

Lab 2

Adriana Sham

Basic R Skills

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.

```
#install.packages("pacman")
#install.packages("testthat")
library(pacman)
library(testthat)
```

- 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.

- 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).

```
#TO-DO
my_reverse = function(vec){
  revvec = rep(NA,length(vec))
  for(i in seq(1, length(vec))){
    revvec[i] = vec[ (length(vec)+1)-i]
  }
  revvec
}
```

Test using the following code:

```
expect_equal(my_reverse(c("1", "2", "3")), c("3", "2", "1"))
expect_equal(my_reverse(v), rev(v))
```

- Let `n = 50`. Create a `nxn` matrix `R` of exactly 50% entries 0's, 25% 1's 25% 2's in random locations.

```
n = 50
numbers = c( rep( 0, n * n * .5 ), rep( 1, n * n * .25 ), rep( 2, n * n * .25 ) )
R = matrix( sample ( numbers ) , nrow = n , ncol = n )
```

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

```
expect_equal(dim(R), c(n, n))
#test that the only unique values are 0, 1, 2
expect_equal(n*n, sum(c(R)==0) + sum(c(R)==1) + sum(c(R)==2))
#test that there are exactly 625 2's
expect_equal(625,sum(c(R)==2))
```

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

```
#TO-DO
R[sample (n*n, n * n * .3)] = NA
```

Use the `testthat` library to test that this worked correctly by ensuring the number of missing entries is between the 0.5%ile and 99.5%ile of the appropriate binomial.

```
NA_in_R = sum(is.na(R))
expect_lt(NA_in_R, qbinom(0.995, n*n, 0.3))
expect_gt(NA_in_R, qbinom(0.005, n*n, 0.3))
```

- Sort the rows matrix R by the largest row sum to lowest. Be careful about the NA's!

```
r_names = c()
for (i in 1:n){
  r_names = c(r_names, sum(R[i,], na.rm = TRUE)) }
row.names(R) = r_names
R = R[order(row.names(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))
}
```

- We will now learn the `apply` function. This is a handy function that saves writing for loops which should be eschewed in R. Use the `apply` function to compute a vector whose entries are the standard deviation of each row. Use the `apply` function to compute a vector whose entries are the standard deviation of each column. Be careful about the NA's!

```
#standard_deviation= sd(R, na.rm = TRUE)

standard_deviation_of_row = apply( R , 1, sd, na.rm = TRUE )
standard_deviation_of_column = apply( R, 2, sd, na.rm = TRUE )
```

- Use the `apply` function to compute a vector whose entries are the count of entries that are 1 or 2 in each column. Try to do this in one line.

```
apply( R >= 1, 2, sum, na.rm=TRUE)
```

```
## [1] 16 16 13 11 27 16 19 16 18 21 19 17 18 17 14 18 16 14 16 16 19 19 27
## [24] 20 17 17 20 21 16 15 13 21 24 23 21 15 14 12 15 19 18 18 17 19 17 21
## [47] 24 13 22 16
```

- Use the `split` function to create a list whose keys are the column number and values are the vector of the columns. Look at the last example in the documentation `?split`.

```
l=split(R, col(R), drop = TRUE)
```

- In one statement, use the `lapply` function to create a list whose keys are the column number and values are themselves a list with keys: “min” whose value is the minimum of the column, “max” whose value is the maximum of the column, “pct_missing” is the proportion of missingness in the column and “first_NA” whose value is the row number of the first time the NA appears. Use the `which` function.

```
list=lapply(l, function(R) {
  minimum = min (R,na.rm = TRUE)
  pct_missing = sum(is.na(R))/ length(R)*100
  first_NA=min(which(is.na(R)))
})
```

```
maximum=max(R,na.rm=TRUE)
c(minimum,maximum,pct_missing,first_NA))
```

- 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)
standard_error = sd(v)/(sqrt(length(v)))
```

- Find the 5%ile of `v` and use the `qnorm` function to compute what it theoretically should be.

```
quant_1 = quantile(v,0.05)
quant_2 = qnorm(p = .05, mean = -10, sd = sqrt(10))
expect_equal(as.numeric(quant_1), as.numeric(quant_2),tolerance = standard_error)
```

- 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.

```
letters = c("A", "B" , "C" , "D" , "E" , "F" , "G","H")
my_list = list()
for(i in 1:8){
  key = letters[i]
  my_list[[key]] = array(seq(1, i) , dim = c(rep(i) , i)))
}
```

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
## 224 bytes
##
## $B
## 232 bytes
##
## $C
## 352 bytes
##
## $D
## 1248 bytes
##
## $E
## 12744 bytes
##
## $F
## 186864 bytes
##
## $G
## 3294416 bytes
```

```
##
## $H
## 67109104 bytes
```

Use `?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?

Answer here in English.

`Object.size` provides an estimate of the memory size that is being used in each key of the `my_list` object. It makes sense in this case because as the dimensions increase, the memory that is being used also becomes larger and larger.

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

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

Basic Binary Classification Modeling

- 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.

```
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.

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

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

- 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?

```
X1 = as.matrix(cbind(drop_level[, 1, drop = FALSE]))
MAX_ITER = 100
w_vec = 0
for (iter in 1 : MAX_ITER){
  for (i in 1 : nrow(X1)){
    x_i = X1[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
  } }

#total error in model
yhat = ifelse(X1 %*% w_vec > 0, 1, 0)
sum(y != yhat) / length(y)
```

```
## [1] 0.5
```

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 == "virginica", "Sepal.Length"])
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      4.900   6.225   6.500   6.588   6.900   7.900
```

Write your answer here in English.

Yes. The summaries provide us enough information to be able to visualize the line between the two sets of `Sepal.Length` data, setosa and versicolor. The min and max’s for each are similar, meaning there would be some error because the data is not linearly separable. Plotting the `Sepal.Length` data will show you that it is not linear separable. Using just `Sepal.Length` is not a good enough way to identify species.

Yes, the summaries has given us enough information to understand that the line we have draw through the two sets of `Sepal.Length` would encounter some errors because the two data sets are very similar and might not be a good way to identify the species.

- What is the total number of errors this model makes (in-sample)?

```
yhat = ifelse(X1 %*% w_vec > 0, 1, 0)
sum(y != yhat) / length(y)
```

```
## [1] 0.5
```