

HW01p

Evangeline Szpylka

February 24, 2018

Welcome to HW01p where the “p” stands for “practice” meaning you will use R to solve practical problems. This homework is due 11:59 PM Saturday 2/24/18.

You should have RStudio installed to edit this file. You will write code in places marked “TO-DO” to complete the problems. Some of this will be a pure programming assignment. The tools for the solutions to these problems can be found in the class practice lectures. I want you to use the methods I taught you, not for you to google and come up with whatever works. You won’t learn that way.

To “hand in” the homework, you should compile or publish this file into a PDF that includes output of your code. Once it’s done, push by the deadline.

R Basics

First, install the package `testthat` (a widely accepted testing suite for R) from <https://github.com/r-lib/testthat> using `pacman`. If you are using Windows, this will be a long install, but you have to go through it for some of the stuff we are doing in class. LINUX (or MAC) is preferred for coding. If you can’t get it to work, install this package from CRAN (still using `pacman`), but this is not recommended long term.

```
if (!require("pacman")){install.packages("pacman")}
```

```
## Loading required package: pacman
```

```
pacman::p_load(testthat)
```

1. Use the `seq` function to create vector `v` consisting of all numbers from -100 to 100.

```
v = seq(-100, 100)
```

Test using the following code:

```
expect_equal(v, -100 : 100)
```

If there are any errors, the `expect_equal` function will tell you about them. If there are no errors, then it will be silent.

2. Create a function `my_reverse` which takes as required input a vector and returns the vector in reverse where the first entry is the last entry, etc. No function calls are allowed inside your function (otherwise that would defeat the purpose of the exercise).

```
my_reverse = function(v){
  vec = c(v)
  for(j in length(v) : 1){
    for(k in length(v) : 1){
      temp = v[j]
      v[j] = vec[k]
      vec[k] = temp
    }
  }
  vec
}
```

Test using the following code:

```
expect_equal(my_reverse(c("A", "B", "C")), c("C", "B", "A"))
expect_equal(my_reverse(v), rev(v))
```

3. Let $n = 50$. Create a $n \times n$ matrix R of exactly 50% entries 0's, 25% 1's 25% 2's in random locations.

```
n = 50
R = matrix(nrow = n, ncol = n, sample(
  c(rep(0, 0.50*n^2), rep(1, 0.25*n^2),
    rep(2, 0.25*n^2)),
  size = n^2,
  replace = FALSE
))
```

Test using the following and write two more tests as specified below:

```
expect_equal(dim(R), c(n, n))
expect_equal(sum(c(R) == 0 | c(R) == 1 | c(R) == 2), n^2 )
expect_equal(sum(c(R) == 2), 0.25 * n^2)
```

4. Randomly punch holes (i.e. NA) values in this matrix so that approximately 30% of the entries are missing.

```
R = replace(R, sample(1 : n^2, 0.3 * n^2), NA)
```

Test using the following code. Note this test may fail 1/100 times.

```
num_missing_in_R = sum(is.na(c(R)))
expect_lt(num_missing_in_R, qbinom(0.995, n^2, 0.3))
expect_gte(num_missing_in_R, qbinom(0.005, n^2, 0.3))
```

5. Sort the rows matrix R by the largest row sum to lowest. See 2/3 way through practice lecture 3 for a hint.

```
row_value = c()

for(j in 1 : n){
  row_value = c(row_value, sum(R[j, ], na.rm = TRUE))
}
rownames(R) = row_value
R = R[order(rownames(R), decreasing = TRUE), ]
```

Test using the following code.

```
for (i in 2 : n){
  expect_gte(sum(R[i - 1, ], na.rm = TRUE), sum(R[i, ], na.rm = TRUE))
}
```

6. Create a vector v consisting of a sample of 1,000 iid normal realizations with mean -10 and variance 10.

```
v = rnorm(1000, mean = -10, sd = sqrt(10))
```

Find the average of v and the standard error of v .

```
average = mean(v)
st_error = sd(v)/sqrt(length(v))
```

Find the 5%ile of v and use the `qnorm` function as part of a test to ensure it is correct based on probability theory.

```
q_1 = quantile(v, probs = 0.05)
q_2 = qnorm(0.05, mean = -10, sd = sqrt(10))

expect_equal(as.numeric(q_1), expected = q_2, tol = 0.05)
```

Find the sample quantile corresponding to the value -7000 of `v` and use the `pnorm` function as part of a test to ensure it is correct based on probability theory.

```
inverse_quantile_obj = ecdf(v[-7000])
iq_1 = inverse_quantile_obj(-7000)
iq_2 = pnorm(-7000, mean = -10, sd = sqrt(10), lower.tail = TRUE, log.p = FALSE)

expect_equal(as.numeric(iq_1), expected = iq_2, tol = 0.05)
```

7. Create a list named `my_list` with keys "A", "B", ... where the entries are arrays of size 1, 2 x 2, 3 x 3 x 3, etc. Fill the array with the numbers 1, 2, 3, etc. Make 8 entries.

```
my_list = list()
n = 8
my_keys = c("A", "B", "C", "D", "E", "F", "G", "H")

for (j in 1:n){
  key = my_keys[j]
  my_list[[key]] = array(seq(1, j), dim = c(rep(j,j)))
}
```

Test with the following uncomprehensive tests:

```
expect_equal(my_list$A[1], 1)
expect_equal(my_list[[2]][, 1], 1 : 2)
expect_equal(dim(my_list[["H"]]), rep(8, 8))
```

Run the following code:

```
lapply(my_list, object.size)
```

```
## $A
## 208 bytes
##
## $B
## 216 bytes
##
## $C
## 336 bytes
##
## $D
## 1232 bytes
##
## $E
## 12728 bytes
##
## $F
## 186848 bytes
##
## $G
## 3294400 bytes
##
```

```
## $H
## 67109088 bytes
```

```
?lapply
```

```
## starting httpd help server ... done
```

```
?object.size
```

Use `?lapply` and `?object.size` to read about what these functions do. Then explain the output you see above. For the later arrays, does it make sense given the dimensions of the arrays?

In the output above, I see the key names of my arrays as well as a corresponding amount of bytes below each key. The `lapply` function returned back a list of the same length as `my_list`, which had 8 elements. The `object.size` function gives us an estimate of how much memory is being used to store an (R) object, which we see above. It does make sense that the later arrays take more memory because they have dimensions, hence more “levels” of data to contain.

Now cleanup the namespace by deleting all stored objects and functions:

```
rm(list = ls())
```

Basic Binary Classification Modeling

8. Load the famous `iris` data frame into the namespace. Provide a summary of the columns and write a few descriptive sentences about the distributions using the code below and in English.

The “iris” dataset has five columns: sepal length, sepal width, petal length, petal width, and species. The three species of flowers are distributed equally, with 50 observations per species for the entire sample. The average petal length is 3.758 units and the average width is 1.199 units. The average sepal (green-looking petal) length is 5.843 units and the average sepal width is 3.057 units. It is interesting that the median values for the petal measurements are larger than the mean values, whereas the means of the sepal measurements are larger than the medians.

```
data(iris)
summary(iris)
```

```
##      Sepal.Length      Sepal.Width      Petal.Length      Petal.Width
##  Min.   :4.300      Min.   :2.000      Min.   :1.000      Min.   :0.100
## 1st Qu.:5.100      1st Qu.:2.800      1st Qu.:1.600      1st Qu.:0.300
## Median :5.800      Median :3.000      Median :4.350      Median :1.300
## Mean   :5.843      Mean   :3.057      Mean   :3.758      Mean   :1.199
## 3rd Qu.:6.400      3rd Qu.:3.300      3rd Qu.:5.100      3rd Qu.:1.800
## Max.   :7.900      Max.   :4.400      Max.   :6.900      Max.   :2.500
##      Species
## setosa      :50
## versicolor:50
## virginica   :50
##
##
##
```

The outcome metric is `Species`. This is what we will be trying to predict. However, we have only done binary classification in class (i.e. two classes). Thus the first order of business is to drop one class. Let’s drop the level “virginica” from the data frame.

```
iris = iris[iris$Species != "virginica", ]
```

Now create a vector `y` that is length the number of remaining rows in the data frame whose entries are 0 if “setosa” and 1 if “versicolor”.

```
y = nrow(iris)

for(i in 1:nrow(iris)){
  if(iris$Species[i] == "versicolor"){
    y[i] = 1
  } else{
    y[i] = 0
  }
}
```

9. Fit a threshold model to `y` using the feature `Sepal.Length`. Try to write your own code to do this. What is the estimated value of the threshold parameter? What is the total number of errors this model makes?

```
Y_1 = as.matrix(cbind(iris[, 1, drop = FALSE]))
MAX_ITER = 100
w_vec = c(0, 0)

for(iter in 1 : MAX_ITER){
  for(i in 1 : nrow(Y_1)){
    x_i = Y_1[i]
    yhat_i = ifelse(sum(x_i * w_vec) > 0, 1, 0)
    y_i = y[i]
    w_vec = w_vec + (y_i - yhat_i) * x_i
  }
}

yhat = ifelse(Y_1 %*% w_vec > 0, 1, 0)
sum(y != yhat) / length(y)
```

```
## [1] 1
```

Does this make sense given the following summaries:

```
summary(iris[iris$Species == "setosa", "Sepal.Length"])
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##  4.300   4.800   5.000   5.006   5.200   5.800
```

```
summary(iris[iris$Species == "versicolor", "Sepal.Length"])
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##  4.900   5.600   5.900   5.936   6.300   7.000
```

Yes, this does make sense given the following summaries. By looking at the summaries above, we can clearly see that the lengths of the two species overlap. In addition, their minimum and maximum lengths are “near” each other. As an example, the flower species setosa has a minimum length of 4.3 units while the versicolor has a minimum length of 4.9 units. Thus, it is reasonable that there would be an error while trying to draw a “correct” line separating the sepal length sets.

10. Fit a perceptron model explaining `y` using all three features. Try to write your own code to do this. Provide the estimated parameters (i.e. the four entries of the weight vector)? What is the total number of errors this model makes?

```
Y_2 = as.matrix(cbind(y, iris[, 1, drop = FALSE], iris[, 2, drop = FALSE],
                      iris[, 3, drop = FALSE], iris[, 4, drop = FALSE]))
```

```

MAX_ITER= 100
w_vec = rep(0, 5)

for(iter in 1 : MAX_ITER){
  for(i in 1 : nrow(Y_2)){
    x_i = Y_2[i,]
    yhat_i = ifelse(sum(x_i * w_vec) > 0, 1, 0)
    y_i = y[i]
    w_vec = w_vec + (y_i - yhat_i) * x_i
  }
}
yhat = ifelse(Y_2 %*% w_vec > 0, 1, 0)
sum(y != yhat) / length(y)

```

```
## [1] 0
```

The four entries of the weight vector are as follows: Sepal length = -1.1 Sepal width = -3.6 Petal length = 5.2 Petal Width = 2.2 In addition, this model makes zero errors.