

# Lab: Week 6

36-350 – Statistical Computing

Week 6 – Fall 2020

Name: Kyle Wagner

Andrew ID: kowagner

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

---

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

---

## Simulation

### Question 1

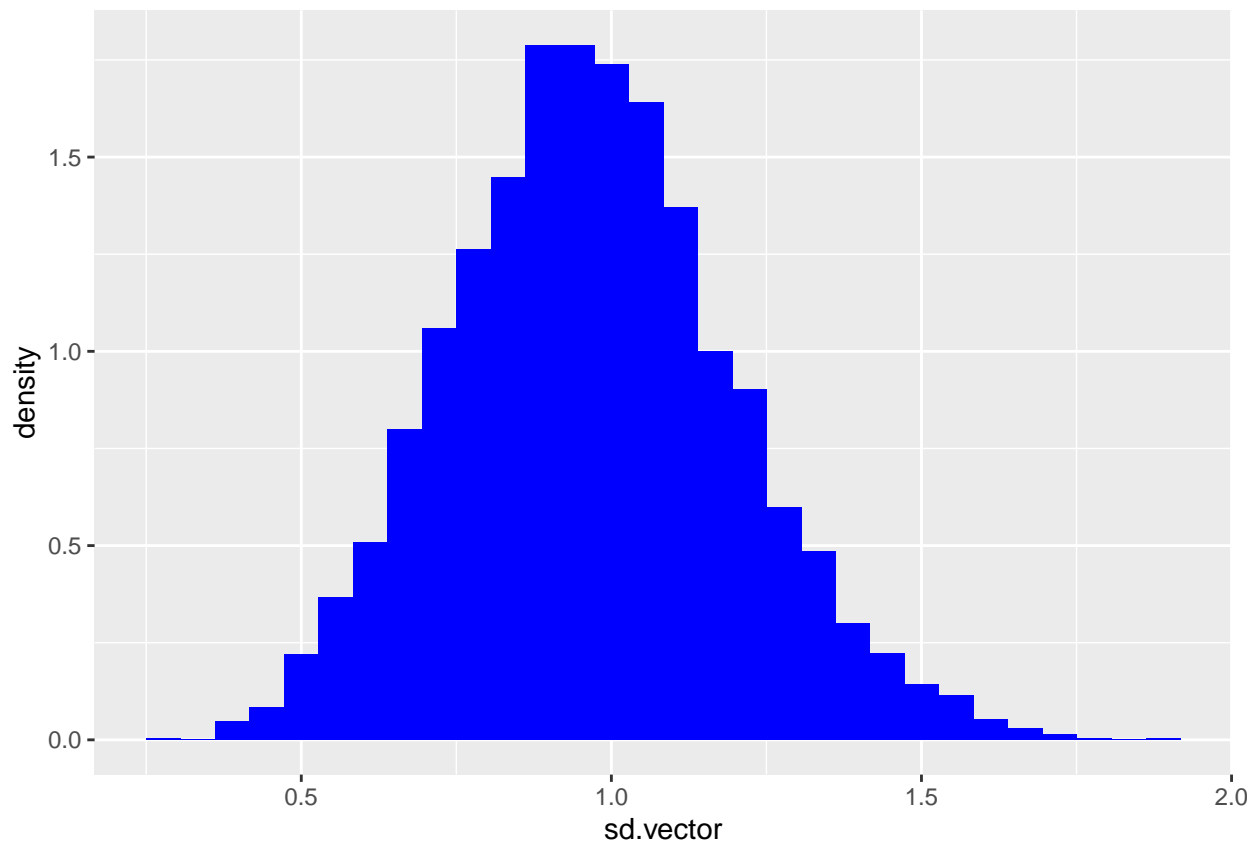
(5 points)

Notes 6A (3-5) and Notes 6C (5)

What is the distribution of the sample standard deviation for 10 draws from a standard normal? Create an empirical distribution by repeating the process of sampling 5000 times. Make a density histogram of your result using `ggplot`. (Pass the argument `aes(y=..density..)` to `geom_histogram()` to change the default frequency histogram into a density histogram). Test the hypothesis that the empirical distribution is itself normal. (Display the p-value.) Google to find an appropriate test of normality. Be sure to set random number seeds here and below for reproducibility. (Also note that `ggplot()` expects a data frame as input, not a vector; you should simply create a one-column data frame from your vector of sample standard deviations and pass that into `ggplot()`.)

```
set.seed(runif(1,0,10000))
n.obs = 5000
n.data = 10
data = matrix(rnorm(n.obs*n.data),ncol = n.data)
sd.vector = apply(data,1,sd)
sd.df = as.data.frame(sd.vector)
ggplot(sd.df,aes(x = sd.vector)) + geom_histogram(aes(y=..density..),fill = "blue")
```

```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```



```
shapiro.test(sd.df$sd.vector)
```

```
##
##  Shapiro-Wilk normality test
##
## data:  sd.df$sd.vector
## W = 0.99616, p-value = 3.937e-10
```

A Shapiro-Wilk normality test gives a small p-value that indicates that sample standard deviation is not normally distributed.

## Question 2

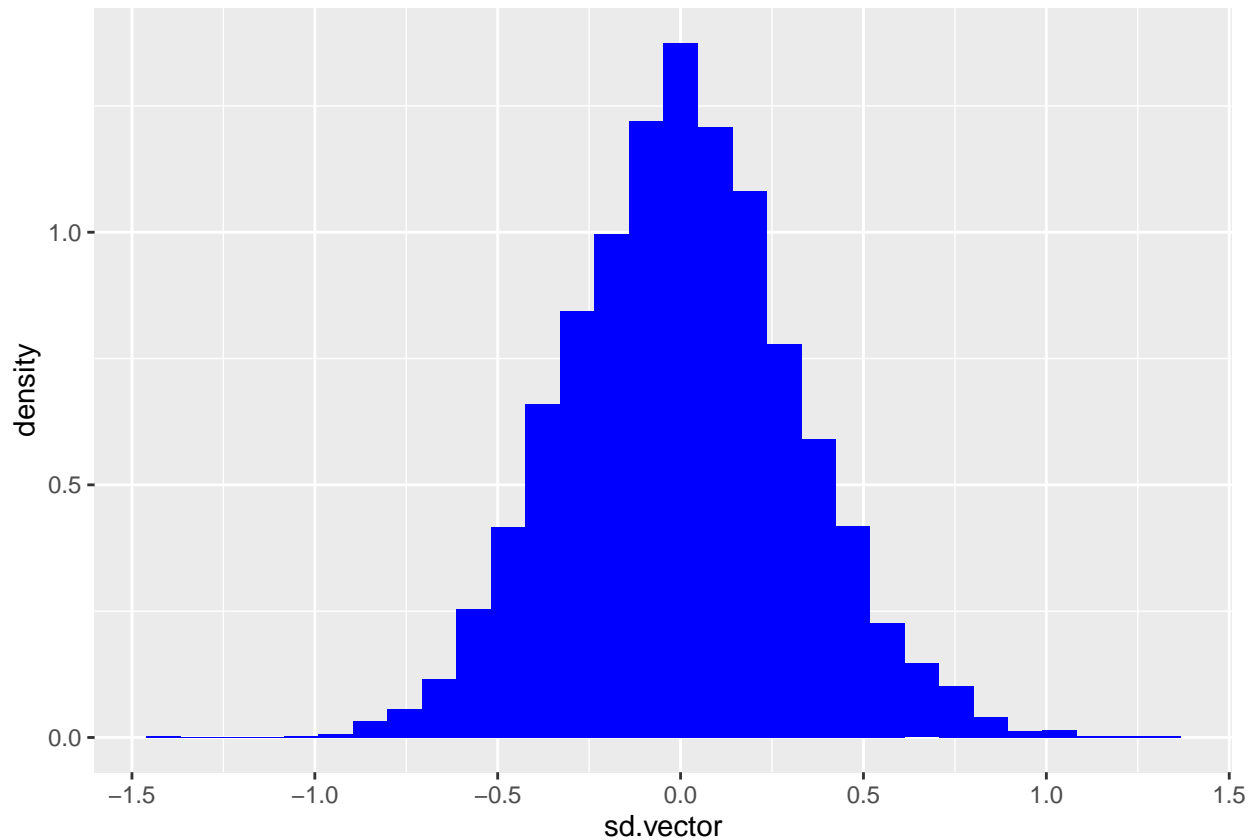
(5 points)

Notes 6A (3-5) and Notes 6C (5)

Repeat the last question, but replace the sample standard deviation with the sample mean. Given the p-value, would you say that the distribution of sample means is closer to being normal than the distribution of sample standard deviations?

```
set.seed(runif(1,0,10000))
n.obs = 5000
n.data = 10
data = matrix(rnorm(n.obs*n.data),ncol = n.data)
sd.vector = apply(data,1,mean)
sd.df = as.data.frame(sd.vector)
ggplot(sd.df,aes(x = sd.vector)) + geom_histogram(aes(y=..density..),fill = "blue")
```

```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```



```
shapiro.test(sd.df$sd.vector)
```

```
##  
##  Shapiro-Wilk normality test  
##  
## data:  sd.df$sd.vector  
## W = 0.99881, p-value = 0.001069
```

Given the p-value I would say that the distribution of means is closer to being normally distributed.

### Question 3

*(5 points)*

*Notes 6A (3-5) and Notes 6C (5)*

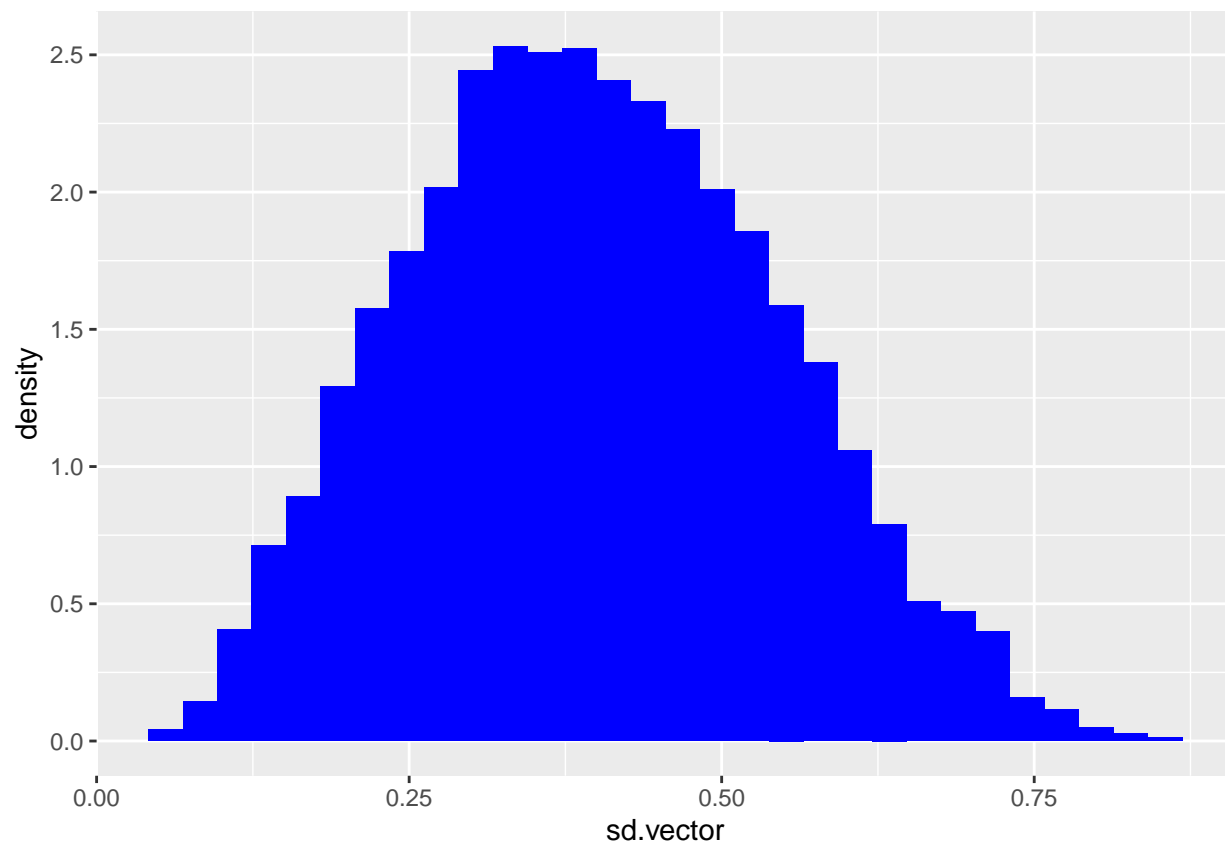
What is the probability distribution function for the difference between the maximum value and the median value when you sample nine data from a Uniform(0,1) distribution? Generate an empirical distribution by repeating the process of sampling nine data 5,000 separate times, and record the differences 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(runif(1,0,10000))
n.obs = 5000
n.data = 9
data = matrix(runif(n.obs*n.data),ncol = n.data)
my.fun = function(x){
  return(max(x) - median(x))
}
sd.vector = apply(data,1,my.fun)
sd.df = as.data.frame(sd.vector)
ggplot(sd.df,aes(x = sd.vector)) + geom_histogram(aes(y=..density..),fill = "blue")

## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.

```



```
mean(sd.df$sd.vector)
```

```
## [1] 0.3981476
```

```
sd(sd.df$sd.vector)
```

```
## [1] 0.1451663
```

## Question 4

(5 points)

*Notes 6A (2-5)*

Estimate the mean and standard error for the sum of three rolled fair dice using 5,000 simulated rolls. Note that there are potentially several ways by which you might attack this problem; how you go about it is up to you, as long as you get a valid final answer. (Which should be close to 10.5.)

```
n.obs = 5000
n.data = 3
data = matrix(sample(1:6,size = n.obs*n.data,replace = TRUE),ncol = n.data)
sums = apply(data,1,sum)
mean(sums)
```

```
## [1] 10.5148
```

```
sd(sums)
```

```
## [1] 3.002463
```

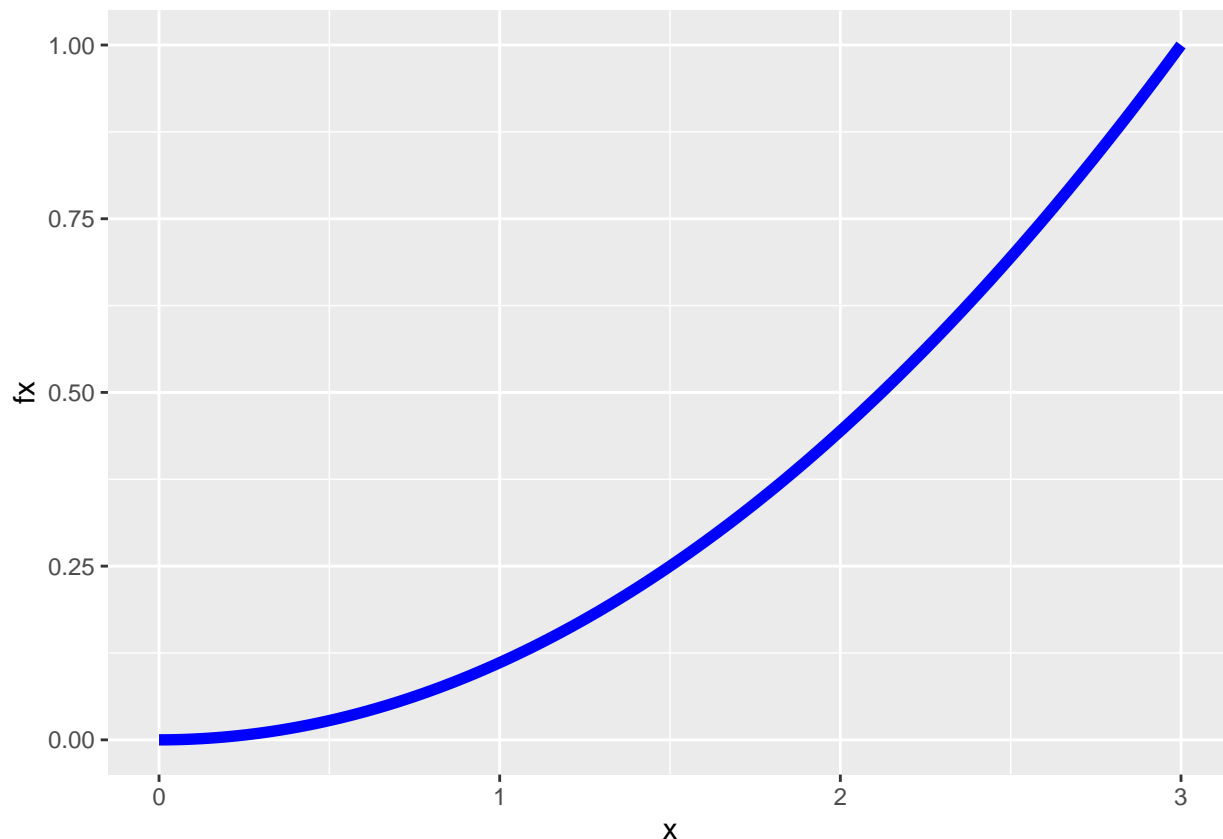
---

You are given the following probability density function:

$$f(x) = \frac{x^2}{9} \quad , \quad x \in [0, 3]$$

It looks like this:

```
x = seq(0,3,by=0.01)
fx = x^2/9
df.pdf = data.frame("x"=x, "fx"=fx)
ggplot(data=df.pdf,mapping=aes(x=x,y=fx)) + geom_line(col="blue",size=2)
```



## Question 5

(10 points)

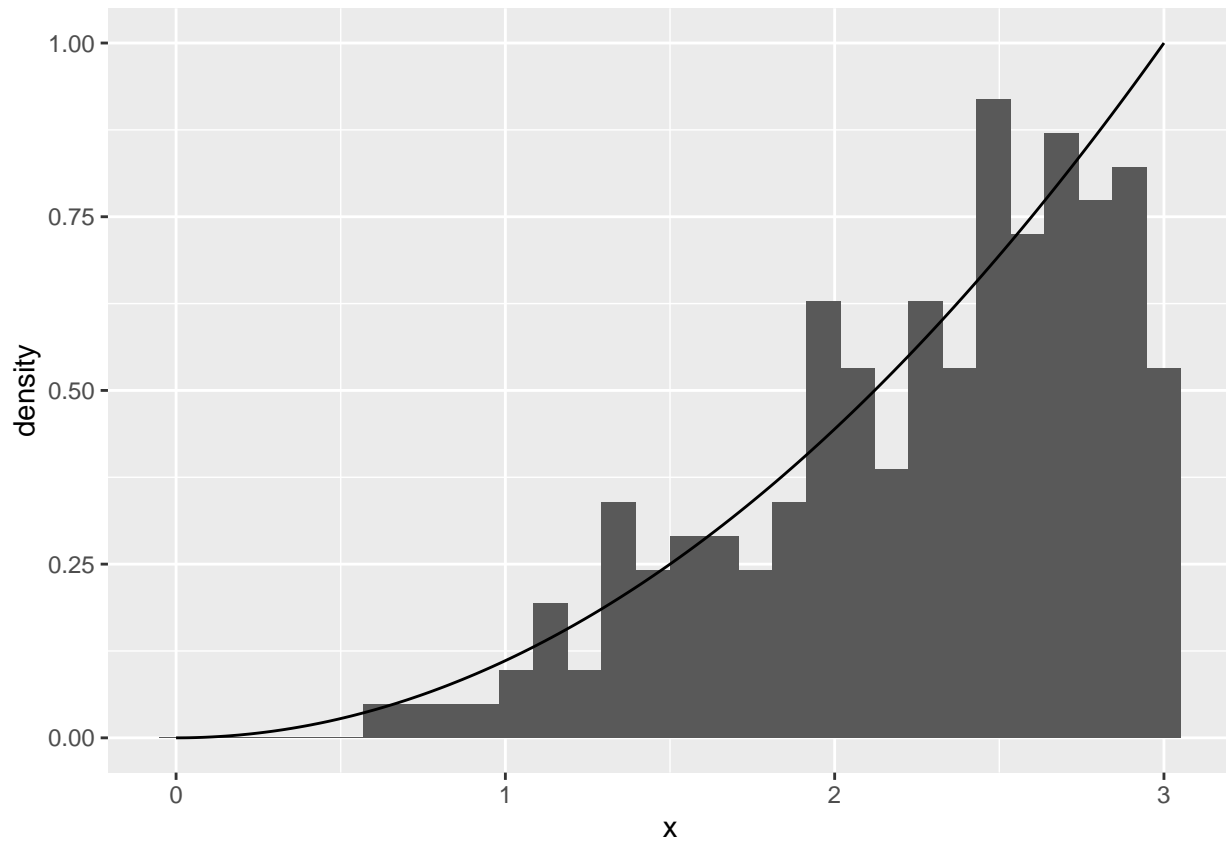
Notes 6A (6-7) and Notes 6C (5)

Code up an inverse transform sampler that allows you to efficiently sample 200 data from this pdf. Histogram your sample, and overlay the line showing the pdf. To do this in `ggplot`, you would use a structure like `ggplot(...) + geom_histogram(...) + geom_line(...)`. The main issue is that the data are histogramming are not the data you would be passing to the line. So: do the `ggplot(...) + geom_histogram(...)` in the same manner you did in previous questions above, where the data frame you point to is the one containing your sampled points, then, when you add on `geom_line(...)`, specify a data argument that points to the `df.pdf` variable defined above and specify a mapping argument that refers to columns of `df.pdf` (and add on a color argument and a size argument like we did above, if you wish).

```
u = runif(n=200)
x = (27*u)^(1/3)
x.df = as.data.frame(x)

ggplot(x.df, aes(x = x)) +
  geom_histogram(aes(y = ..density..)) +
  geom_line(data = df.pdf, aes(x = x, y = fx))
```

## 'stat\_bin()' using 'bins = 30'. Pick better value with 'binwidth'.



## Question 6

(10 points)

Notes 6A (8-9) and Notes 6C (5)

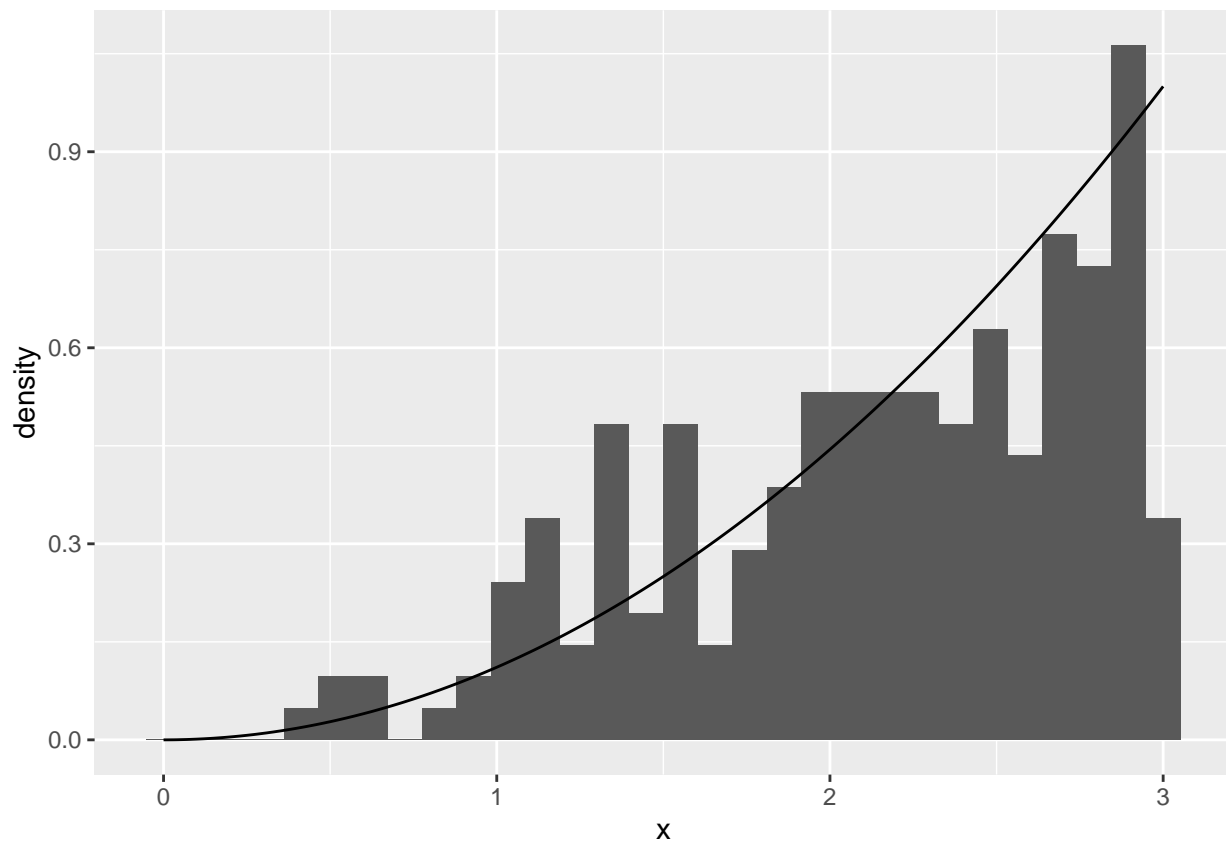
Repeat the previous question, but utilize rejection sampling.

```
k = 200
x = rep(NA,k)
ii = 1

while ( ii <= k) {
  x[ii] = runif(1,min=0,max=3)
  if (runif(1,min=0,max=1) < (x[ii]^2)/9){
    ii = ii + 1
  }
}

x.df1 = as.data.frame(x)
ggplot(x.df1,aes(x = x)) +
  geom_histogram(aes(y = ..density..)) +
  geom_line(data = df.pdf, aes(x = x, y = fx))
```

## 'stat\_bin()' using 'bins = 30'. Pick better value with 'binwidth'.



## Question 7

(5 points)

Notes 6A (3-5)

In 36-225, you learned that an effective rule of thumb regarding the Central Limit Theorem is that if  $n \geq 30$ , the mean of your sample is at least approximately normally distributed. Let's test this out. Construct repeated samples of 30 data from an  $\text{Exponential}(1)$  distribution and record the means: are they normally distributed? Simply show the p-value from an appropriate hypothesis test. If it is  $\ll 0.05$ , we'll conclude the rule of thumb doesn't really hold in this instance.

```
n.obs = 5000
n.data = 30

data = matrix(rexp(n.obs*n.data),ncol = n.data)

res = apply(data,1,mean)
res.df = as.data.frame(res)
shapiro.test(res.df$res)
```

```
##
## Shapiro-Wilk normality test
##
## data:  res.df$res
## W = 0.99025, p-value < 2.2e-16
```



## Question 8

(5 points)

Notes 6A (3-5)

Repeat the last question, but for a Uniform(0,1) distribution.

```
n.obs = 5000
n.data = 30
data = matrix(rexp(n.obs*n.data),ncol = n.data)

res = apply(data,1,mean)
shapiro.test(res)
```

```
##
## Shapiro-Wilk normality test
##
## data:  res
## W = 0.99235, p-value = 8.968e-16
```

## Question 9

(5 points)

Notes 6A (2-5)

You have been called to testify in a trial on the alien planet Slybobia. Prosecutors allege that on this planet, whose intelligent population is 25% blue, 25% green, 25% orange, and 25% fuchsia, it would require bias to construct a panel of 10 Slybobians in which at least 4 are orange *and* at least 4 are fuchsia. But such a panel had been constructed. Your expert testimony is needed: what is the probability that such a panel would be constructed if sampling of the population was truly done randomly? Is it smaller than 0.05? If so, that would indicate that you should reject the null hypothesis that the panel composition is unbiased.

```
n.obs = 5000
n.data = 10
data = matrix(sample(1:4,size = n.obs*n.data,replace = TRUE),ncol = n.data)

my.fun = function(x){
  x = length(which(data[1,] == 1))
  y = length(which(data[1,] == 2))
  if (x >= 4 & y >= 4){
    return(TRUE)
  }
  return(FALSE)
}
fil = apply(data,1,my.fun)
length(fil[which(fil == TRUE)])
```

```
## [1] 0
```

Have a proportion of 0, which indicates zero, therefore would need to reject the null hypothesis.

## Question 10

(5 points)

Notes 6A (8-9)

An old classic: what is the value of  $\pi$ ? Use rejection sampling in two dimensions (and, oh, 10,000,000 samples) to estimate  $\pi$ . (Think of a unit circle inscribed within a  $2 \times 2$  box centered at the origin. The box will have an area of 4, and the circle will have an area of  $\pi$ ... thus the ratio of the total number of random samples inside the unit circle to the total number of samples will approach  $\pi/4$  as the number of samples approaches  $\infty$ .) Also display the percentage error of your estimate, which you can compute using R's built-in value of  $\pi$  (i.e., the R constant `pi`.)

```
k = 10000000
ii = 1
total = 0

while ( ii <= k) {
  x = runif(1,min=-1,max=1) #horizontal
  y = runif(1,min=-1,max=1) #vertical
  #two if statements
  if (x^2 + y^2 <= 1){
    total = total + 1
  }
  ii = ii + 1
}

pi.approx = 4* total / k
pi.approx
```

```
## [1] 3.141104
```

```
(pi.approx - pi)/pi * 100
```

```
## [1] -0.01556706
```

## Question 11

(10 points)

Notes 6A (6-9)

You are given the following bivariate pdf:

$$f(x_1, x_2) = \begin{cases} 3x_1 & 0 \leq x_2 \leq x_1 \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

(This is borrowed from Exercise 5.5 of Wackerly 7.) Code a function for sampling data from this distribution. Sample 5000  $(x_1, x_2)$  pairs, then visualize them with a scatter plot. You should observe a greater density of points as you go from the left to the right.

```
k = 5000
x1 = runif(5000)
x2 = rep(NA,k)
```

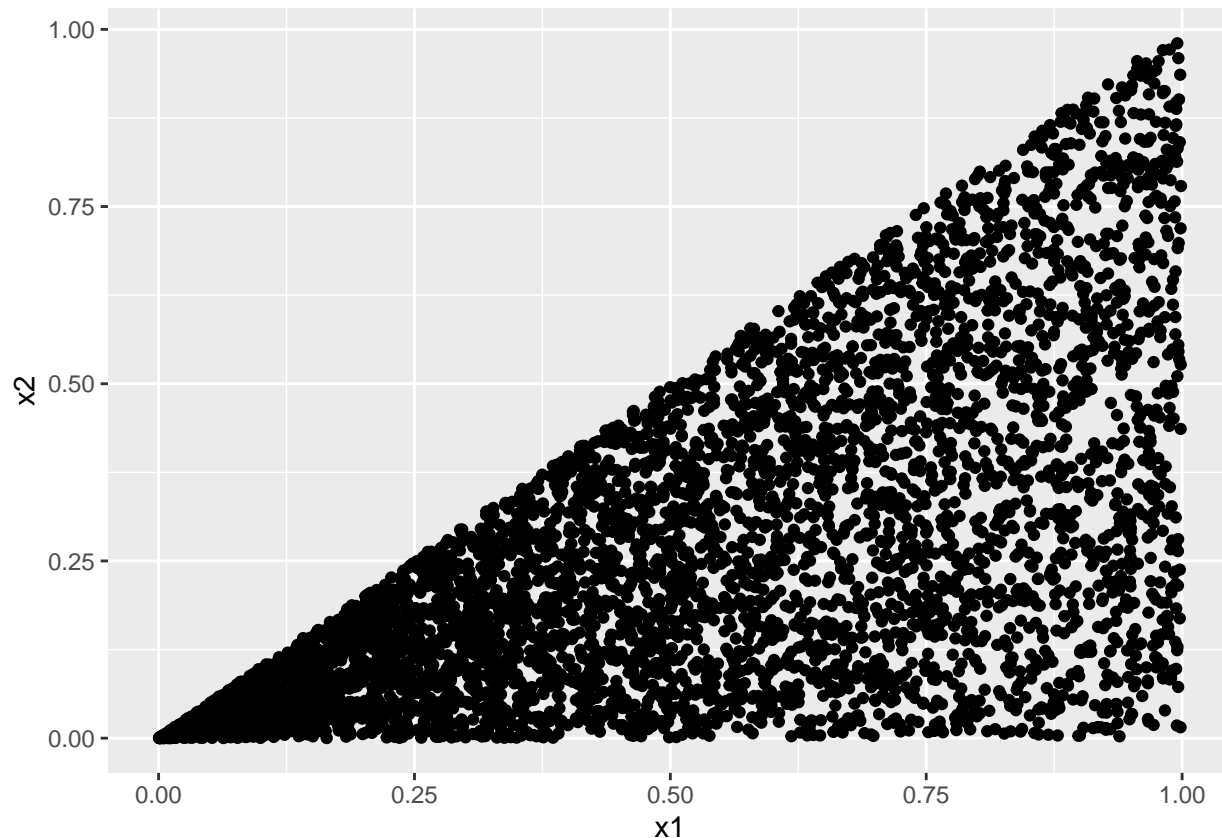
```

ii = 1

while ( ii <= k) {
  x2[ii] = runif(1,min=0,max=x1[ii])
  ii = ii + 1
}

x = data.frame(x1 = x1,x2 = x2)
ggplot(x,aes(x = x1, y= x2)) + geom_point()

```



## Integration

### Question 12

(5 points)

Notes 6B (3-4)

Compute the integral

$$\int_0^3 x^2 e^{-x^4} dx$$

via Monte Carlo integration, with 100,000 points. You should achieve a value close to 0.30635.

```
set.seed(123)
g = function(x) {
  x^2 * exp(-x^4)
}
x = runif(100000,min=0,max=3) # f(x) = 1/3
w = g(x)/(1/3)
print(mean(w))
```

```
## [1] 0.3077858
```

## Question 13

(5 points)

Notes 6B (3-4)

Compute the integral

$$\int_0^1 \int_0^1 \left( \cos\left(\frac{\pi}{2}x_1\right) + x_2^3 \right) dx_1 dx_2$$

via Monte Carlo integration, with 1,000,000 points. Your value should be close to 0.8866.

```
set.seed(123)
g = function(x,y) {
  cos(pi*x/2) + y^3
}
x = runif(1000000)
y = runif(1000000)
w = g(x,y)
print(mean(w))
```

```
## [1] 0.8871224
```

## Question 14

(5 points)

Notes 6A (8-9) and Notes 6C (3-4)

Let's change up that last integral just a bit:

$$\int_0^1 \int_0^{x_2} \left( \cos\left(\frac{\pi}{2}x_1\right) + x_2^3 \right) dx_1 dx_2$$

The region of integration is now the triangle with vertices (0,0), (0,1), and (1,1). (Commence hallucinatory flashbacks to 225, if they haven't started already.) Combine a rejection sampler and MC integration to perform this integral. Sample approximately 1,000,000 points in the region of integration and then use those. (By using the word "approximately," I'm encouraging you to be clever in how you do the sampling... like by sampling 2,000,000 points in a box and keeping the roughly 1,000,000 that lie within the triangular region of integration.) Realize that for a bivariate uniform within the region of integration,  $f(x_1, x_2) = 2$ . (Your answer should be approximately 0.605.)

```

set.seed(123)
g = function(x,y) {
  cos(pi*x/2) + y^3
}
x = runif(1000000)

#rejection sampler for y

k = 1000000
y = rep(NA,k)
ii = 1
while ( ii <= k) {
  y[ii] = runif(n=1)
  if (y[ii] <= x[ii]) {
    ii = ii + 1
  }
}

w = g(x,y)
print(mean(w))

```

```
## [1] 0.6994524
```

## Question 15

(10 points)

Notes 6C (5-8)

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
```

```
##
```

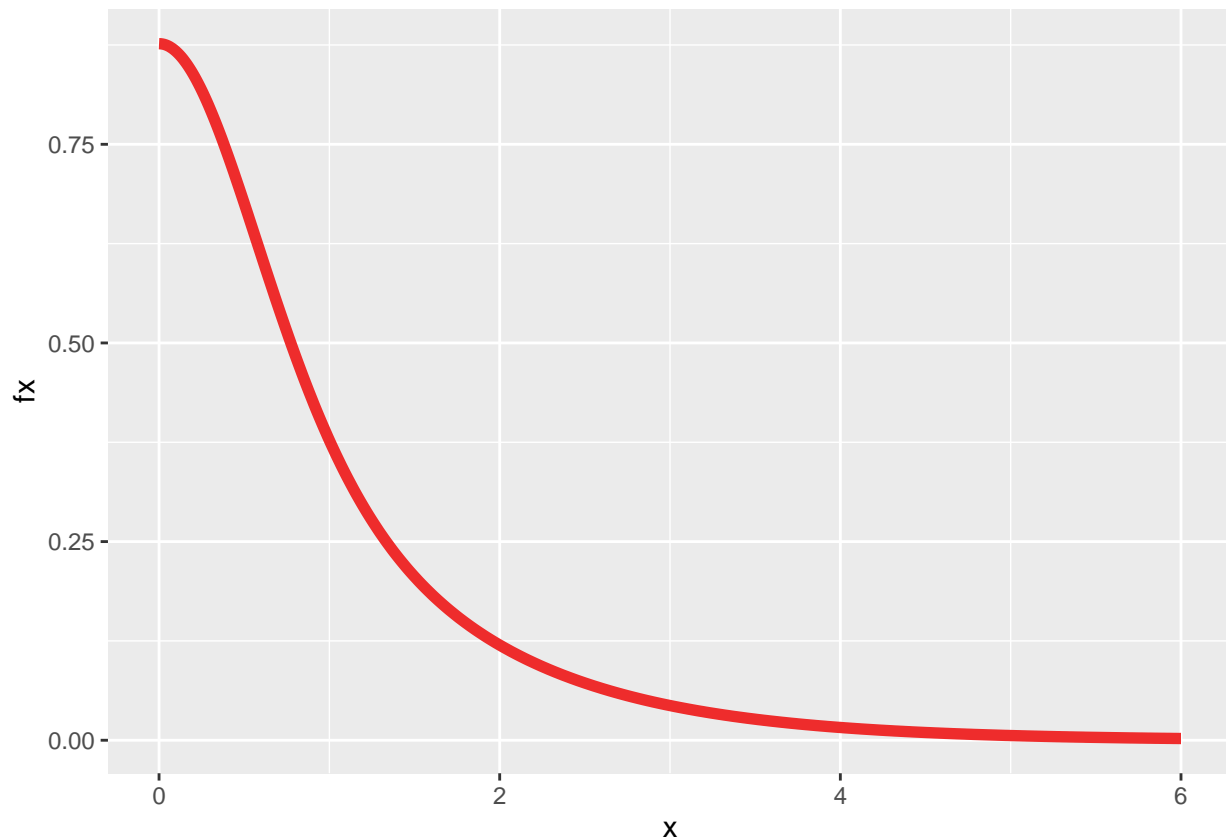
```
## Attaching package: 'VGAM'
```

```
## 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(660)  
N = 100000  
x = rhnorm(N,sigma = 2.5) # Sample x samples from the proposal distribution h(x)  
h = dhnorm(x,sigma = 2.5) # Evaluate h(x)  
g = rep(0,N)  
g = x      # Evaluate g(x)  
  
f = function(x) {  
  f.x = rep(0,length(x))  
  f.x[x>=0] = exp(-x*erf(x))/1.140741  
  return(f.x)  
}  
mean(g*f(x)/h)
```

```
## [1] 0.9485657
```

## Question 16

(5 points)

Notes 6C (7)

And now: estimate the standard deviation of  $f(x)$ . Note: you can reuse many of the variables from above! Your value should again be around 0.95.

```
sqrt(mean(g^2 * f(x)/h) - mean(g*f(x)/h)^2)
```

```
## [1] 0.9570094
```