

Stat 243: Problem Set 4

Eugene Yedvabny

10/15/2014

1

The appearance of the 3 over a 2 is due to lazy evaluation. The value for `x` is only evaluated when `x` is called by the `+` operator within `f1`. Since `x` assigns `y` as part of its evaluation block, lazy eval ensures that default `y=0` is overwritten and thus `x + y = 2 + 1 = 3`.

2

The answer to this question will assume that the time it takes to *create* the index vector or the boolean vector is irrelevant. As such I am going to create vectors of varying sizes and types, as well as their indexing vectors, and only benchmark the subsetting operation itself. I've created the following helper function to generate a dataframe for the each length of the vector I want to subset.

```
library(dplyr)
library(microbenchmark)

# The following function will benchmark list subsetting for numeric and string lists
# The input parameter is the length of the vector to be subset
benchSubsetting <- function(length){
  # Create a random vector of numbers
  test.nums <- runif(length)

  # Create a random vector of strings (of 10 characters)
  test.strs <- replicate(length,
    paste(sample(letters, 10, replace = T),
      collapse=''))

  # Create an accessing boolean vector (skewed probs favor smaller subsets)
  test.bool <- sample(c(T,F), length, replace=T, prob=c(0.3,0.7))

  # Create an interger index corresponding to the same items as above
  test.index <- which(test.bool)

  # Run the test
  test.timing <- microbenchmark(
    nums.ints = test.nums[test.index],
    nums.bools = test.nums[test.bool],
    strs.ints = test.strs[test.index],
    strs.bools = test.strs[test.bool]
  )

  # Summarize on the mean timing?
  test.timing %>%
    group_by(expr) %>%
    summarize(median.time = median(time)) %>%
```

```

    mutate(length = length)
  }

```

The **microbenchmark** package executes each expression 100 times and returns the range of timings. Since there are often high-exec-time outliers, best metric to look at is the median of the timing distribution.

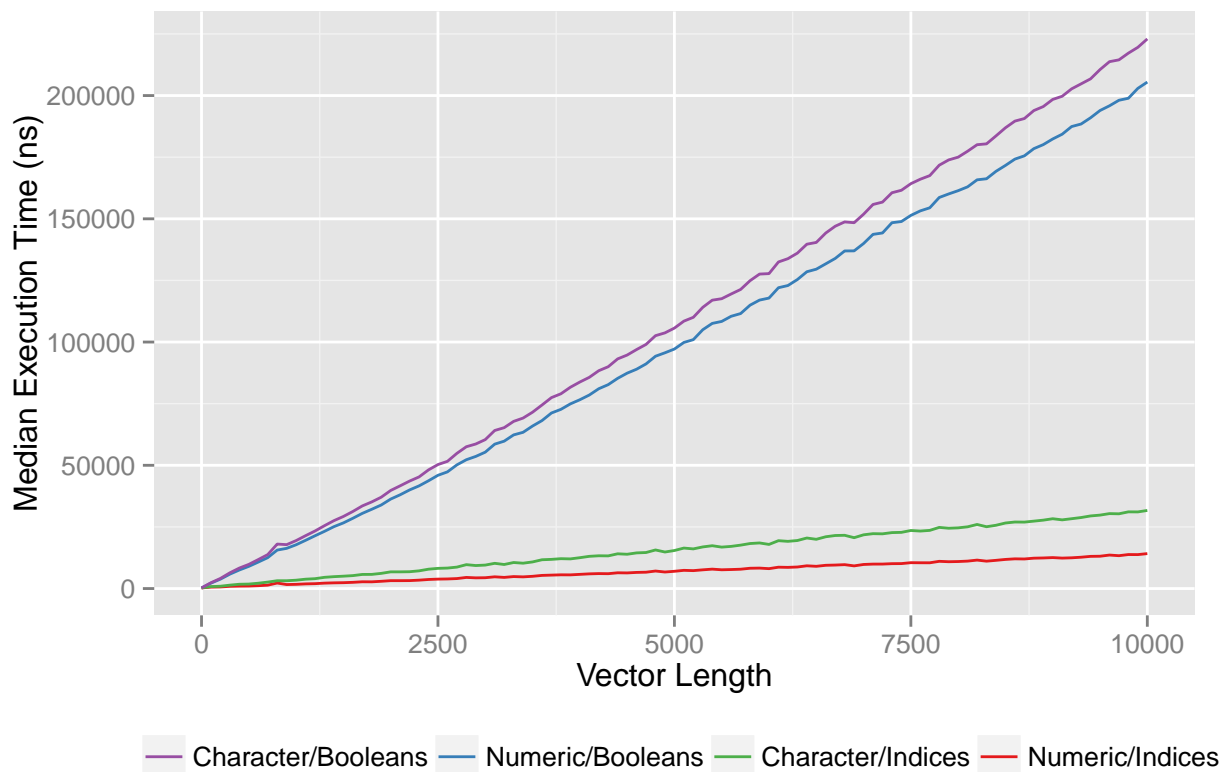
```

source("benchSubsetting.R")

# Enough lengths to get a good plot going
bench.data <- seq(0,10000,100) %>% lapply(benchSubsetting) %>% rbind_all

ggplot(bench.data,aes(x=length, y=median.time, color=expr)) +
  geom_line() +
  scale_color_brewer(type="qual", palette=6, name="",
    breaks=c("strs.bools","nums.bools",
      "strs.ints","nums.ints"),
    labels=c("Character/Booleans","Numeric/Booleans",
      "Character/Indices","Numeric/Indices")) +
  xlab("Vector Length") +
  ylab("Median Execution Time (ns)") +
  theme(legend.position = "bottom")

```



As is evident from the figure above, subsetting by index is *a lot* faster than subsetting by boolean. The performance costs are linear with vector length and are slightly higher for character vectors than numeric vectors. Since numeric vectors are stored continuously while character vectors are essentially arrays of pointers to individual strings, this likely yields the observed difference in subsetting time.

The difference between numeric and character vectors, however, pales in comparison to the gap between index-subsetting and boolean-subsetting. I've generated indexing boolean vectors of the same size as the data vector, so recycling should not be a factor in the timing. The issue really comes down to the number of elements to scan. I've set the parameters to recall about 30% of a vector's elements. In the case of a 100-element data vector, index subsetting would pass in 30 indices and underlying code would know exactly which memory addresses to compile into the new vector. In the case of boolean values, an entire 100-element vector is passed in, and the underlying code has to scan it all to check which of the data elements should be returned. So while the boolean subsetting is *convenient*, as it allows the use of vectorized logic expressions to generate the keys, it scales very poorly with vector size.

What's more telling is that this performance difference holds even when the number of passed-in elements and returned elements is the same for both methods. In the following example I am requesting all 1000 of the vector elements, in one instance by passing a 1000-entry logical vector of T, and in another a 1000-entry vector of indices.

```
test.data <- runif(1000)
test.bool <- rep(T,1000)
test.ind <- c(1:1000)

microbenchmark(
  with.bool = test.data[test.bool],
  with.ind = test.data[test.ind]
)

## Unit: microseconds
##      expr      min       lq      mean  median       uq      max neval
##  with.bool 19.633 20.187 20.38534 20.3650 20.5665 22.960   100
##   with.ind  4.808  4.986  5.23649  5.0505  5.1285 12.995   100
```

Based off the median timing, indexing with booleans is six times slower than indexing with indices, regardless of the vector length or number of elements requested.

3

For the purposes of answering this question, I am going to rely on the **pryr** package as it makes checking memory locations a lot easier. The package is up on CRAN now, so I figure it is fair game. **dplyr** also has convenience functions **location** and **changes** for tracking data frame memory addresses and they are a lot cleaner in presentation than `.Internal(inspect())`.

```
# Let's start off by creating heterogeneous data
strings <- replicate(100, paste(sample(letters, 10, replace = T), collapse=''))
numbers <- runif(100)
booleans <- sample(c(T,F),100,replace=T)
factors <- factor(sample(c("one","two","three"),100,replace=T))
raw.addr <- c(address(strings),address(numbers),address(booleans),address(factors))

# Merge the data into a dataframe
data.f <- data_frame(strings,numbers,booleans,factors)
frame.addr <- location(data.f)[[2]]

# Merge the data into a list
data.l <- list(strings,numbers,booleans,factors)
list.addr <- sapply(data.l,address)
```

```
# We can verify that the data frame is indeed just
# metadata over the underlying vectors
addresses <- data_frame(raw = raw.addr, data.frame = frame.addr, list = list.addr)
knitr::kable(addresses)
```

raw	data.frame	list
0x7f882e1a8100	0x7f882e1a8100	0x7f882e1a8100
0x7f882e1a9390	0x7f882e1a9390	0x7f882e1a9390
0x7f882e1ab140	0x7f882e1ab140	0x7f882e1ab140
0x7f882e1b0970	0x7f882e1b0970	0x7f882e1b0970

Ok, so far so good. The data frame and list are just metadata over the underlying vectors.

```
# Metadata prior to the change
location(data.f)
```

```
## <0x7f88313e72d0>
## Variables:
## * strings: <0x7f882e1a8100>
## * numbers: <0x7f882e1a9390>
## * booleans: <0x7f882e1ab140>
## * factors: <0x7f882e1b0970>
## Attributes:
## * names: <0x7f88313e71b0>
## * row.names: <0x7f882bd40c50>
## * class: <0x7f88313e7318>
```

```
# Change a value in one of the underlying vectors
data.f[10,1]<-"Hello"

location(data.f)
```

```
## <0x7f882f48c160>
## Variables:
## * strings: <0x7f882bd418c0>
## * numbers: <0x7f882e1a9390>
## * booleans: <0x7f882e1ab140>
## * factors: <0x7f882e1b0970>
## Attributes:
## * names: <0x7f88313e71b0>
## * row.names: <0x7f882bd420f0>
## * class: <0x7f88313e7318>
```

The modification has replaced the entire column of the data frame. The address of the data frame itself has changed too, but since a data frame is just a wrapper around lists, that's not a costly operation. The untouched columns have remained the same. Furthermore the change is local to the data-frame; the raw string vector and the list are still pointing into the same original location. Changing one of the two makes another copy and now the raw vector, the list, and the data frame all point in different locations.

```
c(address(strings),address(data.f$strings),address(data.l[[1]]))
```

```
## [1] "0x7f882e1a8100" "0x7f882bd418c0" "0x7f882e1a8100"
```

```
strings[10] <- "Hello"  
c(address(strings),address(data.f$strings),address(data.l[[1]]))
```

```
## [1] "0x7f882bd55eb0" "0x7f882bd418c0" "0x7f882e1a8100"
```

R always makes a copy if the underlying data is referred from other locations. But now that we have three different string vectors, does the copying still take place?

```
strings[11] <- "Goodbye"  
data.f$strings[11] <- "Goodbye"  
data.l[[1]][11] <- "Goodbye"  
c(address(strings),address(data.f$strings),address(data.l[[1]]))
```

```
## [1] "0x7f882bd10ad0" "0x7f882bd0d740" "0x7f882bd22b70"
```

It appears the answer is *yes*. Even though each object has their own version of the string vector, an element change results in an entirely new vector being created. The unmodified elements, however, remain the same across all three objects as they haven't been modified since creation.

```
c(address(numbers),address(data.f$numbers),address(data.l[[2]]))
```

```
## [1] "0x7f882e1a9390" "0x7f882e1a9390" "0x7f882e1a9390"
```

```
numbers[2] <- -10;  
data.f$numbers[2] <- -10;  
data.l[[2]][2] <- -10;  
c(address(numbers),address(data.f$numbers),address(data.l[[2]]))
```

```
## [1] "0x7f882bd45850" "0x7f882bd46080" "0x7f882bd468b0"
```

```
numbers[2] <- -15;  
data.f$numbers[2] <- -15;  
data.l[[2]][2] <- -15;  
c(address(numbers),address(data.f$numbers),address(data.l[[2]]))
```

```
## [1] "0x7f882bd47700" "0x7f882bd47f30" "0x7f882bd48760"
```

Same issue as before: even when only one pointer is accessing the vector, a new copy is still being made. I know for a fact that this should not be the case in newer versions of R, hence the following disclaimer.

Disclaimer: I've ran these codes through RStudio and knitr, which are known to yield different results than raw R. Since I cannot verify 100% what's going in the background with these packages, I am going to default to Hadley Wickham's extensive explanation of R's memory management, available at <http://adv-r.had.co.nz/memory.html>.

I am still a bit confused on how to have the S3 object maintain a list of attributes, such as the starting point, but only pass back a single value. Right now the printing of the result is well-formatted, but the actual returned object has access to all attributes of the rw, including the raw path. In any case, the function is fully vectorized and accomplished the random-walk task as an S3 class

```
# Compute a 2D random walk
# Input: number of steps, whether to return the full path
# Output: an (x,y) vector or list of such vector if full path

rw <- function(num.steps = 0, full.path = F){

  # Sanity check on the input
  if(!is.numeric(num.steps) | num.steps <= 0 ){
    stop("Invalid number of steps")
  }

  # Since this is a random walk, just generate the moves
  moves <- list(c(1,0),c(-1,0),c(0,1),c(0,-1))
  steps <- sample(moves,num.steps,replace=T)
  steps.matrix <- matrix(unlist(steps),ncol=2,byrow=T)

  # Sum up the moves and add them
  steps.matrix[,1] <- cumsum(steps.matrix[,1])
  steps.matrix[,2] <- cumsum(steps.matrix[,2])
  colnames(steps.matrix) <- c("x","y")

  walk <- list(origin = c(0,0), full.path = full.path)
  class(walk)<-"rw"
  walk$path <- steps.matrix

  return(walk)
}

# Print a well-formatted representation
print.rw <- function(walk, ...){
  if(walk$full.path){
    print(walk$path)
  }else{
    print(walk$path[length(walk$path)/2,])
  }
}

# Plot the walk on a 2D chart
plot.rw <- function(walk, ...){
  data <- data_frame(x=walk$path[,1],y=walk$path[,2])
  ggplot(data,aes(x=x,y=y))+
    geom_path()
}

# Return the ith step of the random walk
`[.rw` <- function(walk,i,...){
  if(!is.numeric(i) || i <= 0){
```

```

    return(walk$origin)
  }else{
    return(walk$path[i,])
  }
}

# Return the starting coordinate
start <- function(object, ...) UseMethod("start")
start.rw <- function(object){
  return(object$origin)
}

# Set the starting coordinate
`start<-` <- function(object, ...) UseMethod("start<-")
`start<-.rw` <- function(object,value){
  object$origin <- value
  return(object)
}

```

The code below will run through all the requested features for the class `rw`.

```

source("randomWalk.R")

walker <- rw(1000,T)
walker2 <- rw(1000,F)

# First walker will only return the end
# Second walker will print the full path
walker2

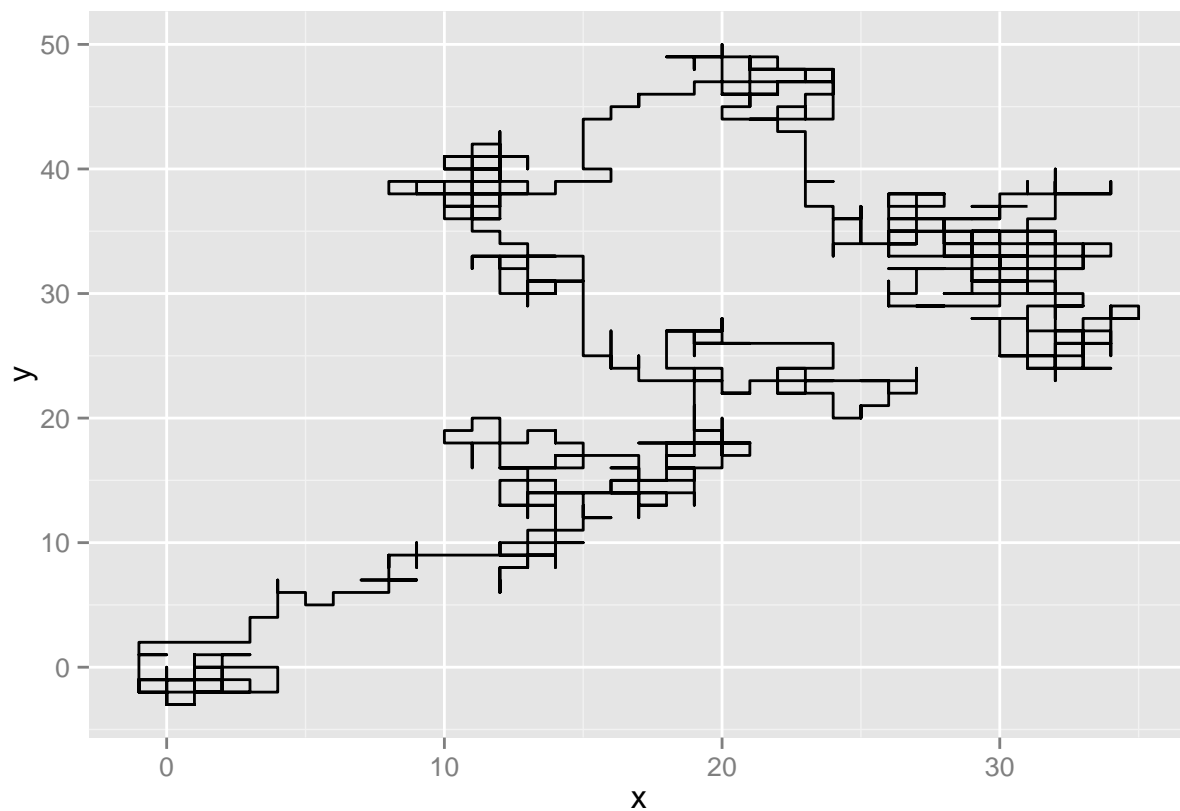
##      x      y
##    -6    -12

head(walker)

##      x      y
## [1,]  1      0
## [2,]  1     -1
## [3,]  2     -1

# Use ggplot2 path geom to chart out the walk
plot(walker)

```



```
# The origin is set to 0,0
# Calling start() allows resetting it
start(walker)
```

```
## [1] 0 0
```

```
start(walker) <- c(5,7)
start(walker)
```

```
## [1] 5 7
```

```
# We can of course get random elements of the walk
walker[10]
```

```
## x y
## 2 -2
```

```
walker[25]
```

```
## x y
## 3 -2
```


5

The **pryr** package again comes in very handy for this part of the assignment. I will be specifically relying on the `object_size` and `mem_used` functions to tell me just how much of the RAM the code is using up.

```
# Create the two variable arrays
object_size(temp.v1 <- sample(seq(1,20,1),1e7,replace=T))
```

```
## 80 MB
```

```
object_size(temp.v2 <- sample(seq(1,20,1),1e7,replace=T))
```

```
## 80 MB
```

```
# Record the initial memory usage
mem_used()
```

```
## 198 MB
```

```
# Load in the provided code
source("fastCount.R")

mem_change(counted <- fastcount(temp.v1,temp.v2))
```

```
## 11.4 kB
```

The actual application of `fastCountOld` does not significantly use up more memory *post execution* since the intermediate objects are all discarded after computation. The big memory costs are in the intermediate arrays; every cast creates a copy, and every logical operation creates a new logical vector.

The execution of `src` and `dummyFun` is negligible as that code is compiled down to C. I do not know the underlying optimizations behind the **inline** package, but goign to assume that `dummyFun` operates with in-memory replacement rather than copying entire arrays. That leaves the code in `fastCountOld`.

The inputs are 10-million element long numeric vectors. At 8 bytes/element, that translates to 80 million bytes, or ~80 Mb as predicted. At this point the function is not modifying them, so no new copies are created. `is.na()` then creates *new* logical vectors with a T/F for each element. Booleans are stored in 4 bytes, so this new logical vector is ~ 40 Mb. Memory usage just spiked up by 80 Mb.

The subsequent two lines are innocuous but very costly, memory wise. `nalineX|nalineY` is a new logical vector, so in addition to computation costs, now there are 40 extra Mb of memory used in subsetting. Once the call ends the memory is recycled, so the two instances of `nalineX|nalineY` don't actually contribute anything to mem usage in the long run. `yvar` and `xvar`, however, now need to be copied from their old versions since there are at least two references to them (in global env and in function's env) so in-memory replacement is not allowed. Combined lines 3&4 thus spike the memory usage by 160 Mb + 40 Mb of the intermediate logic vector.

The logic negation on line 5 is another new vector, and now it's permanent for the duration of the function environment. Since the previous ORs are garbage collected, `useline` requires another recalculation. Final tally is 40 Mb for `useline` + 40 Mb for the intermediate logic vector. At this point the function has generated 280 Mb not including the garbage-collected temporaries.

`tablex` and `tabley` are short numeric vector, so as mentioned in the question prompt, they are negligible compared to the large vectors. For that same reason lines 6,7,10, and 11 (the ones setting `tablex` and `tabley` and the ones accessing `xuse` and `xnames`) are negligible as well.

The line that executes `dummyFun`, however, is very costly. Every array is cast with `as.integer()`, and since integers are only 4 bytes long, this downcasting results in a copy made for every array. `xvar`, `yvar` and `useline` are the large vectors, so this call creates three additional 40 Mb vectors, bringing the total to ~ 400 Mb.

The last three lines starting with `rbind` operate on 21-element vectors, so while there is copying involved in subsetting `tablex` and `tabley`, this impact is negligible as well.

Below is a copy-paste of `fastcount` contents with interleaved memory checks.

```
xvar <- temp.v1
yvar <- temp.v2

start_mem <- mem_used()

nalineX <- is.na(xvar)
nalineY <- is.na(yvar)

mem_used() - start_mem
```

80 MB

```
xvar[nalineX | nalineY] <- 0
yvar[nalineX | nalineY] <- 0

mem_used() - start_mem
```

240 MB

```
useline <- !(nalineX | nalineY);

mem_used() - start_mem
```

280 MB

```
tablex <- numeric(max(xvar)+1)
tabley <- numeric(max(yvar)+1)

mem_used() - start_mem
```

280 MB

```
stopifnot(length(xvar) == length(yvar))
res <- dummyFun(
  tablex = as.integer(tablex), tabley = as.integer(tabley),
  as.integer(xvar), as.integer(yvar), as.integer(useline),
  as.integer(length(xvar)))

mem_used() - start_mem
```

400 MB

```
xuse <- which(res$tablex > 0)
xnames <- xuse - 1
resb <- rbind(res$tablex[xuse], res$tabley[xuse])
colnames(resb) <- xnames

mem_used() - start_mem
```

400 MB

Since casting is so expensive and dummyFun takes in only integers, it makes sense to keep everything as an integer from the get-go. This creates an immediate copy of xvar and yvar but halves the memory for each one and doesn't require casting later down the line.

```
# This really should be the first check
stopifnot(length(temp.v1) == length(temp.v2))

start_mem <- mem_used()

# Coerse both variables into integers early on
xvar <- as.integer(temp.v1)
yvar <- as.integer(temp.v2)

mem_used() - start_mem
```

-80 MB

```
# Multiplying two logical vectors yields an int vector
# If we are using useline to select lines to process
# There is technically no need to set the values to 0
useline <- (!is.na(xvar)) * (!is.na(yvar))

mem_used() - start_mem
```

-80 MB

```
# Not exactly certain why we have different tables
# If they end up different lengths, subsetting by
# xuse might break later on
tablex <- integer(max(xvar) + 1)
tabley <- integer(max(yvar) + 1)

mem_used() - start_mem
```

-80 MB

```
# Everything is already a conditioned integer, so can remove checks
res <- dummyFun(tablex, tabley, xvar, yvar, useline, length(xvar))

mem_used() - start_mem
```

-80 MB

```

# The lower section was actually already pretty good
# Use cbind instead since R is column major
xuse <- which(tablex > 0)
resb <- cbind((xuse-1),res$tablex[xuse],res$tabley[xuse])

mem_used() - start_mem

```

```
## -80 MB
```

6

First things first: let's benchmark how long this code takes to complete.

```
microbenchmark(source("oldMaxLikelihood.R"),times=1)
```

```

## Unit: seconds
##              expr      min       lq      mean   median      uq
## source("oldMaxLikelihood.R") 65.12758 65.12758 65.12758 65.12758 65.12758
##      max neval
## 65.12758     1

```

Yikes, slow. Here's the new and improved version:

```

load('ps4prob6.Rda') # should have A, n, K

ll <- function(Theta, A) {
  sum.ind <- which(A==1, arr.ind=T)
  logLik <- sum(log(Theta[sum.ind])) - sum(Theta)
  return(logLik)
}

oneUpdate <- function(A, n, K, theta.old, thresh = 0.1) {

  Theta.old <- theta.old %*% t(theta.old)
  L.old <- ll(Theta.old, A)

  # Allocate q
  q <- as.numeric(NA)
  length(q) <- n*n*K
  dim(q) <- c(n,n,K)

  #
# Unfortunately didn't have enough time to resolve this :-(
#

  for (i in 1:n) {
    for (j in 1:n) {
      for (z in 1:K) {
        if (theta.old[i, z]*theta.old[j, z] == 0){
          q[i, j, z] <- 0
        } else {

```

```

        q[i, j, z] <- theta.old[i, z]*theta.old[j, z] / Theta.old[i, j]
      }
    }
  }
}

for (z in 1:K) {
  theta.old[,z] <- rowSums(A*q[, ,z])/sqrt(sum(A*q[, ,z]))
}

#
#
#

L.new <- ll(theta.old %*% t(theta.old), A)

converge.check <- abs(L.new - L.old) < thresh

return(list(theta = theta.old/rowSums(theta.old), loglik = L.new, converged = converge.check))
}

# initialize the parameters at random starting values
temp <- runif(n*K)
dim(temp) <- c(n,K)
theta.init <- temp/rowSums(temp)

# do single update
out <- oneUpdate(A, n, K, theta.init)

microbenchmark(source("newMaxLikelihood.R"),times=1)

## Unit: seconds
##           expr      min       lq     mean  median      uq
##  source("newMaxLikelihood.R") 65.0127 65.0127 65.0127 65.0127 65.0127
##      max neval
## 65.0127      1

```