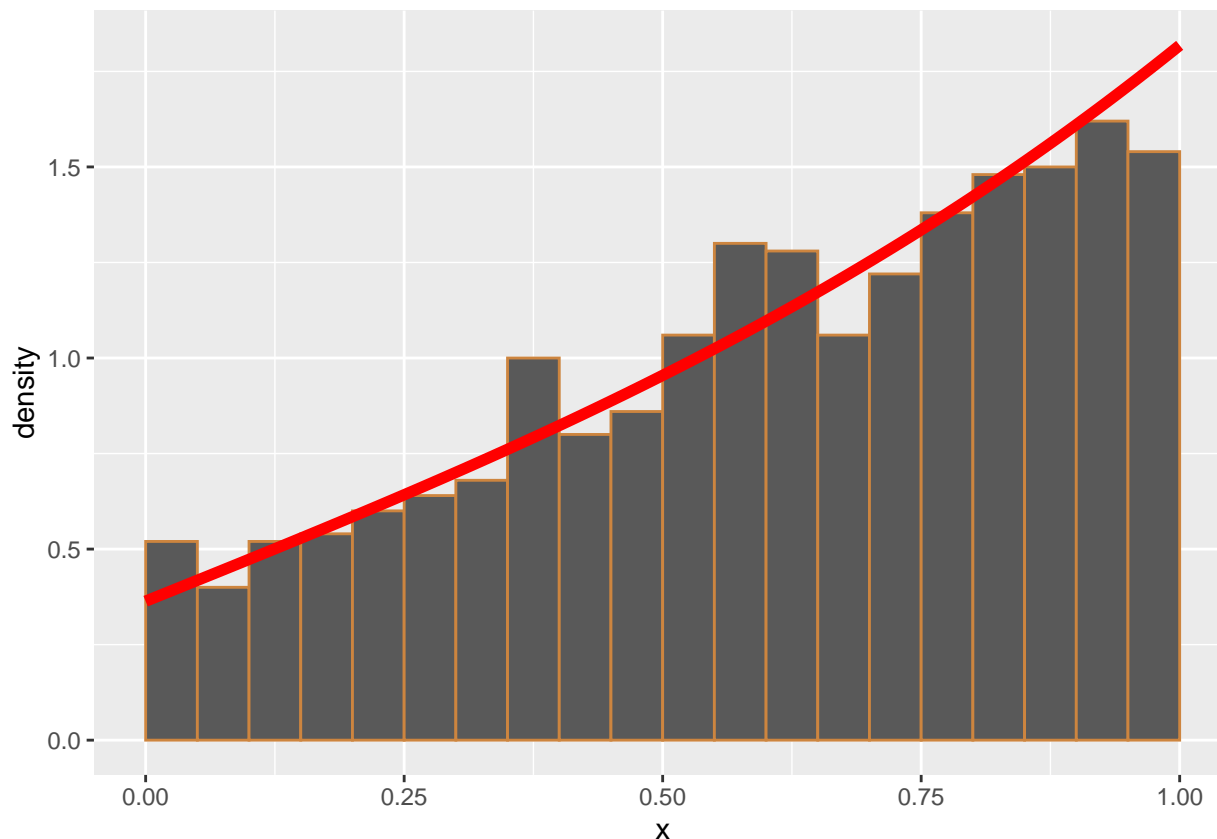


## Question 5

(10 points)

Code up a rejection sampler that allows you to also sample 1000 data from this pdf given in the last question. Again, histogram your sample and overlay the pdf. Save your sample as `sample.rs`.

```
set.seed(555)
sample.rs = rep(NA, 1000)
i = 1
y.max = 4*(1^3+3*1+1)/11
while(i <= 1000) {
  sample.rs[i] = runif(1, min = 0, max = 1)
  if (runif(1, min = 0, max = y.max) < 4*(sample.rs[i]^3 + 3*sample.rs[i]+1)/11) {
    i = i + 1
  }
}
ggplot(data=data.frame(x = as.numeric(sample.rs)), mapping=aes(x=x)) + geom_histogram(col="peru", aes(y
```



## Question 6

(10 points)

Test the hypothesis that `sample.it` and `sample.rs` are both sampled from the same parent population. (I mean they are, but...) Either recall how you would do a two-sample test or Google how you would do it. (Note: I'm not talking about a two-sample t-test here! We are not testing the hypothesis that the distribution means are the same. We are testing the hypothesis that both samples are drawn from the same underlying population.) There are various options for doing this; pick one, and display the p-value. If it less than 0.05, we reject the null. (Hint: it shouldn't be.) In addition, plot the empirical cdfs for both samples; see the documentation for `ecdf()` for help. (To be clear: use the base R function `plot()` here, and not `ggplot()`.) Note that to plot a second ecdf on top of the first, you need to call `plot()` a second time, with the argument `add=TRUE`.

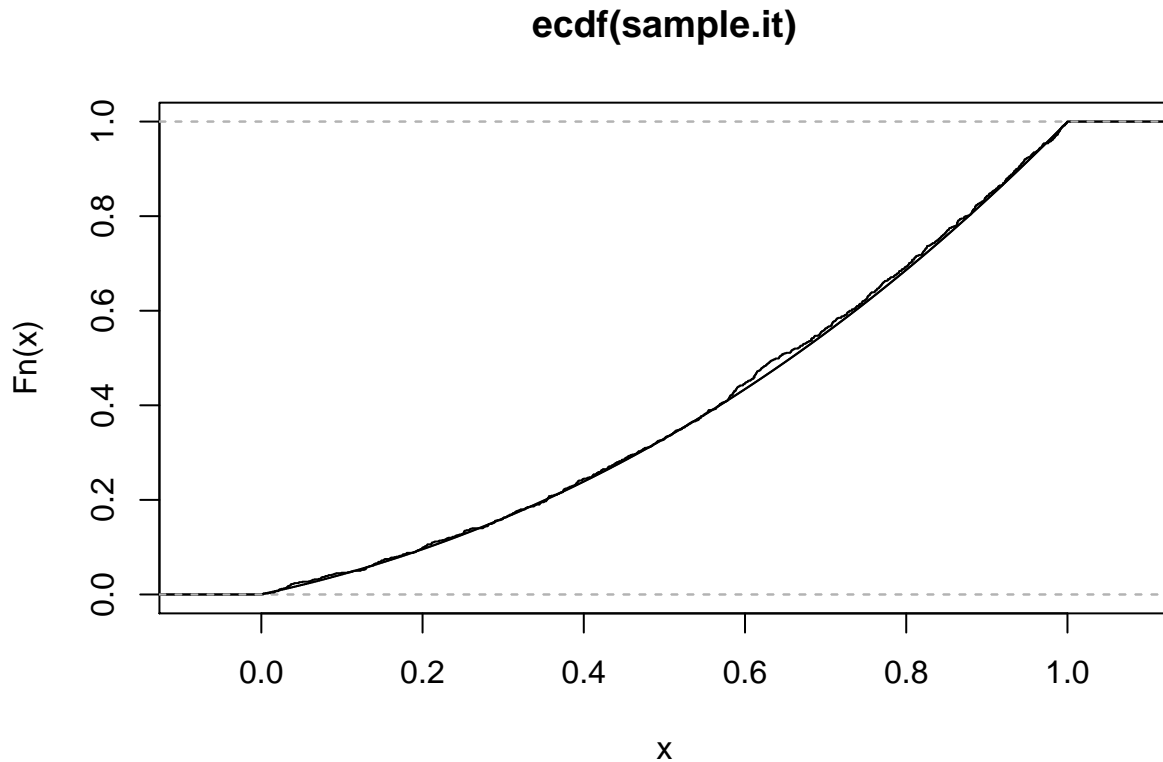
```
set.seed(666)
ks.test(sample.it, sample.rs)$p.value
```

```
## [1] 0.9456963
```

```
plot(ecdf(sample.it))
```

```
## Warning in xy.coords(x, y, setLab = FALSE): imaginary parts discarded in
## coercion
```

```
plot(ecdf(sample.rs), add = TRUE)
```



## Question 7

(10 points)

Write an inverse transform sampler that samples 10,000 data from a exponential distribution with rate parameter 1, but it only keeps data sampled over the ranges  $[0.5, 1]$  and  $[2, 4]$ . Make a probability histogram of your result. This time, tweak the call to `geom_histogram()` by adding the argument `breaks=seq(0,5,by=0.1)`.

This is a bit tricky. (Note: this is an inverse transform sampler, so every single randomly sampled uniform random variable has to get mapped to a valid value of  $x$ .) You might want to start by computing the probabilities  $P[0.5 \leq X \leq 1]$  and  $P[2 \leq X \leq 4]$ . Call these two quantities  $u_{lo}$  and  $u_{hi}$ , and sample random numbers from a  $\text{Uniform}(0, u_{lo} + u_{hi})$  distribution. If the number is  $< u_{lo}$ , the sampled number should be mapped to a sample from the lower range, whereas if the number is  $> u_{lo}$ , it should be mapped to a sample from the upper range. Note: to map from your uniform random variables to exponentially distributed ones, pass your uniform r.v.'s into `qexp()`.

```
set.seed(777)
u_lo = 1-0.5
u_hi = 4-2
sample = rep(NA, 10000)
y <- runif(10000, 0, u_lo + u_hi)
for (i in 1:1000) {
  if (y[i] < u_lo) {
```



```

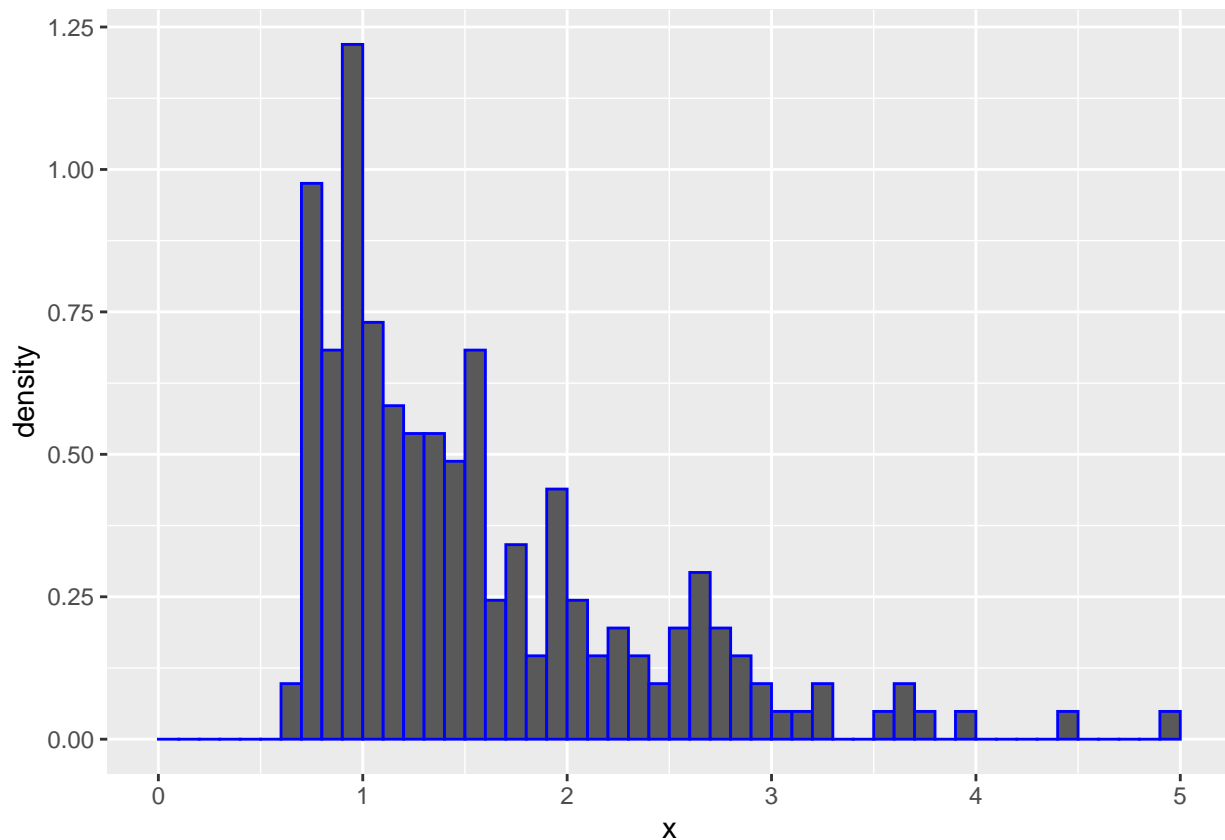
    sample[i] <- 0.5 + y[i]
  }else{
    sample[i] <- 2 + y[i] - u_lo
  }
}

ggplot(data = data.frame(x = qexp(sample)), mapping=aes(x=x)) + geom_histogram(col="blue", aes(y = ..den

```

```
## Warning in qexp(sample): NaNs produced
```

```
## Warning: Removed 9792 rows containing non-finite values (stat_bin).
```



## Question 8

(10 points)

And now for something completely different: randomly sampling a name for your new baby, given 1930 Social Security Administration data. (You did know you were having a baby, right?)

You'll see the `babynames` tibble has five columns. Using `dplyr` techniques, extract the rows for which the year is 1930, and then create a new data frame with just the columns `name` and `prop`. `prop` is a vector which gives the proportion of children given a particular name, but for some strange reason it does not sum to 1. Anyway, use a normalized version of this vector to appropriately sample one name from the `name` column. Show that name. Done. (If you set the random number seed to 111, you should get the name "Bernard".)

```
set.seed(111)
if ( require(babynames) == FALSE ) {
  install.packages("babynames",repos="https://cloud.r-project.org")
  library(babynames)
}
```

```
## Loading required package: babynames
```

```
b = filter(babynames, year == 1930) %>% select(name, prop)
b
```

```
## # A tibble: 9,789 x 2
##   name      prop
##   <chr>    <dbl>
## 1 Mary     0.0550
## 2 Betty    0.0328
## 3 Dorothy  0.0261
## 4 Helen    0.0171
## 5 Margaret 0.0157
## 6 Barbara  0.0157
## 7 Patricia 0.0135
## 8 Joan     0.0133
## 9 Doris    0.0133
## 10 Ruth    0.0128
## # ... with 9,779 more rows
```

```
sample(pull(b, name), 1, prob = pull(b, prop)/sum(pull(b, prop)))
```

```
## [1] "Bernard"
```

## Question 9

(10 points)

Numerically estimate the median of the pdf

$$f(x) = \frac{2.92959}{\sqrt{2\pi}} e^{-x^2/2} \quad x \in [0, 1].$$

(This is a truncated normal distribution.) The median is the value  $y$  such that

$$\int_0^y f(x) dx = 0.5.$$

A not-elegant way to do this is to use `integrate()` over and over again until you hone in on an integral value of 0.5. Don't do this. A more elegant solution is to determine, via `uniroot()`, the root of the function

$$g(y) = \left( \int_0^y f(x) dx \right) - 0.5,$$

i.e., the value of  $y$  such that  $g(y) = 0$ . Do do this.

```
f = function(y) {(2.92959/sqrt(2*pi))*exp((-y^2)/2)}
g = function(y) {integrate(f, 0, y)$value - 0.5}
uniroot(f = g, interval = c(0, 1))$root
```

```
## [1] 0.4417502
```

## Question 10

(10 points)

The ratio of the area of a circle to the area of a square into which the circle is inscribed is  $\pi/4$ . Does this ratio increase or decrease with dimensionality? For instance, what is the ratio of volume of a sphere to the volume of a cube into which the sphere is inscribed? Is it less than  $\pi/4$ ? Compute (and display) the ratio for dimensions 3, 4, ..., 10. The result that you see has manifestations for, e.g., algorithms based on nearest neighbors, etc. Curse of dimensionality, n'tat. (Hint: to do this calculation succinctly, consider putting samples from your uniform distribution into a  $k \times d$  matrix, where  $k$  is the number of sampled points, and  $d$  is the dimensionality. Then you can use `apply()` to determine the distance of the points from the origin and you can easily finish the calculation from there...)

```
set.seed(101)
mat = matrix(nrow = 100, ncol = 8)
i = 1
for (dim in 3:8) {
  mat[,i] = runif(100)
  i = i + 1
  a = apply(mat, MARGIN = 2, FUN = function(x){which(sqrt(dim*(x^2)) <= 1)})
  print(length(a)/100)
}
```

```
## [1] 0.08
## [1] 0.08
## [1] 0.08
## [1] 0.08
## [1] 0.08
## [1] 0.08
```

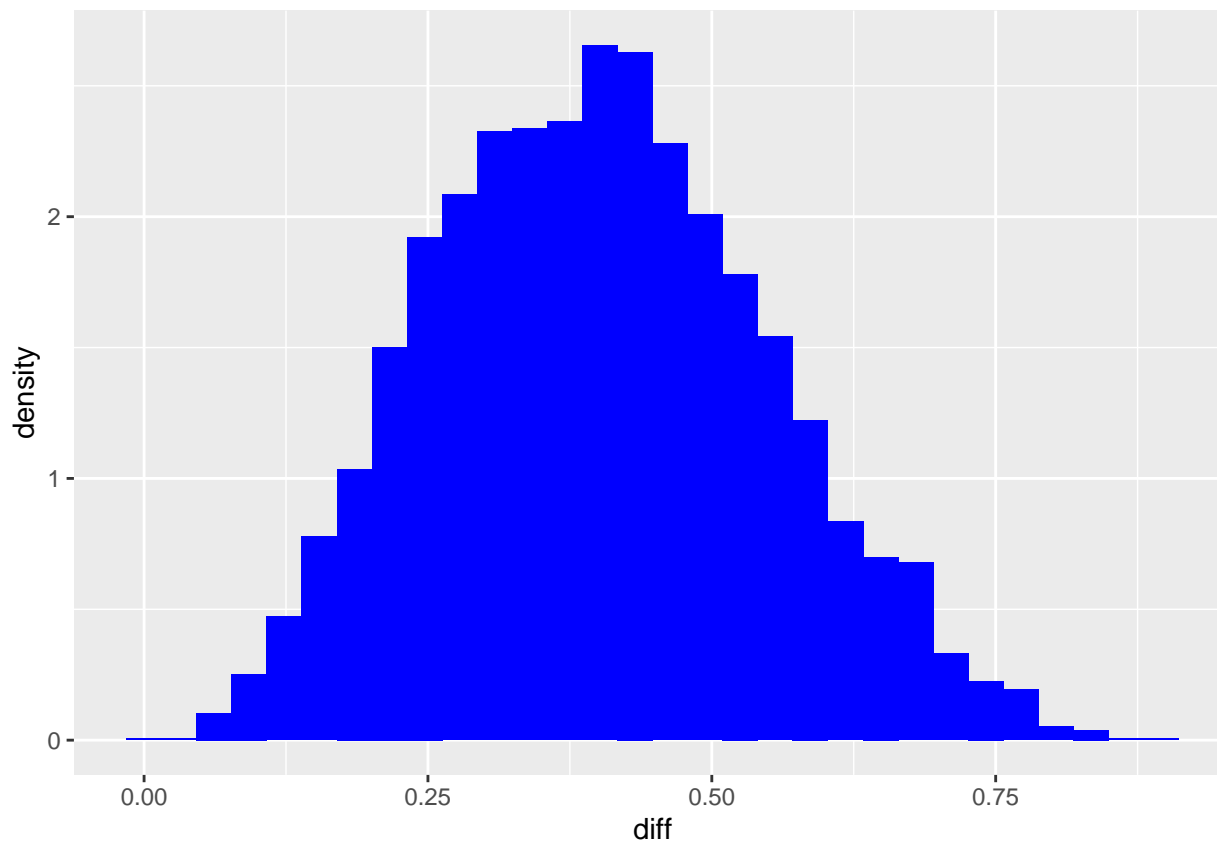
These are slightly more than  $\pi/4$ .

## Question 11

(10 points)

What is the probability distribution function for the difference between the maximum value and the median value when you sample nine data from a `Uniform(0,1)` distribution? Generate an empirical distribution by repeating the process of sampling nine data 5,000 separate times, and record the differences from the maximum and median values. Display a histogram of your result; in addition, display the mean and standard error. The mean should be approximately 0.4.

```
set.seed(101)
data = matrix(runif(9*5000), ncol=9)
df = data.frame("diff"=apply(data,1,max)-apply(data,1,median))
ggplot(data=df, mapping=aes(x=diff)) + geom_histogram(fill="blue", aes(y=..density..), bins = 30)
```



```
mean(df$diff)
```

```
## [1] 0.4015588
```

```
sd(df$diff)/sqrt(5000)
```

```
## [1] 0.002078052
```

## Question 12

(10 points)

You are given the following distribution:

$$f(x) = e^{-x\text{erf}(x)}/1.140741 \quad x > 0.$$

“erf(x)” == the error function with input  $x$ . (You’ll need to install and load the **VGAM** package to be able to compute the error function.) Here’s a plot of  $f(x)$ :

```
if ( require(VGAM) == FALSE ) {
  install.packages("VGAM",repos="https://cloud.r-project.org")
  library(VGAM)
}
```

```
## Loading required package: VGAM
```

```
## Loading required package: stats4
```

```
## Loading required package: splines
```

```
##
```

```
## Attaching package: 'VGAM'
```

```
## The following object is masked from 'package:tidyr':
```

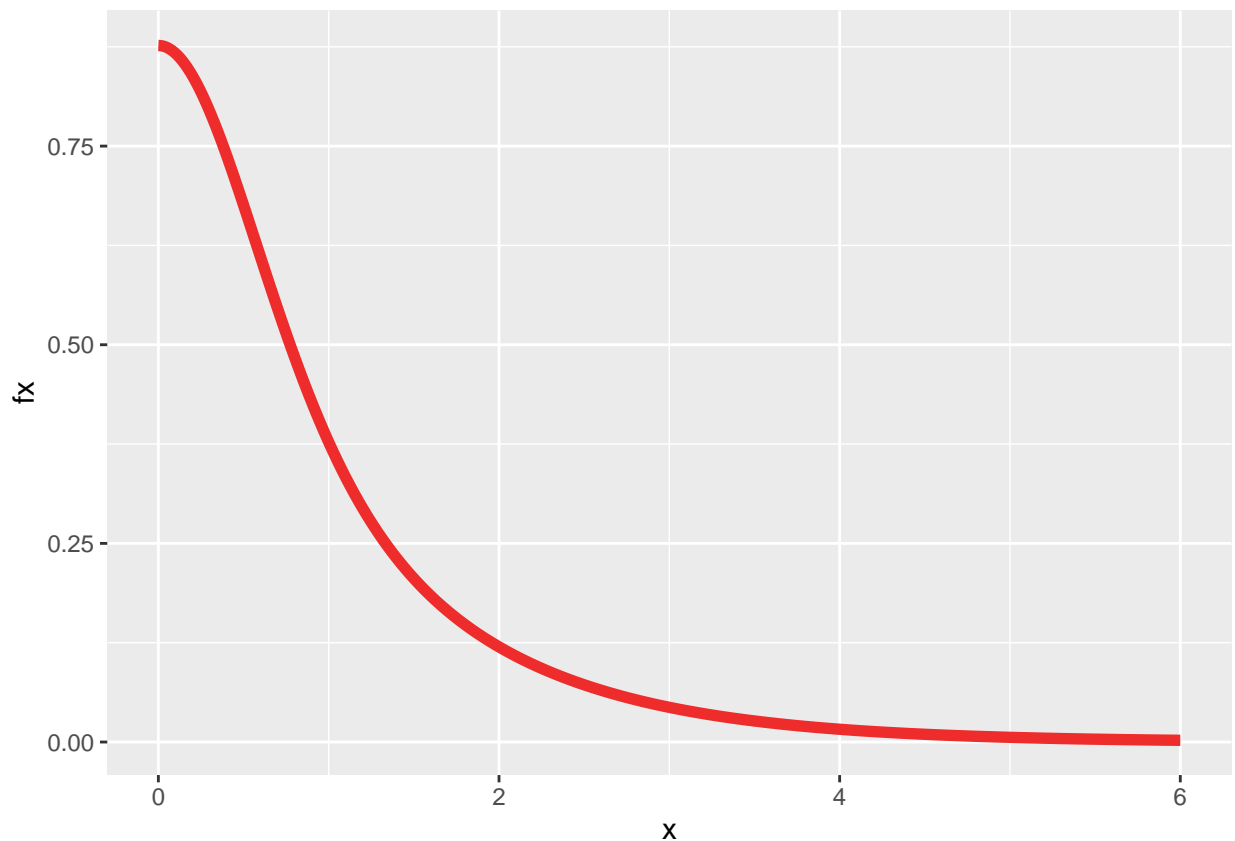
```
##
```

```
## fill
```

```
x = seq(0,6,by=0.01)
```

```
fx = exp(-x*erf(x))/1.140741
```

```
ggplot(data=data.frame(x=x,fx=fx),mapping=aes(x=x,y=fx)) + geom_line(col="firebrick2",size=2)
```



Use importance sampling to estimate the mean of  $f(x)$ . Use 100,000 data points. Your result should be approximately 0.95. (Hint: a half-normal distribution with **sigma** about 2.5 makes a nice proposal distribution here. See the **extraDistr** package.)

```
if ( require(extraDistr) == FALSE ) {  
  install.packages("extraDistr",repos="https://cloud.r-project.org")  
  library(extraDistr)  
}
```

```
## Loading required package: extraDistr
```

```
##
## Attaching package: 'extraDistr'

## The following objects are masked from 'package:VGAM':
##
##      dfrechet, dgev, dgompertz, dgpdp, dgumbel, dhuber, dkumar, dlaplace,
##      dlomax, dpareto, drayleigh, dskellam, dslash, pfrechet, pgev,
##      pgompertz, pgpd, pgumbel, phuber, pkumar, plaplace, plomax,
##      ppareto, prayleigh, pslash, qfrechet, qgev, qgompertz, qgpdp,
##      qgumbel, qhuber, qkumar, qlaplace, qlomax, qpareto, qrayleigh,
##      rfrechet, rgev, rgompertz, rgpd, rgumbel, rhuber, rkumar, rlaplace,
##      rlomax, rpareto, rrayleigh, rskellam, rslash

## The following object is masked from 'package:purrr':
##
##      rdunif

set.seed(333)
k = 100000
x = rhnorm(k, sigma=2.5)
h = dhnorm(x, sigma=2.5)
g = x
f = exp(-x*erf(x))/1.140741
mean(f*g/h)

## [1] 0.9503821
```