

Branch: master ▾

Find file

Copy path

cleaningdataR / week2_quiz_Baker.R

 aurquhart Adam B answers

8da54df on Oct 20, 2014

1 contributor

Raw Blame History



111 lines (80 sloc) 3.89 KB

```
1 # Remove everything from the workspace
2 rm(list = ls())
3
4 # Windows
5 setwd('C://Users//ABaker//Documents//GitHub//Coursera//Getting and Cleaning Data')
6 # OS X
7 setwd('/Users/adam_baker_1/GitHub/Coursera/Getting and Cleaning Data')
8
9 if (!file.exists("data")) {
10   dir.create("data")
11 }
12
13 # Question 1 - What time was the data sharing repo created?
14
15 library(httr)
16 library(httpuv)
17 library(jsonlite)
18 library(dplyr)
19
20 # 1. Find OAuth settings for github:
21 #   http://developer.github.com/v3/oauth/
22 github <- oauth_endpoints("github")
23
24
25 # 2. Register an application at https://github.com/settings/applications;
26 #   Use any URL you would like for the homepage URL (http://github.com is fine)
27 #   and http://localhost:1410 as the callback url
28 #
29 #   Insert your client ID and secret below - if secret is omitted, it will
30 #   look it up in the GITHUB_CONSUMER_SECRET environmental variable.
31 myapp <- oauth_app("github", "9e975720d681777b66d4", "b2a99834bf1543521472d02fd031ed022c63e983")
32
33 # 3. Get OAuth credentials
34 github_token <- oauth2.0_token(oauth_endpoints("github"), myapp)
35
36 # 4. Use API
37 gtoken <- config(token = github_token)
38 req <- GET("https://api.github.com/users/jtleek/repos", gtoken)
39 stop_for_status(req)
40 content(req)
41
42 json1 = content(req)
```

```

43 json2 = jsonlite::fromJSON(toJSON(json1))
44 head(json2)
45
46 json2[json2$full_name == "jtleek/datasharing",] # 2013-11-07T13:25:07Z
47
48 # Question 2 - Which of the following commands will select only the data for the probability weights pwgtp1 wit
49 # sqldf("select * from acs where AGE < 50")
50 # sqldf("select pwgtp1 from acs where AGE < 50")
51 # sqldf("select pwgtp1 from acs")
52 # sqldf("select * from acs")
53
54 library(sqldf)
55
56 # Using method = "curl" doesn't seem to work on Windows
57 download.file("https://d396qusza40orc.cloudfront.net/getdata%2Fdata%2Fss06pid.csv", destfile = "./data/acs_wk2
58
59 # Using method = "curl" on OS X works
60 download.file("https://d396qusza40orc.cloudfront.net/getdata%2Fdata%2Fss06pid.csv", destfile = "./data/acs_wk2
61
62 acs <- read.csv("./data/acs_wk2.csv")
63 head(ac)
64
65 sqldf("select * from acs where AGE < 50") # NO - selects all data for ages less than 50
66 sqldf("select pwgtp1 from acs where AGE < 50") # YES - selects the correct column and for ages less than 50
67 sqldf("select pwgtp1 from acs") # NO - only selects pwgtp1 column
68 sqldf("select * from acs") # NO - selects everything
69
70 # Question 3 - Using the same data frame you created in the previous problem, what is the equivalent function t
71
72 unique(ac$AGE) # this is the desired result
73
74 sqldf("select unique AGE from acs") # NO - syntax error
75 sqldf("select distinct AGE from acs") # YES? - generates the same list but in a different format
76 sqldf("select distinct pwgtp1 from acs") # NO - selects the distinct values of the wrong column
77 sqldf("select unique * from acs") # NO - syntax error
78
79 # Question 4 - How many characters are in the 10th, 20th, 30th and 1000th lines of HTML from this page
80
81 con = url("http://biostat.jhsph.edu/~jtleek/contact.html ")
82 htmlCode = readLines(con)
83 close(con)
84 htmlCode
85
86 nchar(htmlCode[10]) # 45
87 nchar(htmlCode[20]) # 31
88 nchar(htmlCode[30]) # 7
89 nchar(htmlCode[100]) # 25
90
91 # Answer = 45 31 7 25
92
93 # Question 5 - Read the data set into R and report the sum of the numbers in the fourth of the nine columns
94
95 # Downloading data from the Web
96 fileUrl <- "https://d396qusza40orc.cloudfront.net/getdata%2Fwksst8110.for"
97
98 # Using method = "curl" doesn't seem to work on Windows
99 download.file(fileUrl, destfile = "./data/ac_survey.csv")

```

```
100
101 # Using method = "curl" on OS X works
102 download.file(fileUrl, destfile = "./data/q5_data.for", method = "curl")
103 list.files('./data/')
104
105 ?read.fwf
106
107 q5_df <- read.fwf(file = "./data/q5_data.for", widths = c(15, 4, 1, 3, 5, 4), header = FALSE, sep = "\t", skip
108 head(q5_df)
109
110 # Need to sum up the V6 column
111 sum(q5_df$V6) # 32426.7
```