

# exercises7

Nick Gembs

3/7/2023

## Statistical Computing - Exercises 07 - Performance

We return to two problems and add a third: calculating song statistics in the billboard dataset, calculating distances in the flights dataset, and reformatting a hurricane track dataset called hurdat. Below are correct solutions to all three problems. Your task is to make them run faster, while ensuring that your faster code gives the correct result.

Show how long your revised code takes to run, and demonstrate that it gives the right answer.

### Song Statistics

#### My Answer

```
hot <- read.csv("C:/Users/Nick/Downloads/Hot_100.csv")

weeks_below_k <- function( dat, k ){

  dat = dat[dat$chart_position <= k,]
  h <- tapply( dat$instance, dat$song_id, length )

  h = as.data.frame(h)
  h <- cbind(newColName = rownames(h), h)
  rownames(h) <- NULL
  colnames(h) = (c("song_id", "weeks_below_k"))
  h$weeks_below_k <- as.double(h$weeks_below_k)

  songs = data.frame(unique(hot$song_id))
  colnames(songs) = "song_id"

  h = merge(h, songs, by = "song_id", all.y = T)

  h[is.na(h)] = 0

  h = h[
    with(h, order(weeks_below_k, song_id, decreasing = c(T,F),
      method = "radix")),
  ]
}
```

```

]

rownames(h) = NULL

return(h)
}

t1 = proc.time()
res <- weeks_below_k( hot, 1 )
t2 <- proc.time()
print(t2-t1)

```

```

##      user  system elapsed
##    0.46    0.02    0.47

```

## Old Answer

```

weeks_below_k <- function( dat, k ){
  # computes the number of weeks on chart at or below ranking k
  # for each songid in dat
  # get the vector of unique song ids
  unique_songs <- unique( dat$song_id )
  # initialize vector for number of weeks at or below ranking k
  weeks_below_k <- c()
  # loop over all the songs
  for(j in 1:length( unique_songs ) ){
    # find the indices for this song
    inds_song <- which( dat$song_id == unique_songs[j] )
    # initialize a count of number of weeks at or below k
    count <- 0
    # loop over the indices
    for( i in inds_song ){
      # add to the count if the week_position is at or below k
      if( dat$chart_position[i] <= k ){
        count <- count + 1
      }
    }
    # update weeks_below_k with this song's count
    weeks_below_k <- c( weeks_below_k, count )
  }
  # rank the songs
  ord <- order( weeks_below_k, decreasing = TRUE )
  return( data.frame( song_id = unique_songs[ord],
                     weeks_below_k = weeks_below_k[ord] ) )
}

t1 <- proc.time()
r <- weeks_below_k( hot, 1 )
t2 <- proc.time()

```

```
print(t2-t1)
```

```
##      user  system elapsed  
##  104.96    14.87 39605.43
```

```
identical(r,res)
```

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

## Distance matrix for airline data

### My Answer

```
dat <- read.csv("C:/Users/Nick/Documents/GitHub/statcomp2023/datasets/airline_2019-07-01.csv")
```

```
get_distance_matrix <- function( dat ){  
  mat = matrix(data = 0, nrow = length(unique(dat$Origin)),  
               ncol = length(unique(dat$Origin)),  
               dimnames = list(unique(dat$Origin), unique(dat$Origin)))  
  
  mat1 = tapply(dat$Distance , list(dat$Origin, dat$Dest) , mean )  
  
  return(mat1)  
}
```

```
t1 <- proc.time()  
distmat <- get_distance_matrix( dat )  
t2 <- proc.time()  
print(t2-t1)
```

```
##      user  system elapsed  
##    0.19    0.00    0.22
```

```
v <- c("ATL", "ORD", "LGA", "JFK", "DEN", "LAX", "STL", "SEA")  
distmat[v,v]
```

```
##      ATL  ORD  LGA  JFK  DEN  LAX  STL  SEA  
## ATL   NA  606  762  760 1199 1947  484 2182  
## ORD  606   NA  733  740  888 1744  258 1721  
## LGA  762  733   NA   NA 1620   NA  888   NA  
## JFK  760  740   NA   NA 1626 2475   NA 2422  
## DEN 1199  888 1620 1626   NA  862  770 1024  
## LAX 1947 1744   NA 2475  862   NA 1592  954  
## STL  484  258  888   NA  770 1592   NA 1709  
## SEA 2182 1721   NA 2422 1024  954 1709   NA
```

```
new = distmat
```

## Old Answer

```
# distance matrix

get_distance_matrix <- function( dat ){

  # get the set of unique airports
  ports <- sort( unique( c(dat$Origin, dat$Dest) ) )
  nports <- length(ports)

  # set up a distance matrix
  distmat <- matrix(NA, nports, nports )
  rownames(distmat) <- ports
  colnames(distmat) <- ports

  # loop over all possible origins and destinations
  for(p1 in ports){
    for(p2 in ports){

      # get row and column indices for this pair
      j1 <- which( rownames(distmat) == p1 )
      j2 <- which( colnames(distmat) == p2 )

      # get rows of data frame for this pair, and subset
      ii <- which( dat$Origin == p1 & dat$Dest == p2 )
      subdat <- dat[ii,]

      # find the distance
      distmat[j1,j2] <- subdat$Distance[1]

    }
  }
  return(distmat)
}
```

```
t1 <- proc.time()
distmat <- get_distance_matrix( dat )
t2 <- proc.time()
print(t2-t1)
```

```
##      user  system elapsed
## 189.95    6.56   218.47
```

```
v <- c("ATL", "ORD", "LGA", "JFK", "DEN", "LAX", "STL", "SEA")
distmat[v,v]
```

```
##      ATL  ORD  LGA  JFK  DEN  LAX  STL  SEA
```

```
## ATL    NA   606   762   760 1199 1947   484 2182
## ORD   606    NA   733   740   888 1744   258 1721
## LGA   762   733    NA    NA 1620    NA   888    NA
## JFK   760   740    NA    NA 1626 2475    NA 2422
## DEN 1199   888 1620 1626    NA   862   770 1024
## LAX 1947 1744    NA 2475   862    NA 1592   954
## STL   484   258   888    NA   770 1592    NA 1709
## SEA 2182 1721    NA 2422 1024   954 1709    NA
```

```
all.equal(new,distmat)
```

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

## Hurricane data

### My Answer

```
# read in the raw data
dat <- read.csv("C:/Users/Nick/Documents/GitHub/statcomp2023/datasets/hurdat2-1851-2021-041922.txt", header=TRUE)

t1<- proc.time()

date = dat[startsWith(dat$V1, "AL"),]
hurdata = dat[!(startsWith(dat$V1, "AL")),]
date = rep(date$V1, date$V3)
hurdata = cbind(date,hurdata)
rownames(hurdata) <- NULL

t2 <- proc.time()

print(t2-t1)
```

```
##      user   system elapsed
##    0.11    0.02    0.14
```

```
head(hurdata)
```

```
##      date      V1      V2 V3  V4      V5      V6 V7  V8  V9  V10  V11  V12
## 1 AL011851 18510625 0000    HU 28.0N  94.8W 80 -999 -999 -999 -999 -999
## 2 AL011851 18510625 0600    HU 28.0N  95.4W 80 -999 -999 -999 -999 -999
## 3 AL011851 18510625 1200    HU 28.0N  96.0W 80 -999 -999 -999 -999 -999
## 4 AL011851 18510625 1800    HU 28.1N  96.5W 80 -999 -999 -999 -999 -999
## 5 AL011851 18510625 2100 L  HU 28.2N  96.8W 80 -999 -999 -999 -999 -999
## 6 AL011851 18510626 0000    HU 28.2N  97.0W 70 -999 -999 -999 -999 -999
##      V13  V14  V15  V16  V17  V18  V19  V20  V21
## 1 -999 -999 -999 -999 -999 -999 -999 -999 -999
## 2 -999 -999 -999 -999 -999 -999 -999 -999 -999
## 3 -999 -999 -999 -999 -999 -999 -999 -999 -999
## 4 -999 -999 -999 -999 -999 -999 -999 -999 -999
## 5 -999 -999 -999 -999 -999 -999 -999 -999 -999
## 6 -999 -999 -999 -999 -999 -999 -999 -999 -999
```

## Old Answer

```
t1<- proc.time()
# initialize the processed dataset
hurdat <- data.frame( matrix(NA, 0, ncol(dat)+1) )
colnames(hurdat) <- c("date", colnames(dat))

# counter for the row of hurdat
k <- 0

# loop over rows of raw dataset
for(j in 1:nrow(dat)){

  # extract the current row of raw data
  this_row <- dat[j,]

  # check whether this is a code row
  if( substr( this_row[1,1], 1, 2 ) == "AL" ){

    # if so, update the hurricane code
    hur_code <- this_row[1,1]

  } else {

    # otherwise update the counter and write to the next row
    k <- k + 1
    hurdat[k,] <- cbind( hur_code, this_row )

  }
}

row.names(hurdat) = NULL
t2 <- proc.time()

print(t2-t1)
```

```
##      user  system elapsed
## 700.23   17.39   770.40
```

```
identical(hurdat,hurdata)
```

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