# Lab 2

*Fan Yang (fy2232)*

*September 25, 2017*

## Instructions

Before you leave lab today make sure that you upload a .pdf file to the canvas page (this should have a .pdf extension). This should be the PDF output after you have knitted the file, we don't need the .Rmd file (don't upload the one with the .Rmd extension). The file you upload to the Canvas page should be updated with commands you provide to answer each of the questions below. You can edit this file directly to produce your final solutions. Note, however, in the file you upload you should the above header to have the date, your name, and your UNI. Similarly, when you save the file you should replace **UNI** with your actualy UNI.

## Tasks

The work we do today, will build on the work we did in the previous lab.

### Sample distribution of the sample mean

1) Generate 1 million random draws from a normal distribution with $\mu = 3$ and $\sigma^2 = 4$ and save them in a vector named **normal1mil**. Calculate the mean and standard deviation of **normal1mil**.

```
normal1mil <- rnorm(1000000,mean=3,sd=2)
mean(normal1mil)
```

```
## [1] 3.000094
```

```
sd(normal1mil)
```

```
## [1] 2.000371
```

2) Find the mean of all the entries in **normal1mil** that are greater than 3. You may want to generate a new vector first which identifies the elements that fit the criteria.

```
mil3 <- normal1mil[normal1mil>3]
mean(mil3)
```

```
## [1] 4.595736
```

3) Create a matrix **normal1mil_mat** from the vector **normal1mil** that has 10,000 columns (and therefore should have 100 rows).

```
normal1mil_mat <- matrix(normal1mil, nrow = 100, ncol = 10000)
```

4) Calculate the mean of the $1234^{th}$ column.

```
mean(normal1mil_mat[,1234])
```

```
## [1] 3.131722
```

5) Use the **colSums()** functions to calculate the *means* of each column of **normal1mil_mat**. Save the vector of column means as **col_means** as it will be used in the next task. Use the **apply()** function
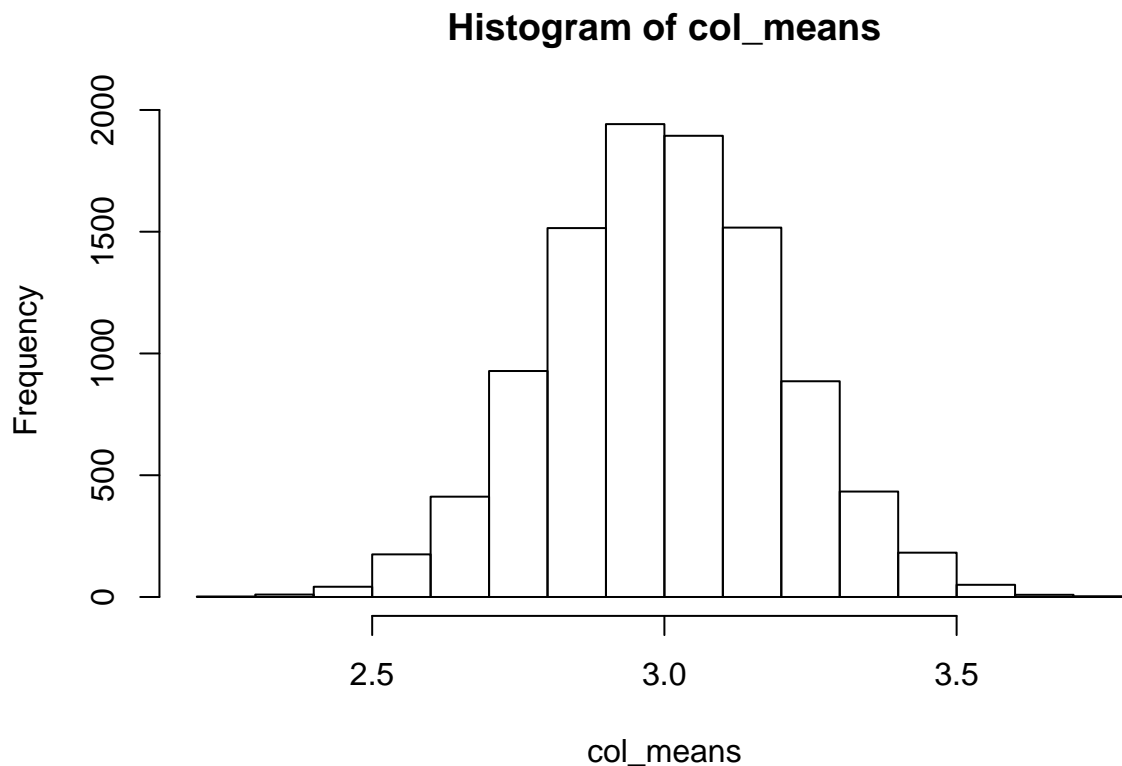
to do the same thing, saving the results in a new vector **col__means2**, then compare the two with
**all(col__means == col__means2)**.

```r
col_means <- c()
col_means[1:10000] <- colSums(normal1mil_mat[,1:10000]) / 100
col_means2 <- apply(normal1mil_mat,2,sum) / 100
all(col_means == col_means2)
```

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

6) Finally, produce a histogram of the column means you calculated in task (5). What is the distribution
   that this histogram approximates (i.e. what is the distribution of the sample mean in this case)?

```r
hist(col_means)
```



**Histogram of col_means**

```r
mean(col_means)
```

```
## [1] 3.000094
```

```r
sd(col_means)
```

```
## [1] 0.1992806
```

According to the histogram, the distribution of the sample mean approximates a normal distribution.

7) Only complete the following questions if you have time. We'd like to count the number of values
   beyond 3 standard deviations from the mean in each column of our matrix and store them in a vector
   called **num__3devs** (that should have length 10,000). (a) Write code that loops over the columns and
   calculates this value for each column filling in the vector **num__3devs** as it iterates. (b) Use a smart
   call to the **apply()** function to do the same thing. (c) Make a table of these values.

(a)

```r
num_3devs <- c()
for (i in 1:10000){
  num_3devs[i]<-sum((abs(normal1mil_mat[,i] - col_means[i]) > 3*sd(normal1mil_mat[,i]))*1)
}
```

(b)

```r
num_3devs2 <- c()
num_3devs2<-apply(normal1mil_mat,2,function(vec) sum((abs(vec - mean(vec)) > 3*sd(vec))*1))
```

(c)

```r
table(num_3devs)
```

```
## num_3devs
##    0    1    2    3
## 7990 1899  110    1
```