# Parellel Computing in R
## from `parallel` to `foreach` and `future`

**Chao Cheng**

*School of Statistics and Management*
**Shanghai University of Finance and Economics, China**

November 22, 2022

# Outline

# Extention of Rcpp

## RcppProgress

○ Article: Using RcppProgress to control the long computations in C++

```
# Rcpp::sourceCpp("rcpp_progress.cpp")
# long_computation(3000)
# 0%   10   20   30   40   50   60   70   80   90   100%
# [----|----|----|----|----|----|----|----|----|----|
# **************************************************|
# [1] 3002.32
```

# Extention of Rcpp

## RcppParallel

# Introduction

## A toy example

- $1 + 2 + \cdots + 100$
- Compute the sum of a vector

## Abstraction

- The whole job can be break into small parts and they can be done independently of each other.
- Map + Reduce

## Useful cases

- Simulation
- Bootstrap, MCMC and cross validation, etc.
- Elementwisely update an vector in ADMM algorithm (Parallel in C++)

# Basic paralle computation for simulation

- Start multiple R sessions
- Preparation: load necessary packages, etc.
- Run simulation scripts, possibly according to session ID.
- Collect and summary the results by hand.

## Abstraction
- Create workers
- Prepare workers
- Run script in parallel and collect the results.

# The `parallel` package

- It's derived from `snow` and `multicore` packages.
- Useful reference:
    - Parallel R. (This book is a bit old.)
    - `parallel`'s documentations.
    - `parallel`'s vignettes.

# A simple template

```r
library(parallel)
# use all the cores of this machine
cls <- makeCluster(detectCores())

# initializing workers
clusterEvalQ(cls, fun)
# pass VARLIST from master to all the workers
clusterExport(cls, VARLIST)


# split full index of all tasks to workers
idx_split <- clusterSplit(cls, idx_full)
# carry out the task parFUNCTION parallely
res <- parLapply(cls, idx_split, parFUNCTION)

# stop workers
stopCluster(cls)
```

# A simple example

```r
library(parallel)
a <- rnorm(12)
slow_function <- function(invec){
    ...    # a slow function
}

cls <- makeCluster(4)
ind_seq <- clusterSplit(cls, a)
clusterExport(cls, varlist = "slow_function")
res_par <- parSapply(cls, ind_seq, slow_function)
res <- sum(res_par)
```

# PSOCK vs FORK

```
a <- rnorm(100)
```

### PSOCK
```
cls <- makeCluster(4)     # default type is PSOCK
```

### FORK
```
cls <- makeCluster(4, type = "FORK")     # NOT available on Windows
```

```
parSapply(cls, 1 : 10, function(id){
    return(a[id])
})
```

# PSOCK vs FORK

## PSOCK

- Pros:
  - Use socket connection, a general approach.
  - All system, locally or remotely with suitable setup such as MPI
- Cons:
  - Might be hard to configure.
  - Manually transport the data.

## FORK

- Pros: use FORK mechanism, no worry about variable transportation.
- Cons: Only for one machine, not available on Windows.

# Parallel random number generation

- Manually use `set.seed()` on every worker.
- Use 'L'Ecuyer-CMRG' multiple RNG stream.
    1. `RNGkind("L'Ecuyer-CMRG")` on your main session.
    2. `set.seed()` on your main session.
    3. `clusterSetupRNGstream()` to set your workers' seed.

- Explicit functions and variables will always be transported.
- FORK will copy the main session at creation.
- Others should be taken care of by hand.
- Additional configuration of `clusterExport` when nested in a function call.

# Dark time of parallel computation

There are so many different parallel backends:

- snow
- multicore
- parallel
- MPI
- Redis
- Hadoop
- Spark
- Slurm
- ...

How to support them? How to maintain code?

# Foreach

`foreach` defines a simple but powerful framework for map/reduce parallel computation.

## Package author/code writer

Decide which part of code can run in parallel.

## End user

Decide how to run in parallel based on their available resources.

`foreach` is syntactically structured in the form of a for loop.

# Foreach

```r
library(foreach)
# library(doParallel)
# registerDoParallel()
a <- 10
foreach(i = 1 : 12, j = 12 : 1, .combine = rbind) %dopar%{
    Sys.sleep(0.5)
    # print will be dropped when run in parallel
    print(paste("i = ", i, ", j = ", j, sep = ""))
    data.frame(i, j, a)
}
```

# future and future.apply

- `future` provides a simple and uniform way of evaluating R expressions asynchronously using various resources available to the user.
- `future.apply` provides worry-free parallel alternatives to base-R "apply" functions.

```r
library(future.apply)    # default plan is sequential
# plan(cluster)
x <- rnorm(16)
future_lapply(1 : 5, function(id){
    print(paste("id = ", id, sep = ""))    # normal print kept
    Sys.sleep(0.5)
    sum(x[1 : id])
})
```

# future

- Asynchronous computation. Not constrained by a for-loop or apply syntax.
- Available extensions:
  - `future.apply`
  - `doFuture`: backends for `foreach`, `BiocParallel` and `plyr`.
  - `furrr`

# future

```r
library(future)
plan(cluster)

x <- future({
    x <- matrix(rnorm(10 ^ 6), nrow = 10 ^ 3)
    for(i in 1 : 5){
        print(paste("i = ", i))
        res <- eigen(x)
    }
    return(res)
}, seed = T)      # not block the main session

resolved(x)      # check whether the future is resolved
a <- rnorm(10)      # we can do other stuff at the main session
```

# Personal suggestions

- Nested parallel is NOT recommended. At least it should be done with careful configuration.
- `future.apply` vs `foreach`
  - Familiar with `foreach`: just use the `doFuture` backends.
  - New to parallel: `future.apply` is a good start point for your code.
  - `future` backend will relay the printed messages.
  - Performance in parallel are close, so-called.
  - Performance for sequential are slower than for-loop.

# I want progress bars

- `RcppProgress` allows to display a progress bar in the R console for long running computations taking place in c++ code, supports OpenMP.
- `pbapply` is a lightweight package that adds progress bar to vectorized R functions ('*apply'). It supports several parallel backends.
- `progress` shows ASCII progress bars.
- `progressr` provides a minimal API for reporting progress updates in R.
  - Developer is responsible for providing progress updates.
  - End user decides if, when, and how progress should be presented.

# progressr

```r
library(progressr)
slow_sum <- function(x) {
    p <- progressr::progressor(along = x)
    sum <- 0
    for (kk in seq_along(x)) {
        Sys.sleep(0.5)
        sum <- sum + x[kk]
        p(message = sprintf("Added %g", x[kk]))
    }
    sum
}
# handlers("default")    # default handler is "txtprogressbar"
with_progress(y <- slow_sum(1:10))
handlers("progress")
with_progress(y <- slow_sum(1:10))
```

# Acknowledgement



**Parallel computing with R using foreach, future, and other packages**

BRYAN LEWIS



**Future: Simple Async, Parallel & Distributed Processing in R - What's Next?**

HENRIK BENGTSSON