

Cran download logs aggregation time summary

A.Birek, M.Kosiński, N.Ryciak, W.Ryciuk

May 11, 2014

CONTENTS

1	Downloading data	3
2	SAS Path	4
3	Traditional \mathcal{R} Path	6
3.1	Unmerged \mathcal{R} files Path	6
3.2	Merged \mathcal{R} files Path	6
4	Rcpp Path	7
A	Preparing report	8
A.1	Data download, unzipt, conversion syntax	8
B	SAS Syntax	9



Figure 1: Powered by ! <https://github.com/MarcinKosinski/AlmostBigData>

CHAPTER
ONE

DOWNLOADING DATA

Syntax used for downloading, unzipping and merging data is available in section [A.1](#). More or less downloading looked like this and took about:

```
start <- as.Date("2012-10-01")
today <- as.Date("2014-05-10")
all_days <- seq(start, today, by = "day")
year <- as.POSIXlt(all_days)$year + 1900
urls <- paste0("http://cran-logs.rstudio.com/", year, "/", all_days, ".csv.gz")

destdir <- "D:/bd1/AlmostBigData/cran-logs/"
n <- length(urls)
i = 1
for (i in 1:n) {
  destfile <- stri_paste(destdir, as.character(all_days[i]))
  download.file(urls[i], destfile)
}
```

Unzipping files syntax looked like this and took:

```
lok <- "D:/bd1/AlmostBigData"
gzpath <- character(n)
i <- 1
for (i in 1:n) {
  gzpath[i] <- paste(lok, "/cran-logs", all_days[i], sep = "")
}
install.packages("R.utils")
library(R.utils)
for (i in 1:n) {
  gunzip(gzpath[i], destname = paste(gzpath[i], ".csv"), remove = TRUE)
}
```

Converting CSV files with proper delimiter syntax looked like this and time spent was:

```
for (i in 1:n) {
  write.csv2(read.csv2(paste(gzpath[i], ".csv"), sep = ","), paste(gzpath[i], "_new.csv"))
}
```

CHAPTER

TWO

SAS PATH

Syntax used for importing, merging and summarizing data is available in chapter [B](#). Importing csv files into **SAS** syntax looked like this and took:

```
proc import datafile='D:/bd1/AlmostBigData/cran-logs2012-10-01 _new.csv'
out=CR.cran1 dbms=csv replace;
    delimiter = ';';
    getnames=yes;
    run;

...

proc import datafile='D:/bd1/AlmostBigData/cran-logs2014-05-09 _new.csv'
out=CR.cran586 dbms=csv replace;
    delimiter = ';';
    getnames=yes;
    run;
```

Merging all those files syntax looked like this and time expired was:

```
data Cr.DANE;
set
CR.cran1,
CR.cran2,
....
CR.cran586;
run;
```

Summaries of each variable syntax looked like this and time expired was:

```
proc summary data=Cr.DANE print;
class package;
run;

proc summary data=Cr.DANE print;
class version;
run;

proc summary data=Cr.DANE print;
class r_arch;
run;

proc summary data=Cr.DANE print;
class r_os;
run;

proc summary data=Cr.DANE print;
class r_version;
run;
```

```
proc summary data=Cr.DANE print;  
class country;  
run;
```

CHAPTER
THREE

TRADITIONAL \mathcal{R} PATH

3.1 Unmerged \mathcal{R} files Path

3.2 Merged \mathcal{R} files Path

CHAPTER
FOUR

RCPD PATH

PREPARING REPORT

A.1 Data download, unzipt, conversion syntax

```
start <- as.Date("2012-10-01")
today <- as.Date("2014-05-10")
all_days <- seq(start, today, by = "day")
year <- as.POSIXlt(all_days)$year + 1900
urls <- paste0("http://cran-logs.rstudio.com/", year, "/", all_days, ".csv.gz")

destdir <- "D:/bd1/AlmostBigData/cran-logs/"
n <- length(urls)
i = 1
for (i in 1:n) {
  destfile <- stri_paste(destdir, as.character(all_days[i]))
  download.file(urls[i], destfile)
}

lok <- "D:/bd1/AlmostBigData"
gzpath <- character(n)
i <- 1
for (i in 1:n) {
  gzpath[i] <- paste(lok, "/cran-logs", all_days[i], sep = "")
}
install.packages("R.utils")
library(R.utils)
for (i in 1:n) {
  gunzip(gzpath[i], destname = paste(gzpath[i], ".csv"), remove = TRUE)
}

for (i in 1:n) {
  write.csv2(read.csv2(paste(gzpath[i], ".csv"), sep = ","), paste(gzpath[i], "_new.csv"))
}
```


SAS SYNTAX

```
proc import datafile='D:/bd1/AlmostBigData/cran-logs2012-10-01 _new.csv'
out=CR.cran1 dbms=csv replace;
    delimiter = ';';
    getnames=yes;
    run;

...

proc import datafile='D:/bd1/AlmostBigData/cran-logs2014-05-09 _new.csv'
out=CR.cran586 dbms=csv replace;
    delimiter = ';';
    getnames=yes;
    run;

data Cr.DANE;
set
CR.cran1,
CR.cran2,
....
CR.cran586;
run;

proc summary data=Cr.DANE print;
class package;
run;

proc summary data=Cr.DANE print;
class version;
run;

proc summary data=Cr.DANE print;
class r_arch;
run;

proc summary data=Cr.DANE print;
class r_os;
run;

proc summary data=Cr.DANE print;
class r_version;
run;

proc summary data=Cr.DANE print;
class country;
run;
```