

DSC520_Week2_Guruprasad_Velikadu_Assignment02

Guruprasad Velikadu Krishnamoorthy

2022-12-11

Contents

1	Assignment 02	2
1.1	Check your current working directory using <code>getwd()</code>	2
1.2	List the contents of the working directory with the <code>dir()</code> function	2
1.3	If the current directory does not contain the <code>data</code> directory, set the	2
1.4	Load the file <code>data/tidynomicon/person.csv</code> to <code>person_df1</code> using <code>read.csv</code>	2
1.5	R interpreted names as factors, which is not the behavior we want	3
1.6	Read the file <code>data/scores.csv</code> to <code>scores_df</code>	3
1.7	Load the <code>readxl</code> library	3
1.8	Using the <code>excel_sheets()</code> function from the <code>readxl</code> package,	3
1.9	Using the <code>read_excel</code> function, read the Voter Turnout sheet	4
1.10	Using the <code>read_excel()</code> function, read the Voter Turnout sheet	5
1.11	Load the DBI library	5
1.12	Create a database connection to <code>data/tidynomicon/example.db</code> using the <code>dbConnect()</code> function	5
1.13	Query the Person table using the <code>dbGetQuery</code> function and the	5
1.14	List the tables using the <code>dbListTables()</code> function	6
1.15	Read all of the tables at once using the <code>lapply</code> function and assign the result to the <code>tables</code> variable	6
1.16	Use the <code>dbDisconnect</code> function to disconnect from the database	7
1.17	Import the <code>jsonlite</code> library	7
1.18	Convert the <code>scores_df</code> dataframe to JSON using the <code>toJSON()</code> function	7
1.19	Convert the <code>scores</code> dataframe to JSON using the <code>toJSON()</code> function with the <code>pretty=TRUE</code> option	8
2	Session Info	11

1 Assignment 02

1.1 Check your current working directory using `getwd()`

```
getwd()
```

```
## [1] "C:/Users/Gurup/GURU/Learning/Masters/Term_2/DSC520_T302_Statistics_for_Data_Science/Week_2/Assignments"
```

1.2 List the contents of the working directory with the `dir()` function

```
dir()
```

```
## [1] "DSC520_Week2_Guruprasad_Velikadu_Assignment00.pdf"
## [2] "DSC520_Week2_Guruprasad_Velikadu_Assignment00_Markdown.pdf"
## [3] "DSC520_Week2_Guruprasad_Velikadu_Assignment00_Markdown0000000000000000.pdf"
## [4] "DSC520_Week2_Guruprasad_Velikadu_Assignment01.pdf"
## [5] "DSC520_Week2_Guruprasad_Velikadu_Assignment01_Markdown.pdf"
## [6] "DSC520_Week2_Guruprasad_Velikadu_Assignment01_Markdown.Rmd"
## [7] "DSC520_Week2_Guruprasad_Velikadu_Assignment02.pdf"
## [8] "DSC520_Week2_Guruprasad_Velikadu_Assignment02_Markdown.pdf"
## [9] "DSC520_Week2_Guruprasad_Velikadu_Assignment02_Markdown.Rmd"
## [10] "Final"
## [11] "Markdown_sample1.pdf"
## [12] "Markdown_sample1.Rmd"
```

1.3 If the current directory does not contain the data directory, set the

1.3.1 working directory to project root folder (the folder should contain the data directory)

1.3.2 Use `setwd()` if needed

```
knitr::opts_knit$set(root.dir = "C:/Users/Gurup/GURU/Learning/Masters/Term_2/DSC520_T302_Statistics_for_Data_Science/Week_2/Assignments")
```

1.4 Load the file `data/tidynomicon/person.csv` to `person_df1` using `read.csv`

1.4.1 Examine the structure of `person_df1` using `str()`

```
person_df1 <- read.csv(file="data/tidynomicon/person.csv",header=TRUE,stringsAsFactors=TRUE)
str(person_df1)
```

```
## 'data.frame': 5 obs. of 3 variables:
## $ person_id : Factor w/ 5 levels "danforth","dyer",...: 2 4 3 5 1
## $ personal_name: Factor w/ 4 levels "Anderson","Frank",...: 4 2 1 3 2
## $ family_name : Factor w/ 5 levels "Danforth","Dyer",...: 2 4 3 5 1
```

1.5 R interpreted names as factors, which is not the behavior we want

1.5.1 Load the same file to `person_df2` using `read.csv` and setting `stringsAsFactors` to `FALSE`

1.5.2 Examine the structure of `person_df2` using `str()`

```
person_df2 <- read.csv(file="data/tidynomicon/person.csv",header=TRUE,stringsAsFactors=FALSE)
str(person_df2)
```

```
## 'data.frame': 5 obs. of 3 variables:
## $ person_id : chr "dyer" "pb" "lake" "roe" ...
## $ personal_name: chr "William" "Frank" "Anderson" "Valentina" ...
## $ family_name : chr "Dyer" "Pabodie" "Lake" "Roerich" ...
```

1.6 Read the file `data/scores.csv` to `scores_df`

1.6.1 Display summary statistics using the `summary()` function

```
scores_df <- read.csv(file="data/scores.csv",header=TRUE,stringsAsFactors=TRUE)
summary(scores_df)
```

```
##      Count      Score      Section
## Min.   :10.00  Min.   :200.0  Regular:19
## 1st Qu.:10.00  1st Qu.:300.0  Sports :19
## Median :10.00  Median :322.5
## Mean   :14.47  Mean   :317.5
## 3rd Qu.:20.00  3rd Qu.:357.5
## Max.   :30.00  Max.   :395.0
```

1.7 Load the `readxl` library

```
library(readxl)
```

1.8 Using the `excel_sheets()` function from the `readxl` package,

1.8.1 list the worksheets from the file `data/G04ResultsDetail2004-11-02.xls`

```
excel_sheets("data/G04ResultsDetail2004-11-02.xls")
```

```
## [1] "Instructions"      "Voter Turnout"    "President"
## [4] "House of Rep"      "Co Clerk"         "Co Reg Deeds"
## [7] "Co Public Defender" "Co Comm 1"        "Co Comm 3"
## [10] "Co Comm 5"         "Co Comm 7"        "St Bd of Ed 2"
## [13] "St Bd of Ed 4"     "Legislature 5"    "Legislature 7"
## [16] "Legislature 9"     "Legislature 11"   "Legislature 13"
```

## [19] "Legislature 23"	"Legislature 31"	"Legislature 39"
## [22] "MCC 1"	"MCC 2"	"MCC 3"
## [25] "MCC 4"	"OPPD"	"MUD"
## [28] "NRD 3"	"NRD 5"	"NRD 7"
## [31] "NRD 9"	"OPS 2"	"OPS 4"
## [34] "OPS 6"	"OPS 8"	"OPS 10"
## [37] "OPS 11"	"OPS 12"	"ESU 2"
## [40] "ESU 3"	"Arlington Sch 24"	"Bennington Sch 59"
## [43] "Elkhorn Sch 10"	"Fremont Sch 1"	"Ft Calhoun Sch 3"
## [46] "Gretna Sch 37"	"Millard Sch 17"	"Ralston Sch 54"
## [49] "Valley Sch 33"	"Waterloo Sch 11"	"Bennington Mayor"
## [52] "Elkhorn Mayor"	"Valley Mayor"	"Ralston Mayor"
## [55] "Ralston Library Bd"	"Bennington City Cnc 1"	"Bennington City Cnc 2"
## [58] "Elkhorn City Cnc A"	"Elkhorn City Cnc B"	"Elkhorn City Cnc C"
## [61] "Ralston City Cnc 1"	"Ralston City Cnc 2"	"Ralston City Cnc 6"
## [64] "Waterloo Bd Trustees"	"Valley City Cnc"	"Amendment 1"
## [67] "Amendment 2"	"Amendment 3"	"Amendment 4"
## [70] "Initiative 417"	"Initiative 418"	"Initiative 419"
## [73] "Initiative 420"		

1.9 Using the read_excel function, read the Voter Turnout sheet

1.9.1 from the data/G04ResultsDetail2004-11-02.xls

1.9.2 Assign the data to the voter_turnout_df1

1.9.3 The header is in the second row, so make sure to skip the first row

1.9.4 Examine the structure of voter_turnout_df1 using str()

```
voter_turnout_df1 <- read_excel("data/G04ResultsDetail2004-11-02.xls", sheet="Voter Turnout", skip=1)
str(voter_turnout_df1)
```

```
## tibble [342 x 4] (S3: tbl_df/tbl/data.frame)
## $ Ward Precinct : chr [1:342] "01-01" "01-02" "01-03" "01-04" ...
## $ Ballots Cast : num [1:342] 421 443 705 827 527 323 358 410 440 500 ...
## $ Registered Voters: num [1:342] 678 691 1148 1308 978 ...
## $ Voter Turnout : num [1:342] 0.621 0.641 0.614 0.632 0.539 ...
```

1.10 Using the `read_excel()` function, read the Voter Turnout sheet

1.10.1 from `data/G04ResultsDetail2004-11-02.xls`

1.10.2 Skip the first two rows and manually assign the columns using `col_names`

1.10.3 Use the names “ward_precint”, “ballots_cast”, “registered_voters”, “voter_turnout”

1.10.4 Assign the data to the `voter_turnout_df2`

1.10.5 Examine the structure of `voter_turnout_df2` using `str()`

```
voter_turnout_df2 <- read_excel("data/G04ResultsDetail2004-11-02.xls", sheet="Voter Turnout", skip=2, col_types="text")
str(voter_turnout_df2)
```

```
## tibble [342 x 4] (S3: tbl_df/tbl/data.frame)
## $ ward_precint      : chr [1:342] "01-01" "01-02" "01-03" "01-04" ...
## $ ballots_cast      : num [1:342] 421 443 705 827 527 323 358 410 440 500 ...
## $ registered_voters: num [1:342] 678 691 1148 1308 978 ...
## $ voter_turnout     : num [1:342] 0.621 0.641 0.614 0.632 0.539 ...
```

1.11 Load the DBI library

```
library(DBI)
library(RSQLite)
```

1.12 Create a database connection to `data/tidynomicon/example.db` using the `dbConnect()` function

1.12.1 The first argument is the database driver which in this case is `RSQLite::SQLite()`

1.12.2 The second argument is the path to the database file

1.12.3 Assign the connection to db variable

```
driver1 <- RSQLite::SQLite()
db <- dbConnect(driver1, "data/tidynomicon/example.db")
```

1.13 Query the Person table using the `dbGetQuery` function and the

1.13.1 `SELECT * FROM PERSON;` SQL statement

1.13.2 Assign the result to the `person_df` variable

1.13.3 Use `head()` to look at the first few rows of the `person_df` dataframe

```
person_df <- dbGetQuery(db,"SELECT * FROM PERSON;",stringsAsFactors=FALSE)
head(person_df)
```

```
##   person_id personal_name family_name
## 1      dyer      William      Dyer
## 2       pb       Frank    Pabodie
## 3     lake    Anderson      Lake
## 4      roe    Valentina    Roerich
## 5 danforth      Frank    Danforth
```

1.14 List the tables using the dbListTables() function

1.14.1 Assign the result to the table_names variable

```
table_names <- dbListTables(db)
```

1.15 Read all of the tables at once using the lapply function and assign the result to the tables variable

1.15.1 Use table_names, dbReadTable, and conn = db as arguments

1.15.2 Print out the tables

```
tables <- lapply(table_names, dbReadTable, conn = db)
```

```
## Warning in result_fetch(res@ptr, n = n): Column 'reading': mixed type, first
## seen values of type real, coercing other values of type string
```

```
tables
```

```
## [[1]]
##   visit_id person_id quantity reading
## 1      619      dyer      rad    9.82
## 2      619      dyer      sal    0.13
## 3      622      dyer      rad    7.80
## 4      622      dyer      sal    0.09
## 5      734       pb      rad    8.41
## 6      734     lake      sal    0.05
## 7      734       pb      temp -21.50
## 8      735       pb      rad    7.22
## 9      735    <NA>      sal    0.06
## 10     735    <NA>      temp -26.00
## 11     751       pb      rad    4.35
## 12     751       pb      temp -18.50
## 13     751     lake      sal    0.00
## 14     752     lake      rad    2.19
## 15     752     lake      sal    0.09
```

```
## 16      752      lake      temp -16.00
## 17      752      roe       sal   41.60
## 18      837      lake      rad    1.46
## 19      837      lake      sal    0.21
## 20      837      roe       sal   22.50
## 21      844      roe       rad   11.25
##
## [[2]]
##   person_id personal_name family_name
## 1      dyer      William      Dyer
## 2      pb       Frank       Pabodie
## 3      lake      Anderson     Lake
## 4      roe       Valentina    Roerich
## 5  danforth      Frank       Danforth
##
## [[3]]
##   site_id latitude longitude
## 1    DR-1   -49.85   -128.57
## 2    DR-3   -47.15   -126.72
## 3   MSK-4   -48.87   -123.40
##
## [[4]]
##   visit_id site_id visit_date
## 1      619    DR-1 1927-02-08
## 2      622    DR-1 1927-02-10
## 3      734    DR-3 1930-01-07
## 4      735    DR-3 1930-01-12
## 5      751    DR-3 1930-02-26
## 6      752    DR-3      <NA>
## 7      837   MSK-4 1932-01-14
## 8      844    DR-1 1932-03-22
```

1.16 Use the `dbDisconnect` function to disconnect from the database

```
dbDisconnect(db)
```

1.17 Import the `jsonlite` library

```
library(jsonlite)
```

1.18 Convert the `scores_df` dataframe to JSON using the `toJSON()` function

```
toJSON(scores_df)
```

```
## [{"Count":10,"Score":200,"Section":"Sports"}, {"Count":10,"Score":205,"Section":"Sports"}, {"Count":20
```

1.19 Convert the scores dataframe to JSON using the toJSON() function with the pretty=TRUE option

```
toJSON(scores_df,pretty = TRUE)
```

```
## [  
##   {  
##     "Count": 10,  
##     "Score": 200,  
##     "Section": "Sports"  
##   },  
##   {  
##     "Count": 10,  
##     "Score": 205,  
##     "Section": "Sports"  
##   },  
##   {  
##     "Count": 20,  
##     "Score": 235,  
##     "Section": "Sports"  
##   },  
##   {  
##     "Count": 10,  
##     "Score": 240,  
##     "Section": "Sports"  
##   },  
##   {  
##     "Count": 10,  
##     "Score": 250,  
##     "Section": "Sports"  
##   },  
##   {  
##     "Count": 10,  
##     "Score": 265,  
##     "Section": "Regular"  
##   },  
##   {  
##     "Count": 10,  
##     "Score": 275,  
##     "Section": "Regular"  
##   },  
##   {  
##     "Count": 30,  
##     "Score": 285,  
##     "Section": "Sports"  
##   },  
##   {  
##     "Count": 10,  
##     "Score": 295,  
##     "Section": "Regular"  
##   },  
##   {
```



```

##      "Count": 10,
##      "Score": 300,
##      "Section": "Regular"
##    },
##    {
##      "Count": 20,
##      "Score": 300,
##      "Section": "Sports"
##    },
##    {
##      "Count": 10,
##      "Score": 305,
##      "Section": "Sports"
##    },
##    {
##      "Count": 10,
##      "Score": 305,
##      "Section": "Regular"
##    },
##    {
##      "Count": 10,
##      "Score": 310,
##      "Section": "Regular"
##    },
##    {
##      "Count": 10,
##      "Score": 310,
##      "Section": "Sports"
##    },
##    {
##      "Count": 20,
##      "Score": 320,
##      "Section": "Regular"
##    },
##    {
##      "Count": 10,
##      "Score": 305,
##      "Section": "Regular"
##    },
##    {
##      "Count": 10,
##      "Score": 315,
##      "Section": "Sports"
##    },
##    {
##      "Count": 20,
##      "Score": 320,
##      "Section": "Regular"
##    },
##    {
##      "Count": 10,
##      "Score": 325,
##      "Section": "Regular"
##    },
##    },

```

```

## {
##   "Count": 10,
##   "Score": 325,
##   "Section": "Sports"
## },
## {
##   "Count": 20,
##   "Score": 330,
##   "Section": "Regular"
## },
## {
##   "Count": 10,
##   "Score": 330,
##   "Section": "Sports"
## },
## {
##   "Count": 30,
##   "Score": 335,
##   "Section": "Sports"
## },
## {
##   "Count": 10,
##   "Score": 335,
##   "Section": "Regular"
## },
## {
##   "Count": 20,
##   "Score": 340,
##   "Section": "Regular"
## },
## {
##   "Count": 10,
##   "Score": 340,
##   "Section": "Sports"
## },
## {
##   "Count": 30,
##   "Score": 350,
##   "Section": "Regular"
## },
## {
##   "Count": 20,
##   "Score": 360,
##   "Section": "Regular"
## },
## {
##   "Count": 10,
##   "Score": 360,
##   "Section": "Sports"
## },
## {
##   "Count": 20,
##   "Score": 365,
##   "Section": "Regular"

```

```
## },
## {
##   "Count": 20,
##   "Score": 365,
##   "Section": "Sports"
## },
## {
##   "Count": 10,
##   "Score": 370,
##   "Section": "Sports"
## },
## {
##   "Count": 10,
##   "Score": 370,
##   "Section": "Regular"
## },
## {
##   "Count": 20,
##   "Score": 375,
##   "Section": "Regular"
## },
## {
##   "Count": 10,
##   "Score": 375,
##   "Section": "Sports"
## },
## {
##   "Count": 20,
##   "Score": 380,
##   "Section": "Regular"
## },
## {
##   "Count": 10,
##   "Score": 395,
##   "Section": "Sports"
## }
## ]
```

2 Session Info

```
sessionInfo()
```

```
## R version 4.2.2 (2022-10-31 ucrt)
## Platform: x86_64-w64-mingw32/x64 (64-bit)
## Running under: Windows 10 x64 (build 22621)
##
## Matrix products: default
##
## locale:
## [1] LC_COLLATE=English_United States.utf8
## [2] LC_CTYPE=English_United States.utf8
```

```

## [3] LC_MONETARY=English_United States.utf8
## [4] LC_NUMERIC=C
## [5] LC_TIME=English_United States.utf8
##
## attached base packages:
## [1] stats      graphics  grDevices  utils      datasets  methods   base
##
## other attached packages:
## [1] jsonlite_1.8.3 RSQLite_2.2.19 DBI_1.1.3      readxl_1.4.1
##
## loaded via a namespace (and not attached):
## [1] Rcpp_1.0.9      knitr_1.41      magrittr_2.0.3  bit_4.0.4
## [5] rlang_1.0.6     fastmap_1.1.0   fansi_1.0.3     blob_1.2.3
## [9] stringr_1.4.1   tools_4.2.2     xfun_0.34       utf8_1.2.2
## [13] cli_3.4.1       htmltools_0.5.3 yaml_2.3.6      bit64_4.0.5
## [17] digest_0.6.30   tibble_3.1.8    lifecycle_1.0.3 vctrs_0.5.0
## [21] cachem_1.0.6    memoise_2.0.1   glue_1.6.2      evaluate_0.18
## [25] rmarkdown_2.18  stringi_1.7.8   compiler_4.2.2  pillar_1.8.1
## [29] cellranger_1.1.0 pkgconfig_2.0.3

```