

# (A bit of) Advanced R

Efficient R programming

Julien Chiquet

Université Paris Dauphine

Juin 2018

<http://github/jchiquet/CourseAdvancedR>

# Resources

- Gillespie & Lovelace (2016): efficient R programming
- Wickham (2014)

## Part 0: Prerequisites

- xply family, do.call, Reduce

# Outline

- 1 Benchmark your code
- 2 Use all your cores when needed
- 3 Part 2: Use multiple cores in your simulations
- 4 Part 3: Be aware of what R is good at
- 5 Part 4: Remember that R is object oriented
- 6 Part 4: Interface with lower-level languages

## Quick (and dirty) benchmarking with `system.time()`

One usually relies on the command `system.time(expr)` to evaluate the timings:

```
func.one <- function(n) {return(rnorm(n,0,1))}  
func.two <- function(n) {return(rpois(n,1))}  
  
n <- 1000  
system.time(replicate(100, func.one(n)))
```

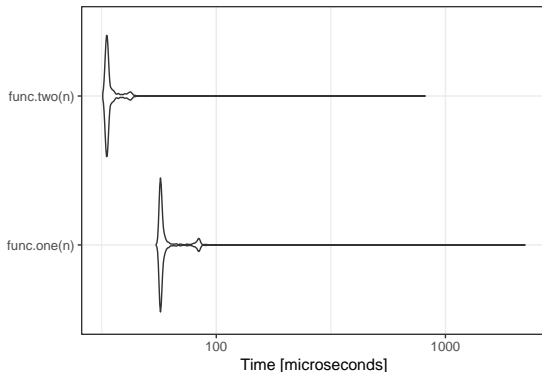
```
##      user  system elapsed  
##  0.012   0.000   0.010
```

```
system.time(replicate(100, func.two(n)))
```

```
##      user  system elapsed  
##  0.008   0.000   0.007
```

# Quick benchmarking with microbenchmark

```
func.one <- function(n) {return(rnorm(n,0,1))}  
func.two <- function(n) {return(rpois(n,1))}  
  
library(microbenchmark)  
  
n <- 1000  
res <- microbenchmark(func.one(n), func.two(n), times=1000)  
ggplot2::autoplot(res)
```



# Profile your code

Suppose you want to evaluate which part of the following function is hot:

```
## generate data, center/scale and perform ridge regression
my_func <- function(n,p) {

  require(MASS)

  ## draw data
  x <- matrix(rnorm(n*p),n,p)
  y <- rnorm(n)

  ## center/scale
  xs <- scale(x)
  ys <- y - mean(y)

  ## return ridge's coefficients
  ridge <- lm.ridge(ys~xs+0,lambda=1)

  return(ridge$coef)
}
```

# Profile your code with base Rprof |

One can rely on the default Rprof function, with somewhat technical outputs

```
Rprof(file="profiling.out", interval=0.05)
res <- my_func(1000,500)
Rprof(NULL)
```

```
summaryRprof("profiling.out")$by.self
```

##	self.time	self.pct	total.time	total.pct
## "La.svd"	1.05	77.78	1.10	81.48
## "matrix"	0.15	11.11	0.15	11.11
## "aperm.default"	0.05	3.70	0.05	3.70
## "apply"	0.05	3.70	0.05	3.70
## "[.data.frame"	0.05	3.70	0.05	3.70

```
summaryRprof("profiling.out")$by.total
```



# Profile your code with base Rprof II

##	total.time	total.pct	self.time	self.pct
## "block_exec"	1.35	100.00	0.00	0.00
## "call_block"	1.35	100.00	0.00	0.00
## "eval"	1.35	100.00	0.00	0.00
## "evaluate"	1.35	100.00	0.00	0.00
## "evaluate_call"	1.35	100.00	0.00	0.00
## "evaluate::evaluate"	1.35	100.00	0.00	0.00
## "FUN"	1.35	100.00	0.00	0.00
## "handle"	1.35	100.00	0.00	0.00
## "in_dir"	1.35	100.00	0.00	0.00
## "knit"	1.35	100.00	0.00	0.00
## "knitr::knit"	1.35	100.00	0.00	0.00
## "lapply"	1.35	100.00	0.00	0.00
## "my_func"	1.35	100.00	0.00	0.00
## "process_file"	1.35	100.00	0.00	0.00
## "process_group"	1.35	100.00	0.00	0.00
## "process_group.block"	1.35	100.00	0.00	0.00
## "rmarkdown::render"	1.35	100.00	0.00	0.00
## "timing_fn"	1.35	100.00	0.00	0.00
## "withCallingHandlers"	1.35	100.00	0.00	0.00
## "withVisible"	1.35	100.00	0.00	0.00
## "lm.ridge"	1.15	85.19	0.00	0.00
## "La.svd"	1.10	81.48	1.05	77.78
## "svd"	1.10	81.48	0.00	0.00
## "matrix"	0.15	11.11	0.15	11.11
## "scale"	0.10	7.41	0.00	0.00
## "scale.default"	0.10	7.41	0.00	0.00

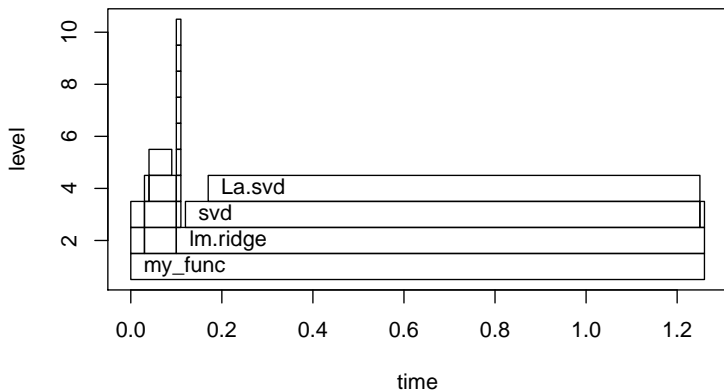
# Profile your code with base Rprof III

## "aperm.default"	0.05	3.70	0.05	3.70
## "apply"	0.05	3.70	0.05	3.70
## "[.data.frame"	0.05	3.70	0.05	3.70
## "["	0.05	3.70	0.00	0.00
## "aperm"	0.05	3.70	0.00	0.00
## "eval.parent"	0.05	3.70	0.00	0.00
## ".External2"	0.05	3.70	0.00	0.00
## "model.frame.default"	0.05	3.70	0.00	0.00
## "na.omit"	0.05	3.70	0.00	0.00
## "na.omit.data.frame"	0.05	3.70	0.00	0.00
## "stats::model.frame"	0.05	3.70	0.00	0.00
## "sweep"	0.05	3.70	0.00	0.00

# Profile your code with profr

The *profr* package is maybe a little easier to understand...

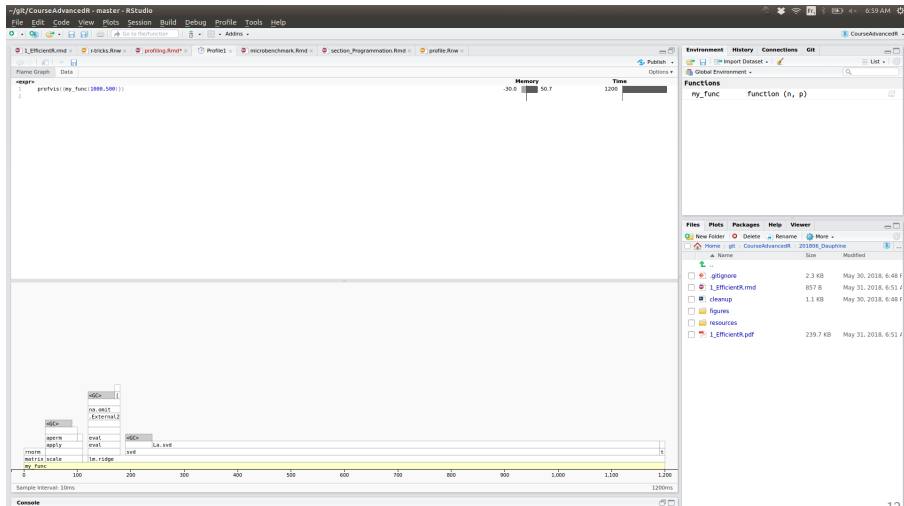
```
library(profr)
profiling <- profr({my_func(1000,500)}), interval = 0.01)
plot(profiling)
```



# Profile your code within R Studio with profvis

Profvis integrates the profiling to the Rstudio API

```
library(profvis)
profvis({my_func(1000,500)})
```



# Outline

- 1 Benchmark your code
- 2 Use all your cores when needed
- 3 Part 2: Use multiple cores in your simulations
- 4 Part 3: Be aware of what R is good at
- 5 Part 4: Remember that R is object oriented
- 6 Part 4: Interface with lower-level languages

# Parallel computing

## Usual Roadmap

- ① Start up and initialize  $M$  'worker' processes
- ② Send data required for each task to the workers
- ③ Split the task into  $M$  roughly equally-sized chunks and send them (including the R code needed) to the workers
- ④ Wait for all the workers to complete their tasks, and ask them for their results
- ⑤ Repeat steps (2–4) for any further tasks
- ⑥ Shut down the worker processes

# Socketing vs Forking

Two approaches achieving the same goal

## The socket approach

- launches a new version of R on each core
- connection is done via networking all happening on your own computer

## The forking approach

- copies the entire current version of R and moves it to a new core
- several processes acheive the same task resulting in different outputs

~> Forking is only possible on Unix systems (Linux, Mac OS)

# Parallel computing with parallel

## Package parallel

- merge of packages multicore and snow
- included in base R and maintained by the R Core team

## Check your computer

```
library(parallel) ## embedded with R since version 2.9 or something
cores <- detectCores() ## How many cores do I have?
print(cores)
```

```
## [1] 4
```

↪ parallel features both socketing (parLapply) and forking (mclapply)



# Forking approach with `parallel::mclapply`

Very easy: use parallel features as soon as you do simulations !

Example: estimates the test error from ridge regression

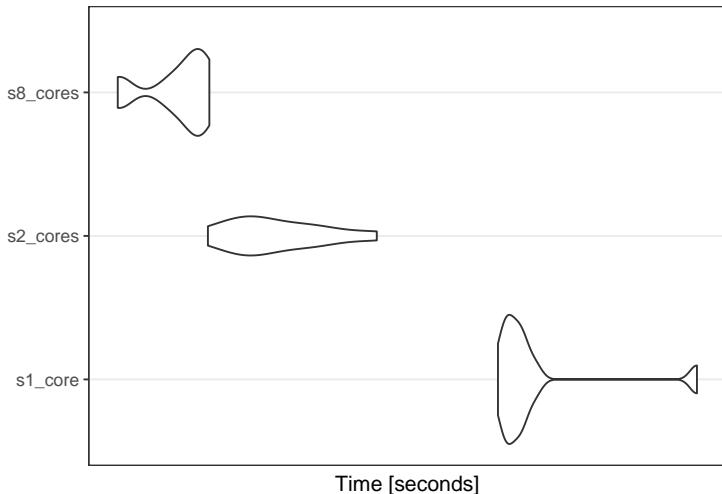
```
one.simu <- function(i) {  
  ## draw data  
  n <- 1000; p <- 500  
  x <- matrix(rnorm(n*p),n,p) ; y <- rnorm(n)  
  ## return ridge's coefficients  
  train <- 1:floor(n/2)  
  test  <- setdiff(1:n,train)  
  ridge <- MASS::lm.ridge(y~x+0,lambda=1,subset=train)  
  err <- (y[test] - x[test, ] %*% ridge$coef )^2  
  return(list(err = mean(err), sd = sd(err)))  
}
```

```
head(do.call(rbind, mclapply(1:8, one.simu, mc.cores = cores)), n = 3)
```

```
##      err      sd  
## [1,] 14.45796 19.79301  
## [2,] 10.3222  14.86089  
## [3,] 12.15729 16.8608
```

# Forking approach with `parallel::mclapply` (cont'd)

```
library(microbenchmark)
res <- microbenchmark(s1_core = mclapply(1:8, one.simu, mc.cores = 1),
                      s2_cores = mclapply(1:8, one.simu, mc.cores = 2),
                      s8_cores = mclapply(1:8, one.simu, mc.cores = 8), times = 10)
```



# Socket approach with `parallel::parLapply`

Windows users need a bit more code to make it work

A possible option: export from base workspace

```
cl <- makeCluster(4)
clusterExport(cl, "one.simu")
res <- parSapply(cl, 1:8, one.simu) # several parLapply call are possible
stopCluster(cl)
res
```

```
##      [,1]      [,2]      [,3]      [,4]      [,5]      [,6]      [,7]
## err 11.86539 11.26563 8.845916 14.37104 12.91292 13.25386 12.50689
## sd  16.6455  15.63073 12.32548 18.74685 18.49328 17.62274 16.61861
##      [,8]
## err 8.508118
## sd  12.48525
```

# Parallel computing with parallel: final remarks

- Parallelize piece of code complex enough
- Do not choose stupidly the number of cores
- Screen outputs are lost in Rstudio: use pbmcapply (progress bar)

```
pbmcapply::pbmcapply(1:8, FUN = one.simu, mc.cores = 2)
```

```
##      [,1]      [,2]      [,3]      [,4]      [,5]      [,6]      [,7]      [,8]  
## err 10.72632 14.25714 9.77985  9.589137 10.90891 11.8773  10.73228 11.2432  
## sd  13.99144 18.25551 13.08587 14.26683 15.5244  15.03568 16.37883 15.9956
```

# Outline

- 1 Benchmark your code
- 2 Use all your cores when needed
- 3 Part 2: Use multiple cores in your simulations
- 4 Part 3: Be aware of what R is good at
- 5 Part 4: Remember that R is object oriented
- 6 Part 4: Interface with lower-level languages

# Outline

- 1 Benchmark your code
- 2 Use all your cores when needed
- 3 Part 2: Use multiple cores in your simulations
- 4 Part 3: Be aware of what R is good at
- 5 Part 4: Remember that R is object oriented
- 6 Part 4: Interface with lower-level languages

# Outline

- 1 Benchmark your code
- 2 Use all your cores when needed
- 3 Part 2: Use multiple cores in your simulations
- 4 Part 3: Be aware of what R is good at
- 5 Part 4: Remember that R is object oriented
- 6 Part 4: Interface with lower-level languages

# Outline

- 1 Benchmark your code
- 2 Use all your cores when needed
- 3 Part 2: Use multiple cores in your simulations
- 4 Part 3: Be aware of what R is good at
- 5 Part 4: Remember that R is object oriented
- 6 Part 4: Interface with lower-level languages



# References

Gillespie, C., & Lovelace, R. (2016). *Efficient r programming*. " O'Reilly Media, Inc." Retrieved from <https://bookdown.org/csgillespie/efficientR/>

Wickham, H. (2014). *Advanced r*. CRC Press. Retrieved from <http://adv-r.had.co.nz/>