

02 Data Import and Manipulation

June 19, 2020

Contents

1	Introduction	1
2	Read and Write Data Files	1
2.1	CSV Files	1
2.1.1	Read CSV Files	1
2.1.2	Write CSV Files	3
2.2	XLSX Files	3
2.2.1	Read XLSX Files	3
2.2.2	Write XLSX Files	4
3	Database	5
3.1	MySQL	5
3.1.1	Connect to MySQL Server	5
3.1.2	Reading Tables by SQL Commands	6
3.2	SQLite	7
4	Packages for Data Handling	7

1 Introduction

In this notes, examples concerning data processing are presented

2 Read and Write Data Files

2.1 CSV Files

2.1.1 Read CSV Files

The Syntax

```
read.csv(file, header = TRUE, sep = ",", quote = "\"",  
         dec = ".", fill = TRUE, comment.char = "", ...)
```

Read a CSV File with Default Options

- `read.csv` returns a `data.frame`

The comma-separated values (CSV) file used here for demonstration has a content:

```
"", "BPchange", "Dose", "Run", "Treatment", "Animal"
"1", 0.5, 6.25, "C1", "Control", "R1"
"2", 4.5, 12.5, "C1", "Control", "R1"
"3", 10, 25, "C1", "Control", "R1"
"4", 26, 50, "C1", "Control", "R1"
"5", 37, 100, "C1", "Control", "R1"
"6", 32, 200, "C1", "Control", "R1"
```

```
[1]: # Getting data from file "rabbit.csv"
rabbit_sample <- read.csv("datasets/rabbit.csv")

# Print the class of the variable rabbit_sample
print(class(rabbit_sample))

# Printing first few lines of the dataframe
head(rabbit_sample)
```

```
[1] "data.frame"
```

A data.frame: 6 × 6

	X	BPchange	Dose	Run	Treatment	Animal
	<int>	<dbl>	<dbl>	<chr>	<chr>	<chr>
1	1	0.5	6.25	C1	Control	R1
2	2	4.5	12.50	C1	Control	R1
3	3	10.0	25.00	C1	Control	R1
4	4	26.0	50.00	C1	Control	R1
5	5	37.0	100.00	C1	Control	R1
6	6	32.0	200.00	C1	Control	R1

Not Assuming the First Row in the CSV File is Labels

- The column labels will be “V1”, “V2”, etc...

```
[2]: # Getting data from file "rabbit.csv"
rabbit_sample <- read.csv("datasets/rabbit.csv", header = FALSE)

# Printing first few lines of the dataframe
head(rabbit_sample)
```

A data.frame: 6 × 6

	V1	V2	V3	V4	V5	V6
	<int>	<chr>	<chr>	<chr>	<chr>	<chr>
1	NA	BPchange	Dose	Run	Treatment	Animal
2	1	0.5	6.25	C1	Control	R1
3	2	4.5	12.5	C1	Control	R1
4	3	10	25	C1	Control	R1
5	4	26	50	C1	Control	R1
6	5	37	100	C1	Control	R1

Using Custom Column Names

- The rule is the same as rows.

```
[3]: # Getting data from file "rabbit.csv"
rabbit_sample <- read.csv("datasets/rabbit.csv", col.names = c("A", "B", "C", "D", "E", "F"))

# Printing first few lines of the dataframe
head(rabbit_sample)
```

A data.frame: 6 × 6

	A	B	C	D	E	F
	<int>	<dbl>	<dbl>	<chr>	<chr>	<chr>
1	1	0.5	6.25	C1	Control	R1
2	2	4.5	12.50	C1	Control	R1
3	3	10.0	25.00	C1	Control	R1
4	4	26.0	50.00	C1	Control	R1
5	5	37.0	100.00	C1	Control	R1
6	6	32.0	200.00	C1	Control	R1

2.1.2 Write CSV Files

The Syntax

```
write.csv(x, file = "", quote = TRUE, eol = "\n", na = "NA", row.names = TRUE, fileEncoding = "UTF-8")
```

- A more general implementation is `write.table`. Check `?write.table` for more detail.

Simple Use of `write.csv`

```
[4]: # Write the data.frame to "testing.csv"
write.csv(rabbit_sample, "datasets/testing.csv")
```

The file “testing.csv” contains:

```
"", "A", "B", "C", "D", "E", "F"
"1", 1, 0.5, 6.25, "C1", "Control", "R1"
"2", 2, 4.5, 12.5, "C1", "Control", "R1"
"3", 3, 10, 25, "C1", "Control", "R1"
"4", 4, 26, 50, "C1", "Control", "R1"
"5", 5, 37, 100, "C1", "Control", "R1"
"6", 6, 32, 200, "C1", "Control", "R1"
```

2.2 XLSX Files

2.2.1 Read XLSX Files

The Syntax

```
read.xlsx(
  file,
  sheetIndex,
  sheetName = NULL,
  rowIndex = NULL,
  startRow = NULL,
  endRow = NULL,
```

```

colIndex = NULL,
as.data.frame = TRUE,
header = TRUE,
colClasses = NA,
keepFormulas = FALSE,
encoding = "unknown",
password = NULL,
...
)

```

Here ... are other arguments to 'data.frame', for example 'stringsAsFactors'

Read a CSV File with Default Options

- `xlsx::read.xlsx` returns a data.frame
- The xlsx file used for demonstration contains the following data:

	A	B	C	D	E	F	
1		Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species	
2	1	5.1	3.5	1.4	0.2	setosa	
3	2	4.9	3	1.4	0.2	setosa	
4	3	4.7	3.2	1.3	0.2	setosa	
5	4	4.6	3.1	1.5	0.2	setosa	
6	5	5	3.6	1.4	0.2	setosa	
7	6	5.4	3.9	1.7	0.4	setosa	
8	7	4.6	3.4	1.4	0.3	setosa	

```

[5]: # Loading the xlsx library
library(xlsx)

# Get the iris dataset from iris.xlsx, the second argument is the index of the
↪ worksheet in the xlsx file.
iris_table <- xlsx::read.xlsx("datasets/iris.xlsx", 1)

# Print first few lines of the table
head(iris_table)

```

		NA.	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
		<chr>	<dbl>	<dbl>	<dbl>	<dbl>	<chr>
A data.frame: 6 × 6	1	1	5.1	3.5	1.4	0.2	setosa
	2	2	4.9	3.0	1.4	0.2	setosa
	3	3	4.7	3.2	1.3	0.2	setosa
	4	4	4.6	3.1	1.5	0.2	setosa
	5	5	5.0	3.6	1.4	0.2	setosa
	6	6	5.4	3.9	1.7	0.4	setosa

2.2.2 Write XLSX Files

The Syntax

```
write.xlsx(
  x,
  file,
  sheetName = "Sheet1",
  col.names = TRUE,
  row.names = TRUE,
  append = FALSE,
  showNA = TRUE,
  password = NULL
)
```

Write a CSV File with Default Options

```
[6]: # Staff table to export
staff_table = data.frame(
  ID = c(1L, 2L, 3L, 4L),
  Name = c("Tom", "Ann", "Peter", "Kelly"),
  Phone = c(73490245L, 77990904L, 47876737L, 35146136L)
)

# Write the xlsx file to the file namely staff_table.xlsx
xlsx::write.xlsx(staff_table, "datasets/staff_table.xlsx", append = FALSE)
```

- The output xlsx file:

	A	B	C	D	E
1		ID	Name	Phone	
2	1		1 Tom	73490245	
3	2		2 Ann	77990904	
4	3		3 Peter	47876737	
5	4		4 Kelly	35146136	
6					
7					

3 Database

3.1 MySQL

3.1.1 Connect to MySQL Server

To connection to MySQL servers, we need to include two libraries:

```
[7]: # Include libraries for MySQL connection
library(DBI)
library(RMySQL)
```

Then we connect to database namely “classicmodels” on the MySQL server at 127.0.0.1 using function DBI::dbConnect:

```
[8]: con <- DBI::dbConnect(RMySQL::MySQL(),
                           dbname="classicmodels",
                           host="127.0.0.1",
                           user="alan",
                           password="password")
```

Now the connection pipe is stored in object `con`. To list tables, we could use `DBI::dbListTables`.

```
[9]: DBI::dbListTables(conn = con)
```

1. 'customers' 2. 'employees' 3. 'offices' 4. 'orderdetails' 5. 'orders' 6. 'payments' 7. 'productlines'
8. 'products'

3.1.2 Reading Tables by SQL Commands

Using `DBI::dbGetQuery`

```
[10]: select_result <- DBI::dbGetQuery(conn = con, statement = "
      select customerNumber,customerName,phone from customers;
    ")

cat("\nThe type of the output object:", class(select_result) ,". \n")

head(select_result)
```

The type of the output object: `data.frame` .

	customerNumber	customerName	phone
	<int>	<chr>	<chr>
1	103	Atelier graphique	40.32.2555
2	112	Signal Gift Stores	7025551838
3	114	Australian Collectors, Co.	03 9520 4555
4	119	La Rochelle Gifts	40.67.8555
5	121	Baane Mini Imports	07-98 9555
6	124	Mini Gifts Distributors Ltd.	4155551450

A data.frame: 6 × 3

Using `DBI::dbSendQuery` and `DBI::dbFetch`

```
[11]: select_result_raw <- DBI::dbSendQuery(conn = con, statement = "
      select customerNumber,customerName,state from customers;
    ")

select_result <- DBI::dbFetch(select_result_raw)

cat("\nThe type of the output object:", class(select_result) ,". \n")

head(select_result)
```

The type of the output object: `data.frame` .

		customerNumber	customerName	state
		<int>	<chr>	<chr>
A data.frame: 6 × 3	1	103	Atelier graphique	NA
	2	112	Signal Gift Stores	NV
	3	114	Australian Collectors, Co.	Victoria
	4	119	La Rochelle Gifts	NA
	5	121	Baane Mini Imports	NA
	6	124	Mini Gifts Distributors Ltd.	CA

3.2 SQLite

4 Packages for Data Handling