

ARCHIWIZACJA ARTEFAKTÓW Z PAKIETEM **archivist**

Przemysław Biecek, Marcin Kosiński
przemyslaw.biecek@gmail.com, m.p.kosinski@gmail.com

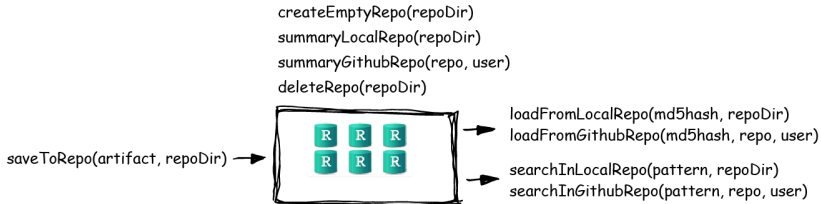
<https://github.com/pbiecek/archivist>

16 Października, 2014

DO CZEGO SŁUŻY **archivist**?

- pozwala **magazynować i archiwizować** obiekty w repozytoriach przechowywanych na dysku bądź na Githubie
- udostępnia proste w obsłudze funkcje uprawniające **szukanie i odzyskiwanie** obiektów
- wspiera filozofię **powtarzalnych badań** (*reporoducible research*)
- idealnie sprawuje się jako **pamięć podręczna**

JAK DZIAŁA archivist?



Each repository contains a database with objects metadata.

Objects are stored as binary files.

Each object has a unique key - md5 hash.

Metadata, like object class, name, creation date, relations with other objects are useful when searching for an object in a repository.

CASE STUDY: PAMIĘĆ PODRĘCZNA

Funkcja `getMaxDistribution()` podsumowuje rozkład maksimum z N niezależnych obserwacji z rozkładu D dla R replikacji.

```
getMaxDistribution <- function(D = rnorm, N = 10,  
                               R = 1000000) {  
  res <- replicate(R, max(D(N)))  
  summary(res)  
}  
system.time( getMaxDistribution(rnorm, 10) )
```

user	system	elapsed
49.11	0.41	50.53

```
system.time( getMaxDistribution(rexp, 20) )
```

user	system	elapsed
49.02	0.13	50.06

CASE STUDY: PAMIĘĆ PODRĘCZNA

Spróbuj sam. Funkcję można łatwo pobrać z repozytorium przechowywanego na Githubie.

```
install.packages("archivist")
library("archivist")
# library(devtools)
# install_github("pbiecek/archivist")

loadFromGithubRepo("c", user="MarcinKosinski",
                    repo="Museum")
```

CASE STUDY: PAMIĘĆ PODRĘCZNA

```
cache <- function(cacheRepo, FUN, ...) {  
  tmp1 <- list(...)  
  tmp1$.FUN <- FUN  
  outputHash <- digest(tmp1)  
  isInRepo <- searchInLocalRepo(paste0("cacheId:",  
                                       outputHash), cacheRepo)  
  if (length(isInRepo) > 0)  
    return(loadFromLocalRepo(isInRepo[1], repoDir = cacheRepo))  
  output <- do.call(FUN, list(...))  
  attr( output, "tags") <- paste0("cacheId:", outputHash)  
  attr( output, "call") <- ""  
  saveToRepo(output, repoDir = cacheRepo, archiveData = TRUE,  
             archiveMiniature = FALSE, rememberName = FALSE)  
  output  
}  
loadFromGithubRepo("7", user="MarcinKosinski",  
                  repo="Museum")
```

CASE STUDY: PAMIĘĆ PODRĘCZNA

Używając pakietu **archivist** można przygotować repozytorium przechowujące wywołania funkcji `cache()`, aby uniknąć ich powtórnego wywołania w przyszłości.

```
cacheRepo <- tempdir()  
createEmptyRepo( cacheRepo )  
library( "digest" )  
system.time( cache(cacheRepo, getMaxDistribution,  
                    rnorm, 10) )
```

user	system	elapsed
46.85	0.17	47.16

```
system.time( cache(cacheRepo, getMaxDistribution,  
                    rexp, 10) )
```

user	system	elapsed
44.29	0.08	46.52

CASE STUDY: PAMIĘĆ PODRĘCZNA

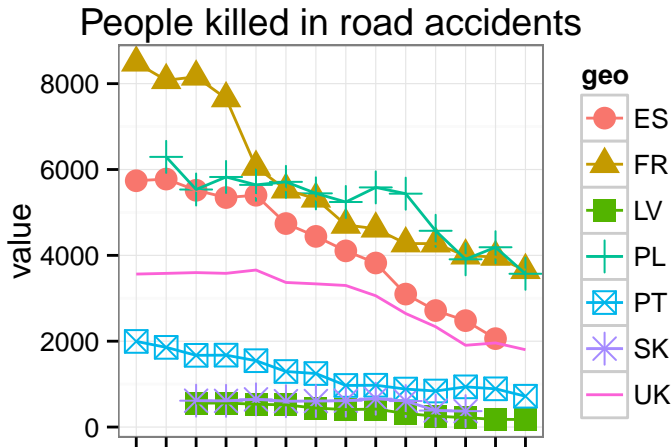
Pierwsze wywołanie nie różniło się prędkością od zwykłego wywołania samej funkcji `getMax`., jednak drugie wykorzystuje wyniki z pierwszego wywołania i jest o wiele szybsze.

```
system.time( cache( cacheRepo, getMaxDistribution,  
                    rnorm, 10) )
```

user	system	elapsed
0.03	0.00	0.03

INNE PRZYKŁADY

```
wykres <-loadFromGithubRepo( md5hash = "fcd7" ,  
user = "pbiecek", repo = "graphGallery", value = TRUE)  
print( wykres )
```



INNE PRZYKŁADY

```
hash <- searchInGithubRepo( "name:crime",  
user="MarcinKosinski", repo="Museum", fixed = FALSE)  
getTagsGithub( hash, user="MarcinKosinski",  
                repo="Museum")
```

```
name:crime.by.state %.%  
filter(State == "New York", Year == 2005) %.%  
  arrange(desc(Count)) %.%  
  select(Type.of.Crime, Count) %.%  
  mutate(Proportion = Count/sum(Count)) %.%  
  group_by(Type.of.Crime) %.%  
    summarise(num.types = n(), counts = sum(Count))
```

WIĘCEJ?

Te i inne zastosowania można zobaczyć pod tymi linkami lub na stronie <http://pbiecek.github.io/archivist/>

Cache with the `archivist` package

Retrieving all plots with other github repository (example with flights data from Hadley Wickham `useR!2014` tutorial)

Archiving artifacts with their chaining code

Just get the object

Lazy load with **`archivist`**