

## STAT430 Homework #2: Due Friday, February 8, 2019.

Name: Oliver Shanklin

---

0. We are continuing with **Section 7.2** of the text in this homework: there is a lot to cover in that section.

To get started, let's look at how to approximate the probability of an event via simulation.

In **Example 7.2** (pages 354-355),  $Y_1, Y_2, \dots, Y_9$  are iid  $N(\mu, \sigma^2)$  with  $\sigma^2 = 1$ . The text shows how to compute

$$\begin{aligned} P(|\bar{Y} - \mu| \leq 0.3) &= P\left(\frac{-0.3}{\sqrt{1/9}} \leq Z \leq \frac{0.3}{\sqrt{1/9}}\right) \\ &= P(-0.9 \leq Z \leq 0.9) \\ &= P(Z \leq 0.9) - P(Z \leq -0.9), \end{aligned}$$

which evaluates to

```
pnorm(0.9) - pnorm(-0.9)
```

```
## [1] 0.6318797
```

Let's consider how to approximate this probability via simulation. In R, if we make logical statements like

```
x <- c(3 < 5, 6 < 5) # logical results "combined" into a vector with c()
x
```

```
## [1] TRUE FALSE
```

then R evaluates those statements as TRUE or FALSE, which are treated numerically as TRUE = 1 and FALSE = 0. So,

```
mean(x) # mean of logical vector is the proportion TRUE
```

```
## [1] 0.5
```

In this little example, 1/2 of the statements are true. Now we just want to simulate a large number of standard normals, and check the empirical proportion that are between  $-0.9$  and  $0.9$ :

```
Z <- rnorm(100000)
x <- (abs(Z) < 0.9) # big vector of logical results
mean(x)             # mean of logical vector is the proportion TRUE
```

```
## [1] 0.63066
```

This is quite close to the theoretical calculation. In general, we can approximate the probability of an event by (1) simulating a vector of logical outcomes that are TRUE when the event occurs; and (2) computing the mean of the logical vector. We can use this technique to check an analytical calculation, or to approximate the theoretical answer when the analytical calculation is hard or intractable. For example, suppose we wanted some weird probability like

$$P(\ln|Z| > 0.5).$$

We can get a quick, approximate answer via

```
Z <- rnorm(100000)
x <- (log(abs(Z)) > 0.5) # big vector of logical results
mean(x)                 # mean of logical vector is the proportion TRUE
```

```
## [1] 0.09712
```

---

1. Complete **Exercise 7.11** of the text, using the `pnorm` function in R to compute the exact, theoretical probability. Then check your answer using the probability approximation above, via simulation of 10,000  $N(0, 1)$  random variables. As usual, set your random number seed to 4302019.

---

**Answer:**

```
set.seed(4302019)
```

```
#find the prob that the sample mean will be within 2 sq-in of the population mean
```

```
pnorm(1.5) - pnorm(-1.5)
```

```
## [1] 0.8663856
```

```
Z <- rnorm(10000)
```

```
x <- (abs(Z) <= 1.5)
```

```
mean(x)
```

```
## [1] 0.8712
```

$$\begin{aligned}
 P(|\bar{Y} - \mu| \leq 2) &= P(-2 \leq \bar{Y} - \mu \leq 2) \\
 &= P\left(\frac{-2}{4/3} \leq \frac{\bar{Y} - \mu}{S_{\bar{Y}}} \leq \frac{2}{4/3}\right) \\
 &= P(-1.5 \leq Z \leq 1.5) \\
 &= \text{pnorm}(1.5) - \text{pnorm}(-1.5)
 \end{aligned}$$

So from the simulation, the probability approximation is .8712, and the exact probability using R is .86638, which is pretty close.

2. Complete **Exercise 7.12** of the text, finding the appropriate sample size,  $n$ . Then assess your results by simulation as follows. You know that since  $Y_1, \dots, Y_n$  iid  $N(\mu, \sigma^2)$ , the exact sampling distribution of the sample mean is  $\bar{Y} \sim N(\mu, \sigma^2/n)$ . So reset your random number seed to 4302019 and use the R function `rnorm` with arguments `mean` =  $\mu$  (you can choose any real number that doesn't break your computer) and `sd` =  $\sqrt{\sigma^2/n} = \sqrt{4^2/n}$  to generate 10000 simulated  $\bar{Y}$  values (since we know the exact sampling distribution, we don't have to start by simulating the raw data as in Homework 1: we can actually simulate from the  $\bar{Y}$  distribution directly). With your value of  $n$ , is  $\bar{Y}$  within 1 square inch of your  $\mu$  at least 90% of the time? Convince yourself that changing the value of  $\mu$  does not change your answer.

---

**Answer:**

$$\begin{aligned}
 P(|\bar{Y} - \mu| \leq 1) &= .90 \\
 P(-1 \leq \bar{Y} - \mu \leq 1) &= .90 \\
 P\left(\frac{-1}{4/\sqrt{n}} \leq Z \leq \frac{1}{4/\sqrt{n}}\right) &= .90
 \end{aligned}$$

Since I need to have 5% in the upper and lower tails, I need to get a  $qnorm(0.95)$ .

```
qnorm(.95)
```

```
## [1] 1.644854
```

Now,

$$\frac{\sqrt{n}}{4} = .16448$$
$$n = 43.2887$$

So, I would round up to  $n = 44$ , which means we need at least 44 samples to have the probability of 0.90 for the sample mean to be within 1 square inch of the population mean.

```
set.seed(4302019)
```

```
n = 44
```

```
sd <- sqrt(4^2/n)
```

```
Z <- rnorm(10000, 0, sd)
```

```
x <- (abs(Z) <= 1)
```

```
mean(x)
```

```
## [1] 0.9093
```

So, from the simulation, 90.93% of the samples are within 1 square inch of the mean.

- 
3. In R, it is often useful to write your own functions for computations that you do repeatedly. For example, we can do the general computation

$$P(|\bar{Y} - \mu| \leq \delta) = P\left(\frac{-\delta}{\sqrt{\sigma^2/n}} \leq Z \leq \frac{\delta}{\sqrt{\sigma^2/n}}\right)$$
$$= P\left(Z \leq \frac{\delta}{\sqrt{\sigma^2/n}}\right) - P\left(Z \leq \frac{-\delta}{\sqrt{\sigma^2/n}}\right)$$

by constructing the R function

```
my_prob <- function(delta, sigma, n){ # you can call your function anything you like
  prob <- pnorm(delta / sqrt(sigma ^ 2 / n)) - pnorm(-delta / sqrt(sigma ^ 2 / n))
  return(prob)
}
```

Then we can run our function on a problem like **Example 7.2**:

```
my_prob(0.3, 1, 9)
```

```
## [1] 0.6318797
```

Notice that we did not need to name the arguments, because they came in the order expected by the function. If for some reason we used a different order, we need to use the names:

```
my_prob(sigma = 1, delta = 0.3, n = 9)
```

```
## [1] 0.6318797
```

In R, we can often give a vector argument and get a vector response. For example, if we want to look at sample sizes  $n = 9, 10, 11, 12$ , we can use

```
my_prob(0.3, 1, c(9, 10, 11, 12))
```

```
## [1] 0.6318797 0.6572183 0.6802576 0.7013024
```

Use this new function to complete **Exercise 7.9 and 7.10** of the text. (For 7.9(d), give a better answer than “Yes”!)

---

**Answer:**

```
set.seed(4302019)
qprob <- function(delta, sigma, n){
  prob <- pnorm(delta / sqrt(sigma ^ 2 / n)) - pnorm(-delta / sqrt(sigma ^ 2 / n))
  return(prob)
}
```

Exercise 7.9

```
## a)
```

```
qprob(.3, 1, 9)
```

```
## [1] 0.6318797
```

```
## b)
```

```
qprob(.3,1,c(25,36,49,64))
```

```
## [1] 0.8663856 0.9281394 0.9642712 0.9836049
```

- c) As  $n$  gets larger, you get a higher probability of the sample mean being within 0.3 ounce of the true mean.
- d) So, from Example 7.3, we see that 42 was not enough samples to achieve 95%, but adding one more to the sample increased that probability to above 95%. Indicating, more samples will have a higher probability to being within  $\delta$  of the true population mean.

Exercise 7.10

```
## a)
```

```
## using qprob() from above
```

```
qprob(.3,2,9)
```

```
## [1] 0.3472896
```

```
## b)
```

```
qprob(.3,2,c(25,26,49,64))
```

```
## [1] 0.5467453 0.5556409 0.7062819 0.7698607
```

- c) Again, the values of the probability increase with a larger  $n$ .
  - d) When the standard deviation is higher, the lower the probability is for the mean of the samples to be within 0.3 of the true population mean.
-

4. In class, we denoted by  $z_{\alpha/2}$  the value such that for  $Z \sim N(0, 1)$ ,

$$P(Z > z_{\alpha/2}) = P(Z \leq -z_{\alpha/2}) = \frac{\alpha}{2}.$$

In R, you can compute  $\pm z_{\alpha/2}$  with `qnorm`; for example,

```
round(qnorm(c(0.025, 0.05, 0.95, 0.975)), 3)
```

```
## [1] -1.960 -1.645 1.645 1.960
```

are the values I told you to memorize for  $\alpha = 0.05, 0.10$ . Also in R, the “next largest integer” function is `ceiling`:

```
ceiling(c(3.9, 41.1, 2.7))
```

```
## [1] 4 42 3
```

Use these functions and follow the example above to write your own function that calculates the minimum (integer) value of  $n$  needed to guarantee that if  $Y_1, \dots, Y_n$  iid  $N(\mu, \sigma^2)$ , then  $P(|\bar{Y} - \mu| \leq \delta) \geq 1 - \alpha$ . Check your function by repeating **Example 7.3** of the book, repeating **Exercise 7.12** above, and completing **Exercise 7.14** using your function (we did this one by hand in class).

---

**Answer:**

```
set.seed(4302019)
```

```
nprob <- function(delta, sigma, p){
  n <- ceiling(((qnorm((1-p)/2)*sigma)/delta)^2)
  return(n)
}
```

Example 7.3

```
nprob(.3, 1, .95)
```

```
## [1] 43
```

Exercise 7.12

```
nprob(1, 4, .9)
```

```
## [1] 44
```

Exercise 7.14

```
nprob(.5, sqrt(.4), .95)
```

```
## [1] 7
```

7 tests should be ran in order to satisfy the question.

---

5. Complete **Exercise 7.15** of the text, showing the steps in (a) and (b). For (c), you can use the function you created above.

---

**Answer:**

Exercise 7.15

a)

$$E(\bar{X} - \bar{Y}) = E(\bar{X}) - E(\bar{Y}) = \mu_1 - \mu_2$$

b)

$$V(\bar{X} - \bar{Y}) = V(\bar{X}) + V(\bar{Y}) = \frac{\sigma_1^2}{m} + \frac{\sigma_2^2}{n}$$

c)

```
nprob(1, sqrt(4.5), .95)
```

```
## [1] 18
```

- 
6. In class, I mentioned the paper by Student (1908) and the simulation experiment that involved writing the “height and left middle finger measurements of 3000 criminals” on “3000 pieces of cardboard, which were then very thoroughly shuffled and drawn at random.” Let’s do this experiment without using cardboard! The data are available in R in a slightly strange form, so run the following bit of code to extract the 3000 heights into a vector:

```
require(stats)
tmp <- as.numeric(colnames(crimtab)) / 2.54
height_inches <- as.numeric(rep(tmp, colSums(crimtab)))
```

- (a) Draw a histogram of the criminals’ heights in inches and comment: do they look approximately normal? (b) Set your seed as usual and use the method of simulation from Homework #1 to draw 10,000 simulated samples of size  $n = 4$  with replacement from the 3000 criminals. (c) For each simulated sample, compute the sample variance,  $S^2$ , using the `var()` function in R. (d) Approximate the true variance,  $\sigma^2$ , by `var(height_inches)`. Is the mean of your 10,000 sample variances close to  $\sigma^2$ ? (e) Recall that, if  $Y_1, Y_2, Y_3, Y_4$  are iid  $N(\mu, \sigma^2)$ , then the sample variance satisfies

$$\frac{(4-1)}{\sigma^2} S^2 \sim \chi_3^2.$$

Rescale your 10,000 sample variances by multiplying by  $3/\sigma^2$ , plot the histogram of your rescaled sample variances (but set `breaks = 40` to get more bins in your histogram than the default number), and add the theoretical pdf of  $\chi_3^2$  using the `dchisq` function in R. (e). As an alternative to a histogram, use the *empirical cumulative distribution function*, implemented in R by `plot(ecdf(your_rescaled_variances))`. Add to your ecdf plot the theoretical cumulative distribution function using the `pchisq` function.

---

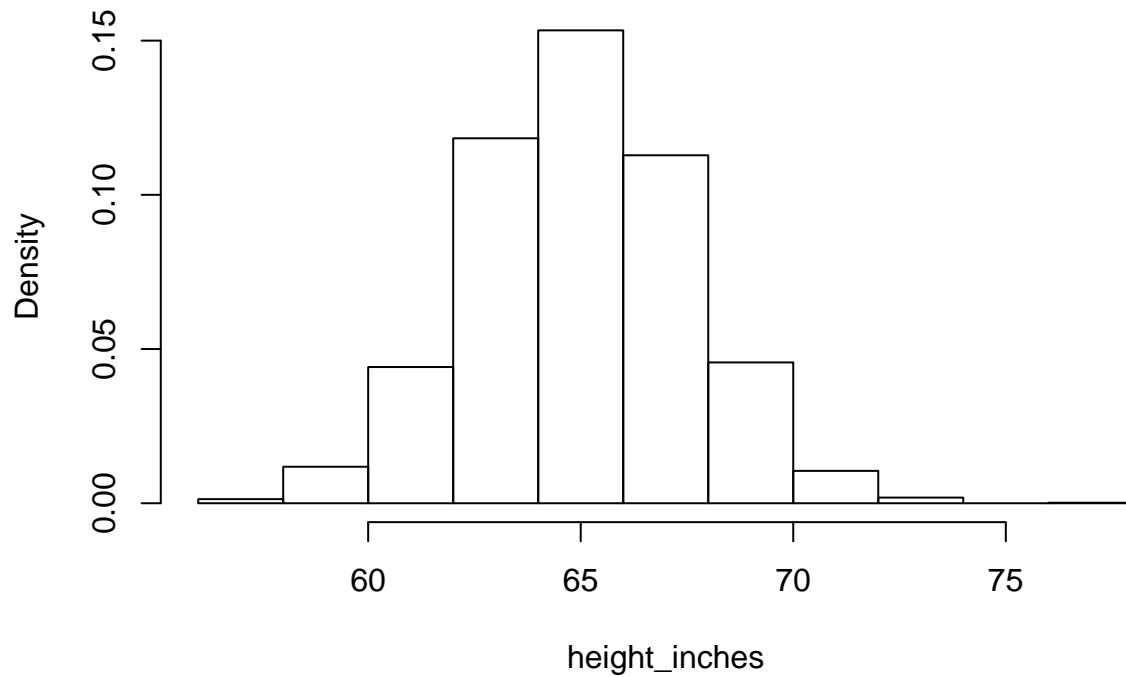
**Answer:**

```
set.seed(4302019)

require(stats)
tmp <- as.numeric(colnames(crimtab)) / 2.54
height_inches <- as.numeric(rep(tmp, colSums(crimtab)))

hist(height_inches, freq = FALSE)
```

## Histogram of height\_inches



```
varVec <- rep(0,10000)
for (i in seq(0,10000,1)) {
  varVec[i]<-var(sample(height_inches, 4, replace = T))
}
mean(varVec)
```

```
## [1] 6.496275
```

```
var(height_inches)
```

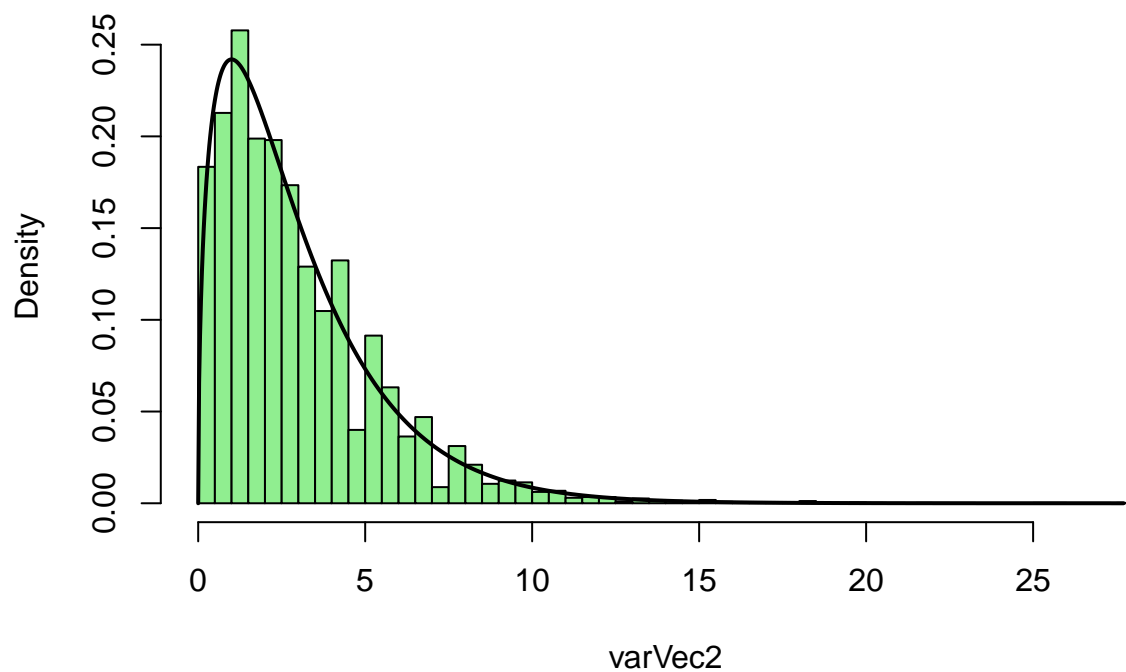
```
## [1] 6.542118
```

d) The sampling variance is fairly close to the variance of the data set.

e)

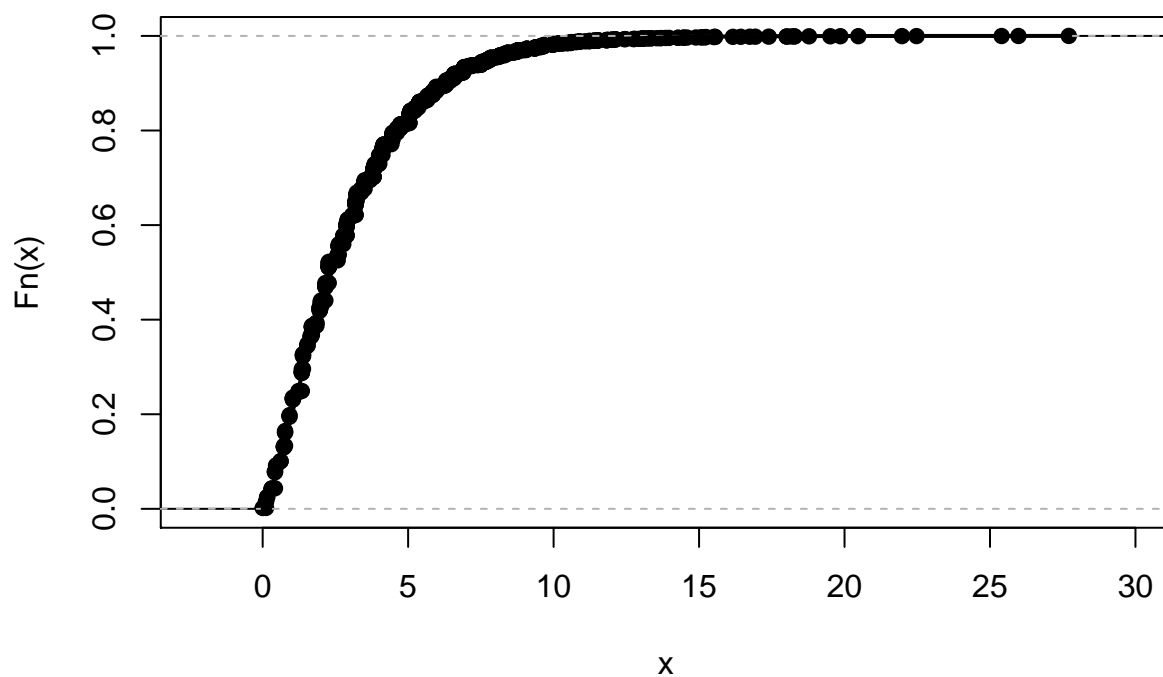
```
varVec2 <- varVec
for (i in seq(0,10000,1)) {
  varVec2[i] <- varVec[i]*3/mean(varVec)
}
hist(varVec2, col="lightgreen", breaks = 40, freq = F)
x_vec <- seq(0,max(varVec2), length = 1000)
lines(x_vec, lwd="2", col="black", dchisq(x_vec, df = 3))
```

## Histogram of varVec2



```
plot(ecdf(varVec2))  
lines(x_vec, pchisq(x_vec, df=3), lwd="2", col="black")
```

## ecdf(varVec2)



7. We can avoid using the tables in the back of the book by using R functions. For this problem, use the R



function `qchisq` to do the following computations: (a). Reproduce the *row* of quantiles (or “percentage points”) in Table 6, pages 850-851, for 10 degrees of freedom. (b). Reproduce the *column* of quantiles in Table 6 on page 851 that corresponds to an upper tail area of  $\alpha = 0.05$ . (Small approximation errors might make the book’s answer slightly different in some cases.)

---

**Answer:**

a)

```
qchisq(c(.005,.01,.025,.05,.1), df= 10 )
```

```
## [1] 2.155856 2.558212 3.246973 3.940299 4.865182
```

b)

```
qchisq(df = seq(1,29,1), .05)
```

```
## [1] 0.00393214 0.10258659 0.35184632 0.71072302 1.14547623
## [6] 1.63538289 2.16734991 2.73263679 3.32511284 3.94029914
## [11] 4.57481308 5.22602949 5.89186434 6.57063138 7.26094393
## [16] 7.96164557 8.67176020 9.39045508 10.11701306 10.85081139
## [21] 11.59130521 12.33801458 13.09051419 13.84842503 14.61140764
## [26] 15.37915658 16.15139585 16.92787504 17.70836618
```

```
qchisq(df= seq(30,100, 10), .05)
```

```
## [1] 18.49266 26.50930 34.76425 43.18796 51.73928 60.39148 69.12603 77.92947
```

- 
8. Complete **Exercise 7.19** in the text. You can use the `pchisq` function to complete the probability. Then check your answer by setting the usual seed, simulating 10,000 sample variances under the stated conditions, and checking the empirical proportion of sample variances that are bigger than 0.065. Also, we know that in theory,  $E(S^2) = \sigma^2$ , so check that the sample mean of your 10,000 simulated sample variances is approximately  $\sigma^2 = 0.04$ , as given in the problem.

---

**Answer:**

Exercise 7.19

$$\text{Let, } X \sim \chi_9^2$$

$$\sigma^2 = 0.04$$

$$n = 10$$

$$P(s^2 > 0.065) = P\left(\frac{\sigma^2 X}{(n-1)} > 0.065\right)$$

$$P(s^2 > 0.065) = 1 - pchisq(0.065, 9) = 0.10176$$

```
1-pchisq(14.625,9)
```

```
## [1] 0.1017651
```

```
set.seed(4302019)
```

```
sample_chi <- (.04/9)*rchisq(10000, df = 9, ncp=0)
mean(sample_chi>0.065)
```

```
## [1] 0.1034
```

```
mean(sample_chi)
```

```
## [1] 0.04011735
```

The proportion of sample variances greater than 0.065 are similar to the actually proportion.

---

9. Suppose that  $Y_1, \dots, Y_n$  are iid  $N(\mu, \sigma^2)$  and let  $S^2$  denote the sample variance, as usual. (a). Use the fact that  $(n-1)S^2/\sigma^2 \sim \chi^2$  with  $n-1$  df to determine the theoretical variance of the sample variance,  $V(S^2)$ . Show your work. (You can use the fact that a  $\chi^2_\nu$  random variable has  $V(\chi^2_\nu) = 2\nu$ .) (b). Compute the value of your theoretical variance formula using the numbers given in **Exercise 7.19**. (c). Check your computed theoretical value in (b) against the empirical variance of your simulated variances from the previous problem, (#8 above).
- 

**Answer:**

a)

$$\text{Let } X \sim \chi^2_{(n-1)}$$
$$X = \frac{S^2(n-1)}{\sigma^2} \implies S^2 = \frac{\sigma^2 X}{(n-1)}$$

b)

```
var(sample_chi)
```

```
## [1] 0.0003541971
```

---