

# MATH 4780 (MSSC 5780) Homework 1

## Probability and Statistics Review

**Due Date: September 12, 2022 11:59 PM**  
**Henrique Medeiros Dos Reis**

### 1 Homework Instruction and Requirement

- Homework 1 covers course materials of Week 1 to 2.
- Please submit your work in **one PDF** file to **D2L > Assessments > Dropbox**. *Multiple files or a file that is not in pdf format are not allowed.*
- In your homework, please number and answer questions **in order**.
- There is no need to submit your work on R programming part. However, I highly recommend reviewing basic R syntax for R data structures and graphics if you haven't got familiar with it.
- It is your responsibility to let me understand what you try to show. If you type your answers, make sure there are no typos. I grade your work based on *what you show, not what you want to show*. If you choose to handwrite your answers, write them neatly. If I can't read your sloppy handwriting, your answer is judged as wrong.

### 2 R Programming

Please sharpen your R skill. **No need** to show your work on this part! : )

1. Please register RStudio Cloud or get R and RStudio installed in your laptop. Read the R and RStudio slides in Week 1 module for installation instruction and basic usage of RStudio.
2. If you are not familiar with basic R syntax and data types, please review my slides of MATH 3570 in Week 1 module.
3. You can test your understanding of R by doing the problems in `basic_r.pdf` that is actually Homework 1 of my MATH 3570 course Spring 2021.

### 3 Probability and Statistics Review

We will use some facts or properties discussed in MATH 4700 and 4710. Here we don't learn why (which should be taught in MATH 4700 and 4710), but just know what they are, and apply them for regression analysis.

1. Use the linearity of expected value  $E(X + Y) = E(X) + E(Y)$  and  $E(aX) = aE(X)$  where  $X$  and  $Y$  are random variables and  $a$  is a constant to show that for a random variable  $Y$ ,

$$\text{Var}(Y) = E(Y^2) - [E(Y)]^2$$

**Solution:**

By definition:

$$\text{Var}(Y) = E[(Y - E(Y))^2]$$

then

$$\text{Var}(Y) = E[(Y - E(Y))(Y - E(Y))]$$

$$\text{Var}(Y) = E[Y^2 - 2E[Y]Y + E[Y]^2]$$

since

$$E(X + Y) = E(X) + E(Y)$$

$$\text{Var}(Y) = E[Y^2] - 2E[Y]E[Y] + E[E[Y]^2]$$

$$\text{Var}(Y) = E[Y^2] - 2E[E[Y]^2] + E[E[Y]^2]$$

$$\text{Var}(Y) = E[Y^2] - E[E[Y]^2]$$

and

$$E(aX) = aE(X)$$

$$\text{Var}(Y) = E(Y^2) - [E(Y)]^2$$

2. Let the two independent random variables be  $Y_1 \sim N(3, 8)$  and  $Y_2 \sim N(1, 4)$ . What is the distribution of the variable  $2Y_1 + 3Y_2$ ?

**Solution:**

Since  $Y_1$  and  $Y_2$  are independents, it is known that

$$aX + bY \sim N(a\mu_X + b\mu_Y, a^2\sigma_X^2 + b^2\sigma_Y^2)$$

then

$$(2Y_1 + 3Y_2) \sim N(2(3) + 3(1), 2^2(8) + 3^2(4))$$

$$(2Y_1 + 3Y_2) \sim N(9, 68)$$

3. Suppose  $Y_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$ ,  $i = 1, 2, \dots, n$ . Show that  $\bar{Y} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$ , where  $\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$ .

**Solution:**

$$E[\bar{Y}] = E\left[\frac{1}{n} \sum_{i=1}^n Y_i\right] = \frac{1}{n} E\left[\sum_{i=1}^n Y_i\right] = \frac{1}{n} \sum_{i=1}^n E[Y_i]$$

$$\frac{1}{n} n\mu = \mu$$

and

$$\text{Var}(\bar{Y}) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n Y_i\right) = \left(\frac{1}{n}\right)^2 \text{Var}\left(\sum_{i=1}^n Y_i\right)$$

$$\frac{1}{n^2} \sum_{i=1}^n \text{Var}(Y_i) = \frac{1}{n^2} \sum_{i=1}^n \sigma^2 = \frac{1}{n^2} n\sigma^2$$

$$\text{Var}(\bar{Y}) = \frac{\sigma^2}{n}$$

therefore

$$\bar{Y} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

4. Let the sample variance be  $S^2 = \frac{\sum_{i=1}^n (Y_i - \bar{Y})^2}{n-1}$ . The two facts are

i.  $\frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2$

ii.  $\frac{\bar{Y} - \mu}{\sigma/\sqrt{n}}$  and  $\frac{(n-1)S^2}{\sigma^2}$  are independent.

Use i. and ii. and the fact that *If  $Z \sim N(0, 1)$ ,  $V \sim \chi_v^2$  and  $Z$  and  $V$  are independent, then  $\frac{Z}{\sqrt{V/v}} \sim t_v$  to show that*

$$\frac{\bar{Y} - \mu}{S/\sqrt{n}} \sim t_{n-1}$$

**Solution:**

$$\text{If } Z = \frac{\bar{Y} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1) \text{ and } V = \frac{vS^2}{\sigma^2} \sim \chi_v^2 = \frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2$$

then

$$\frac{Z}{\sqrt{V/v}} \sim t_v = \frac{\frac{\bar{Y} - \mu}{\sigma/\sqrt{n}}}{\sqrt{\frac{(n-1)S^2/\sigma^2}{n-1}}} \sim t_{n-1}$$

$$\frac{\sqrt{n} \frac{\bar{Y} - \mu}{\sigma}}{\sqrt{\left(\frac{(n-1)S^2}{\sigma^2}\right)/(n-1)}}$$

$$\frac{\sqrt{n} \frac{\bar{Y} - \mu}{\sigma}}{\sqrt{\left(\frac{(n-1)S^2}{\sigma^2}\right)\left(\frac{1}{n-1}\right)}} = \frac{\sqrt{n} \frac{\bar{Y} - \mu}{\sigma}}{\sqrt{\frac{S^2}{\sigma^2}}} = \frac{\sqrt{n} \frac{\bar{Y} - \mu}{\sigma}}{\frac{S}{\sigma}}$$

$$\frac{\sqrt{n}(\bar{Y} - \mu)}{\sigma} \frac{\sigma}{S} = \frac{\sqrt{n}(\bar{Y} - \mu)}{S} = \frac{\bar{Y} - \mu}{S/\sqrt{n}}$$

Therefore

$$\frac{\bar{Y} - \mu}{S/\sqrt{n}} \sim t_{n-1}$$

5. Suppose  $Y_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$ ,  $i = 1, 2, \dots, n$ , with unknown  $\mu$  and  $\sigma$ . The  $100(1 - \alpha)\%$  confidence interval (CI) for the population mean  $\mu$  is  $\left(\bar{y} - t_{\alpha/2, n-1} \frac{s}{\sqrt{n}}, \bar{y} + t_{\alpha/2, n-1} \frac{s}{\sqrt{n}}\right)$ . Use simulation with  $\alpha = 0.1$ ,  $\mu = 4$  and  $\sigma = 2$  to verify that such CIs contain  $\mu$  about  $100(1 - \alpha)\%$  of times. Fill the percentage in the following table, and comment on your results. Attach your code at the end of your homework PDF file.

| Simulation times | $n = 5$ | $n = 30$ | $n = 200$ |
|------------------|---------|----------|-----------|
| 20               | 90%     | 85%      | 90%       |
| 1000             | 88.9%   | 90.4%    | 92.1%     |
| 20000            | 90.27%  | 89.69%   | 89.69%    |

As we can see in the table there is a tendency that the percentage of times where  $\mu$  is inside the CI is getting close and closer to  $100(1 - \alpha) = 90\%$  as the number of simulations and the size of  $n$  increases.

**Solution:**

```
#set seed, so results are consistent
set.seed(1)
#set variables
alpha <- 0.1
mu <- 4
sigma <- 2
numsOfN <- c(5, 30, 200)
simNum <- c(20,1000,20000)

#fist lets loop through all the values of n
for(n in numsOfN){
  #make sure we are starting with 0 for results needed
  inCI <- 0
  percentInCI <- 0
  for(j in 1:3) {# do the process once for each simulation number
    inCI <- 0
    for(i in 1:simNum[j]){#loop from 1 to the max number of simulations
      #compute numbers for the normal distribution, t value and CI
      result.rnorm <- rnorm(n,mu,sigma)
      t <- qt(1-alpha/2, n-1)*sd(result.rnorm)/sqrt(n)
      lower <- mean(result.rnorm)-t
      upper <- mean(result.rnorm)+t
      if(mu <= upper & mu >= lower){#in case it is in the CI
        inCI <- inCI+1#increase the number of numbers in CI
      }
      if(i == simNum[j]){
        #calculate and print formated results in percentages
        percentInCI <- inCI/simNum[j]
        cat('When doing ',simNum[j],' simulations with n = ', n,
          ' the percentage is ',percentInCI*100,'% \n',sep='')
      }
    }
  }
}
```

```
## When doing 20 simulations with n = 5 the percentage is 90%
## When doing 1000 simulations with n = 5 the percentage is 88.9%
## When doing 20000 simulations with n = 5 the percentage is 90.27%
## When doing 20 simulations with n = 30 the percentage is 85%
## When doing 1000 simulations with n = 30 the percentage is 90.4%
## When doing 20000 simulations with n = 30 the percentage is 89.69%
## When doing 20 simulations with n = 200 the percentage is 90%
## When doing 1000 simulations with n = 200 the percentage is 92.1%
## When doing 20000 simulations with n = 200 the percentage is 89.685%
```

6. If  $U_1$  and  $U_2$  are independent and both are uniform variables over  $[0,1]$  interval ([https://en.wikipedia.org/wiki/Continuous\\_uniform\\_distribution](https://en.wikipedia.org/wiki/Continuous_uniform_distribution)), then  $X_1$  and  $X_2$  defined by

$$X_1 = \sqrt{-2\ln(U_1)} \cos(2\pi U_2), \quad X_2 = \sqrt{-2\ln(U_1)} \sin(2\pi U_2)$$

are independent  $N(0,1)$  variables. Draw 10,000 samples for  $U_1$  and  $U_2$  using the `runif()` function, and use the transformation to generate the samples of  $X_1$  and  $X_2$ . Verify

- the standard normality of  $X_1$  and  $X_2$  by plotting their histogram with a superimposed standard normal density.
- the independence of  $X_1$  and  $X_2$  by plotting the scatterplot of  $X_1$  and  $X_2$  and computing their correlation coefficient.

### Solution:

First bullet point:

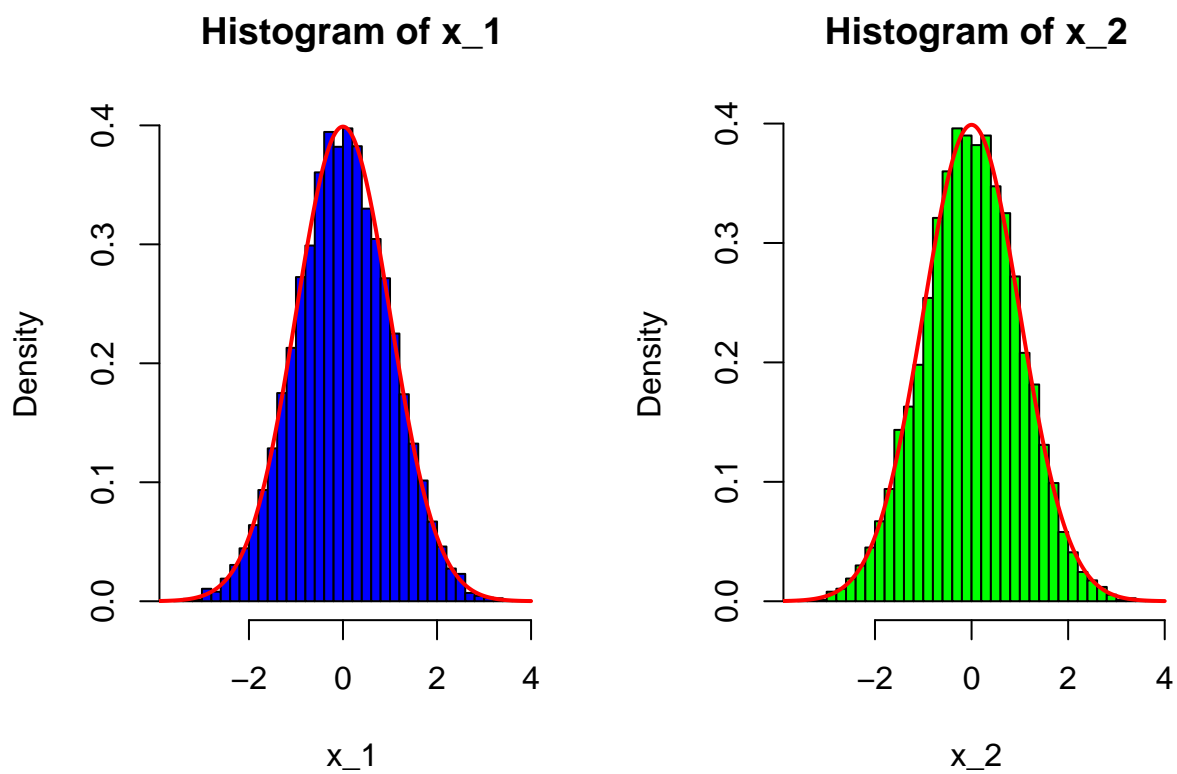
```
#compute random numbers
u_1<-runif(10000)
u_2<-runif(10000)
#do the transformation
x_1 <- sqrt(-2*log(u_1, exp(1)))*cos(2*pi*u_2)
x_2 <- sqrt(-2*log(u_1, exp(1)))*sin(2*pi*u_2)

#compute a length and use that to construct a normal distribution
len <- seq(-4,4, length=1000)
normalDist <- dnorm(len)

#do 2 graphs in the same window
par(mfrow = c(1,2))

hist(x_1, probability = TRUE,breaks = 50, col="Blue")
lines(len,normalDist, col="Red", lwd=2)

hist(x_2, probability = TRUE,breaks = 50, col="Green")
lines(len,normalDist, col="Red", lwd=2)
```

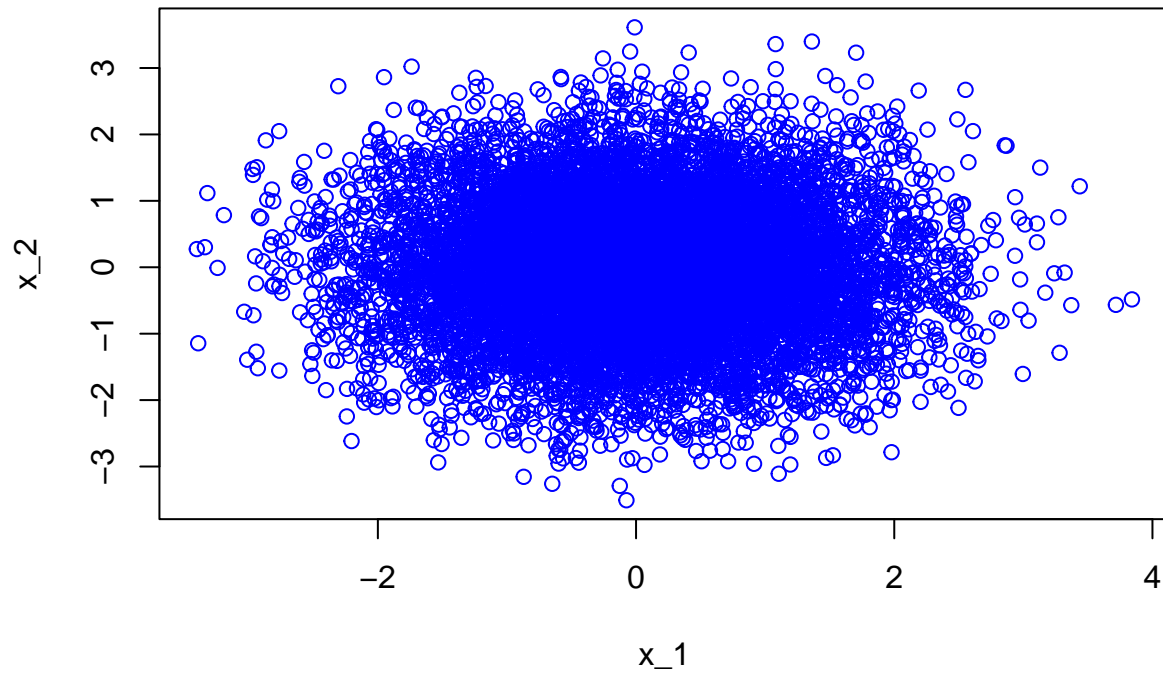


Second Bullet point:

```
#u_1<-runif(10000)
#u_2<-runif(10000)
#x_1 <- sqrt(-2*log(u_1, exp(1)))*cos(2*pi*u_2)
#x_2 <- sqrt(-2*log(u_1, exp(1)))*sin(2*pi*u_2)

plot(x_1,x_2, main="Correlation Between X_1 and X_2", col="Blue")
```

## Correlation Between X\_1 and X\_2



```
#concatenate text and the computation for the correlation  
cat("Correlation coefficient between X_1 and X_2 =", cor(x_1,x_2))
```

```
## Correlation coefficient between X_1 and X_2 = -0.01210106
```