

STAT40780 Data Programming with C (online)

Lab Sheet 1

This week's lab will involve profiling slow R code to identify bottlenecks, writing an improved/more efficient version of the code, and benchmarking the improved version against the original.

1 Profiling R code

The R code on the next page counts the odd numbers in each row of an integer matrix. The function `countOdds()` takes an integer matrix as input, and returns a vector of counts of the odd numbers for each row of the matrix. The function `countOdds()` makes calls to the function `vecOdds()` to count the odd numbers in a given row (vector). The `vecOdds()` function in turn passes each element of the row (vector) to the `isOdd()` function - if the element is odd it returns a value of 1, and returns 0 otherwise. This code contains nested `for` loops - best avoided in R.

You may not have come across the modulus operator in R, `%%`, that returns the remainder from the division of two integer values. For example, `10 %% 2` returns 0, while `9 %% 2` returns 1.

Step 1: Have a look through this code and ensure you understand how it works, before proceeding.

Step 2: Create a new R script in RStudio. Then copy and paste the below R code into your R script. Save your script file as `lab1.R`.

Counts of odd numbers in each row of a matrix

```

1 countOdds <- function( X ){
2
3   #type checking
4   if( !is.matrix(X) ) stop("input is not a matrix")
5   if( !is.integer(X) ) stop("input should be of type integer")
6
7   N <- nrow(X) #number of rows of X
8   numOdds <- rep(0, N ) #vector to store counts; initialize to zero
9
10  for( i in 1:N ){ #for each row i in X
11    numOdds[i] <- vecOdds( X[ i, ] ) #count odds in row i
12  } #end of for
13
14  return(numOdds) #return the vector of counts
15
16 } #end of countOdds()
17
18
19 #count the number of odd numbers in a vector y
20 vecOdds <- function(y){
21
22   K <- length(y) #length of the vectors (row)
23   #keep count of number of odd numbers in y
24   sumOdds <- 0 #initialize counter to 0 (no odd numbers)
25
26   for(j in 1:K){ #for each element of the vector y
27     sumOdds <- sumOdds + isOdd( y[j] ) #update sumOdds
28   }
29   return(sumOdds)
30 } #end of vecOdds()
31
32
33 #function that checks whether its input is an odd number
34 #returns 1 if input is odd and 0 otherwise
35 isOdd <- function(num){
36   if( num %% 2 ){
37     return(1)
38   } else{
39     return(0)
40   }
41 }

```

Step 3 Create an integer matrix to input to the R function.

Generate data from a binomial distribution, with number of events = 30, and probability of success = 0.5. This will generate integer data valued between 0 and 30. Make a barchart of your data to see what it looks like. Then arrange your random data into a matrix with 100000 rows.

Use the following R code to do this:

Generate a matrix of integer values

```
1 #generate 10000000 data points from a binomial distribution
2 simdata <- rbinom(n=10000000, size=30, prob = 0.5)
3
4 #check what the simulated data look like
5 barplot(table(simdata))
6
7 #arrange simulated data into a matrix with 100000 rows
8 X <- matrix(simdata, nrow= 100000)
9 dim(X) #print dimensions of matrix
```

Call the countOdds() function

```
1 countOdds(X) #call to the countOdds() function
```

Step 4 Profiling.

This is the main task of this section of the lab. Profile the R function `countOdds()` to identify bottlenecks.

2 Benchmarking code

Step 1: Using the results of your profiling as a guideline, write a more efficient version of the `countOdds()` function - be sure to give the revised function a different name! Hint: for improving efficiency, make use of the `apply()` function.

Step 2: Benchmark your code. Compare the original `countOdds()` function to your (hopefully) improved version. How much faster is your new function, relative to the original? Hint: use the `Rprof()` function in R. Hint: use the `Rprof` function.

You can experiment with specifying the option `line.profiling=TRUE` to instruct `Rprof` to profile code by line, rather than by function/operation.