

HW 2 More on R

Wylan Gao STAT 5400

Due: Sep 13, 2024 9:30 AM

Submit your solutions as an .Rmd file and accompanying .pdf file. Include all the **relevant** R code and output. With the echo option in R markdown, you opt to hide some R code that is not very relevant (but still run it to generate the document).

Always comment on your result whenever it is necessary. For example, in problem 3, display and also comment on the Q-Q plot that you produce

Reading assignments.

Read Chapters 2-4, 7 of *Using R for Data Analysis and Graphics*. <https://cran.r-project.org/doc/contrib/usingR.pdf>.

Problems

1. Use the `system.time` function in R to time the performance of the same task in two different ways:
 - Generate a vector of 500,000 random variates from a Normal (0, 1) density and use the `sum` function to calculate their sum.
 - Create a variable called `answer` and initialize it to 0. Then, using a `for` loop, do the following steps 500,000 times: generate a single Normal (0, 1) value and add it to sum contained in `answer`.
 - Besides including the R code and output, compare on the relevant timings for both methods and state which one is more efficient.

```
system.time(sum(rnorm(500000, mean = 0, sd = 1)))
```

```
##      user      system elapsed  
##    0.02      0.00      0.03
```

```
system.time({  
  count <- 0  
  answer <- function(x) {  
    for (i in 1:500000) {  
      count <- count + rnorm(1)  
    }  
    return(count)  
  }  
  print(answer())  
})
```

```
## [1] 500.4147
```

```
##      user  system elapsed
##      0.70    0.16    1.33
```

The one that is more efficient is the package for sum running runif.

2. Use R to do the following:

- Create a matrix called M with the following entries:

$$\begin{bmatrix} 1 & 3 & 5 \\ 2 & 4 & 6 \\ 3 & 6 & 9 \end{bmatrix}.$$

- Create a vector called v with the following entries: 17 46 181.
- Compute and display the product Mv produced by matrix multiplication.
- Compute and display the transpose of M .
- Display only those elements of v that have values less than 50.

#I create a matrix

```
m <- matrix(c(1, 3, 5,
              2, 4, 6,
              3, 6, 9),
            nrow = 3, ncol = 3, byrow = TRUE)
```

```
v <- c(17, 46, 181)
```

#this multiple the vector and the colimn

```
m %*% v
```

```
##      [,1]
## [1,] 1060
## [2,] 1304
## [3,] 1956
```

#this is a tranpose

```
t(m)
```

```
##      [,1] [,2] [,3]
## [1,]    1    2    3
## [2,]    3    4    6
## [3,]    5    6    9
```

#display only elements that have Less than 50

```
v[ v< 50]
```

```
## [1] 17 46
```

```
library(moments)
```

#the name of the function is skewness

It is in the package of moments.

- Use either `help.search` function or just google it to locate a package that contains a function to compute the skewness of a vector of numbers. Make sure that it uses the standard definition of skewness. What is the name of the function, and which package is it in?

It is in moments. The function is called skewness.

- Locate an R function that computes the five-number summary of a vector of numbers. What is the name of the function, and which package is it in?

It is in the `fivenum` function. It gives me the minimum, first quartile, median, third quartile, maximum. Summary provides a similar output.

- Write an R function that does the following:
 - Accepts one argument: a vector.
 - Checks whether the vector is numeric.
 - If not, displays the message `Vector must be numeric and exist.`
 - If yes, computes the skewness of the values (after removing any missing values).
 - If the absolute value of skewness is less than 1, returns a list containing two objects: skewness in an object named `skewness`; a vector consisting of the mean and standard deviation in an object named `descstats`.
 - Otherwise, returns a list containing two objects: skewness in an object named `skewness`; a vector consisting of the five-number summary in an object named `descstats`.
 - Run your function in R three times, using the following vectors as arguments:
 - `c("stat", "actuarial", "2022")`
 - `rnorm(100)`
 - `rexp(5)`
 - In the document that you submit for homework, show both the R code and output for the three calls to it.

```
library(moments)
```

```
givemevec <- function(x) {  
  if (is.numeric(x)) {  
    vec <- na.omit(x)
```

```

    if (abs(skewness(vec)) < 1) {
      return(list(skewness = skewness(vec) , descstats = c(mean(vec),
sd(vec))))
    } else {
      element_2 <- fivenum(vec)
      return(list(skewness = skewness(vec), descstats = fivenum(vec)))
    }
  } else {
    return(print('Vector must be numeric'))
  }
}

```

```
test1 <- givemevec(c("stat", "actuarial", "2022"))
```

```
## [1] "Vector must be numeric"
```

```
test2 <- givemevec(rnorm(100))
```

```
test3 <- givemevec(rexp(5))
```

```
test1
```

```
## [1] "Vector must be numeric"
```

```
test2
```

```
## $skewness
```

```
## [1] -0.2802904
```

```
##
```

```
## $descstats
```

```
## [1] 0.04020901 0.96040373
```

```
test3
```

```
## $skewness
```

```
## [1] 0.7265261
```

```
##
```

```
## $descstats
```

```
## [1] 1.270177 1.445190
```

4. Review the example on two-sample t-test on slide 45 of S2P1.pdf. If you are not familiar with t-test, you may google some online tutorials.
 - Imagine we have run 1000 experiments (rows) and each of which collects data on 50 individuals (columns). The first 25 individuals in each experiment are assigned to group 1 and the rest to group 2.

- You may imagine in practice we only have 50 observations, 25 of which belongs to group 1. With simulations, we have luxury to generate 1000 samples, each of which has 50 observations, to explore the distribution of t-statistics.

- We first generate some random data to represent this problem.

```
set.seed(1)
m <- 1000
n <- 50
X <- matrix(rnorm(m * n, mean=10, sd=3), nrow=m)
grp <- rep(1:2, each=n/2)
```

- For the first sample (i.e., the first row), compute the t-statistic manually. Compare your result with the output from the following code.

```
t.test(X[1, grp==1], X[1, grp==2])$stat

##           t
## -0.5284632

#we manually compute the value of T
x_1 <- X[1, grp == 1]
x_2 <- X[1, grp == 2]

top <- mean(x_1) - mean(x_2)

var_x1 <- var(x_1)
var_x2 <- var(x_2)

n_1 <- length(x_1)
n_2 <- length(x_2)

bottom <- sqrt((var_x1 / n_1) + (var_x2 / n_2))

t_stat_manual <- top / bottom

t_stat_manual
## [1] -0.5284632

library(matrixStats)

optimized_t_test_all <- function(X, grp) {
  group_1_indices <- grp == 1
  group_2_indices <- grp == 2
```

```

X1 <- X[, group_1_indices]
X2 <- X[, group_2_indices]

t_stats <- (rowMeans(X1) - rowMeans(X2)) / sqrt(rowVars(X1) / ncol(X1) +
rowVars(X2) / ncol(X2))

return(t_stats)
}

# Measure the time for optimized function
system.time({
  t_stats_optimized <- optimized_t_test_all(X, grp)
})

##      user  system elapsed
##         0         0         0

# Compare with the built-in t.test
system.time({
  t_stats_t_test <- apply(X, 1, function(row_data) {
    t.test(X[1, grp==1], X[1, grp==2])$stat
  })
})

##      user  system elapsed
##    0.16    0.00    0.25

#the authors cpde
rowtstat <- function(X, grp){
  t_stat <- function(X) {
    m <- rowMeans(X)
    n <- ncol(X)
    var <- rowSums((X - m) ^ 2) / (n - 1)

    list(m = m, n = n, var = var)
  }

  g1 <- t_stat(X[, grp == 1])
  g2 <- t_stat(X[, grp == 2])

  se_total <- sqrt(g1$var / g1$n + g2$var / g2$n)
  (g1$m - g2$m) / se_total
}
system.time(t3 <- rowtstat(X, grp))

##      user  system elapsed
##    0.00    0.00    0.02

```

```
#>   user  system elapsed  
#> 0.015   0.001   0.014
```

Their code is faster because it refers used vectorized operations., the calculations are efficient and uses libraries. However, the run time seems similar on the `system.time` function. It is important to use vectors and columns and not store variables when doing calculation and optimizations. R is a line by line language similar to python.

- Use R to compute the t-statistics for all 1000 samples. Try to optimize your code to make it efficient.
- Check the R code in the end of section 24.7 on <https://adv-r.hadley.nz/perf-improve.html#t-test>. Use `system.time` to compare the computing time of the R function provided on the webpage with your own function. Comment on why your code is slower or faster. Don't check their code before you finish writing your own code.