

sqldf - package in R

An 'R' package for running SQL statements on data frames

# Load the package

library(sqldf)

Titanic # data set in R

This data set provides information on the fate of passengers on the fatal maiden voyage of the ocean liner 'Titanic', summarized according to economic status (class), sex, age & survival.

Format

A 4 dimensional array resulting from cross-tabulating 2201 observations on 4 variables. The variables & their levels are as follows.

No	Name	Levels
1	class	1st, 2nd, 3rd, crew
2	sex	Male, Female
3	Age	child, Adult
4	survived	No, yes.

## Examples:

These examples show how to use a variety of data frame manipulations in R without SQL and then again with SQL.

warpbreaks

# dataframe in R

	breaks	wool	tension
1	26	A	L
2	30	A	L
3	54	A	L
4	25	A	L
5	70	A	L
6	52	A	L

# head

returns first 6 observations

```
als ← head(warpbreaks)
```

```
als ← sqldf("select * from warpbreaks  
limit 6")
```

```
identical(als, als)
```

```
[1] TRUE
```

subset function

```
a2r ← subset(co2, grep("Qn", plant))
```

```
a2s ← subset("select * from co2 where  
plant like 'Qn%'")
```

```
all.equal(as.data.frame(a2r), a2s)  
[1] TRUE
```

co2 # dataframe

	Plant	Type	Treat ment	conc	uptake
1	Qn1	Quebec	nonchilled	95	16.0
2	Qn1	Quebec	"	175	30.4
3	Qn1	"	"	250	34.8
4	Qn1	"	"	350	35.9
5	Qn1	"	"	500	35.3
6	Qn1	"	"	675	39.2
7	Qn1	"	"	1000	39.7
8	Qn2	"	"	95	13.6

1/5

5/1 en

1/5.1

③

data(farms, package = "MASS")

a3r ← subset(farms, Manag %in%  
c("BF", "HF"))

a3s ← sqlf("select \* from farms where  
Manag in ('BF', 'HF')")

row.names(a3r) ← NULL

identical(a3r, a3s)

[1] TRUE

a3r ← subset(farms, Manag %in%  
c("BF", "HF"))

a3r

	Mois	Manag	Use	Manure
2	M1	BF	U2	C2
5	M1	HF	U1	C2
6	M1	HF	U2	C2
7	M1	HF	U3	C3
8	M5	HF	U3	C3
9	M4	HF	U1	C1
10	M2	BF	U1	C1
10	M1	BF	U3	C1

to specify  
multiple  
values  
in  
select



asr ← subset (farms, mois = 'M1')

ass ← ~~subset~~ ("select \* from farms where  
sold#

mois = 'M1' ", row.names = TRUE

identical (asr, ass)

[1] FALSE

	asr			
	mois	manag	use	manure
1	M1	SF	U2	C4
2	M1	BF	U2	C2
5	M1	HF	U1	C2
6	M1	HF	U2	C2
7	M1	HF	U3	C3
11	M1	BF	U3	C1
18	M1	NM	U1	C0

# should be identical but getting  
false

$\text{row.names} = \text{TRUE}$   
 $\text{a6r} \leftarrow \text{subset}(\text{farms}, \text{Mois} = 'M2')$   
 $\text{a6s} \leftarrow \text{sql}(\text{"select * from farms where"},$   
 $\text{Mois} = 'M2', \text{row.names} = \text{TRUE})$

$\text{identical}(\text{a6r}, \text{a6s})$

a6r

	Mois	Manag	Use	Manure
3	M2	SF	U2	C4
4	M2	SF	U2	C4
10	M2	BF	U1	C1
17	M2	NM	U1	C0

a6s

	Mois	Manag	Use	Manure
3	M2	SF	U2	C4
4	M2	SF	U2	C4
10	M2	BF	U1	C1
17	M2	NM	U1	C0

(1) FALSE

? should be identical

Rbind & Union all

a78 ← rbind(a58, a68)

a75 ← sqlDF("select \* from a55 Union  
all select \* from a65")

a77

	Mois	Manag	Use	Manure
1	M1	SF	U2	C4
2	M1	BF	U2	C2
5	M1	HF	U1	C2
6	M1	HF	U2	C2
7	M1	HF	U3	C3
11	M1	BF	U3	C3
18	M1	NM BF	U1	C0
3	M2	SF NM	U2	C6
4	M2	SF	U2	C4
10	M2	BF	U1	C4
17	M2	NM	U1	C0

a58 = 7  
 a68 = 4

} 11 rows retrieved

Note - `saldf` drops the unused levels of  
factors but `sbind` does not; however,  
# all data is the same and the  
other columns are identical

```
rownames(a78) ← NULL # set it to  
identical(a78[-1], a78[-1]) NULL  
rownames  
[1] [TRUE] TRUE
```

Aggregate -

```
a88 ← aggregate[1:2], iris[5], mean)  
a88 ← saldf('select species avg("sepal.  
length") 'sepal.Length',  
avg("sepal.Width") 'sepal.Width' from  
iris group by species')  
all.equal(a88, a88)  
[TRUE]
```

a88	Species	sepal.Length	sepal.Width
1	setosa	5.006	3.428
2	versicolor	5.936	2.770
3	virginica	6.588	2.974



485 ← sqrt

a85

Species	sepal.Length	Sepal.Width
setosa	5.006	3.428
Versicolor	5.936	2.770
Virginica	6.588	2.974

all.equal(a85, a88) (a98, a95)

[1] TRUE

[1] "Attributes: <component 'row.names':  
modes=character, numeric >"

[2] "Attributes: <component 'row.names':  
target is character, current is numeric >"

by = avg conc and total capture by plant  
and type

```
aqs ← do.call(rbind, by(iris, iris[5], function(x)  
  with(x, data.frame(species = Species[1],  
    mean.Sepal.Length = mean(sepal.  
      Length),  
    mean.Sepal.Width = mean(sepal.  
      width),  
    mean.Sepal.ratio = mean(sepal.Length /  
      sepal.Width))))))
```

```
row.names(aqs) ← NULL
```

```
aqs ← sqlDF('select species, avg("sepal.Length")  
  'mean.Sepal.Length',  
  avg("sepal.Width") 'mean.Sepal.Width',  
  avg("sepal.Length" / "sepal.Width")  
  'mean.Sepal.ratio' from iris  
  group by species')
```

```
all.equal(aqs, aqs)
```

```
aqs
```

Species	mean.Sepal.Length	m.S.W	Sepal.ratio
setosa	5.006	3.428	1.470188
versicolor	5.936	2.770	2.160402
virginica	6.588	2.974	2.230453

Species	m.S.L	m.S.W	sep
Setosa → same	5.006	3.428	1.47
versicolor "	5.936	2.770	2.16
virginica "	6.588	2.974	2.23

all-equal (a98, a95)

"Attributes!"

"Attributes!"

Reading excel files in R

24-2

3 packages

```
install.packages("XLconnect")
```

```
install.packages("xlsx")
```

```
install.packages("gdata")
```

using XLconnect

```
getwd()
```

setwd() → for setting the path of your file

```
install.packages("rJava")
```

```
any(grepl("rJava", installed.packages()))
```

```
[1] TRUE
```

```
any(grepl("methods", installed.packages()))
```

```
[1] TRUE
```

```
library("rJava")
```

Error: onLoad failed in loadNamespace() for rJava details:

call: fun(libname, pkgname)

error: Java\_HOME cannot be determined from the Registry

Error: package or name space load failed for 'rJava'



```
Sys.setenv(JAVA_HOME = 'C:\\Program Files(x86)\\Java\\jre7')
```

liblaly (xJava)

Error: onload failed in loadNamespace () for 'xJava', details:

call: 'inDL (x, as.logical(local), as.logical(now), ...)

error: unable to load shared object

'C:/users/sidhar/documents/R/win-library/3.3/xJava/libs/i386/xJava.dll':

LoadLibrary failure: The specified module could not be found

Error: package or namespace load failed for 'xJava'

Environment Variable  
Path Setting

```
Sys.setenv(JAVA_HOME = 'C:\Program Files \\  
Java \ jre1.8.0-121')
```

```
library(xJava) = no quotes
```

```
library("xlconnect")
```

```
wb ← loadWorkbook("simple.xlsx")
```

```
myDF ← readWorksheet(wb, sheet = "sheet1",  
                      $  
                      header = TRUE)
```

Example 1:

mtcars xlsx file from demofiles  
subfolder of package XLconnect.

```
demoExcelFile ← system.file("demofiles /  
mtcars.xlsx", package = "XLconnect")
```

```
# Load workbook
```

```
wb ← loadwworkbook(demoExcelFile)
```

```
# Read worksheet 'mtcars' (providing no  
specific area bounds)
```

```
# with default header: TRUE)
```

```
data ← readWorksheet(wb, sheet = "mtcars")
```

Example 2:

```
demoExcelFile ← system.file("demofiles / mtcars.  
xlsx",  
package = XLconnect")
```

```
wb ← loadWorkbook(demoExcelFile)
```

```
data ← readWorksheet(wb, sheet = "mtcars",  
startRow = 1, startCol = 3,  
endRow = 15, endCol = 8)
```

# Providing area bounds using the region argument; with default header: TRUE

```
data ← readWorksheet(wb, sheet = "mtcars",  
                    region = "CH4",  
                    region = "C1:H15")
```

Example 4

```
df ← readWorksheet(wb, sheet = "conversion",  
                  header = TRUE,  
                  colTypes = c(XLC$DATA_TYPE.NUMERIC,  
                              XLC$DATA_TYPE.DATETIME,  
                              XLC$DATA_TYPE.BOOLEAN),  
                  forceConversion = TRUE,  
                  dateTimeFormat = "%Y-%m-%d %H:%M")
```

# Read worksheet 'Conversion' with the  
pre specified column types.

Note. in the worksheet all data was entered  
as strings.

forceConversion = TRUE is used to force  
conversion from string into less  
generic data types Numeric, DateTime  
& Boolean.



Example 5

```
load  
wb ← Workbook(demoExcelFile)
```

```
data ← readWorksheet(wb, sheet = "mtcars",  
keep = c(1, 3, 5))
```

# Read the columns 1, 3 and 5 from the  
sheet 'mtcars' (with default header =  
TRUE)

Eg:

```
excelFile ← system.file("demoFiles/conversion.xlsx",  
package = "XLconnect")
```

```
wb ← loadWorkbook(excelFile)
```



## Package - gdata

```
read.xls(xls, sheet=1, verbose=FALSE, pattern,  
na.strings=c("NA", "#DIV/0!"),  
method=c("csv", "tsv", "tab"), perl="perl")
```

**xls**: Path to the Microsoft Excel file - supports "http://", "https://", and "ftp://" URLs.

**sheet**: name or number of the worksheet to read

**verbose**: logical flag indicating whether details should be printed as the file is processed.

**pattern**: if specified, then skip all lines before the first containing this string

**perl**: name of the perl executable to be called.

**method**: intermediate file format, "csv" for comma-separated and "tab" for tab-separated.

**na.strings**: a character vector of strings which are to be interpreted as 'NA' values.

**blank.lines.skip**: logical flag indicating whether blank lines in the original file should be ignored.

read xls = returns a data frame.

Note: Either a working version of Perl must be present in the executable search path or the exact path of the perl executable must be provided via the perl argument.



Examples:

```
xlsfile ← file.path(path.package("gdata"),  
  "xls", "ixis.xls")
```

# ixis.xls is included in the gdata package for use as an example.

```
xlsfile
```

```
ixis ← read.xls(xlsfile) # defaults to csv format
```

```
ixis ← read.xls(xlsfile, method = "csv")
```

```
# specify csv format
```

```
ixis ← read.xls(xlsfile, method = "tab")
```

```
# specify tab format
```

Examples:

```
mydf ← read.xls("simple.xls", sheet = 2,  
  header = TRUE)
```

	mydf			
	sid	sname	marks	age
1	1	A	89	21
2	2	B	98	22
3	3	C	78	23
4	4	D	88	24
5	5	E	68	25



```
df ← read.xls("some file" sheet = 2,
              perl = "C:/perl/bin/perl.exe")
```

df

Example 2:

```
ixis ← read.xls(xisfile, perl = "C:/
               perl/bin/perl.exe")
```

# Example specifying exact perl path to  
default MS-Windows install of  
ActiveState perl.

Unix systems:

```
ixis ← read.xls(xisfile, perl = "/usr/bin/
               perl")
```

Finding Perl

```
perl ← gdata ::: findperl("perl")
```

```
ixis ← read.xls(xisfile, perl = perl)
```

# read.xls automatically calls findperl so  
this is rarely needed.)

# read xls file from url

```
nba.url ← "http://mgtclass.mgt.unm.edu/bose/  
Excel/Tutorial.05/Cases/NBA.xls"
```

```
nba ← read.xls(nba.url)
```

# read xls file ignoring all lines prior to first containing state.

```
crime.url ← "http://www.jrsainfo.org/jabg/  
state-dat2/Tribal_Data00.xls"
```

```
crime ← read.xls(crime.url, pattern = "state")
```

# use of xls2csv - open con, print two lines, close con.

```
con ← xls2csv(crime.url)
```

```
print(readLines(con, 2))
```

```
file.remove(summary(con) $ description)
```

# Examples demonstrating

```
exampleFile ← file.path(path.package('gdata',  
  'xls', 'ExampleExcelFile.xls'))
```

```
exampleFile2007 ← file.path(path.package('gdata',  
  'xls', 'ExampleExcelFile.xlsx'))
```

# Examples demonstrating selection of specific sheets from the example XLS file 'ExampleExcelFile.xls'

Eg: see the no. of names of sheets

```
sheetCount(exampleFile)  
if('xlsx' %in% xlsFormats())  
sheetCount(exampleFile2007)
```

# Names of the sheets

```
sheetNames(exampleFile)  
if('xlsx' %in% xlsFormats())  
sheetNames(exampleFile2007)
```

Examples:

```
data ← read.xls(exampleFile)
```

```
# default is first worksheet
```

to),

```
data ← read.xls(exampleFile, sheet = 2)
```

```
# second worksheet - by no:
```

```
data ← read.xls(exampleFile, sheet = "sheet  
second",
```

```
v = TRUE)
```

```
## and by name
```

```
# load the third worksheet, skipping the  
first two non-data lines
```

```
if('xlsx' %in% xlsFormats())
```

```
# if xlsx is supported
```

```
data ← read.xls(exampleFile2007, sheet = sheet  
with initial text", skip = 2)
```

load a file containing data & column  
names using latin-1 characters

```
latinFile ← file.path(path.packages('gdata'),  
'xls', 'latin-1.xls')
```

```
latin1 ← read.xls(latinFile, fileEncoding =  
"latin1")
```

```
colnames(latin1)
```



data = read.xls("emp.xlsx")

data

eid	ename	dep
501	Saritha	CSE
502	Haritha	IT
503	Vanitha	CSE
504	Hari	CIVIL
505	Veena	MECH

df = read.xls("emp.xlsx", sheet = 1,  
perl = "C:/perl/bin/perl.exe")

df = read.xls("emp.xlsx",  
perl = "C://perl/bin/perl.exe")

same: data o/p

perl = "C:/perl/bin/perl.exe")

df = read.xls("emp.xlsx", 1,  
perl = " ")

df = read.xls("emp.xlsx", sheet = 1)

## read.xlsx functions

```
data1 ← read.xlsx("emp.xlsx")
```

Error in read.xlsx("emp.xlsx")

Please provide a sheet name or a sheet index

```
data1 ← read.xlsx("emp.xlsx", sheet = 1)
```

Error in

argument 2 matches multiple formal arguments.

```
data1 ← read.xlsx("emp.xlsx", sheet = "sheet 1")
```

same error.

```
data1 ← read.xlsx("emp.xlsx", 1)
```

data1

eid	ename	dep
501	Salitha	CSE
502	Haritha	IT
503	Vanitha	CSE
504	Hari	CIVIL
505	Veena	MECH

Read R output in excel

```
head(mtcars)
```

```
ht ← head(mtcars)
```

```
write.table(ht, "mt.csv", sep = ",",  
            header = TRUE)
```

Error unused argument 'header = TRUE'

```
write.table(ht, "mt.csv", sep = ",")
```

```
write.table(ht, "mt.csv", sep = ",",  
            row.names = FALSE)
```

using csv file:-

```
write.csv(ht, "mt2.csv", row.names =  
          FALSE)
```

Examples:-

```
read.csv(" ")
```

```
write.csv(x, file = "foo.csv", row.names  
         = FALSE)
```

```
write.csv(x, file = "foo.csv")
```

Text book example:-

```
x ← cbind(rnorm(20), runif(20))
```

```
colnames(x) ← c("A", "B")
```

x

```
write.table(x, file = "myfile.csv", sep = ",",  
            row.names = FALSE)
```

a ← 1:10

b ← 10:1

cbind(a, b)

	a	b
1,	1	10
2,	2	9
3,	3	8

```
write.csv(cbind(a, b), file = "mycolumn1.csv",  
          row.names = FALSE)
```

```
write.csv(cbind(a, b), file = "mycolumn2.csv")
```