

Data Science with R

Writing Functions in R

Graham.Williams@togaware.com

13th July 2013

Visit <http://onepager.togaware.com/> for more OnePageR's.

The required packages for this module include:

```
library(rattle)
```

As we work through this module, new R commands will be introduced. Be sure to review the command's documentation and understand what the command does. You can ask for help using the `?` command as in:

```
?read.csv
```

We can obtain documentation on a particular package using the *help=* option of `library()`:

```
library(help=rattle)
```

This module is intended to be hands on. To learn effectively, you are encouraged to have R running (e.g., RStudio) and to run all the commands as they appear here. Check that you get the same output, and you understand the output. Try some variations. Explore.

Copyright © 2013 Graham J Williams. You can freely copy, distribute, transmit, adapt, or make commercial use of this module, as long as the attribution is retained and derivative work is provided under the same license.



1 Functions

```
mult10 <- function(x)
{
  if (is.character(x))
  {
    result <- apply(sapply(x, rep, 10), 2, paste, collapse=" ")
    names(result) <- NULL
  }
  else
  {
    result <- x * 10
  }

  return(result)
}
```

This page
is under
develop-
ment.

2 Function Calls

```
4 + 5
## [1] 9
"+"(4, 5)
## [1] 9
```

```
1 + 2 + 3 + 4 + 5
## [1] 15
Reduce("+", 1:5)
## [1] 15
```

```
cmd <- "1 + 2 + 3 + 4 + 5"
eval(parse(text=cmd))
## [1] 15
```

This page
is under
develop-
ment.

3 Flow Control

```
for (i in 0:4)
  for (j in 5:9)
    print(paste0(i, j))

## [1] "05"
## [1] "06"
## [1] "07"
## [1] "08"
....
```

This page
is under
develop-
ment.

4 Exercise Dataset: WeatherAUS

For our exercises we will use the **weatherAUS** dataset from **rattle** ([Williams, 2013](#)). We will essentially follow the template presented in the Models module.

```
ds <- read.csv(file="data/weatherAUS.csv")
```

```
dim(ds)
```

```
## [1] 66672    24
```

```
head(ds)
```

```
##           Date Location MinTemp MaxTemp Rainfall Evaporation Sunshine
## 1 2008-12-01   Albury    13.4    22.9      0.6          NA         NA
## 2 2008-12-02   Albury     7.4    25.1      0.0          NA         NA
## 3 2008-12-03   Albury    12.9    25.7      0.0          NA         NA
## 4 2008-12-04   Albury     9.2    28.0      0.0          NA         NA
## 5 2008-12-05   Albury    17.5    32.3      1.0          NA         NA
## 6 2008-12-06   Albury    14.6    29.7      0.2          NA         NA
....
```

```
tail(ds)
```

```
##           Date Location MinTemp MaxTemp Rainfall Evaporation Sunshine
## 66667 2012-11-24   Darwin    25.5    34.1      0.0          5.0         6.2
## 66668 2012-11-25   Darwin    24.4    35.7      0.2          4.8        11.7
## 66669 2012-11-26   Darwin    25.0    35.4      0.0          7.4        11.7
## 66670 2012-11-27   Darwin    26.5    35.9      0.0          8.0        10.3
## 66671 2012-11-28   Darwin    27.4    35.0      0.0          7.8         6.5
## 66672 2012-11-29   Darwin    24.8    33.5      3.4          7.4         4.7
....
```

```
str(ds)
```

```
## 'data.frame': 66672 obs. of  24 variables:
## $ Date       : Factor w/ 1826 levels "2007-11-01","2007-11-02",...: 397 ...
## $ Location   : Factor w/ 46 levels "Adelaide","Albany",...: 3 3 3 3 3 ...
## $ MinTemp    : num  13.4 7.4 12.9 9.2 17.5 14.6 14.3 7.7 9.7 13.1 ...
## $ MaxTemp    : num  22.9 25.1 25.7 28 32.3 29.7 25 26.7 31.9 30.1 ...
## $ Rainfall   : num  0.6 0 0 0 1 0.2 0 0 0 1.4 ...
## $ Evaporation: num  NA NA NA NA NA NA NA NA NA NA ...
....
```

```
summary(ds)
```

```
##           Date           Location           MinTemp           MaxTemp
## 2009-01-01:    46   Canberra: 1826   Min.      :-8.5   Min.      :-3.1
## 2009-01-02:    46    Sydney  : 1734   1st Qu.: 7.3   1st Qu.:17.6
## 2009-01-03:    46   Adelaide: 1583   Median :11.7   Median :22.1
## 2009-01-04:    46   Brisbane: 1583   Mean     :11.9   Mean     :22.7
## 2009-01-05:    46    Darwin  : 1583   3rd Qu.:16.6   3rd Qu.:27.5
## 2009-01-06:    46    Hobart  : 1583   Max.      :33.9   Max.      :48.1
....
```

5 Exercises: Prepare for Modelling

Following the template presented in the Models module, we continue with setting up some of the modelling parameters.

```
target <- "RainTomorrow"
risk    <- "RISK_MM"
dsname  <- "weather"

ds[target] <- as.factor(ds[[target]])
summary(ds[target])

## RainTomorrow
## No :50187
## Yes :15259
## NA's: 1226
....

vars      <- colnames(ds)
ignore    <- vars[c(1, 2, if (exists("risk")) which(risk==vars))]
vars      <- setdiff(vars, ignore)
(inputs   <- setdiff(vars, target))

## [1] "MinTemp"      "MaxTemp"      "Rainfall"     "Evaporation"
## [5] "Sunshine"     "WindGustDir"  "WindGustSpeed" "WindDir9am"
## [9] "WindDir3pm"   "WindSpeed9am" "WindSpeed3pm"  "Humidity9am"
## [13] "Humidity3pm"  "Pressure9am"  "Pressure3pm"   "Cloud9am"
....

nobs      <- nrow(ds)
dim(ds[vars])

## [1] 66672    21

(form <- formula(paste(target, "~ .")))
## RainTomorrow ~ .

set.seed(142)

length(train <- sample(nobs, 0.7*nobs))

## [1] 46670

length(test  <- setdiff(seq_len(nobs), train))

## [1] 20002
```

6 Exercise: varWeights()

The first exercise is to write a function to take a dataset and return probabilities associated with each input variable in the dataset, that relate to the correlation between the input variable and the target variable.

```
varw <- varWeights(form, ds)
```

We will use the R correlation functions to calculate the correlation between each column (variable) of the **data** frame and the values of the **target** vector. Below are some hints.

```
n1 <- ds[["Temp3pm"]]
c1 <- ds[["WindGustDir"]]
t1 <- ds[[target]]

cor(as.numeric(n1), as.numeric(t1), use="pairwise.complete.obs")
## [1] -0.1857

cor(as.numeric(c1), as.numeric(t1), use="pairwise.complete.obs")
## [1] 0.04414
```

The template for the function is:

```
varWeights <- function(formula, data)
{
  ...
}
```

The actual solution will produce the following output:

```
varWeights(form, ds)

##      Date      Location      MinTemp      MaxTemp      Rainfall
## 0.0028545 0.0004961 0.0210742 0.0336550 0.0544825
## Evaporation      Sunshine WindGustDir WindGustSpeed WindDir9am
## 0.0249688 0.1006235 0.0096812 0.0518932 0.0067119
....
```

It is time now to write that function.

7 Exercise: selectVars()

The next exercise is to write a function that will return **n** variables chosen at random from all of the variables in a **dataset**, but chosen with a probability proportional to the correlation of the target variable.

```
vars <- selectVars(form, ds, 3)
```

The **sample()** function might come in use for this function. Note the **prob=** argument of **sample**. We will, of course, also make use of the **varWeights()** function we defined previously.

The template for the function is:

```
selectVars <- function(formula, data, n)
{
  ...
}
```

The actual solution will produce the following output:

```
selectVars(form, ds, 3)
## [1] "Pressure3pm" "MaxTemp"      "RISK_MM"
selectVars(form, ds, 3)
## [1] "Cloud9am" "RISK_MM"  "Sunshine"
selectVars(form, ds, 3)
## [1] "WindDir9am" "Humidity3pm" "Cloud9am"
selectVars(form, ds, 3)
## [1] "Cloud3pm"    "Evaporation" "Humidity3pm"
```


8 Exercise: `wsrpart()`

This exercise is to write a function to build a subspace decision tree. The function `wsrpart()` (for weighted subspace `rpart`) will take a dataset (`data`) and return a decision tree (built using `rpart()`) that uses only a subspace of the variables available. The number of variables to use is, by default, $\text{trunc}(\log_2(n+1))$ (overridden by `nvars=`) and the variables are chosen according to the weighted selection implemented through `selectVars()`.

The idea is similar, but not identical to, the concept of random forests developed by Leo Breiman. Note that Breiman's original random forest paper (on which this idea is based) specifies $\text{trunc}(\log_2 n + 1)$ which is ambiguous in terms of being either $\text{trunc}(\log_2(n+1))$ or $\text{trunc}(\log_2(n) + 1)$, although he probably meant the latter.

```
dt <- wsrpart(form, data)
```

The template for the function is:

```
wsrpart <- function(formula, data, nvars, ...)  
{  
  ...  
}
```

Notice the use of “...” in the argument list. This allows us to pass other arguments on through to `rpart()`.

The actual solution will produce the following output:

```
## system.time(model <- wsrpart(form, ds[train, vars]))  
  
##      user  system elapsed  
## 0.164    0.020    2.457  
  
model  
  
## A multiple rpart model with 1 tree.  
##  
## Variables used (11): Humidity3pm, WindGustSpeed, Cloud9am, MaxTemp, Cloud3pm,  
##                      Temp3pm, Humidity9am, Sunshine,  
##                      ....
```

9 Exercise: Multiple Decision Trees

Now extend the function `wsrpart()` to build multiple decision trees. The function will take a formula, a dataset (`data`) and a number of trees to build (`ntrees`), and returns a list of `ntrees` decision trees. Each element of the list will itself be a list, with at least one element named `model`. This is the actual `rpart` model. The result should be of class `mrpart` for “multiple `rpart`.”

The actual solution will produce the following output:

```
system.time(model <- wsrpart(form, ds[train, vars], 4))

##      user  system elapsed
## 118.492    2.932   49.340

class(model)

## [1] "mrpart"

length(model)

## [1] 50

class(model[[1]]$model)

## [1] "rpart"

model[[1]]$model

## n=45794 (876 observations deleted due to missingness)
##
## node), split, n, loss, yval, (yprob)
##      * denotes terminal node
##
## 1) root 45794 10720 No (0.7660 0.2340)
##   2) Humidity3pm< 75.5 39713 6399 No (0.8389 0.1611) *
##   3) Humidity3pm>=75.5 6081 1762 Yes (0.2898 0.7102) *
```

10 Exercise: `predict.mrpart()`

Score a Dataset Using the Forest

Define a function `predict.mrpart()`. Returns the proportion of trees voting for the positive case, assuming binary classification models.

The actual solution will produce the following output:

```
predict(model, ds[test,vars])  
##      8      9     12     16     19     27     36     45     46     51     55     59  
##    No    No   Yes    No    No    No    No    No    No    No    No    No  
##    61    62    66    68    78    83    86    87    88    89    93    95  
##    No    No    No    No    No    No    No    No    No    No    No    No  
....
```

This page
is under
develop-
ment.

11 Further Reading

12 References

R Core Team (2013). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. URL <http://www.R-project.org/>.

Williams GJ (2009). “Rattle: A Data Mining GUI for R.” *The R Journal*, **1**(2), 45–55. URL http://journal.r-project.org/archive/2009-2/RJournal_2009-2_Williams.pdf.

Williams GJ (2011). *Data Mining with Rattle and R: The art of excavating data for knowledge discovery*. Use R! Springer, New York. URL http://www.amazon.com/gp/product/1441998896/ref=as_li_qf_sp_asin_tl?ie=UTF8&tag=togaware-20&linkCode=as2&camp=217145&creative=399373&creativeASIN=1441998896.

Williams GJ (2013). *rattle: Graphical user interface for data mining in R*. R package version 2.6.27, URL <http://rattle.togaware.com/>.