# Speeding up R code using Rcpp and foreach packages.

## Pedro Albuquerque

Universidade de Brasília

## February, 8th 2017

**UnB**

**1** Speeding up R.

**2** foreach package.
      Bootstrap.

**3** Rcpp package.
      isPrime function.

**4** Exercises.

**UnB**

R is a good choice as open source and free software.
But R has some disadvantages:

- Limited memory for big datasets.

- Inefficient standard loops.

There are some solutions for both disadvantages, especially to speed up:

- Using apply family functions and vectorized operations.

- Using parallel processing.

- Using C++ (Rcpp package).

**UnB**

**Speeding up R.**      foreach package.      Rcpp package.      Exercises.

○○○○○○      ○○○

R is a good choice as open source and free software.
But R has some disadvantages:

- Limited memory for big datasets.
- Inefficient standard loops.

There are some solutions for both disadvantages, especially to speed up:

- Using apply family functions and vectorized operations.
- Using parallel processing.
- Using C++ (Rcpp package).

**UnB**

R is a good choice as open source and free software.
But R has some <span style="color:red">disadvantages</span>:

- Limited memory for big datasets.
- Inefficient standard loops.

There are some solutions for both <span style="color:red">disadvantages</span>, especially to speed up:

- Using apply family functions and vectorized operations.
- Using parallel processing.
- Using C++ (Rcpp package).

**UnB**

**Speeding up R.**         foreach package.        Rcpp package.        Exercises.

○○○○○○        ○○○

R is a good choice as open source and free software.

But R has some <span style="color:red">disadvantages</span>:

- Limited memory for big datasets.
- Inefficient standard loops.

There are some solutions for both <span style="color:red">disadvantages</span>, especially to speed up:

- Using apply family functions and vectorized operations.
- Using parallel processing.
- Using C++ (Rcpp package).

**UnB**

Speeding up R.                     foreach package.                     Rcpp package.              Exercises.

○○○○○○                           ○○○

R is a good choice as open source and free software.
But R has some <span style="color:red">disadvantages</span>:

- Limited memory for big datasets.
- Inefficient standard loops.

There are some solutions for both <span style="color:red">disadvantages</span>, especially to speed up:

- Using apply family functions and vectorized operations.
- Using parallel processing.
- Using C++ (Rcpp package).

**UnB**

R is a good choice as open source and free software.
But R has some <span style="color:red">disadvantages</span>:

- Limited memory for big datasets.
- Inefficient standard loops.

There are some solutions for both <span style="color:red">disadvantages</span>, especially to speed up:

- Using apply family functions and vectorized operations.
- Using parallel processing.
- Using C++ (Rcpp package).

**UnB**

## foreach package.

Suppose we are interested in estimate a variance of a complicated function of parameters, for example, the coefficient of variation:
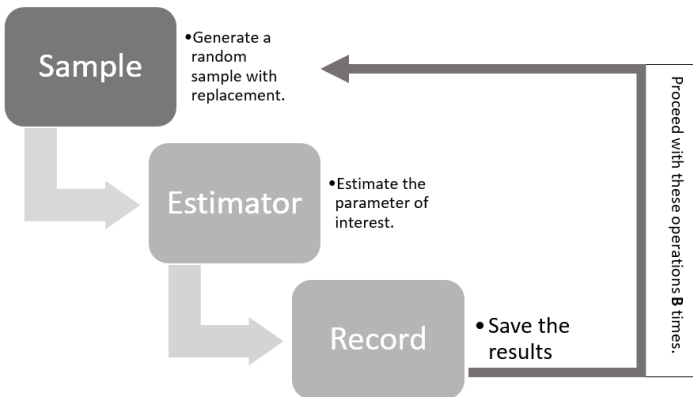
$$\theta = \frac{\mu}{\sigma}$$

where $\mu$ and $\sigma$ are the population mean and standard deviation.

We can estimate the variance of $\theta$ using, for example, the bootstrap estimator.

**UnB**

## foreach package.

Suppose we are interested in estimate a variance of a complicated function of parameters, for example, the coefficient of variation:

$$\theta = \frac{\mu}{\sigma}$$

where $\mu$ and $\sigma$ are the population mean and standard deviation.

We can estimate the variance of $\theta$ using, for example, the bootstrap estimator.

**UnB**

## foreach package.

Suppose we are interested in estimate a variance of a complicated function of parameters, for example, the coefficient of variation:

$$\theta = \frac{\mu}{\sigma}$$

where $\mu$ and $\sigma$ are the population mean and standard deviation.

We can estimate the variance of $\theta$ using, for example, the bootstrap estimator.

**UnB**

Speeding up R.          **foreach package.**          Rcpp package.          Exercises.
                        ●○○○○○                        ○○○

Bootstrap.

# Bootstrap.

Speeding up R.          **foreach package.**          Rcpp package.          Exercises.
                        ○●○○○○○              ○○○

Bootstrap.

# Bootstrap.

| Bootstrap sample (b=1,...B) | Estimator |
| --- | --- |
| 1 | $\hat{\theta}_1 = \hat{\mu}_1/\hat{\sigma}_1$ |
| 2 | $\hat{\theta}_2 = \hat{\mu}_2/\hat{\sigma}_2$ |
| 3 | $\hat{\theta}_3 = \hat{\mu}_3/\hat{\sigma}_3$ |
| . | |
| . | |
| B | $\hat{\theta}_B = \hat{\mu}_B/\hat{\sigma}_B$ |
| **MEAN** | $\widehat{E(\theta)} = \sum \hat{\theta}_b/B$ |
| **VARIANCE** | $\widehat{V(\theta)} = \sum \left(\theta_b - \widehat{E(\theta)}\right)^2/(B-1)$ |

**UnB**

Speeding up R.                    **foreach package.**                    Rcpp package.                    Exercises.
                                   ○○●○○○                                 ○○○

Bootstrap.

## Aplication.

```
1 #Set the seed
2 set.seed(1)
3 #Generate fake data
4 data <- rnorm(n=10000, mean=10, sd=10)
5 #True coefficient of variation
6 CV.true <- 10/10
7 #Estimated coefficient of variation
8 CV.hat <- mean(data)/sd(data)
```

Algorithm 1: Generating fake data.

The estimate is $\hat{\theta} \approx 0.9813$.

**UnB**

Speeding up R.                    **foreach package.**              Rcpp package.              Exercises.
                                  000●00                           000

Bootstrap.

# Aplication.

```
1  #Create the bootstrap vector
2  result <- rep(NA, 100000)
3  #Start the clock
4  ptm <- proc.time()
5  for(b in 1:100000){
6    #Generate the random sample with replacement
7    ids <- sample(x=length(data),size=length(data),replace
       =T)
8    sample <- data(ids)
9    #Calculate the CV
10   result(b) <- mean(sample)/sd(sample)
11 }
12 #Stop the clock - Time elapsed 20.40.
13 proc.time() - ptm
```

Algorithm 2: Time elapsed 20.40.

Speeding up R.                    **foreach package.**                        Rcpp package.                  Exercises.
                                  ○○○○●○                                        ○○○

Bootstrap.

# Aplication.

```
1  library(foreach)
2  library(doParallel)
3  #See how many cores we have
4  ncl<-detectCores()
5  #Register the cores
6  cl <- makeCluster(ncl)
7  registerDoParallel(cl)
```

Algorithm 3: Using foreach package.

**UnB**

Speeding up R.    **foreach package.**    Rcpp package.    Exercises.
○○○○○●    ○○○

Bootstrap.

# Aplication.

```
1  #Start the clock
2  ptm <- proc.time()
3  #Create the bootstrap vector
4  result <- rep(NA, 100000)
5  #Bootstrap using foreach
6  result <- foreach(b=1:100000, .combine=c) %dopar% {
7      #Generate the random sample with replacement
8      ids <- sample(x=length(data),size=length(data),replace=T)
9      sample <- data(ids)
10     #Calculate the CV
11     mean(sample)/sd(sample)
12  }
13  #Stop the clock - Time elapsed 5.32.
14  proc.time() - ptm
15  #Release the cores
16  stopCluster(cl)
```

Algorithm 4: Time elapsed 5.32.

**UnB**

# Rcpp package.

Another solution is to use the Rcpp package:



**UnB**

# Rcpp package.

We need to compile the functions before use in R:

# isPrime function.

Suppose we are interested in find if a number is prime or not.

```
1  #Create the isPrime function
2  isPrime <- function(num){
3    prime <- TRUE
4    den <- num -1
5    while(prime==TRUE & den >1){
6      #Remainder of the division
7      if(num%%den==0){
8        prime <- FALSE
9      }
10     #Decrease the number
11     den <- (den - 1)
12   }
13   return(prime)
14 }
15 #Example isPrime(12)
```

Algorithm 5: isPrime function.

**UnB**

# isPrime function.

Using the same idea with the Rcpp:

```cpp
#include <Rcpp.h>
using namespace Rcpp;
// ((Rcpp::export))
bool isPrimeCpp(int num) {
   bool prime = true;
   int den = num -1;
   while(prime==true & den >1){
   //Test if is divided
     if(num % den==0){
        prime = false;
     }
 //Decrease the number
    den = (den - 1);
   }
   return(prime);
}
```

# isPrime function.

You can save your functions in a *cpp* file and invoke using:

```
1  #Rcpp library
2  library(Rcpp)
3  #Call the functions
4  sourceCpp("MyFile.cpp")
5  #Example
6  isPrimeCpp(12)
```

Algorithm 6: isPrime function.

**UnB**

Speeding up R.        foreach package.        Rcpp package.        **Exercises.**

oooooo        ooo

- Create a Rcpp function to calculate the sample mean of a vector.

- Create a Rcpp function to calculate the $k$ Fibonacci numbers.

- Use foreach and the bootstrap technique to estimate an arbitrary parameter function: $g(\theta) = \sqrt{\theta}/2$.

**UnB**