

Big Data Analytics

Lecture 2: Computation and Memory

*Prof. Dr. Ulrich Matter
(University of St. Gallen)*

27/02/2020

1 Components of a standard computing environment

Figure 1 illustrates the key components of a standard computing environment to process digital data. In our case, these components serve the purpose of computing a statistic, given a large data set as input.

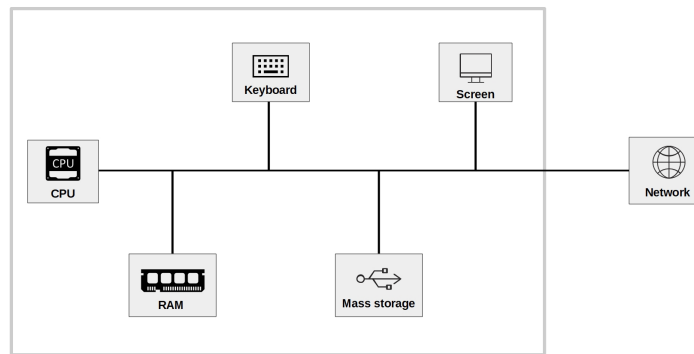


Figure 1: Basic components of a standard computing environment.

- The component actually *processing* data is the Central Processing Unit (CPU). When using R to process data, R commands are translated into complex combinations of a small set of basic operations which the *CPU* then executes.
- In order to work with data (e.g., in R), it first has to be loaded into the *memory* of our computer. More specifically, into the Random Access Memory (*RAM*). Typically, data is only loaded in the RAM as long as we work with it.
- *Mass Storage* refers to the type of computer memory we use to store data in the long run. This is what we call the *hard drive* or *hard disk*. In these days, the relevant hard disk is actually often not the one physically built into our computer but a hard disk ‘in the cloud’ (built into a server to which we connect over the Internet).

Very simply put, the difference between ‘data analytics’ and ‘Big Data analytics’ is that in the latter case, the standard usage of one or several of these components fails or works very inefficiently because the amount of data overwhelms its normal capacity.

From this hardware-perspective, there are two basic strategies to cope with the situation that one of these components is overwhelmed by the amount of data:

- *Scale up* (‘horizontal scaling’): Extend the physical capacity of the affected component by building a system with large RAM shared between applications. This sounds like a trivial solution (‘if RAM is too small, buy more RAM...’), but in practice it can be very expensive.
- *Scale out* (‘vertical scaling’): Distribute the workload over several computers (or separate components of a system).

From a software-perspective, there are many (context-specific) strategies that can help us to use the resources available more efficiently in order to process large amounts of data. In the context of computing statistics based on big data, this can involve:

- Implementing the computation of a given statistical procedure in a more efficient way (make better use of a given programming language or choose another programming language).
- Choosing/implementing a more efficient statistical procedure/algorithm (see, e.g., the *Uluru* algorithm).
- At a lower level, improving how the system allocates resources.

2 Units of information/data storage

The smallest unit of information in computing/digital data is called a *bit* (from *binary digit*; abbrev. ‘b’) and can take one of two (symbolic) values, either a 0 or a 1 (“off” or “on”). Consider, for example, the decimal number 139. Written in the binary system, 139 corresponds to the binary number 10001011. In order to store this number on a hard disk, we require a capacity of 8 bits, or one *byte* (1 byte = 8 bits; abbrev. ‘B’). Historically, one byte encoded a single character of text (i.e., in the ASCII character encoding system). 4 bytes (or 32 bits) are called a *word*. When thinking of a given data set in its raw/binary representation, we can simply think of it as a row of 0s and 1s, as illustrated in the following figure.

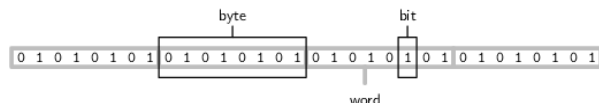


Figure 2: Writing data stored in RAM to a Mass Storage device (hard drive). Figure by Murrell (2009) (licensed under CC BY-NC-SA 3.0 NZ).

Bigger units for storage capacity usually build on bytes:

- 1 kilobyte (KB) = $1000^1 \approx 2^{10}$ bytes
- 1 megabyte (MB) = $1000^2 \approx 2^{20}$ bytes
- 1 gigabyte (GB) = $1000^3 \approx 2^{30}$ bytes
- 1 terabyte (TB) = $1000^4 \approx 2^{40}$ bytes
- 1 petabyte (PB) = $1000^5 \approx 2^{50}$ bytes
- 1 exabyte (EB) = $1000^6 \approx 2^{60}$ bytes
- 1 zettabyte (ZB) = $1000^7 \approx 2^{70}$ bytes

1ZB = 1000000000000000000000 bytes = 1 billion terabytes = 1 trillion gigabytes.

2.1 Example in R: Data types and information storage

Given the fact that computers only understand 0s and 1s, different approaches are taken to map these digital values to other symbols or images (text, decimal numbers, pictures, etc.) that we humans can more easily make sense of. Regarding text and numbers, these mappings involve *character encodings* (in which combinations of 0s and 1s represent a character in a specific alphabet) and *data types*.

Let's illustrate the main concepts with the simple numerical example from above. When we see the decimal number 139 written somewhere, we know that it means ‘one-hundred-and-thirty-nine’. The fact that our computer is able to print 139 on the screen means that our computer can somehow map a sequence of 0s and 1s to the symbols 1, 3, and 9. Depending on what we want to do with the data value 139 on our computer, there are different ways of how the computer can represent this value internally. Inter alia, we could load it into RAM as a *string* (‘text’/‘character’) or as an *integer* (‘natural number’) or *double* (numeric, floating

point number). All of them can be printed on screen but only the latter two can be used for arithmetic computations. This concept can easily be illustrated in R.

We initiate a new variable with the value 139. By using this syntax, R by default initiates the variable as an object of type double. We then can use this variable in arithmetic operations.

```
my_number <- 139
# check the class
typeof(my_number)
```

```
## [1] "double"
```

```
# arithmetic
my_number*2
```

```
## [1] 278
```

When we change the *data type* to ‘character’ (string) such operations are not possible.

```
# change and check type/class
my_number_string <- as.character(my_number)
typeof(my_number_string)
```

```
## [1] "character"
```

```
# try to multiply
my_number_string*2
```

```
## Error in my_number_string * 2: non-numeric argument to binary operator
```

If we change the variable to type integer, we can still use math operators.

```
# change and check type/class
my_number_int <- as.integer(my_number)
typeof(my_number_int)
```

```
## [1] "integer"
```

```
# arithmetics
my_number_int*2
```

```
## [1] 278
```

Having all variables in the right type is relevant for data analytics with all kind of sample sizes. However, given the fact that different data types have to be represented differently internally, different types might take up more or less memory and therefore substantially affect the performance when dealing with massive amounts of data.

We can illustrate this point with `object.size()`:

```
object.size("139")
```

```
## 112 bytes
```

```
object.size(139)
```

```
## 56 bytes
```

3 Big Data econometrics

When thinking about how to approach a data analytics task based on large amounts of data, it is helpful to consider two key aspects concerning the computational burden involved. (i) how is the statistic we have in mind computed. That is, what is the formal definition of the statistic (and how computationally demanding is it)? And (ii), given a definition of the statistic, how does the program/software/language implement the computation thereof?

Regarding the former, we might realize that there is an alternative statistical procedure that would provide essentially the same output but that happens to be more efficient (here: computationally efficient in contrast to statistically efficient). The latter point is a question of how to efficiently implement the given statistical procedure in your computing environment (taking into consideration the computer's available resources: CPU, RAM, etc.).

Below, we look at both of these two aspects by means of illustrative examples.

3.1 Example: Fast least squares regression

As an illustration of how an alternative statistical procedure can speed up our analysis, we look at one such procedure that has recently been developed to estimate linear models when the classical OLS estimator is too computationally intense for very large samples: The *Uhuru* algorithm (Dhillon et al. 2013).

3.1.1 OLS as a point of reference

Recall the OLS estimator in matrix notation, given the linear model $\mathbf{y} = \mathbf{X}\beta + \epsilon$:

$$\hat{\beta}_{OLS} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.$$

In order to compute $\hat{\beta}_{OLS}$, we have to compute $(\mathbf{X}^\top \mathbf{X})^{-1}$, which implies a computationally expensive matrix inversion.¹ If our data set is large, \mathbf{X} is large and the inversion can take up a lot of computation time. Moreover, the inversion and matrix multiplication to get $\hat{\beta}_{OLS}$ needs a lot of memory. In practice, it might well be that the estimation of a linear model via OLS with the standard approach in R (`lm()`) brings a computer to its knees, as there is not enough RAM available.

To further illustrate the point, we implement the OLS estimator in R.

```
beta_ols <-  
  function(X, y) {  
  
    # compute cross products and inverse  
    XXi <- solve(crossprod(X,X))  
    Xy <- crossprod(X, y)  
  
    return( XXi  %*% Xy )  
  }
```

Now, we will test our OLS estimator function with a few (pseudo) random numbers in a Monte Carlo study. First, we set the sample size parameters `n` (how many observations shall our pseudo sample have?) and `p` (how many variables shall describe these observations?) and initiate the data set `X`.

```
# set parameter values  
n <- 10000000  
p <- 4
```

¹The computational complexity of this is larger than $O(n^2)$. That is, for an input of size n , the time needed to compute (or the number of operations needed) is n^2 .

```

# Generate sample based on Monte Carlo
# generate a design matrix (~ our 'dataset') with four variables and 10000 observations
X <- matrix(rnorm(n*p, mean = 10), ncol = p)
# add column for intercept
X <- cbind(rep(1, n), X)

```

Now we define how the real linear model looks like that we have in mind and compute the output y of this model, given the input X .²

```

# MC model
y <- 2 + 1.5*X[,2] + 4*X[,3] - 3.5*X[,4] + 0.5*X[,5] + rnorm(n)

```

Finally, we test our `beta_ols` function.

```

# apply the ols estimator
beta_ols(X, y)

```

```

##           [,1]
## [1,]  2.0031121
## [2,]  1.4999626
## [3,]  3.9996541
## [4,] -3.5000947
## [5,]  0.5001843

```

3.1.2 The Uluru algorithm as an alternative to OLS

Following Dhillon et al. (2013), we implement a procedure to compute $\hat{\beta}_{Uluru}$:

$$\hat{\beta}_{Uluru} = \hat{\beta}_{FS} + \hat{\beta}_{correct}$$

, where

$$\hat{\beta}_{FS} = (\mathbf{X}_{subs}^T \mathbf{X}_{subs})^{-1} \mathbf{X}_{subs}^T \mathbf{y}_{subs}$$

, and

$$\hat{\beta}_{correct} = \frac{n_{subs}}{n_{rem}} \cdot (\mathbf{X}_{subs}^T \mathbf{X}_{subs})^{-1} \mathbf{X}_{rem}^T \mathbf{R}_{rem}$$

, and

$$\mathbf{R}_{rem} = \mathbf{Y}_{rem} - \mathbf{X}_{rem} \cdot \hat{\beta}_{FS}$$

.

The key idea behind this is that the computational bottleneck of the OLS estimator, the cross product and matrix inversion, $(\mathbf{X}^T \mathbf{X})^{-1}$, is only computed on a sub-sample (\mathbf{X}_{subs} , etc.), not the entire data set. However, the remainder of the data set is also taken into consideration (in order to correct a bias arising from the sub-sampling). Again, we implement the estimator in R to further illustrate this point.

```

beta_uluru <-
  function(X_subs, y_subs, X_rem, y_rem) {

    # compute beta_fs (this is simply OLS applied to the subsample)
    XXi_subs <- solve(crossprod(X_subs, X_subs))
    Xy_subs <- crossprod(X_subs, y_subs)
    b_fs <- XXi_subs %*% Xy_subs

    # compute \mathbf{R}_{rem}

```

²In reality we would not know this, of course. Acting as if we knew the real model is exactly point of Monte Carlo studies to analyze the properties of estimators by simulation.

```

    R_rem <- y_rem - X_rem %*% b_fs

    # compute \hat{\beta}_{correct}
    b_correct <- (nrow(X_subs)/(nrow(X_rem))) * Xxi_subs %*% crossprod(X_rem, R_rem)

    # beta uluru
    return(b_fs + b_correct)
}

```

Test it with the same input as above:

```

# set size of subsample
n_subs <- 1000
# select subsample and remainder
n_obs <- nrow(X)
X_subs <- X[1L:n_subs,]
y_subs <- y[1L:n_subs]
X_rem <- X[(n_subs+1L):n_obs,]
y_rem <- y[(n_subs+1L):n_obs]

# apply the uluru estimator
beta_uluru(X_subs, y_subs, X_rem, y_rem)

```

```

##           [,1]
## [1,]  1.9558444
## [2,]  1.4988402
## [3,]  3.9987101
## [4,] -3.4963467
## [5,]  0.5032868

```

This looks quite good already. Let's have a closer look with a little Monte Carlo study. The aim of the simulation study is to visualize the difference between the classical OLS approach and the *Uluru* algorithm with regard to bias and time complexity if we increase the sub-sample size in *Uluru*. For simplicity, we only look at the first estimated coefficient β_1 .

```

# define subsamples
n_subs_sizes <- seq(from = 1000, to = 500000, by=10000)
n_runs <- length(n_subs_sizes)
# compute uluru result, stop time
mc_results <- rep(NA, n_runs)
mc_times <- rep(NA, n_runs)
for (i in 1:n_runs) {
  # set size of subsample
  n_subs <- n_subs_sizes[i]
  # select subsample and remainder
  n_obs <- nrow(X)
  X_subs <- X[1L:n_subs,]
  y_subs <- y[1L:n_subs]
  X_rem <- X[(n_subs+1L):n_obs,]
  y_rem <- y[(n_subs+1L):n_obs]

  mc_results[i] <- beta_uluru(X_subs, y_subs, X_rem, y_rem)[2] # the first element is the intercept
  mc_times[i] <- system.time(beta_uluru(X_subs, y_subs, X_rem, y_rem))[3]
}

```

```
# compute ols results and ols time
ols_time <- system.time(beta_ols(X, y))
ols_res <- beta_ols(X, y)[2]
```

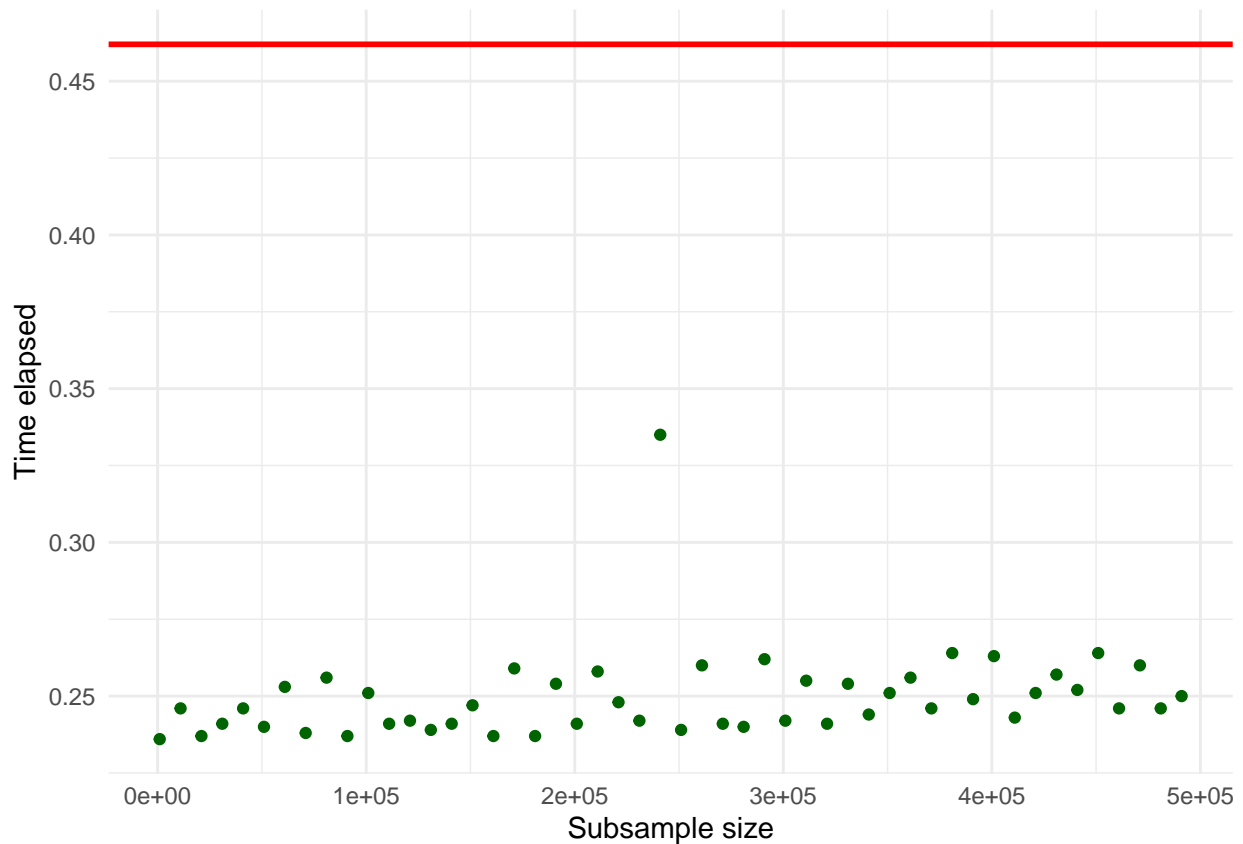
Let's visualize the comparison with OLS.

```
# load packages
library(ggplot2)

# prepare data to plot
plotdata <- data.frame(beta1 = mc_results,
                        time_elapsed = mc_times,
                        subs_size = n_subs_sizes)
```

First, let's look at the time used estimate the linear model.

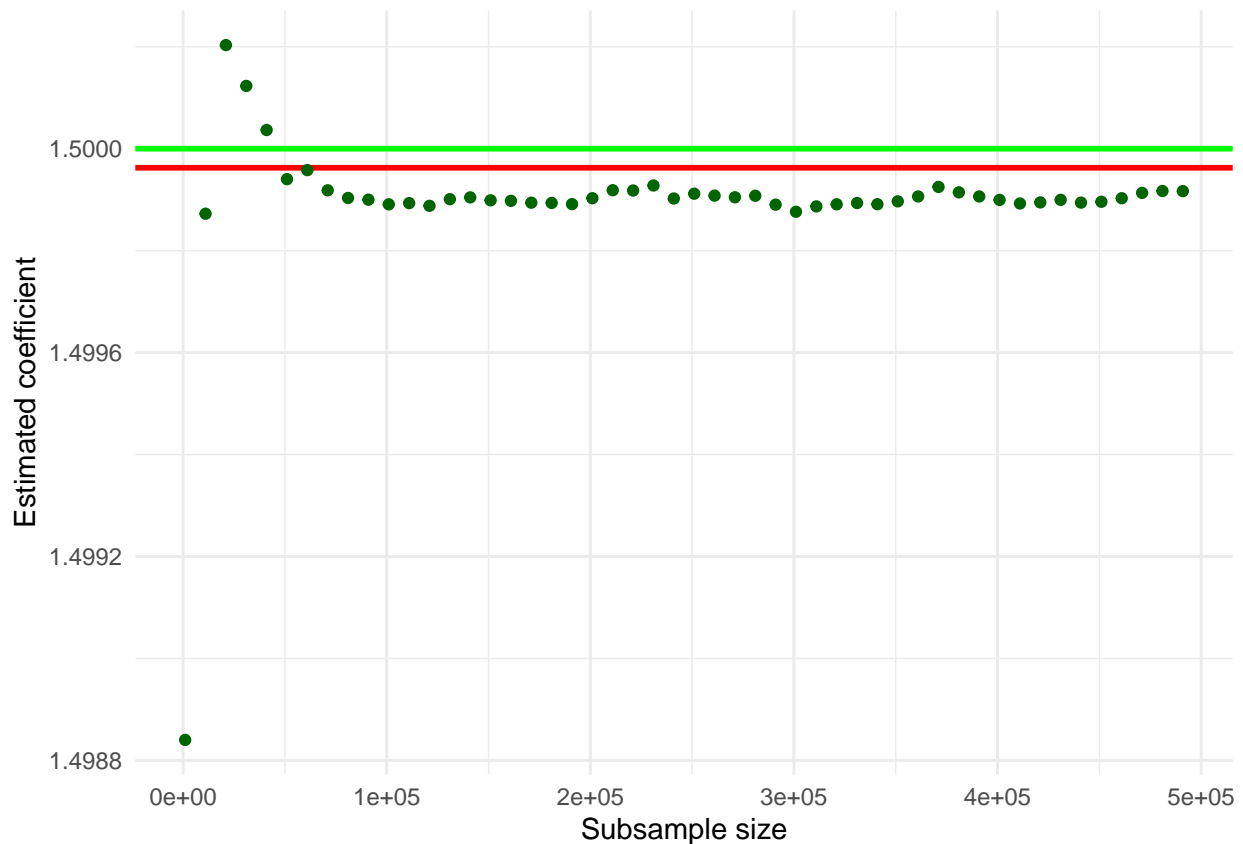
```
ggplot(plotdata, aes(x = subs_size, y = time_elapsed)) +
  geom_point(color="darkgreen") +
  geom_hline(yintercept = ols_time[3],
             color = "red",
             size = 1) +
  theme_minimal() +
  ylab("Time elapsed") +
  xlab("Subsample size")
```



The horizontal red line indicates the computation time for estimation via OLS, the green points indicate the computation time for the estimation via the Uluru algorithm. Note that even for large sub-samples, the computation time is substantially lower than for OLS.

Finally, let's have a look at how close the results are to OLS.

```
ggplot(plotdata, aes(x = subs_size, y = beta1)) +  
  geom_hline(yintercept = ols_res,  
             color = "red",  
             size = 1) +  
  geom_hline(yintercept = 1.5,  
             color = "green",  
             size = 1) +  
  geom_point(color="darkgreen") +  
  
  theme_minimal() +  
  ylab("Estimated coefficient") +  
  xlab("Subsample size")
```



The horizontal red line indicates the size of the estimated coefficient, when using OLS. The horizontal green line indicates the size of the actual coefficient. The green points indicate the size of the same coefficient estimated by the Uluru algorithm for different sub-sample sizes. Note that even relatively small sub-samples already deliver estimates very close to the OLS estimates.

4 Resource allocation

When optimizing the performance of an analytics program processing large amounts of data, it is useful to differentiate between the efficient allocation of computational (CPU) power, and the allocation of RAM (and mass storage).[^]In many data analysis tasks the two are, of course, intertwined. However, keeping both aspects in mind when optimizing an analytics program helps to choose the right tools.] Below, we will look

at both aspects in turn.

4.1 Case study: Parallel processing

In this example, we estimate a simple regression model that aims to assess racial discrimination in the context of police stops.³ The example is based on the ‘Minneapolis Police Department 2017 Stop Dataset’, containing data on nearly all stops made by the Minneapolis Police Department for the year 2017.

We start with importing the data into R.

```
stopdata <- read.csv("https://vincentarelbundock.github.io/Rdatasets/csv/carData/MplsStops.csv")
```

We specify a simple linear probability model that aims to test whether a stopped person identified as ‘white’ is less likely to have her vehicle searched when stopped by the police. In order to take into account level-differences between different police precincts, we add precinct-indicators to the regression specification

First, let’s remove observations with missing entries (NA) and code our main explanatory variable and the dependent variable.

```
# remove incomplete obs
stopdata <- na.omit(stopdata)
# code dependent var
stopdata$vsearch <- 0
stopdata$vsearch[stopdata$vehicleSearch=="YES"] <- 1
# code explanatory var
stopdata$white <- 0
stopdata$white[stopdata$race=="White"] <- 1
```

We specify our baseline model as follows.

```
model <- vsearch ~ white + factor(policePrecinct)
```

And estimate the linear probability model via OLS (the `lm` function).

```
fit <- lm(model, stopdata)
summary(fit)
```

```
##
## Call:
## lm(formula = model, data = stopdata)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.13937 -0.06329 -0.05473 -0.04227  0.97729
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    0.054733   0.005154  10.619 < 2e-16 ***
## white         -0.019553   0.004465  -4.380 1.19e-05 ***
## factor(policePrecinct)2  0.008556   0.006757   1.266  0.2054
## factor(policePrecinct)3  0.003409   0.006483   0.526  0.5990
## factor(policePrecinct)4  0.084639   0.006232  13.582 < 2e-16 ***
## factor(policePrecinct)5 -0.012465   0.006371  -1.956  0.0504 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

³Note that this example aims to illustrate a point about computation in an applied econometrics context. It does not make any argument about identification or the broader research question whatsoever.

```
##
## Residual standard error: 0.254 on 19078 degrees of freedom
## Multiple R-squared:  0.02502,    Adjusted R-squared:  0.02476
## F-statistic: 97.92 on 5 and 19078 DF,  p-value: < 2.2e-16
```

A potential problem with this approach (and there might be many more in this simple example) is that observations stemming from different police precincts might be correlated over time. If that is the case, we likely underestimate the coefficient's standard errors. There is a standard approach to compute estimates for so-called *cluster-robust* standard errors which would take the problem of correlation over time within clusters into consideration (and deliver a more conservative estimate of the SEs). However this approach only works well if the number of clusters in the data is roughly 50 or more. Here we only have 5.

The alternative approach is typically to compute bootstrapped standard errors. That is, we apply the bootstrap resampling procedure at the cluster level. Specifically, we draw B samples (with replacement), estimate and record for each bootstrap-sample the coefficient vector, and then estimate SE_{boot} based on the standard deviation of all respective estimated coefficient values.

```
# load packages
library(data.table)
# set the 'seed' for random numbers (makes the example reproducible)
set.seed(2)

# set number of bootstrap iterations
B <- 10
# get selection of precincts
precincts <- unique(stopdata$policePrecinct)
# container for coefficients
boot_coefs <- matrix(NA, nrow = B, ncol = 2)
# draw bootstrap samples, estimate model for each sample
for (i in 1:B) {

  # draw sample of precincts (cluster level)
  precincts_i <- sample(precincts, size = 5, replace = TRUE)
  # get observations
  bs_i <- lapply(precincts_i, function(x) stopdata[stopdata$policePrecinct==x,])
  bs_i <- rbindlist(bs_i)

  # estimate model and record coefficients
  boot_coefs[i,] <- coef(lm(model, bs_i))[1:2] # ignore FE-coefficients
}
```

Finally, let's compute SE_{boot} .

```
se_boot <- apply(boot_coefs,
  MARGIN = 2,
  FUN = sd)
se_boot
```

```
## [1] 0.004042725 0.004689611
```

Note that even with a very small B , computing SE_{boot} takes up some time to compute. When setting B to over 500, computation time will be substantial. Also note that running this code does hardly use up more memory than the very simple approach without bootstrapping (after all, in each bootstrap iteration the data set used to estimate the model is approximately the same size as the original data set). There is little we can do to improve the script's performance regarding memory. However we can tell R how to allocate CPU resources more efficiently to handle that many regression estimates.

Particularly, we can make use of the fact that most modern computing environments (such as a laptop) have

CPU's with several *cores*. We can exploit this fact by instructing the computer to run the computations *in parallel* (simultaneously computing in several cores.). The following code is a parallel implementation of our bootstrap procedure which does exactly that.

```
# install.packages("doSNOW", "parallel")
# load packages for parallel processing
library(doSNOW)

# get the number of cores available
ncores <- parallel::detectCores()
# set cores for parallel processing
ctemp <- makeCluster(ncores) #
registerDoSNOW(ctemp)

# set number of bootstrap iterations
B <- 10
# get selection of precincts
precincts <- unique(stopdata$policePrecinct)
# container for coefficients
boot_coefs <- matrix(NA, nrow = B, ncol = 2)

# bootstrapping in parallel
boot_coefs <-
  foreach(i = 1:B, .combine = rbind, .packages="data.table") %dopar% {

    # draw sample of precincts (cluster level)
    precincts_i <- sample(precincts, size = 5, replace = TRUE)
    # get observations
    bs_i <- lapply(precincts_i, function(x) stopdata[stopdata$policePrecinct==x,])
    bs_i <- rbindlist(bs_i)

    # estimate model and record coefficients
    coef(lm(model, bs_i))[1:2] # ignore FE-coefficients

  }

# be a good citizen and stop the snow clusters
stopCluster(cl = ctemp)
```

As a last step, we compute again SE_{boot} .

```
se_boot <- apply(boot_coefs,
  MARGIN = 2,
  FUN = sd)
se_boot
```

```
## (Intercept)      white
## 0.003708425 0.003390722
```

4.2 Case study: Memory allocation

Consider the first steps of a data pipeline in R. The first part of our script to import and clean the data looks as follows.

```
#####
# Big Data Statistics: Flights data import and preparation
#
# U. Matter, January 2019
#####

# SET UP -----

# fix variables
DATA_PATH <- "../data/flights.csv"

# DATA IMPORT -----
flights <- read.csv(DATA_PATH)

# DATA PREPARATION -----
flights <- flights[,-1:-3]
```

When running this script, we notice that some of the steps take a while. Moreover, while none of these steps obviously involves a lot of computation (such as a matrix inversion or numerical optimization), it quite likely involves memory allocation. We first read data into RAM (allocated to R by our operating system). It turns out that there are different ways to allocate RAM when reading data from a CSV file. Depending on the amount of data to be read in, one or the other approach might be faster. We first investigate the RAM allocation in R with `mem_change()` and `mem_used()`.

```
# SET UP -----

# fix variables
DATA_PATH <- "../data/flights.csv"
# load packages
library(pryr)

# check how much memory is used by R (overall)
mem_used()

## 1.08 GB

# check the change in memory due to each step

# DATA IMPORT -----
mem_change(flights <- read.csv(DATA_PATH))

## -6.64 MB

# DATA PREPARATION -----
flights <- flights[,-1:-3]

# check how much memory is used by R now
mem_used()

## 1.07 GB
```

The last result is kind of interesting. The object `flights` must have been larger right after importing it than at the end of the script. We have thrown out several variables, after all. Why does R still use that much memory? R does by default not ‘clean up’ memory unless it is really necessary (meaning no more memory is available). In this case, R has still way more memory available from the operating system, thus there is no need to ‘collect the garbage’ yet. However, we can force R to collect the garbage on the spot with `gc()`. This

can be helpful to better keep track of the memory needed by an analytics script.

```
gc()
```

```
##           used (Mb) gc trigger (Mb) limit (Mb) max used (Mb)
## Ncells   1025186 54.8   2079574 111.1         NA   2079574 111.1
## Vcells 126561998 965.6  256314384 1955.6       16384 253107984 1931.1
```

Now, let's see how we can improve the performance of this script with regard to memory allocation. Most memory is allocated when importing the file. Obviously, any improvement of the script must still result in importing all the data. However, there are different ways to read data into RAM. `read.csv()` reads all lines of a csv file consecutively. In contrast, `data.table::fread()` first 'maps' the data file into memory and only then actually reads it in line by line. This involves an additional initial step, but the larger the file, the less relevant is this first step with regard to the total time needed to read all the data into memory. By switching on the `verbose` option, we can actually see what `fread` is doing.

```
# load packages
library(data.table)
```

```
# DATA IMPORT -----
```

```
flights <- fread(DATA_PATH, verbose = TRUE)
```

```
##   omp_get_num_procs()           4
##   R_DATATABLE_NUM_PROCS_PERCENT unset (default 50)
##   R_DATATABLE_NUM_THREADS      unset
##   omp_get_thread_limit()       2147483647
##   omp_get_max_threads()        4
##   OMP_THREAD_LIMIT             unset
##   OMP_NUM_THREADS              unset
##   RestoreAfterFork             true
##   data.table is using 2 threads. See ?setDTthreads.
## Input contains no \n. Taking this to be a filename to open
## [01] Check arguments
##   Using 2 threads (omp_get_max_threads()=4, nth=2)
##   NAstrings = [<<NA>>]
##   None of the NAstrings look like numbers.
##   show progress = 0
##   0/1 column will be read as integer
## [02] Opening the file
##   Opening file ../data/flights.csv
##   File opened, size = 29.53MB (30960660 bytes).
##   Memory mapped ok
## [03] Detect and skip BOM
## [04] Arrange mmap to be \0 terminated
##   \n has been found in the input and different lines can end with different line endings (e.g. mixed
## [05] Skipping initial rows if needed
##   Positioned on line 1 starting: <<year,month,day,dep_time,sched_>>
## [06] Detect separator, quoting rule, and ncolumns
##   Detecting sep automatically ...
##   sep=', ' with 100 lines of 19 fields using quote rule 0
##   Detected 19 columns on line 1. This line is either column names or first data row. Line starts as:
##   Quote rule picked = 0
##   fill=false and the most number of columns found is 19
## [07] Detect column types, good nrow estimate and whether first row is column names
##   Number of sampling jump points = 100 because (30960659 bytes from row 1 to eof) / (2 * 8882 jump0s
##   Type codes (jump 000)      : 555555555A5AAA5555A Quote rule 0
```

```
## Type codes (jump 100) : 55555555A5AAA5555A Quote rule 0
## 'header' determined to be true due to column 1 containing a string on row 1 and a lower type (int32)
## =====
## Sampled 10048 rows (handled \n inside quoted fields) at 101 jump points
## Bytes from first data row on line 2 to the end of last row: 30960501
## Line length: mean=92.03 sd=3.56 min=68 max=98
## Estimated number of rows: 30960501 / 92.03 = 336403
## Initial alloc = 370043 rows (336403 + 9%) using bytes/max(mean-2*sd,min) clamped between [1.1*estn
## =====
## [08] Assign column names
## [09] Apply user overrides on column types
## After 0 type and 0 drop user overrides : 55555555A5AAA5555A
## [10] Allocate memory for the datatable
## Allocating 19 column slots (19 - 0 dropped) with 370043 rows
## [11] Read the data
## jumps=[0..30), chunk_size=1032016, total_size=30960501
## Read 336776 rows x 19 columns from 29.53MB (30960660 bytes) file in 00:00.106 wall clock time
## [12] Finalizing the datatable
## Type counts:
##      14 : int32      '5'
##       5 : string     'A'
## =====
## 0.001s ( 1%) Memory map 0.029GB file
## 0.004s ( 3%) sep=', ' ncol=19 and header detection
## 0.000s ( 0%) Column type detection using 10048 sample rows
## 0.005s ( 4%) Allocation of 370043 rows x 19 cols (0.033GB) of which 336776 ( 91%) rows used
## 0.097s (92%) Reading 30 chunks (0 swept) of 0.984MB (each chunk 11225 rows) using 2 threads
## + 0.047s (44%) Parse to row-major thread buffers (grown 0 times)
## + 0.048s (45%) Transpose
## + 0.003s ( 3%) Waiting
## 0.000s ( 0%) Rereading 0 columns due to out-of-sample type exceptions
## 0.106s      Total
```

Let's put it all together and look at the memory changes and usage. For a fair comparison, we first have to delete `flights` and collect the garbage with `gc()`.

```
# SET UP -----
```

```
# fix variables
```

```
DATA_PATH <- "../data/flights.csv"
```

```
# load packages
```

```
library(pryr)
```

```
library(data.table)
```

```
# housekeeping
```

```
flights <- NULL
```

```
gc()
```

```
##           used (Mb) gc trigger (Mb) limit (Mb) max used (Mb)
## Ncells  1014203 54.2   2079574 111.1      NA   2079574 111.1
## Vcells 124227066 947.8  256314384 1955.6   16384 253107984 1931.1
```

```
# check the change in memory due to each step
```

```
# DATA IMPORT -----
```

```
mem_change(flights <- fread(DATA_PATH))
```

36.4 MB

Note that `fread()` uses up more memory. By default, `fread` does not parse strings as `factors` (and `read.csv()` does). Storing strings in factors is more memory efficient.

5 Beyond memory

In the previous example we have inspected how RAM is allocated to store objects in the R computing environment. But what if all RAM of our computer is not enough to store all the data we want to analyze?

Modern operating systems have a way to dealing with such a situation. Once all RAM is used up by the currently running programs, the OS allocates parts of the memory back to the hard-disk which then works as *virtual memory*. The following figure illustrates this point.

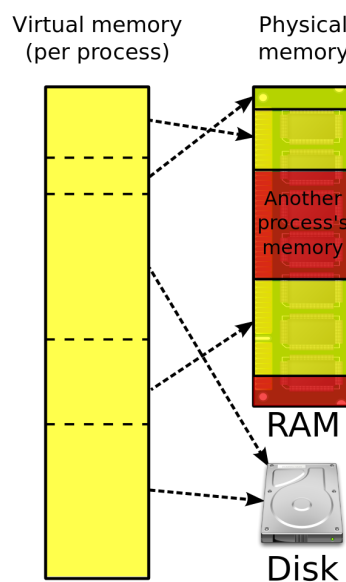


Figure 3: Virtual memory. Figure by Ehamberg (CC BY-SA 3.0).

For example, when we implement an R-script that imports one file after the other into the R environment, ignoring the RAM capacity of our computer, the OS will start *paging* data to the virtual memory. This happens ‘under the hood’ without explicit instructions by the user. When this happens, we quite likely notice that the computer slows down a lot.

While this default usage of virtual memory by the OS is helpful to run several applications at the same time, each taking up a moderate amount of memory, it is not a really useful tool for processing large amounts of data in one application (R). However, the underlying idea of using both RAM and Mass storage simultaneously in order to cope with a lack of memory, is very useful in the context of big data statistics.

Several R packages have been developed that exploit the idea behind virtual memory explicitly for big data analytics. The basic idea behind these packages is to map a data set to the hard disk when loading it into R. The actual data values are stored in chunks on the hard-disk, while the structure/metadata of the data set is loaded into R. See Walkowiak (2016), Chapter 3 for more details and example code.

References

- Dhillon, Paramveer, Yichao Lu, Dean P. Foster, and Lyle Ungar. 2013. “New Subsampling Algorithms for Fast Least Squares Regression.” In *Advances in Neural Information Processing Systems 26*, 360–68.
- Murrell, Paul. 2009. *Introduction to Data Technologies*. London, UK: CRC Press.
- Walkowiak, Simkon. 2016. *Big Data Analytics with R*. Birmingham, UK: PACKT Publishing.